



## ARTICLE

DOI: 10.1038/s42003-018-0094-7

OPEN

# Evidence of non-tandemly repeated rDNAs and their intragenomic heterogeneity in *Rhizophagus irregularis*

Taro Maeda <sup>1</sup>, Yuuki Kobayashi<sup>1</sup>, Hiromu Kameoka <sup>1</sup>, Nao Okuma<sup>1,2</sup>, Naoya Takeda<sup>3</sup>, Katsushi Yamaguchi<sup>4</sup>, Takahiro Bino<sup>4</sup>, Shuji Shigenobu<sup>2,4</sup> & Masayoshi Kawaguchi<sup>1,2</sup>

Arbuscular mycorrhizal fungus (AMF) species are some of the most widespread symbionts of land plants. Our much improved reference genome assembly of a model AMF, *Rhizophagus irregularis* DAOM-181602 (total contigs = 210), facilitated a discovery of repetitive elements with unusual characteristics. *R. irregularis* has only ten or 11 copies of complete 45S rDNAs, whereas the general eukaryotic genome has tens to thousands of rDNA copies. *R. irregularis* rDNAs are highly heterogeneous and lack a tandem repeat structure. These findings provide evidence for the hypothesis that rDNA heterogeneity depends on the lack of tandem repeat structures. RNA-Seq analysis confirmed that all rDNA variants are actively transcribed. Observed rDNA/rRNA polymorphisms may modulate translation by using different ribosomes depending on biotic and abiotic interactions. The non-tandem repeat structure and intragenomic heterogeneity of AMF rDNA/rRNA may facilitate successful adaptation to various environmental conditions, increasing host compatibility of these symbiotic fungi.

<sup>1</sup>Division of Symbiotic Systems, National Institute for Basic Biology, Myodaiji Nishigonaka, Okazaki, Aichi 444-8585, Japan. <sup>2</sup>The Graduate University for Advanced Studies [SOKENDAI], Hayama, Miura, Kanagawa 240-0193, Japan. <sup>3</sup>School of Science and Technology, Kwansai Gakuin University, Gakuen, Mita, Hyogo 669-1337, Japan. <sup>4</sup>Functional Genomics Facility, National Institute for Basic Biology, Myodaiji Nishigonaka, Okazaki, Aichi 444-8585, Japan. Correspondence and requests for materials should be addressed to S.S. (email: [shige@nibb.ac.jp](mailto:shige@nibb.ac.jp)) or to M.K. (email: [masayosi@nibb.ac.jp](mailto:masayosi@nibb.ac.jp))

The arbuscular mycorrhizal fungus (AMF) is an ancient fungus with origins at least as old as the early Devonian period<sup>1,2</sup>. AMF colonizes plant roots and develops highly branched structures called arbuscules in which soil nutrients (phosphate and nitrogen) are efficiently delivered to the host plant<sup>3</sup>. AMF forms symbiotic networks with most land plant species<sup>4,5</sup>, and the mycelial network formed by various AMF species contributes to plant biodiversity and productivity within the terrestrial ecosystem<sup>6</sup>. The distinctive features of AMF have made it an important model in ecology and evolution<sup>7,8</sup>; these features include coenocytic mycelia<sup>5</sup>, nutrition exchange with plant, classification as an obligate biotroph<sup>9</sup>, signal crosstalk during mycorrhiza development<sup>9,10</sup>, and extremely high symbiotic ability<sup>9,11</sup>.

Recently, multiple genome projects have advanced the understanding of AMF species. Genomic data have been provided for *Rhizophagus irregularis* DAOM-181602 (=DOAM-197198)<sup>12–14</sup>, *Gigaspora rosea*<sup>12</sup>, *Rhizophagus clarus*<sup>15</sup>, and other isolates of *R. irregularis*<sup>14,16</sup>. These studies revealed potential host-dependent biological pathways<sup>12,17</sup> and candidate genes for plant infection and sexual reproduction<sup>15,16,17</sup>. However, fragmented genome sequences limit the ability to analyze repetitive structures and to distinguish between orthologous and paralogous genes<sup>14</sup>. The first published genome sequence of *R. irregularis* DAOM-181602 (JGI\_v1.0)<sup>17</sup> contained 28,371 scaffolds and an N50 index of 4.2 kb (Supplementary Table 1). The second sequence by Lin et al.<sup>13</sup> (Lin14) contained 30,233 scaffolds with an N50 of 16.4 kb (Supplementary Table 1). Recently published assemblies by Chen et al.<sup>14</sup> (JGI\_v2.0) contained 1123 scaffolds with an N50 of 336.4 kb (Supplementary Table 1). The quality of genomic sequence data for other AMF species did not surpass that of DAOM-181602<sup>12,15</sup>. In contrast, many fungi that are not AMF species contain less than several hundred scaffolds and N50 lengths over 1 Mb<sup>18</sup>. For example, a genomic sequence of an asymbiotic fungus closely related to AMF, *Rhizopus delemar* (GCA000149305.1), was constructed from 83 assemblies with an N50 of 3.1 Mb<sup>19</sup>. Thus, we here present an improved whole-genome sequence of *R. irregularis* DAOM-181602 to facilitate examination of the genomics underlying specific features of AMF species. Taking an advantage of the highly contiguous assembly with little ambiguous regions, we focus on the investigation of the repetitive structures including transposable elements (TE), highly duplicated genes, and rDNA gene copies.

A general eukaryotic genome has tens to thousands of rDNA copies<sup>20</sup> (Supplementary Figure 1a), and the sequences of the copies are identical or nearly identical. However, since Sanders et al.<sup>21</sup>, many studies have indicated intracellular polymorphisms of rDNA (ITS) in various AMF species<sup>22–24</sup>, and the sequencing of isolated nuclei from *Claroideoglossum etunicatum* and *R. irregularis* DAOM-181602 suggested sequence variation among the paralogous rDNAs, i.e., intragenomic heterogeneity<sup>13,25</sup>. This heterogeneity has potentially high impact of studying AMF species, because the rDNA is a fundamental marker of the AMF phylogeny and ecology<sup>8,26–28</sup>, and studies have assumed that these rDNAs have no intragenomic sequence variation<sup>29</sup>. Hence, determining the variation degree could cause a reevaluation of the previous understanding of geographic distribution<sup>8</sup>, species identification<sup>28</sup>, and evolutionary processes of AMF. However, the degree of the variation among the 48S rDNA paralogs has been ambiguous because previous studies by Sanger or Illumina sequencing were unable to distinguish each rDNA paralog in a genome. Moreover, the number of rDNA genes in an AMF genome has never been investigated.

The tandem repeat structure (TRS) of the rDNAs is also an attractive topic for evolutionary studies. General organisms

require many rDNA copies to make a sufficient amount of rRNA for protein translation<sup>30,31</sup>. However, in the evolutionary time-scale, multicopy genes reduce in number due to homologous recombination (Supplementary Figure 1b)<sup>32,33</sup> and single-strand annealing (Supplementary Figure 1c)<sup>34</sup>. To maintain the number of rDNAs, eukaryotes increase the number of copies by unequal sister chromatid recombination (USCR) using the rDNA TRS (Supplementary Figure 1d, e)<sup>32</sup>. Because this rDNA replacement causes a bottleneck effect in the genome, almost all eukaryotes have homogenous rDNAs in their genomes<sup>20</sup>. This process, termed concerted evolution is an essential system for maintaining eukaryotic protein translation by ribosomes<sup>30</sup>. The heterogeneous rDNAs observed in AMF species implies the collapse of their concerted evolution, and suggest the unique maintenance system of rDNA copy number.

In this study, we built an improved reference genome assembly of *R. irregularis* DAOM-181602, which allowed us to discover repetitive elements with unique characteristics of the AMF genome. We identified an unusually small number of rDNA genes in the *R. irregularis* genome. We also found that the rDNA copies are highly heterogeneous and lack a TRS.

## Results

**A contiguous DAOM-181602 genome generated by PacBio data.** We primarily used single-molecule, real-time (SMRT) sequencing technology for sequencing and assembling the *R. irregularis* genome. We generated a 76-fold whole-genome shotgun sequence (11.7 Gb in total) (Supplementary Table 2) from genome DNA isolated from a spore suspension of a commercial strain of *R. irregularis* DAOM-181602 using the PacBio SMRT sequencing platform. A total of 766,159 reads were generated with an average length of 13.1 kb and an N50 length of 18.8 kb (Supplementary Table 2). We assembled these PacBio reads using the HGAP3 program<sup>35</sup> (149.9 Mb composed of 219 contigs). To detect erroneous base calls, we generated 423 Mb of 101 bases-paired-end Illumina whole-genome sequence data (Supplementary Table 2) and aligned them to the HGAP3 assembly. Through variant calling, we corrected 3032 single base call errors and 10,841 small indels in the HGAP3 assembly. Nine contigs were almost identical to carrot DNA sequences deposited in the public database (Supplementary Table 3), and these were removed as contaminants derived from a host plant used by the manufacturer. We evaluated the completeness of the final assembly using CEGMA<sup>36</sup>; of the 248 core eukaryotic genes, 244

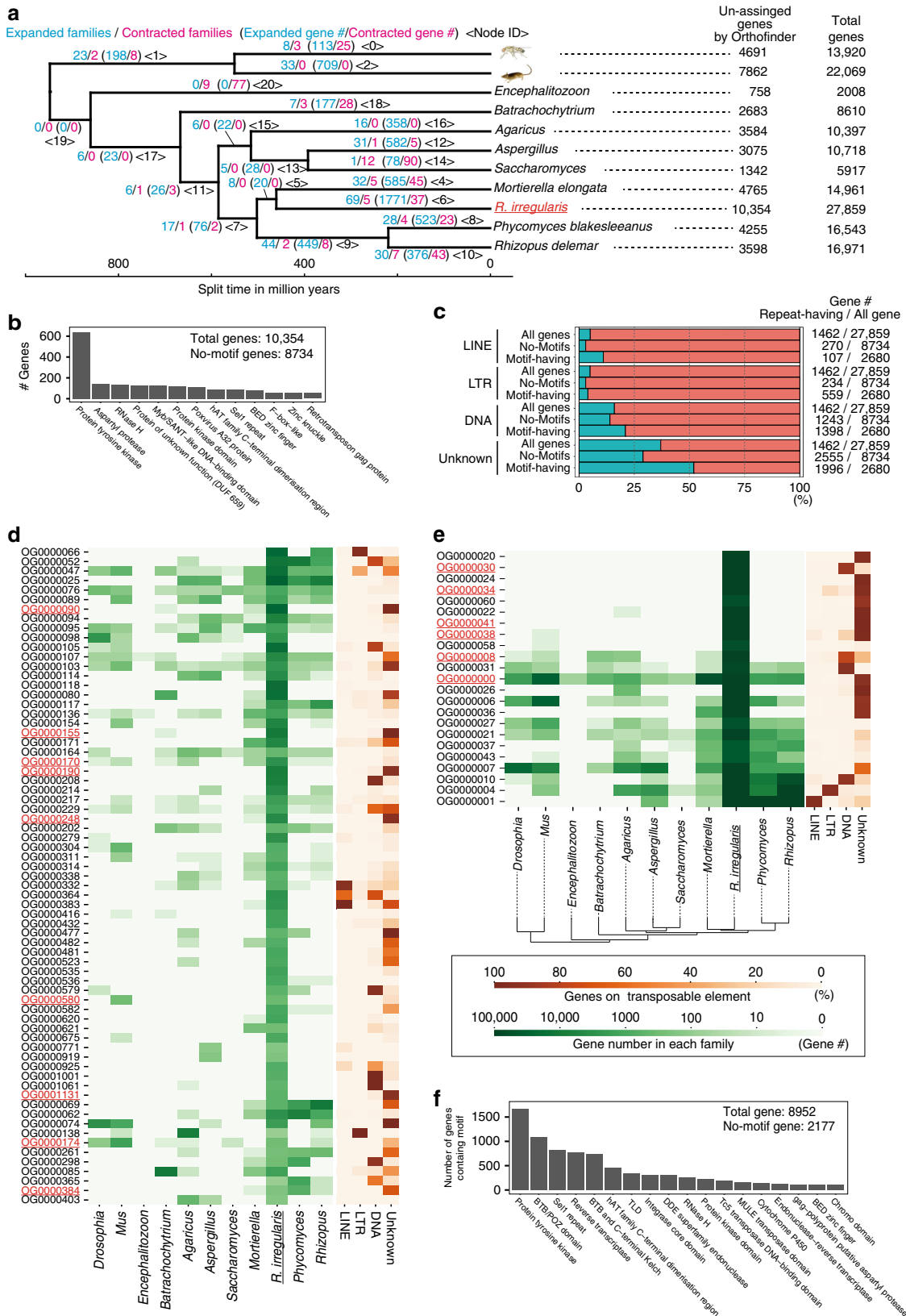
**Table 1 Assembly statistics of *R. irregularis* genome**

	RIR17
Accession number	BDIQ01000000
Predicted genome size by flow cytometry	154 Mb
Total length of contigs (% of genome)	149,750,837 bases (97%)
# Contigs	210
# N bases	0
Longest contig (bp)	5,727,599
Contig N50 (bases)	2,308,146
L50	23
GC %	27.9%
CEGMA completeness for genome contigs	98.4%
# of predicted genes	41,572
BUSCO completeness for gene models (DB; fungi_odb9)	94.1% (273/290)
Complete single copy	83.8% (243/290)
Complete duplicated	10.3% (30/290)
Fragmented	3.8% (11/290)
Missing	2.1% (6/290)

genes (98.4%) were completely assembled (Table 1 and Supplementary Table 1). Consequently, we obtained a high-quality reference genome assembly of *R. irregularis* DAOM-181602, which is referred to as RIR17.

Compared with previous assemblies<sup>13,14,17</sup>, RIR17 represents a decrease in assembly fragmentation (1123 to 210) and an

improvement in contiguity using the N50 contig length as a metric (Table 1 and Supplementary Table 1). The total size of the assembly was 9–59 Mb greater than that of previous versions, reaching 97.24% coverage of the whole genome (154 Mb)<sup>17</sup> (Table 1 and Supplementary Table 1). The new assembly contains no ambiguous bases (N-bases), whereas previous assemblies had



30,115–6,925,426 N-bases (Table 1 and Supplementary Table 1). Approximately 1–7 Mb of sequences from previous assemblies were not contained in RIR17, and JGI\_v2.0 has one more conserved gene family than RIR17 (Supplementary Table 1), indicating that a few genomic parts remain to be uncovered by our improvement with continuous sequences. On the other hand, RIR17 was aligned with 95–99.2% of previous assemblies (Supplementary Table 4), suggesting that RIR17 covers the majority of the previously sequenced areas with high sequence contiguity. Moreover, RIR17 contained 8–47 Mb of regions unassigned in previous genomes (Supplementary Table 4). These regions are newly revealed by our improvement.

RIR17 contains a greater extent of repetitive regions than JGI\_v2.0. The RepeatModeler<sup>37</sup> and RepeatMasker<sup>37</sup> pipeline identified 64.4 Mb (43.03%) of RIR17 as repetitive regions (Supplementary Data 1). These regions total 18.9 Mb more than those of JGI\_v2.0 (Supplementary Data 1). Previous fosmid sequences predicted that DAOM-181602 contains ~55 Mb of repetitive regions<sup>17</sup>, suggesting that RIR17 covers the majority of the repetitive regions of DAOM-181602.

We confirmed a unique repeat profile in the AMF genome. The majority of the interspersed repeats (62.83%) could not be categorized with known repeat classes (Supplementary Data 1), indicating that the AMF genome accumulated novel classes of interspersed repeats. Moreover, DAOM-181602 lacks short interspersed nuclear elements (SINEs), which are abundant in closely related fungi (Supplementary Data 1). Several types of SINEs proliferate using transposases on long interspersed nuclear elements (LINEs)<sup>38</sup>. Although the AMF has 23 LINEs containing the transposase gene (Supplementary Data 2), SINEs have never been found in previous genomes<sup>13,14,17</sup> or RIR17. DAOM-181602 may have a system to suppress the invasion and proliferation of SINEs (e.g., a high number of very active Argonaute proteins, as predicted by Tisserant et al.<sup>17</sup>).

#### New gene annotation details gene family expansion in AMF.

Using the RIR17 assembly together with strand-specific RNA-Seq data (Rir\_RNA\_SS in Supplementary Table 2), we built a set of 41,572 gene models (Supplementary Data 1 and Supplementary Table 5). Of the genes predicted, 27,859 (67.0%) had either RNA-Seq expression support, homology support, or protein motif support (Supplementary Data 1 and Supplementary Table 5). The gene models having any support were submitted to the DDBJ as standard genes and were used in downstream analyses. The models having no support were assigned as PROVISIONAL gene models (Supplementary Data 1 and Supplementary Table 5). Using Orthofinder with previous genomic gene sets indicated that our gene models cover the majority of previously provided genes (Supplementary Figure 2). Although new models showed more coverage of Benchmarking Universal Single-Copy Orthologs (BUSCOs)<sup>39</sup> (Supplementary Table 5) than JGI\_v1.0 and Lin14, their gene completeness was slightly lower than that of JGI v2.0 (nine BUSCO families overlooked, Supplementary Table 5), indicating the advantage of using the JGI annotation pipeline to discuss the gene variety in DAOM-181602<sup>14</sup>. However, we considered our model set suitable for the analysis of the repetitive

region and highly paralogous genes because our model is based on highly continuous assemblies, and the number of genes on repetitive regions was increased to 2349–12,559 genes from the number in JGI\_2.0 (Supplementary Table 6).

*R. irregularis* has one of the largest numbers of genes in fungi (Supplementary Figure 3). Our ortholog analyses indicate that the gene number inflation was caused by lineage-specific expansions of gene families and not by whole-genome duplications. An Orthofinder analysis of nine fungal genomes and two animal data sets (Supplementary Table 7) showed that many of the single-copy genes in other fungi were also single copies in RIR17 (216/239 families, Supplementary Data 3), negating the possibility of whole-genome duplication in *R. irregularis*. The large number of species-specific single-copy genes in DAOM-181602 (10,354 genes, Fig. 1a, Supplementary Data 3) suggests that the AMF genes inflated by new gene constructions through gene fusion and mutation accumulation. Moreover, several common gene families in Opisthokonta also contributed to gene inflation; the *R. irregularis* lineage had 92 rapidly expanded families, containing 8952 genes (Fig. 1a, d, e, Supplementary Data 3 and 4), suggesting that *R. irregularis* has also acquired many genes by the duplication of particular gene families.

The motif annotation indicates that inflated genes may contribute to signaling pathways of AMF species. Our Pfam search annotated 1620 species-specific single-copy genes and 6755 rapidly expanded genes (Fig. 1b, f, Supplementary Tables 6 and Supplementary Data 5). The most frequently observed motif was protein tyrosine kinase; PF07714 (Fig. 1b, f, Supplementary Data 6), which is often found in signaling proteins in multicellular organisms<sup>40</sup>, consistently with previous genome studies<sup>14</sup>. Other signal-related motifs (e.g., Sell repeat and BED zinc finger) were also found in the inflated genes (Fig. 1b, f, Supplementary Data 6). AMF has developed a unique signal pathway for symbiosis (e.g., establishments of symbiosis with pathways via SIS1<sup>41</sup> and lyso-phosphatidylcholine<sup>42</sup>). This inflation of signaling-related genes may have led the development of a complex signaling pathway in AMF.

We then investigated the contribution of the TEs to gene inflation based on the overlapping of highly paralogous genes and the TEs. Previous studies hypothesized that the gene inflation in *R. irregularis* relates with the expansion of TEs<sup>14</sup>. Our analysis showed that in several rapidly expanded families (e.g., OG0000090 and OG0000020), over 90% of the genes were located with TEs (Fig. 1d, e, Supplementary Data 7), suggesting that TEs accelerated the gene expansion in these families. However, some of the families had no correspondence with TEs (e.g., OG0000025 and OG0000058 in Fig. 1c, e). In species-specific single-copy genes, TEs were slightly more frequently found with motifs than in all gene sets but were less frequently found in species-specific single-copy genes without motifs (Fig. 1c). This detailed analysis supports the contribution of TEs to gene inflation in several gene families but also clarified that several families show TE-independent expansion. Although more genome data for AMF species and sister groups are required to reveal the gene expansion process and its contribution to AM symbiosis, our data provide a fundamental dataset to reveal the evolution of gene redundancy in AMF species.

**Fig. 1** Gene inflation in *R. irregularis*. **a** Rapidly expanded/contracted ortholog groups based on CAFE analysis. Total gene number of analyzed species and unassigned genes by Orthofinder analysis (species-specific single-copy genes) are described on the right side of the tree. Illustrations were modified from the resources distributed in the Togo picture gallery (licensed under CC-BY 4.0 ©Togo picture gallery). **b** The number of *R. irregularis*-specific single-copy genes having protein motifs. Minor motifs (<50 genes) were omitted from the figure (raw-data; Supplementary Data 6). **c** The proportion of genes among the species-specific single-copy genes having each repeat element. **d** Sixty-nine rapidly expanded orthologous groups (OGs). Green heat map shows the number of genes in each OG. Orange heat map indicates the proportion of genes with each repeat element. The OGs containing the protein tyrosine kinase domain are marked in red. **e** Rapidly expanded OGs based on z-score analysis. The colors have the same meaning as in (d). **f** The number of rapidly expanded ortholog genes having protein motifs. Minor motifs (<100 genes) are omitted from the figure (raw-data; Supplementary Data 6)

**Losing conserved fungal genes.** Previous AMF studies suggested the loss of several categories of genes by symbiosis with plant<sup>12,13,17</sup>. Our RIR17 genome assembly confirmed the loss of genes involved in the degradation of plant cell walls such as cellobiohydrolases (GH6 and GH7 in the CAZy database), polysaccharide lyases (PL1 and PL4), proteins with cellulose-binding motif 1 (CBM1), and lytic polysaccharide monoxygenases (Supplementary Data 2 and Supplementary Table 8) and nutritional biosynthetic genes, including fatty acid synthase (FAS) and the thiamine biosynthetic pathway (Supplementary Data 8). Given that fatty acids and thiamine are essential nutrients for fungi<sup>43,44</sup>, *R. irregularis* should take up those essential nutrients from a host plant without digestion of the plant cell wall. Several recent papers have described the transport of lipids from plants to AMF<sup>45–47</sup>.

***R. irregularis* has an exceptionally low rDNA copy number.** The general eukaryotic genome has tens to thousands of rDNA copies<sup>20</sup> (Supplementary Figure 1a). However, the RIR17 genome assembly contained only ten copies of the complete 45S rDNA cluster, which was composed of 18S rRNA, ITS1, 5.8S rRNA, ITS2, and 28S rDNAs (Fig. 2a, Supplementary Data 9). To confirm that no rDNA clusters were overlooked, we also estimated the rDNA copy number based on the read depth of coverage. Mapping the Illumina reads of the genomic sequences (Rir\_DNA\_PE180) onto the selected reference sequences indicated that the coverage depth of the consensus rDNA was 8–11 times deeper than the average coverage depth of the single-copy genes (Fig. 2b, Supplementary Table 9), the number of rDNA copies is approximately 10, and the RIR17 assembly covers almost all of the rDNA copies. This AMF rDNA copy number is the lowest among eukaryotes<sup>48</sup> other than pneumonia-causing *Pneumocystis* (one rDNA)<sup>49</sup> and malaria-causing *Plasmodium* (seven rDNAs)<sup>50</sup>.

This low copy number suggests a unique ribosome synthesis system in AMF. The rDNA copy number has relevance for the efficiency of translation because multiple rDNAs are required to synthesize sufficient rRNA. For instance, an experimental decrease in rDNA copy number in yeast (approximately 150 rDNAs in wild type) resulted in no isolated strain having <20 copies, which is considered the minimum number to allow yeast growth<sup>30</sup>. The doubling time of yeast with 20 rDNA copies (TAK300) was 20% longer than that of the wild type<sup>30</sup>. In DAOM-181602, successive cultivation in an infected state with a plant has been widely observed, suggesting that this exceptionally small rDNA copy number is enough to support growth. The multinucleate feature of AMF would increase the rDNA copy number per cell and thereby perhaps supply enough rRNA to support growth. A similar trend in rDNA reduction is observed in the organellar DNA (e.g., mitochondria and plastids)<sup>51</sup>. Revealing the details of translation in AMF requires a future tracking study of the rRNA production and degradation process in AMF. Elucidation of the mechanism to produce mass rRNAs from a few rDNAs may contribute to the understanding of not only AMF evolution but also other polynuclear cells (e.g., striated muscle and Ulvophyceae green algae) and symbiont-derived organelles.

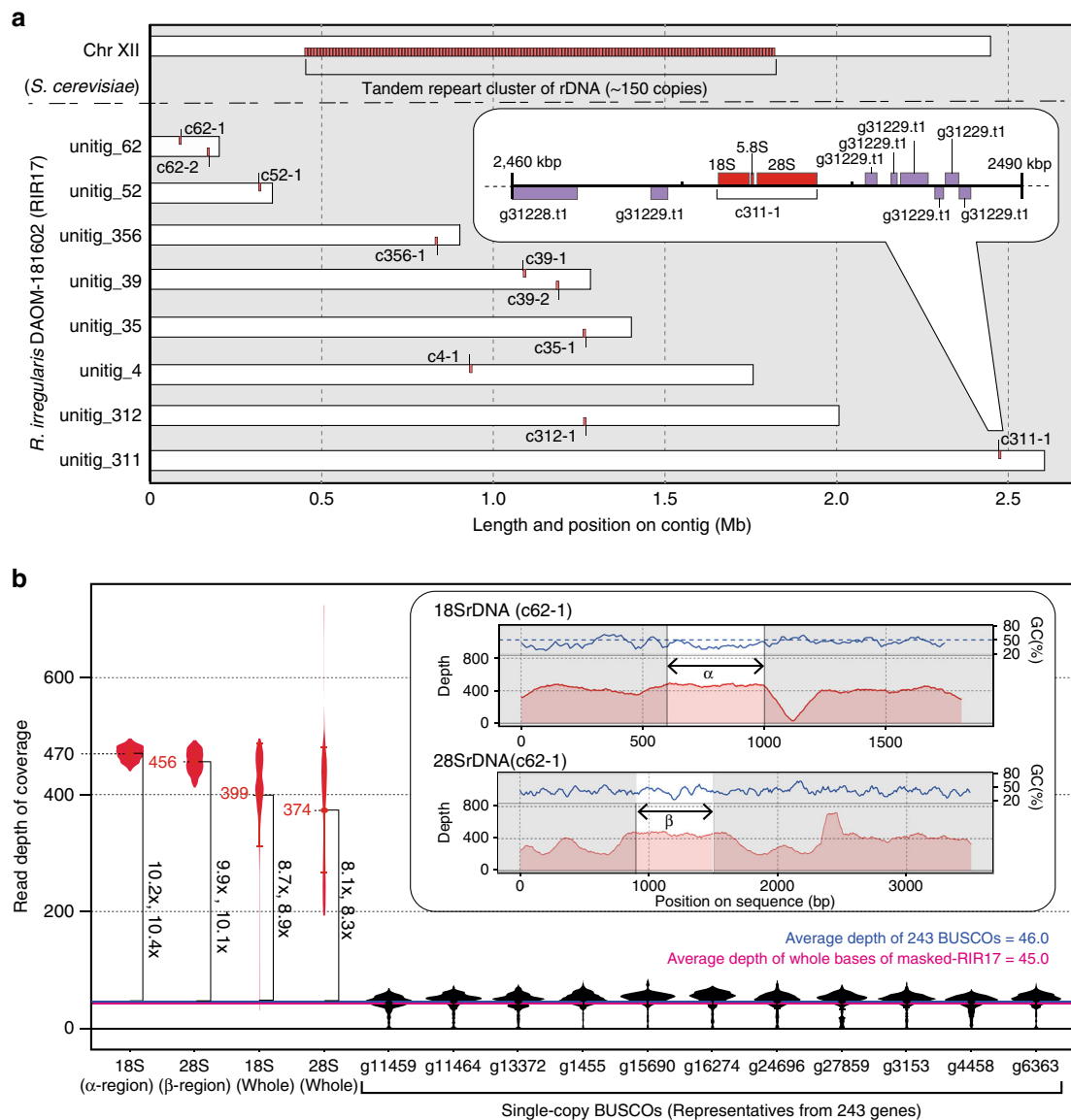
**rDNAs are heterogeneous and completely lack a TRS.** Interestingly, none of the RIR17 rDNAs form a TRS, in contrast to most eukaryotic rDNAs, which comprise tens to hundreds of tandemly repeated units<sup>20</sup>. Most of the rDNA clusters in RIR17 were placed on different contigs; a single copy of rDNA was found in unitig\_311, \_312, \_35, \_356, \_4, and \_52, and two copies were found in unitig\_39 and \_62 (Fig. 2b, Supplementary Data 9). In the cases where two rDNA clusters were found, the two copies resided apart from each other and did not form a tandem repeat;

the distances between the clusters were over 70 kb (76,187 bases in unitig\_62 and 89,986 bases in unitig\_39, Fig. 2b, Supplementary Data 9), the internal regions contained 31 and 42 protein-coding genes, respectively, and the two clusters were located on reverse strands from each other (Fig. 2a, Supplementary Data 9). Since all rDNA copies are located over 28 kb away from the edge of each contig (Fig. 2a, Supplementary Data 9), the lack of TRSs is unlikely to be an artifact derived from an assembly problem often caused by highly repetitive sequences.

The lack of tandem rDNA structure was also supported by mapping our PacBio reads to RIR17 and searching for rDNA on JGI\_v2.0 assemblies. BWA-MEM<sup>52</sup> mapping showed multiple PacBio reads across the 5' non-coding region, 48S rDNA, and 3' non-coding regions of each rDNA contig (Fig. 3a, Supplementary Figure 4). Because our PacBio analysis directly sequenced the DNA molecules in AMF, this syntenic structure is not due to chimeric fragments from DNA amplification but reflects the natural sequence. The 5' and 3' non-coding regions of each rDNA have sequences that are not similar other than the highly similar 5' regions on c62-1 and c62-2 (Fig. 3b and Supplementary Figure 4), negating the possibility of mapping confusion due to sequence similarity. We reproducibly obtained the PacBio reads passing the rDNA regions from our three PacBio datasets, which had been constructed from different spore suspensions. Furthermore, our rDNA searching by RNAmmer detected a non-tandem 48S rDNA region from three JGI\_v2.0 scaffolds (Fig. 3c, Supplementary Figure 5 and Supplementary Table 10). Although the seven rDNAs cannot be reconstructed from JGI\_v2.0, two partial rDNA sequences on JGI\_v2.0 had corresponding down- or upstream sequences that matched our RIR17 rDNAs (Supplementary Figure 5 and Supplementary Table 10), indicating that our assembly around the rDNA genes is consistent with previous assemblies.

We then examined polymorphism among the 45S rDNA clusters on RIR17. rDNA heterogeneity has been reported in various AMF species, including DAOM-181602<sup>13,17,25,29</sup>. However, the distribution and degree of the variation among the rDNA paralogs were unclear. Pairwise comparisons of the ten rDNA copies detected 27.3 indels and 106.1 sequence variants with 98.18% identity on average (Supplementary Data 10 and Supplementary Table 11), whereas the sequences of rDNA clusters at c311-1 and c52-1 were identical. Polymorphisms were distributed unevenly throughout the rDNA; percent identities were 99.91% in 18S rDNA, 97.93% in 28S rDNA, 96.65% in 5.8S rDNA, 93.45% in ITS1, and 90.28% in ITS2 (Fig. 4, Supplementary Data 10 and Supplementary Table 11). The number of polymorphic sites in *R. irregularis* rDNAs reached 4.07 positions per 100 bases, much higher than in other fungi, which have polymorphic sites at 0.04–1.97 positions per 100 bases (Table 2). The rDNA polymorphisms observed in RIR17 covered most of the polymorphisms previously reported in this species (Fig. 5), providing incentive to review the molecular ecology of AMF. The degree of intragenomic variation was not high enough to disrupt species-level identification (i.e., all rDNA genotypes from a genome are jointed to the *R. irregularis* clade), but was sufficient to cause erroneous identification of *R. irregularis* strains (Fig. 5). These findings pose a caution that previous studies on geographic distribution<sup>8</sup>, species identification<sup>28</sup>, and evolutionary processes of AMF assuming rDNA homogeneity require reevaluation, considering the high-level intra-genomic heterogeneity of rDNA sequences in AMFs. Recently developed population genomic approach with ddRAD-Seq<sup>53</sup> may also help us understand the biodiversity of AMF.

**A model for the relaxation of rDNA homogeneity.** The revealed non-tandem structure of AMF rDNA led to a model for the

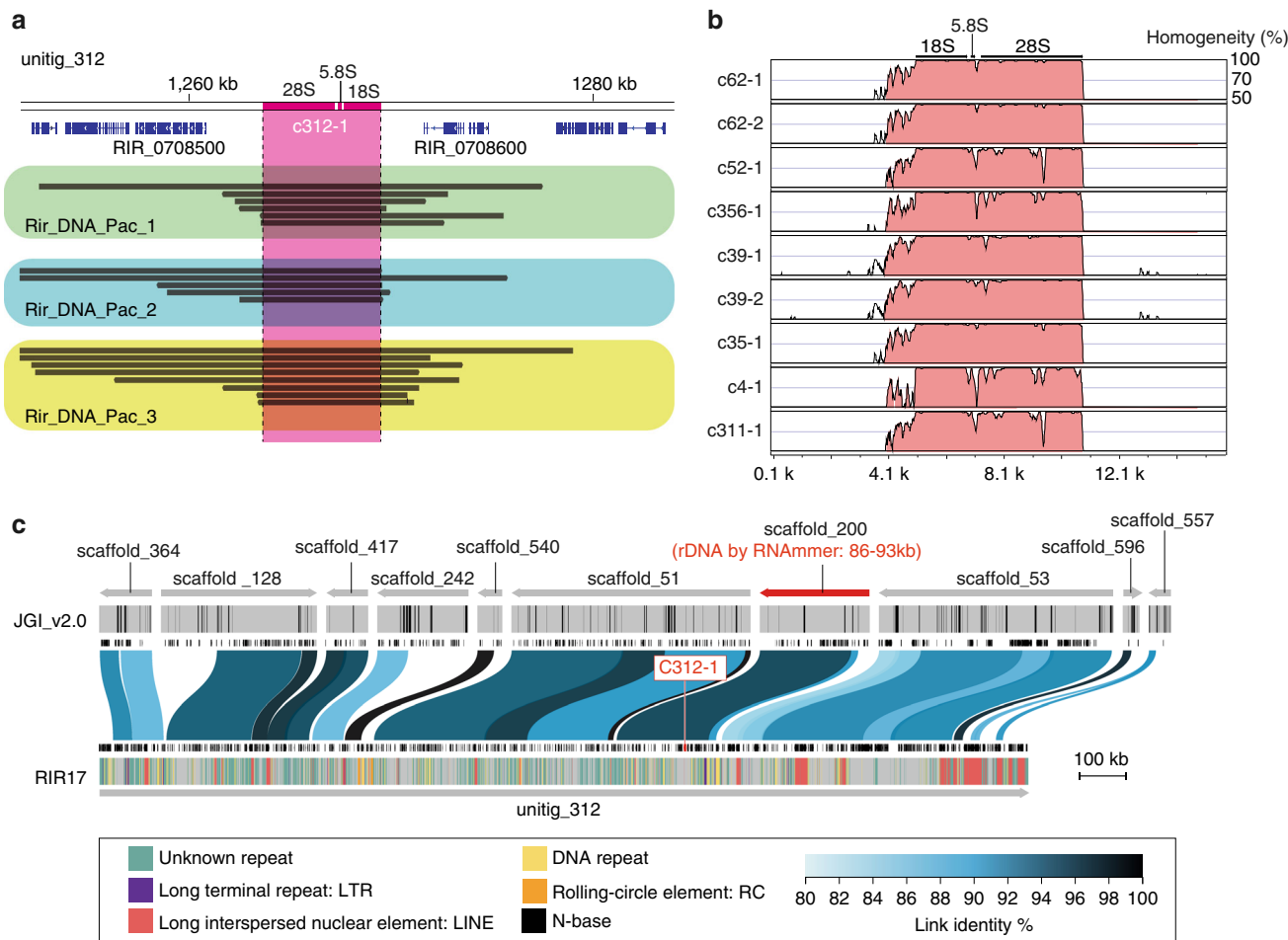


**Fig. 2** Physical maps of rDNA structures and copy numbers in RIR17. **a** Distribution of *R. irregularis* rDNA units in the genome. Each 48S rDNA unit is represented as a red box. For comparison, rDNA clusters on *Saccharomyces cerevisiae* chromosome XII are shown<sup>90,91</sup>. Inset is a magnified view of a 48S rDNA unit (c311-1) with nearby protein-encoding genes (purple boxes). Genes encoded by the plus-strand genome are depicted on the top side, and those encoded by the minus strand are shown on the bottom side. **b** Copy number of rDNA in DAOM-181602 based on the read depth of coverage. Averages of the read depth of coverage are represented as dots and with italic labels. Error bars and violin plots show standard deviations and normalized coverage distribution. The depths of rDNA regions are marked in red. For comparison, the data from representative single-copy BUSCO genes on RIR17 are shown in black. The mean depth of means from 243 BUSCOs is marked with a horizontal blue line, and the mean depth of all RIR17 bases is marked with a magenta line. The changes in the depth of rDNA regions are in vertical bold labels and square brackets. rDNA regions adapted for the copy number estimation ( $\alpha$ - and  $\beta$ -regions) are marked in the inset with the depth of coverage and the GC content of each sequence position

mechanism responsible for its intragenomic heterogeneity. Kuhn et al.<sup>54</sup> and Pawlowska and Taylor<sup>25</sup> predicted that rDNA heterogeneity is caused by relaxation of concerted rDNA evolution in Glomerales including *Rhizophagus*. However, details of the relaxation have been unclear. Here, we propose a hypothetical mechanism: the loss of TRSs precludes the presence of DNA conformations associated with rDNA amplification and the maintenance of its homogeneity. The standard model of concerted rDNA evolution needs two or more tandemly repeated rDNA segments because the rDNA duplicates using tandemly repeated rDNAs as binding sites and templates for replication (Supplementary Figure 1c)<sup>55</sup>. Although non-tandem rDNAs are rare in eukaryotes, this trend of heterogeneity in non-tandem rDNAs has been detected; *Arabidopsis*

*thaliana* has one pseudogenic rDNA (lacking 270 bases of an important helix as rRNA) besides the main tandem repeat rDNA arrays<sup>56,57</sup>, and the lack of rDNA tandem repeats in malaria-causing *Plasmodium* parasites<sup>50,58</sup> indicates intragenomic rDNA polymorphisms. These observations support our hypothesis that rDNA heterogeneity in AMF is related to their lack of TRSs. AMF species may not amplify their rDNA by the general eukaryotic rDNA amplification system (USCR), which may increase their rDNA heterogeneity.

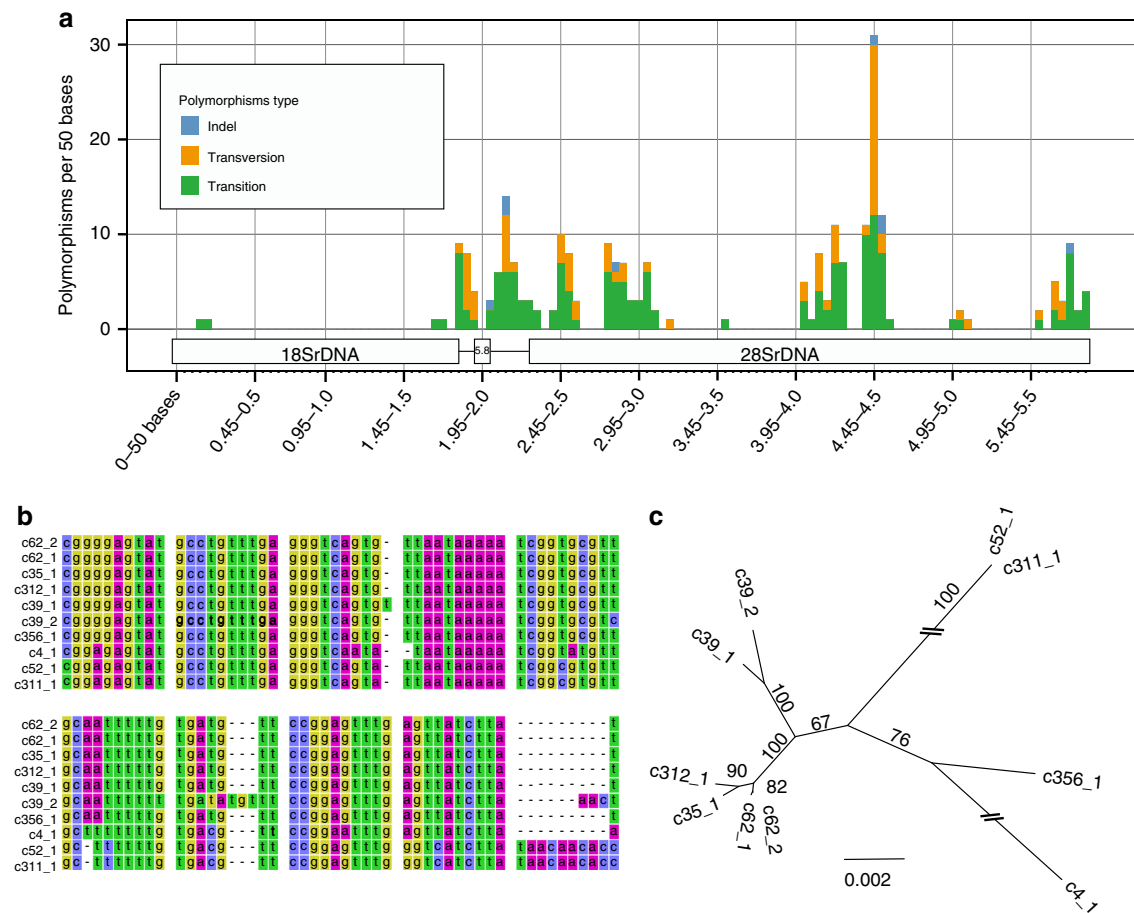
On the other hand, our phylogenetic analysis suggests that AMF has a system to maintain weak similarity among the paralogs without TRSs. Previously observed rDNA heterogeneity in Glomerales suggests that concerted evolution was relaxed



**Fig. 3** Evidence for the lacking of tandem repeat structures of rDNA. **a** Mapped PacBio read for the rDNA regions on unitig\_312 contig in RIR17. The top bar and tick marks indicate sequence positions on the contig. The rDNA region (c312-1) is indicated in magenta. Blue boxes show the predicted protein-coding genes. The mapped read indicates black bar, and reads from different DNA samples and libraries (Supplementary Table 2) are boxed with green, light-blue, and yellow colors, in each. Mapped reads for the other rDNA regions are summarized in Supplementary Figure 4. **b** Sequence similarity of c312-1 rDNAs with other rDNA regions on RIR17. The 5 kb upstream and downstream sequences of each rDNA region are separated from each contig. Alignment and similarity were calculated with mVISTA<sup>92</sup>. Red color shows the sequence regions with similarity over the threshold (>70% similarity for 100 b). **c** Positions and identities of JGI\_v2.0 scaffold aligned against unitig\_312 contigs of RIR17. Scale bar represents the length of sequence assemblies. The top area indicates aligned scaffolds and their strands. A scaffold containing the predicted rDNA gene is marked in red. The positions of N-base on JGI\_v2.0 are marked with black bars in the next line. Predicted protein-coding genes from Chen et al.<sup>14</sup> are indicated with the next black boxes. Aligned positions and their similarity are marked with blue or black bands on the next line. The area below the black boxes shows the predicted genes in the present study. Repetitive regions are marked with colored lines on the bottom band. Types of repetitive elements and the legend of similarity coloration are indicated in the bottom box

before the diversification of *Rhizophagus* species<sup>25,29</sup>. When the observed ten rDNAs duplicated before speciation and evolved independently, each of the duplicated genes formed a clade with orthologs in other species. However, we found no orthologous rDNA genes from other *Rhizophagus* species (Fig. 5). Our tree suggests that the observed rDNAs in *R. irregularis* either expanded or were assimilated after speciation. One hypothetical mechanism that would cause this similarity is homologous recombination via synthesis-dependent strand annealing (Supplementary Figure 6)<sup>59</sup>. This conserved system to repair double-strand breaks (DSBs) results in non-crossover recombination and gene conversion wherein nonreciprocal genetic transfer occurs between two homologous sequences (Supplementary Figure 6). Decreases in divergence by gene conversion are widely observed in duplicated genes. RIR17 showed that two rDNA pairs on the same contigs (c39-1 and c39-2, c62-1 and c62-2) had higher similarity than other paralogs (Fig. 4c). This similarity may be caused by the high gene conversion rate between these loci.

Our model raises a new question about the mechanism that maintains the number of rDNAs without gene duplication by USCR. Even if rDNA lacks TRSSs, crossover recombination and single-strand annealing delete paralogous genes. Observed inverted repeat structures between rDNAs in proximity may contribute to inhibiting single-strand annealing between them and prevent copy number reduction. Plastidial rDNAs of land plants also make inverted repeat structures and conserve two rDNA copies on their plastidial DNA. Another probable system is the suppression of crossover recombination. When Holliday junctions dissociate without crossover, DSBs are repaired without gene number reduction. AMF species may keep their rDNA copy number by their highly controlled Holliday junction dissociation. Although the detail of AMF reproduction system and its contribution to the recombination still have many unresolved issues<sup>60</sup>, it should be noted that the majority of these crossovers arise during meiosis in eukaryotes<sup>59</sup> and *R. irregularis* can asexually make the spore without meiotic stage.



**Fig. 4** Polymorphisms of 48S rDNA paralogs in RIR17. **a** The distribution of rDNA sequence variants within the 48S rDNA of RIR17. The position and types of polymorphisms were called based on the paralog c62-1. **b** Alignment of a heterogeneous region among the 48S rDNA paralogs. Partial sequences of MAFFT-aligned 48S rDNAs (corresponding to 2049–2136 base positions on c62-1). **c** Neighbor-joining tree for phylogenetic relationships among the ten rDNA paralogs based on 5847 aligned positions. Bootstrap values are described at each node. Scale bar represents the branch length (substitution per site)

**RNA-level impact and probable fitness of TRS-lacking.** To confirm the transcriptional activity of each rDNA, we conducted total-RNA-Seq (RNA sequencing without poly-A tail selection; see Methods section). Illumina sequencing of a modified library for rRNA sequencing (Rir\_RNA\_rRNA in Supplementary Table 2) produced 18,889,290 reads (read length = 100–301 bases) from DAOM-181602. We mapped the reads to all gene models from RIR17 (43,675 protein-encoding isoforms and ten 48S rDNA paralogs) and estimated the expression levels of each gene by eXpress software<sup>61</sup>. All rDNA paralogs had over 5000 Fragments Per Kilobase of exon per Million mapped fragments (FPKM) (Table 3), and multiple reads were matched to the specific region of each paralog, indicating that the ten rDNA copies are transcriptionally active. In general, eukaryotes silence a part of the rDNA copies<sup>62</sup>, and some eukaryotes change the transcribed rRNA sequences by RNA editing<sup>63</sup>. These editing and silencing processes were not detected in the AMF, and the rRNA was as polymorphic as the rDNA. These results show that DAOM-181602 has multiple types of ribosomes, each containing different rRNAs. Additionally, we detected highly duplicated ribosomal protein genes (e.g., ribosomal protein S17/S11) (Supplementary Tables 6 and Supplementary Data 5) and tRNA genes, indicating unknown amino acid isotypes, which may also account for the heterogeneity of ribosomes (Supplementary Table 12).

The evolutionary significance of the non-tandemly repeated heterogeneous rDNAs is unclear. One of the probable factors is a

reduction in the need to maintain numerous rDNAs in a genome. As described in the above sections, the AMF rDNA copy number suggests a system that efficiently produces rRNA from a few rDNAs, and the inverted repeats structure of rDNAs will also reduce the deletion rate of rDNAs. AMF may thus no longer need to rapidly amplify rDNA copies using TRSs, and the slowed replacement rate of rDNA may then cause the heterogeneity as a side effect. Another possibly adaptive effect is the enhancement of phenotypic plasticity by ribosomal heterogeneity (Fig. 6). Recent studies have started to reveal that various eukaryotes (e.g., yeast, mice, and *Arabidopsis*) produce heterogeneous ribosomes and subsequently alter phenotypes via proteins translated by particular ribosomes<sup>64</sup>. Accelerated accumulation of AMF rDNA mutations by the lack of TRSs may lead to functional variety in produced ribosomes and increases in the rate of adaptation by different translation activities within the same species. Previous studies have reported that large variations in the fungal phenotype were observed among single spore lines derived from one parent AMF<sup>65–67</sup>. This might be not only due to genetic variation but also due to variations in each rDNA expression. Although the functional effects of observed rDNA mutations remain to be determined, the middle area of our 28S rDNA (4450–4500 bases on c62-1) had a higher mutation rate than ITS regions (Fig. 4a). Because the ITS regions (encoding non-functional RNA) vary under neutral mutation rates, the accumulated variants in the middle-28S region may have



**Table 2 Numbers of intragenomic polymorphic sites in fungal rDNAs**

Species	# Polymorphic sites	Repeat unit length (bp)	# of units in genome	# of polymorphic sites/100 bases
<i>Rhizoglyphus irregularis</i>	238	5847	10	4.07
<i>Rhizoglyphus irregularis</i> <sup>13</sup>	38	1563	-	2.43
<i>Ashbya gossypii</i> <sup>30</sup>	3	8147	50	0.04
<i>Saccharomyces paradoxus</i> <sup>30</sup>	13	9103	180	0.14
<i>Saccharomyces cerevisiae</i> <sup>30</sup>	4	9081	150	0.04
<i>Aspergillus nidulans</i> <sup>30</sup>	11	7651	45	0.14
<i>Cryptococcus neoformans</i> <sup>30</sup>	37	8082	55	0.46
<i>Phoma exigua</i> var. <i>exigua</i> <sup>93</sup>	27	1672	-	1.61
<i>Mycosphaerella punctiformis</i> <sup>93</sup>	26	1669	-	1.56
<i>Teratosphaeria microspora</i> <sup>93</sup>	16	1671	-	0.96
<i>Davidiella tassiana</i> <sup>93</sup>	33	1672	-	1.97

functional effects favored by natural selection (via diversifying selection). This region is thus a useful target for the future functional analyses of AMF rRNA.

AMF species are similar to the malaria parasite in that they both have heterogeneous non-tandem rDNAs and infect distantly related host species<sup>50</sup>. In the malaria parasite, changes in the ribosome properties depend on the host (human or mosquito), which is likely able to alter the rate of translation, either globally or of specific messenger RNAs, thereby changing the rate of cell growth or altering patterns of cell development<sup>50</sup>. The relationship between the diversity of host organisms and rDNA polymorphisms will be an important area for further research. The phenotypic plasticity caused by heterogeneous translation machinery may allow adaptation for various host species having slightly different symbiotic systems. Previous studies have proposed that the heterokaryosity in AMF species drives variable genetic combinations of mycelia<sup>68</sup>. Recent genomic studies, furthermore, discovered signatures of sexual reproduction within the dikaryon-like stage<sup>16,69</sup>. Our hypothesis does not exclude current theories for the genetic and phenotypic plasticity of AMF species (heterokaryosis and sexual reproduction) but proposes a multilayered diversification mechanism leading to their widespread distribution.

## Conclusion

We here reported an improved genome assembly of *R. irregularis* DAOM-181602. Improved genome revealed that common concepts of eukaryotic rDNA are not applicable to AMF. Its rDNA copies are highly heterogeneous, reduced in number, and lack TRS. This frequently used ecological and phylogenetic marker gene should be adopted cautiously for AMF. The sequence diversity and reduced copy number of rDNA may result from the collapse of the concerted evolution system due to the lack of TRS. Although the adaptive significance of the TRS lacking in AMF remains to be determined, future investigations on the functional differences among the heterogeneous rRNAs may reveal mechanisms that facilitate adaptation to various environmental conditions in AMF.

## Methods

**PacBio-based assembling.** DNA preparation: The DNA sample for the PacBio and Illumina sequencing was extracted from a commercial strain of *R. irregularis* DAOM-181602 (MYCORISE® ASP, Premier Tech Biotechnologies, Canada). The DNA extraction followed the method of Fulton et al.<sup>70</sup> with some modifications described below. Purchased spore suspensions (including approximately 1,000,000 spores) were centrifuged (4500 rpm, 20 min), and washed three times with distilled water. Precipitated spores were frozen with liquid nitrogen, ground with pestle, and dispersed in extraction buffer (100 mM Tris-HCl pH 8.0, 20 mM EDTA, 0.75% Sarkosyl, 0.1% PVP, 0.75% cetyl trimethylammonium bromide (CTAB), 0.13 M sorbitol, 0.75 M NaCl, and 0.1 mg ml<sup>-1</sup> proteinase K). After incubation at 37 °C for 10 min, the aqueous phase was centrifuged (15,000 rpm, 4 min), and the pellet was discarded. An equal volume of phenol/chloroform (1:1,

vol:vol) was added, and the sample was gently mixed and centrifuged (15,000 rpm, 2 min). The aqueous phase was collected, and an equal volume of chloroform was added to the sample, which was then mixed and centrifuged (15,000 rpm, 2 min). The aqueous phase was collected again, and 1:10 vol of sodium acetate and 0.7 vol of isopropanol were added. The sample was then mixed and centrifuged (12,000 rpm, 20 min). The resulting pellet was washed twice with 70% EtOH and eluted with TE buffer. Extracted DNA was purified with Genomic-tip (Qiagen, Germany) following the manufacturer's instructions.

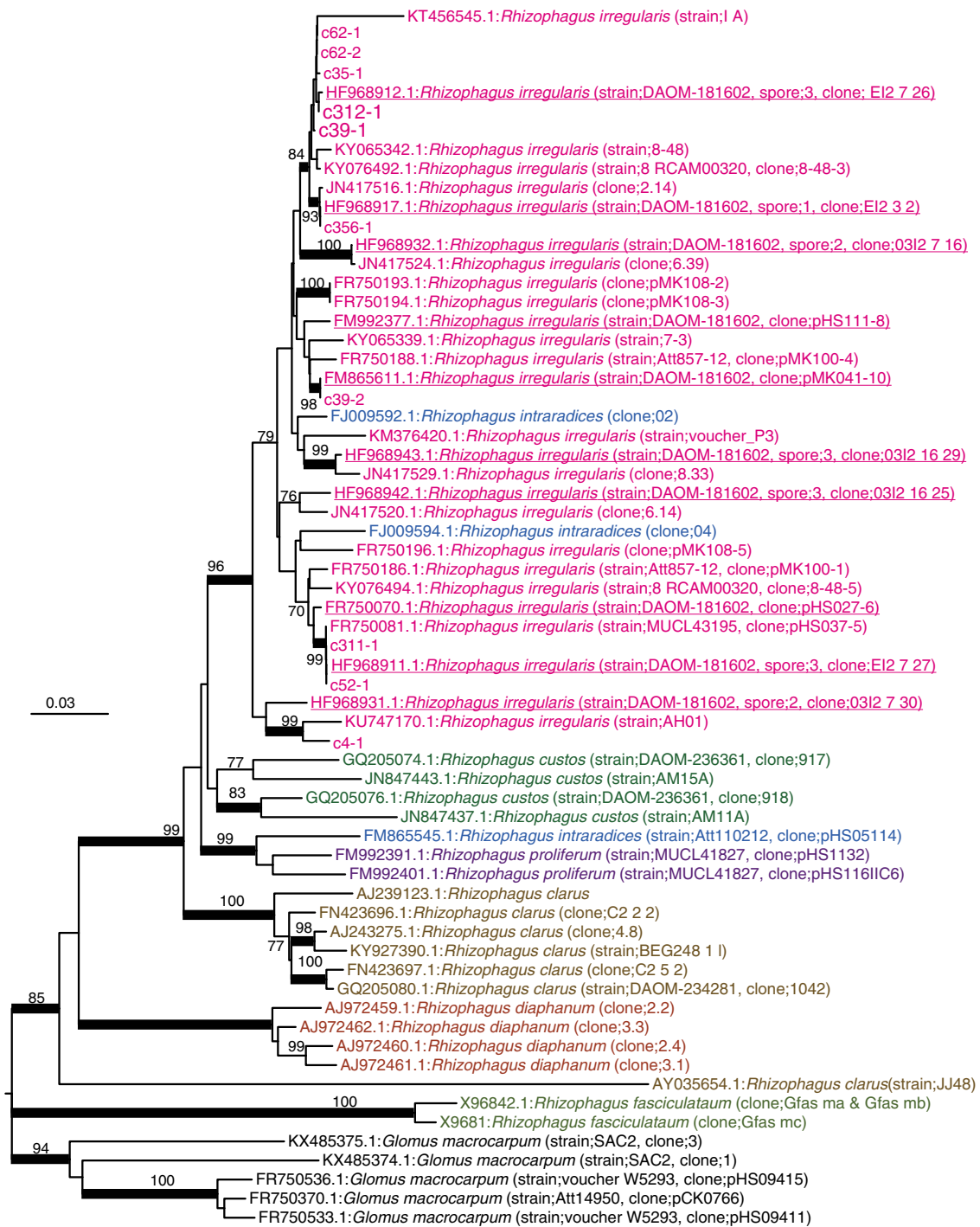
**PacBio sequencing:** Long-read sequences were generated with a PacBio RS II sequencer (Pacific Biosciences, Menlo Park, CA, USA) using a DNA/Polymerase Binding Kit P6 v2 (Pacific Biosciences) and a DNA Sequencing Reagent Kit 4.0 (Pacific Biosciences). The library was prepared according to the 20-kb Template Preparation Using BluePippin™ Size-Selection System (Sage Science, MA, USA). A total sequence of 11.7 Gb in 955,841 reads (76× coverage of the genome, assuming a genome size of 154 Mb) was obtained from 29 SMRT cells (Supplementary Table 2). The N50 length of the raw reads was 13,107 bases.

**PacBio-based genome assembly:** The *R. irregularis* genome was assembled using the RS\_HGAP\_Assembly.3 protocol for assembly and Quiver for genome polishing in SMRT Analysis v2.3.0 (Pacific Biosciences). The procedure consisted of three parts, involving (1) generation of preassembled reads with improved consensus accuracy; (2) assembly of the genome through overlap consensus accuracy using Celera Assembler; and (3) one round of genome polishing with Quiver. For HGAP, the following parameters were used: PreAssembler Filter v1 (minimum subread length = 500 bases, minimum polymerase read quality = 0.80, minimum polymerase read length = 100 bases); PreAssembler v2 (minimum seed length = 6000 bases, number of seed read chunks = 6, alignment candidates per chunk = 10, total alignment candidates = 24, minimum coverage for correction = 6, and BLASR options = "noSplitSubreads, minReadLength = 200, maxScore = 1000, and maxLCPLength = 16"); AssembleUnitig v1 (genome size = 150 Mb, target coverage = 25, overlapper error rate = 0.06, overlapper min length = 40 bases and overlapper k-mer = 14); and BLASR v1 mapping of reads for genome polishing with Quiver (maximum divergence = 30, minimum anchor size = 12). Assembly polishing with PacBio reads was carried out with Quiver using only unambiguously mapped reads. The statistics of the PacBio-only assembly set and previously sequenced data (Lin14, JGI\_v1.0, JGI\_v2.0) were evaluated using QUAST ver. 4.3<sup>71</sup>. The percentage of genome coverage was estimated assuming the genome size to be 154 Mb based on Tisserant et al.<sup>17</sup>.

**Error correction and identification of host plant contamination.** After polishing using Illumina data, we eliminated the sequences derived from contaminated DNAs during the sample preparation. BLASTn search of the polished assemblies against the refseq\_genomic database detected nine assemblies showing similarity with sequences from carrot (BLAST ver. 2.2.31+, query coverage per subject >95%, percentages of identical matches >90%, bit score >1000) (Supplementary Table 2), which might be used as a host plant by the manufacturer for the cultivation of *R. irregularis* samples. After elimination of the nine contaminated contigs, we submitted the assemblies to the DDBJ as whole-genome shotgun sequence data (RIR17) of *R. irregularis* DAOM-181602 (BDIQ01).

**Genomic alignment with previous genome assemblies.** The quality of our genome assembly was evaluated by alignment with previously available *R. irregularis* DAOM-181602 genome assemblies. A one-by-one genome alignment was constructed by MUMmer ver. 4.0.0beta2<sup>72</sup> between RIR17, JGI\_v2.0, Lin14, and JGI\_v1.0 assemblies. Each genome set was aligned by the nucmer function in MUMmer, and the statistics of the alignments were extracted by the dnadiff wrapper with the default setting.

**Gene prediction and annotation.** De novo repeat motifs were identified using RepeatModeler ver. 1.0.8, which combines RECON and RepeatScout programs<sup>37</sup>. Based on the identified motif, the repetitive region in the assemblies was masked



**Fig. 5** NJ tree based on 586 positions of 48S rDNA. The numbers at the nodes and the scale bar have the same meanings as in Fig. 4c. Partial 18S, ITS1, 5.8S, ITS2, and partial 28S rDNAs were used. The ten rDNA paralogs from RIR17 and 58 *Rhizophagus* sequences from the DDBJ were chosen as operational taxonomic units (OTUs). The 58 *Rhizophagus* sequences were selected from 329 OTUs in the DDBJ (209 OTUs for DAOM-181602, 57 OTUs for other *R. irregularis* strains, and 63 OTUs for other *Rhizophagus* species) using CD-Hit clustering ( $-c$  0.98  $-n$  5). Five *Glomus* sequences were used as outgroup OTUs. Red underlined OTUs are sequences from *R. irregularis* DAOM-181602, and other red OTUs are data from other strains of *R. irregularis*. Nodes supported by over 80 bootstrap values are marked by a bold line. All *R. irregularis* OTUs made a single clade with *Rhizophagus intraradices* that is a morphologically non-distinct sister group of *R. irregularis*

with RepeatMasker ver. 4.0.5<sup>37</sup>. We used the default parameters for the identification and the masking.

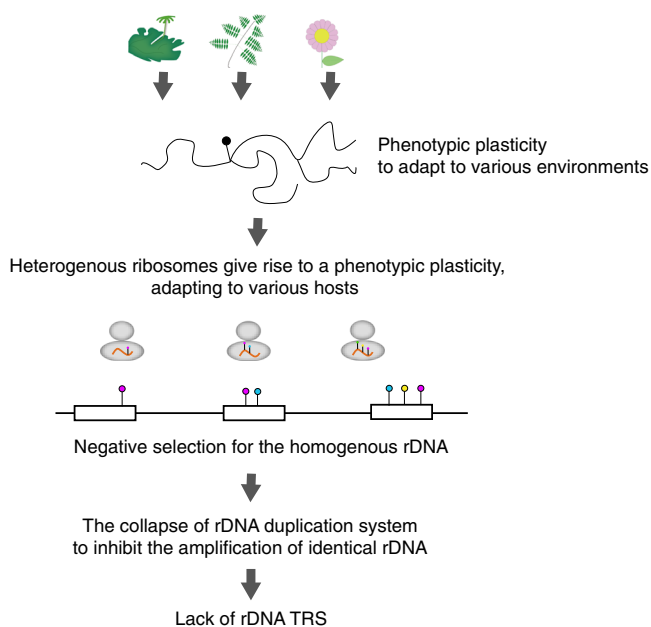
For the gene models constructed from RIR17 assemblies, standard RNA-Seq data were obtained from *R. irregularis* spores and hyphae. The RNA was extracted with an RNeasy Plant Mini kit (Qiagen) after incubation of the purchased spores (MYCORISE® ASP) in a minimum nutrient medium for 1 day. An Illumina RNA-Seq library was constructed with a TruSeq Stranded mRNA Library prep kit (Illumina). The library was sequenced (101 bases from each end) on a HiSeq 1500

platform (Illumina). A total of 16,122,964 raw reads (3.2 Gb) were obtained from the library (Supplementary Table 2). After filtering low-quality and adapter sequences, RNA-Seq data were mapped to RIR17 assemblies with TopHat ver. 2.1.1<sup>73</sup> with the default setting.

Then, the RIR17 assemblies were processed through the RNA-Seq-based gene model construction pipeline using AUGUSTUS ver. 3.2.1 software<sup>74</sup>. We constructed *R. irregularis*-specific probabilistic models of the gene structure based on 495 manually constructed gene models from the longest unigntg\_392 sequence in

**Table 3** Transcription activity of the rDNA paralogs

Target ID	FPKM	Confidence interval (95%)	
		Low	High
c312_1	28,888	28,672	29,103
c39_1	20,719	20,537	20,901
c39_2	20,358	20,177	20,538
c62_2	19,431	19,254	19,608
c4_1	19,430	19,255	19,605
c52_1	19,054	18,879	19,228
c311_1	19,054	18,879	19,228
c356_1	16,151	15,990	16,311
c35_1	10,053	9927	10,180
c62_1	7656	7546	7766



**Fig. 6** Hypothetical model for the evolution of unique rDNAs/rRNAs in AMF. Evolutionary model for the lack of TRSs in AMF and its sequence heterogeneity. The various environmental conditions (e.g., various host species) may lead to the evolution of phenotypic plasticity via multiple types of ribosomes in AMF. If the rDNA is exposed to disruptive selection, rDNA duplication by TRSs and USCR may be non-adaptive because the duplication of particular rDNA types reduces the variety of rDNA types

RIR17. Manual gene models were made with ab initio AUGUSTUS analysis based on probabilistic models for *Rhizopus oryzae* and by manual refinement using the homology data with already-known genes and mapped RNA-Seq data. Then, with the trained probabilistic models and the intron-hints data from the mapped RNA-Seq read, 37,639 optimal gene models were constructed using the AUGUSTUS pipeline. We then confirmed whether the AUGUSTUS pipeline overlooked the called genes in previous genome studies. We mapped all transcript sequences obtained from previous gene modeling on Lin14 and JGI\_v1.0 against our RIR17 genomic sequences with Exonerate<sup>75</sup> (ver. 2.2.0, option --model est2genome --bestn 1), resulting in the recruitment of 3933 overlooked genes. The completeness of the constructed gene model was evaluated with BUSCO ver. 2.0<sup>39</sup>. The BUSCO analysis used the Fungi odb9 gene set (<http://buscocodev.ezlab.org/datasets/fungiodb9.tar.gz>) as a benchmark and employed the -m proteins option to analyze the preconstructed protein data without the ab initio gene modeling step.

The confidences of the obtained 41,572 gene models were estimated based on (1) RNA-Seq expression support, (2) homology evidence, and (3) protein motif evidence. For the calculation of gene expression levels, we mapped our Rir\_RNA\_SS data and 32 RNA-Seq data submitted to the sequence read archive (SRA) database (24 data sets from DRP002784 and eight data sets from DRP003319) and calculated the gene expression levels (FPKM) using FeatureCounts<sup>76</sup> with the default setting (Supplementary Data 2). Homology with previously known genes was determined by BLAST searches against the orthoDB

(odb9) (Supplementary Tables 6 and Supplementary Data 5). The protein motif was searched using Pfam analysis in InterProScan ver. 5.23-62.0<sup>77</sup> (Supplementary Data 2).

Constructed gene models were annotated by several in-silico searches. Gene functions were predicted based on BLASTp (Database = nr, RefSeq, and UniProt), and Pfam in InterProScan (Supplementary Data 2). We manually selected the descriptive nomenclatures from those four searches and submitted to the DDBJ. Orthologous relationships were classified with Orthofinder (ver. 1.1.2)<sup>78</sup>, and rapidly expanded/contracted families were analyzed with CAFE (ver. 4.1)<sup>79</sup> from Orthofinder results. Phylogenetic trees for the CAFE analysis were constructed with IQ-tree (ver. 1.6.1)<sup>80</sup> for maximum likelihood (ML) analysis and r8s (v1.81) for a conversion for an ultrametric tree. An ML tree was made from 159 single-copy genes from the Orthofinder results (Supplementary Data 2) and was converted to an ultrametric tree based on the divergence times of AMF-Mortierellales (460 Myr)<sup>27</sup> and Deuterostomia-Protostomia (550 Myr)<sup>81</sup>. Overlapping genes with TEs were extracted from AUGUSTUS and RepeatMasker results using bedtools (ver. 2.26.0, bedtools intersect with -wa option)<sup>82</sup>.

The missing ascomycete core gene (MACG) orthologs were sought using BLAST with the -evalue 0.0001 option, and the reference sequences for the MACG search were selected from protein data from an S288C reference in the *Saccharomyces* genome data base (SGD) (Supplementary Data 8). Genes involved in the degradation of plant cell walls were sought by BLAST with the same settings as the MACG search, and the reference sequences were selected from *Aspergillus niger* CBS 513.88 data in GenBank based on CAZY classification (Supplementary Data 8). Other gene annotations based on the CAZY database were performed with the dbCAN HMMs 6.0 web service<sup>83</sup> (Supplementary Data 2).

**Detection of ribosomal DNA and intragenomic polymorphisms.** Ribosomal DNA regions were detected by RNAmmer ver. 1.2<sup>84</sup> from whole RIR17 assemblies and were manually refined based on the MAFFT v7.294b<sup>85</sup> alignment to the 48S rRNA in *Saccharomyces cerevisiae* S288C. The genomic positions of rDNAs were visualized with Python ver. 3.4.0 (BasicChromosome ver. 1.68, and GenomeDiagram ver. 0.2 modules) (Fig. 2a).

The number of rDNA paralogs in the genome was estimated by mean depth of coverage. We masked repetitive regions (based on RepeatModeler analysis) and all rDNA regions on RIR17 except one rDNA copy (c62-1). Then, trimmed R1 Illumina reads from Rir\_DNA\_PE180 library were mapped to the repeat-masked RIR17 using bowtie2 ver. 2.2.9<sup>86</sup>. The coverage depth of the rDNA region and 243 single-copy BUSCOs were obtained using bedtools (bedtools coverage command with -d option), and the statistics of each region were calculated and visualized by R software ver. 3.4.2 with the ggplot2 library (Fig. 2b). To prevent copy number estimation from depth fluctuation due to the intragenomic heterogeneity, we confirmed the coverage depth using the consensus sequences of all ten rDNA paralogs; the joined Illumina reads (from Rir\_DNA\_PE180 library) were mapped back to a consensus rDNA sequences and ten single-copy BUSCO genes from RIR17, and the depth of coverages was then counted by bedtools (genomeCoverageBed) (Supplementary Table 9).

The syntenic structure around rDNA genes was confirmed by the mapping of PacBio raw reads and comparison with JGI\_v2.0 assemblies. All of the filtered-subreads from SMART Analysis software were mapped to RIR17 assemblies by BWA-MEM (ver. 0.7.15-r1140) with the -x pacbio option. Mapped reads were visualized with Integrative Genomics Viewer (ver. 2.4), and the reads covering the rDNA regions were selected by eye. Alignment between JGI\_v2.0 and RIR17 was done by a combination of MUMmer, LASTz (ver. 1.04.00), and AliTV<sup>87</sup> (ver. 1.0.4) software. JGI\_v2.0 scaffolds having regions corresponding with RIR17 sequences were selected by the nucmer and delta-filter (with -l option) functions in MUMmer. Then, we extracted the JGI\_v2.0 scaffolds corresponding to RIR17 contigs with rDNAs (unitig\_311, \_312, \_35, \_356, \_4, and \_52). Selected scaffolds were aligned to the corresponding RIR17 contigs by alitv.pl scripts (with alignment: lastz and --ambiguous = n settings) and alitv-filter.pl (with --min-link-identity 80 and --min-link-length 10000 option) in the AliTV package and visualized with the AliTV web service (<http://alitvteam.github.io/Alitv/d3/Alitv.html>).

The difference among the rDNA paralogs was calculated from the aligned sequences by MAFFT ver. 7.309 (options: --localpair, --op 5, --ep 3, --maxiterate 1000) using the pairwise comparison with CLC Main Workbench 7.8.1 (Qiagen). The mutation type was called by eye from the alignment, and we chose the c62-1 paralog as a reference sequence for mutation calling (Fig. 4a). Phylogenetic trees (Figs. 4c and 5) were constructed from the MAFFT alignment by the neighbor-joining method with MEGA<sup>88</sup> ver. 7.0.21 under the maximum composite likelihood model and were tested for robustness by bootstrapping (500 pseudoreplicates).

**Heterogeneity of translation machineries.** The expression levels of the rDNA paralogs were examined with modified Illumina sequencing of *R. irregularis* spores and hyphae. Total RNA was extracted with an RNeasy Plant Mini kit (Qiagen) after the incubation of the MYCORISE<sup>®</sup> spores in a minimum nutrient medium for 7 days. An Illumina RNA-Seq library was constructed with a TruSeq Stranded mRNA Library prep kit (Illumina). To skip the poly-A tailing selection step in the

library construction, we started from the fragmentation step of the standard manufacturer's instructions. The library was sequenced (301 bases from each end) on a MiSeq platform (Illumina). A total of 16,122,964 raw reads (3.2 Gb) were obtained from the library (Supplementary Table 2). After filtering low-quality and adapter sequences, RNA-seq data were mapped to the RIR17 assembly with TopHat with the default settings. FPKMs for each gene were calculated with eXpress ver. 1.5.1 with the --no-bias-correct option. Transfer RNAs were identified with tRNAscan-SE<sup>89</sup> ver. 1.3.1.

**Data availability.** Raw reads, genome assemblies, and annotations were deposited at INSDC under the accessions as follows: Sequence read archive: DRA004849, DRA004878, DRA004889, DRA004835, DRA005204, and DRA006039; Whole genome assembly: BDIQ01000001-BDIQ01000210; Annotations: GBC10881-GBC54553. All the other data generated or analyzed during this study are included in this published article and its supplementary information.

Received: 16 November 2017 Accepted: 12 June 2018

Published online: 10 July 2018

## References

- Remy, W., Taylor, T. N., Hass, H. & Kerp, H. Four hundred-million-year-old vesicular arbuscular mycorrhizae. *Proc. Natl. Acad. Sci. U.S.A.* **91**, 11841–11843 (1994).
- Redecker, D., Kodner, R. & Graham, L. E. Glomalean fungi from the Ordovician. *Science* **289**, 1920–1921 (2000).
- Smith, S. E., Jakobsen, I., Grønlund, M. & Andrew Smith, F. Roles of arbuscular mycorrhizas in plant phosphorus nutrition: interactions between pathways of phosphorus uptake in arbuscular mycorrhizal roots have important implications for understanding and manipulating plant phosphorus acquisition. *Plant Physiol.* **156**, 1050–1057 (2011).
- Bougoure, J., Ludwig, M., Brundrett, M. & Grierson, P. Identity and specificity of the fungi forming mycorrhizas with the rare mycoheterotrophic orchid *Rhizanthella gardneri*. *Mycol. Res.* **113**, 1097–1106 (2009).
- Parniske, M. Arbuscular mycorrhiza: the mother of plant root endosymbioses. *Nat. Rev. Microbiol.* **6**, 763–775 (2008).
- van der Heijden, M. G. A. et al. Mycorrhizal fungal diversity determines plant biodiversity, ecosystem variability and productivity. *Nature* **396**, 69–72 (1998).
- Johnson, N. C., Graham, J. H. & Smith, F. A. Functioning of mycorrhizal associations along the mutualism–parasitism continuum. *New Phytol.* **135**, 575–586 (1997).
- Davison, J. et al. Global assessment of arbuscular mycorrhizal fungus diversity reveals very low endemism. *Science* **349**, 970–973 (2015).
- Fellbaum, C. R. et al. Fungal nutrient allocation in common mycorrhizal networks is regulated by the carbon source strength of individual host plants. *New Phytol.* **203**, 646–656 (2014).
- Lee, J. The distribution of cytoplasm and nuclei within the extra-radical mycelia in *Glomus intraradices*, a species of arbuscular mycorrhizal fungi. *Mycobiology* **39**, 79–84 (2011).
- Zhang, Y. & Guo, L. D. Arbuscular mycorrhizal structure and fungi associated with mosses. *Mycorrhiza* **17**, 319–325 (2007).
- Tang, N. et al. A survey of the gene repertoire of *Gigaspora rosea* unravels conserved features among Glomeromycota for obligate biotrophy. *Front. Microbiol.* **7**, 233 (2016).
- Lin, K. et al. Single nucleus genome sequencing reveals high similarity among nuclei of an endomycorrhizal fungus. *PLoS Genet.* **10**, e1004078 (2014).
- Chen, E. C. H. et al. High intraspecific genome diversity in the model arbuscular mycorrhizal symbiont *Rhizophagus irregularis*. *New Phytol.* <https://doi.org/10.1111/nph.14989> (2018).
- Toro, K. S. & Brachmann, A. The effector candidate repertoire of the arbuscular mycorrhizal fungus *Rhizophagus clarus*. *BMC Genomics* **17**, 101 (2016).
- Ropars, J. et al. Evidence for the sexual origin of heterokaryosis in arbuscular mycorrhizal fungi. *Nat. Microbiol.* **1**, 16033 (2016).
- Tisserant, E. et al. Genome of an arbuscular mycorrhizal fungus provides insight into the oldest plant symbiosis. *Proc. Natl. Acad. Sci. U.S.A.* **110**, 20117–20122 (2013).
- Mohanta, T. K. & Bae, H. The diversity of fungal genome. *Biol. Proced. Online* **17**, 8 (2015).
- Ma, L. J. et al. Genomic analysis of the basal lineage fungus *Rhizopus oryzae* reveals a whole-genome duplication. *PLoS Genet.* **5**, e1000549 (2009).
- Eickbush, T. H. & Eickbush, D. G. Finely orchestrated movements: evolution of the ribosomal RNA genes. *Genetics* **175**, 477–485 (2007).
- Sanders, I. R., Alt, M., Groppe, K., Boller, T. & Wiemken, A. Identification of ribosomal DNA polymorphisms among and within spores of the Glomales—application to studies on the genetic diversity of arbuscular mycorrhizal fungal communities. *New Phytol.* **130**, 419–427 (1995).
- LloydMacgilp, S. A. et al. Diversity of the ribosomal internal transcribed spacers within and among isolates of *Glomus mosseae* and related mycorrhizal fungi. *New Phytol.* **133**, 103–111 (1996).
- Hosny, M., Hijri, M., Passerieux, E. & Duliue, H. rDNA units are highly polymorphic in *Scutellospora castanea* (Glomales, Zygomycetes). *Gene* **226**, 61–71 (1999).
- Hijri, M. & Sanders, I. R. The arbuscular mycorrhizal fungus *Glomus intraradices* is haploid and has a small genome size in the lower limit of eukaryotes. *Fungal Genet. Biol.* **41**, 253–261 (2004).
- Pawlowska, T. E. & Taylor, J. W. Organization of genetic variation in individuals of arbuscular mycorrhizal fungi. *Nature* **427**, 733–737 (2004).
- Rosendahl, S. & Stukenbrock, E. H. Community structure of arbuscular mycorrhizal fungi in undisturbed vegetation revealed by analyses of LSU rDNA sequences. *Mol. Ecol.* **13**, 3179–3186 (2004).
- Schussler, A., Schwarzott, D. & Walker, C. A new fungal phylum, the Glomeromycota: phylogeny and evolution. *Mycol. Res.* **105**, 1413–1421 (2001).
- Krüger, M., Krüger, C., Walker, C., Stockinger, H. & Schüssler, A. Phylogenetic reference data for systematics and phylotaxonomy of arbuscular mycorrhizal fungi from phylum to species level. *New Phytol.* **193**, 970–984 (2012).
- VanKuren, N. W., den Bakker, H. C., Morton, J. B. & Pawlowska, T. E. Ribosomal RNA gene diversity, effective population size, and evolutionary longevity in asexual Glomeromycota. *Evolution* **67**, 207–224 (2013).
- Ganley, A. R. D. & Kobayashi, T. Highly efficient concerted evolution in the ribosomal DNA repeats: total rDNA repeat variation revealed by whole-genome shotgun sequence data. *Genome Res.* **17**, 184–191 (2007).
- Milo, R. et al. *Cell Biology by the Numbers*. <http://book.bionumbers.org> (2015).
- Ganley, A. R. D., Ide, S., Saka, K. & Kobayashi, T. The effect of replication initiation on gene amplification in the rDNA and its relationship to aging. *Mol. Cell* **35**, 683–693 (2009).
- Kobayashi, T. Ribosomal RNA gene repeats, their stability and cellular senescence. *Proc. Jpn. Acad. Ser. B, Phys. Biol. Sci.* **90**, 119–129 (2014).
- Bhargava, R., Onyango, D. O. & Stark, J. M. Regulation of single-strand annealing and its role in genome maintenance. *Trends Genet.* **32**, 566–575 (2016).
- Chin, C. S. et al. Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing data. *Nat. Methods* **10**, 563–569 (2013).
- Parra, G., Bradnam, K. & Korf, I. CEGMA: a pipeline to accurately annotate core genes in eukaryotic genomes. *Bioinformatics* **23**, 1061–1067 (2007).
- Smit, A. F. A., Hubley, R. & Green, P. *RepeatMasker* <http://repeatmasker.org>.
- Kajikawa, M. & Okada, N. LINEs mobilize SINEs in the eel through a shared 3' sequence. *Cell* **111**, 433–444 (2002).
- Simao, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* **31**, 3210–3212 (2015).
- Radha, V., Nambirajan, S. & Swarup, G. Association of Lyn tyrosine kinase with the nuclear matrix and cell-cycle-dependent changes in matrix-associated tyrosine kinase activity. *Eur. J. Biochem.* **236**, 352–359 (1996).
- Tsuzuki, S., Handa, Y., Takeda, N. & Kawaguchi, M. Strigolactone-induced putative secreted protein 1 is required for the establishment of symbiosis by the arbuscular mycorrhizal fungus *Rhizophagus irregularis*. *Mol. Plant Microbe Interact.* **29**, 277–286 (2016).
- Drissner, D. et al. Lyso-phosphatidylcholine is a signal in the arbuscular mycorrhizal symbiosis. *Science* **318**, 265–268 (2007).
- Tehlivets, O., Scheuringer, K. & Kohlwein, S. D. Fatty acid synthesis and elongation in yeast. *Biochim. Biophys. Acta* **1771**, 255–270 (2007).
- Li, M. G. et al. Thiamine biosynthesis in *Saccharomyces cerevisiae* is regulated by the NAD(+)-dependent histone deacetylase Hst1. *Mol. Cell. Biol.* **30**, 3329–3341 (2010).
- Wewer, V., Brands, M. & Dormann, P. Fatty acid synthesis and lipid metabolism in the obligate biotrophic fungus *Rhizophagus irregularis* during mycorrhization of *Lotus japonicus*. *Plant J.* **79**, 398–412 (2014).
- Keymer, A. et al. Lipid transfer from plants to arbuscular mycorrhizal fungi. *eLife* **6**, e29107 (2017).
- Bravo, A., Brands, M., Wewer, V., Dormann, P. & Harrison, M. J. Arbuscular mycorrhiza-specific enzymes FatM and RAM2 fine-tune lipid biosynthesis to promote development of arbuscular mycorrhiza. *New Phytol.* **214**, 1631–1645 (2017).
- Prokopowich, C. D., Gregory, T. R. & Crease, T. J. The correlation between rDNA copy number and genome size in eukaryotes. *Genome* **46**, 48–50 (2003).
- Cushion, M. T. & Keely, S. P. Assembly and annotation of *Pneumocystis jirovecii* from the human lung microbiome. *mBio* **4**, e00224 (2013).
- Gardner, M. J. et al. Genome sequence of the human malaria parasite *Plasmodium falciparum*. *Nature* **419**, 498–511 (2002).
- Timmis, J. N., Ayliffe, M. A., Huang, C. Y. & Martin, W. Endosymbiotic gene transfer: organelle genomes forge eukaryotic chromosomes. *Nat. Rev. Genet.* **5**, 123–135 (2004).

52. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. Preprint at arXiv:1303.3997 [q-bio.GN] (2013).
53. Savary, R. et al. A population genomics approach shows widespread geographical distribution of cryptic genomic forms of the symbiotic fungus *Rhizophagus irregularis*. *ISME J.* **12**, 17–30 (2018).
54. Kuhn, G., Hijri, M. & Sanders, I. R. Evidence for the evolution of multiple genomes in arbuscular mycorrhizal fungi. *Nature* **414**, 745–748 (2001).
55. Gibbons, J. G., Branco, A. T., Godinho, S. A., Yu, S. K. & Lemos, B. Concerted copy number variation balances ribosomal DNA dosage in human and mouse genomes. *Proc. Natl. Acad. Sci. U.S.A.* **112**, 2485–2490 (2015).
56. Kaul, S. et al. Analysis of the genome sequence of the flowering plant *Arabidopsis thaliana*. *Nature* **408**, 796–815 (2000).
57. Mentewab, A. B., Jacobsen, M. J. & Flowers, R. A. Incomplete homogenization of 18S ribosomal DNA coding regions in *Arabidopsis thaliana*. *BMC Res. Notes* **4**, 93 (2011).
58. Vembar, S. S., Droll, D. & Scherf, A. Translational regulation in blood stages of the malaria parasite *Plasmodium* spp.: systems-wide studies pave the way. *Wiley Interdiscip. Rev. RNA* **7**, 772–792 (2016).
59. Andersen, S. L. & Sekelsky, J. Meiotic versus mitotic recombination: two different routes for double-strand break repair: the different functions of meiotic versus mitotic DSB repair are reflected in different pathway usage and different outcomes. *BioEssays* **32**, 1058–1066 (2010).
60. Sanders, I. R. Sex, plasticity, and biologically significant variation in one Glomeromycotina species: a response to Bruns et al. (2017) ‘Glomeromycotina: what is a species and why should we care?’. *New Phytol.* <https://doi.org/10.1111/nph.15049> (2018).
61. Cappé, O. & Moulines, E. On-line expectation-maximization algorithm for latent data models. *J. R. Stat. Soc. Ser. B, Stat. Methodol.* **71**, 593–613 (2009).
62. Birch, J. L. & Zomerdijk, J. C. Structure and function of ribosomal RNA gene chromatin. *Biochem. Soc. Trans.* **36**, 619–624 (2008).
63. Brennicke, A., Marchfelder, A. & Binder, S. RNA editing. *FEMS Microbiol. Rev.* **23**, 297–316 (1999).
64. Xue, S. F. & Barna, M. Specialized ribosomes: a new frontier in gene regulation and organismal biology. *Nat. Rev. Mol. Cell Biol.* **13**, 355–369 (2012).
65. Ehinger, M. O., Croll, D., Koch, A. M. & Sanders, I. R. Significant genetic and phenotypic changes arising from clonal growth of a single spore of an arbuscular mycorrhizal fungus over multiple generations. *New Phytol.* **196**, 853–861 (2012).
66. Angelard, C., Colard, A., Niculita-Hirzel, H., Croll, D. & Sanders, I. R. Segregation in a mycorrhizal fungus alters rice growth and symbiosis-specific gene transcription. *Curr. Biol.* **20**, 1216–1221 (2010).
67. Angelard, C. & Sanders, I. R. Effect of segregation and genetic exchange on arbuscular mycorrhizal fungi in colonization of roots. *New Phytol.* **189**, 652–657 (2011).
68. Sanders, I. R. & Croll, D. Arbuscular mycorrhiza: the challenge to understand the genetics of the fungal partner. *Annu. Rev. Genet.* **44**, 271–292 (2010).
69. Corradi, N. & Brachmann, A. Fungal mating in the most widespread plantsymbionts? *Trends Plant Sci.* **22**, 175–183 (2017).
70. Fulton, T. M., Chunwongse, J. & Tanksley, S. D. Microprep protocol for extraction of DNA from tomato and other herbaceous plants. *Plant Mol. Biol. Report.* **13**, 207–209 (1995).
71. Gurevich, A., Saveliev, V., Vyahhi, N. & Tesler, G. QUAST: quality assessment tool for genome assemblies. *Bioinformatics* **29**, 1072–1075 (2013).
72. Marçais, G. et al. MUMmer4: a fast and versatile genome alignment system. *PLoS Comput. Biol.* **14**, e1005944 (2018).
73. Trapnell, C., Pachter, L. & Salzberg, S. L. TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics* **25**, 1105–1111 (2009).
74. Keller, O., Kollmar, M., Stanke, M. & Waack, S. A novel hybrid gene prediction method employing protein multiple sequence alignments. *Bioinformatics* **27**, 757–763 (2011).
75. Slater, G. S. & Birney, E. Automated generation of heuristics for biological sequence comparison. *BMC Bioinformatics* **6**, 31 (2005).
76. Liao, Y., Smyth, G. K. & Shi, W. FeatureCounts: an efficient general purpose program for assigning sequence reads to genomic features. *Bioinformatics* **30**, 923–930 (2014).
77. Jones, P. et al. InterProScan 5: genome-scale protein function classification. *Bioinformatics* **30**, 1236–1240 (2014).
78. Emms, D. M. & Kelly, S. OrthoFinder: solving fundamental biases in whole genome comparisons dramatically improves orthogroup inference accuracy. *Genome Biol.* **16**, 157 (2015).
79. De Bie, T., Cristianini, N., Demuth, J. P. & Hahn, M. W. CAFE: a computational tool for the study of gene family evolution. *Bioinformatics* **22**, 1269–1271 (2006).
80. Nguyen, L.-T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* **32**, 268–274 (2015).
81. Smith, A. B. Cambrian problematica and the diversification of deuterostomes. *BMC Biol.* **10**, 79 (2012).
82. Quinlan, A. R. BEDTools: the swiss-army tool for genome feature analysis. *Curr. Protoc. Bioinformatics* **47**, 11.12.1–11.12.34 (2014).
83. Yin, Y. et al. dbCAN: a web resource for automated carbohydrate-active enzyme annotation. *Nucleic Acids Res.* **40**, W445–W451 (2012).
84. Lagesen, K. et al. RNAMmer: consistent and rapid annotation of ribosomal RNA genes. *Nucleic Acids Res.* **35**, 3100–3108 (2007).
85. Katoh, K., Misawa, K., Kuma, K. & Miyata, T. MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* **30**, 3059–3066 (2002).
86. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–359 (2012).
87. Ankenbrand, M. J., Hohlfeld, S., Hackl, T. & Förster, F. AliTV—interactive visualization of whole genome comparisons. *PeerJ Comput. Sci.* **3**, e116 (2017).
88. Kumar, S., Stecher, G. & Tamura, K. MEGA7: molecular evolutionary genetics analysis version 7.0 for bigger datasets. *Mol. Biol. Evol.* **33**, 1870–1874 (2016).
89. Schattner, P., Brooks, A. N. & Lowe, T. M. The tRNAscan-SE, snoscan and snoGPS web servers for the detection of tRNAs and snoRNAs. *Nucleic Acids Res.* **33**, W686–W689 (2005).
90. Kwan, E. X., Wang, X. B. S., Amemiya, H. M., Brewer, B. J. & Raghuraman, M. K. rDNA copy number variants are frequent passenger mutations in *Saccharomyces cerevisiae* deletion collections and *de novo* transformants. *G3: Genes, Genomes, Genet.* **6**, 2829–2838 (2016).
91. Goffeau, A. et al. Life with 6000 genes. *Science* **274**, 546 (1996).
92. Frazer, K. A., Pachter, L., Poliakov, A., Rubin, E. M. & Dubchak, I. VISTA: computational tools for comparative genomics. *Nucleic Acids Res.* **32**, W273–W279 (2004).
93. Simon, U. K. & Weiss, M. Intragenomic variation of fungal ribosomal genes is higher than previously thought. *Mol. Biol. Evol.* **25**, 2251–2254 (2008).

## Acknowledgments

This work was supported by JST ACCEL Grant Number JPMJAC1403, Japan. We thank Miwako Matsumoto, the Functional Genomics Facility and the Data Integration and Analysis Facility at the National Institute for Basic Biology for technical support; Katsuharu Saito, Kohki Akiyama; and present and past members of the Kawaguchi Lab and the Shigenobu Lab.

## Author contributions

T.M., S.S. and M.K. conceived of and designed the experiments; T.M., Y.K., H.K., N.T., K.Y., and T.B. performed the experiments; T.M., N.O., S.S. and T.B. analyzed the data; T.M., Y.K., H.K., K.Y., T.B., S.S. and M.K. wrote the manuscript.

## Additional information

**Supplementary information** accompanies this paper at <https://doi.org/10.1038/s42003-018-0094-7>.

**Competing interests:** The authors declare no competing interests.

**Reprints and permission** information is available online at <http://npg.nature.com/reprintsandpermissions/>

**Publisher's note:** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2018