

1 **Mechanisms of surface antigenic variation in the**
2 **human pathogenic fungus *Pneumocystis jirovecii***

3
4
5
6

7 **Emanuel Schmid-Siegert**¹, **Sophie Richard**², **Amanda Luraschi**²,
8 **Konrad Mühlethaler**³, **Marco Pagni**¹, **Philippe M. Hauser**²

9
10

11 ¹ Vital-IT Group, SIB Swiss Institute of Bioinformatics, Lausanne, Switzerland

12 ² Institute of Microbiology, Lausanne University Hospital, Lausanne, Switzerland

13 ³ Institut für Infektionskrankheiten, Universität Bern, Bern, Switzerland

14
15

16 Corresponding author :

17 P. Hauser

18 Av. Bugnon 48

19 IMUL CHUV

20 1011 Lausanne

21 Phone +41 21 314 40 84

22 Philippe.Hauser@chuv.ch

23
24

25 **Running title** : Antigenic variation in *Pneumocystis jirovecii*

26
27

28 **Words text**: 5180

29
30

31 **Methods text**: 2028

32
33

34 **Abstract**: 268

35
36

37 **Keywords**: *Pneumocystis carinii*, PCP, major surface glycoprotein, adhesin, subtelomere, mosaicism, PacBio sequencing.

38
39
40

38 **Abstract**

39 **Background:** Microbial pathogens commonly escape the human immune system by varying
40 surface proteins. Here we investigated the mechanisms used for that purpose by *Pneumocystis*
41 *jirovecii*. This uncultivable fungus is an obligate pulmonary pathogen which causes pneumonia
42 in immuno-compromised individuals, a major life-threatening infection.

43 **Results:** Long-read PacBio sequencing was used to assemble a set of subtelomeres of a single *P.*
44 *jirovecii* strain from a bronchoalveolar lavage fluid specimen of a single patient. A total of 113
45 genes encoding surface proteins were identified, including 28 pseudogenes. These genes formed
46 a subtelomeric gene superfamily which included five families encoding adhesive GPI-anchored
47 glycoproteins, and one family encoding excreted glycoproteins. Numerical analyses suggested
48 that diversification of the glycoproteins relies on mosaic genes created by ectopic recombination,
49 and occurs only within each family. DNA motifs suggested that all genes are expressed
50 independently, except those of the family encoding the most abundant surface glycoproteins
51 which are subject to mutually exclusive expression. PCR analyses showed that exchange of the
52 expressed gene of the latter family occurs frequently, possibly favoured by the location of the
53 genes proximal to the telomere because this allows concomitant telomere exchange.

54 **Conclusions:** Our observations suggest that (i) the structure of *P. jirovecii* cell surface is made
55 of a complex mixture of different glycoproteins, (ii) genetic mosaicism ensures variation of the
56 glycoproteins, and (iii) the strategy of the fungus consists in the continuous production of new
57 subpopulations composed of cells which are antigenically different. This strategy is unique
58 among human pathogens and may be associated to the particular niche within lungs which
59 tolerates the presence of low abundant fungi within the natural microbiota.

60 **Introduction**

61 *Pneumocystis jirovecii* is a fungus colonizing specifically human lungs. It has developed
62 strategies to survive in healthy human lungs, at least transiently, and can turn into a deadly
63 pathogen causing pneumonia in individuals with debilitated immune system (Cushion et al.
64 2007; Cushion and Stringer 2010; Hauser 2014; Ma et al. 2016). This disease is the second most
65 frequent life-threatening invasive fungal infection with ca. 400'000 cases per year worldwide
66 (Brown et al. 2012). However, the biology of this pest remains difficult to study in the lab
67 because of the lack of any established methods for *in vitro* culture. Recent progresses in
68 understanding *P. jirovecii* biology strongly benefitted from the publication of two assemblies of
69 its genome from two different clinical samples (Ma et al. 2016; Cissé et al. 2012).

70 In contrast to other pathogenic fungi, the cells of *P. jirovecii* lack chitin as well as glucans
71 during part of the cell cycle, which may avoid eliciting innate and acquired immune responses
72 (Ma et al. 2016). Moreover, a mechanism of surface antigenic variation, to which ca 5% of the
73 genome is dedicated (Ma et al. 2016), seems crucial to escape from the human immune system
74 during colonisation, although it has not been understood in details so far. Surface antigenic
75 variation is a common strategy among major microbial human pathogens, for example in
76 *Plasmodium*, *Trypanosoma*, *Candida*, *Neisseria*, and *Borrelia*. It relies on various genetic and/or
77 epigenetic mechanisms aimed at expressing only one or few of them at once (Deitsch et al.
78 2009). Such systems often involve gene families encoding surface antigens localized at
79 subtelomeres, presumably because these regions of the genome are prone to gene silencing,
80 which is used for mutually exclusive expression, and possibly to enhanced mutagenesis (Barry et
81 al. 2003). Moreover, the formation of clusters of telomeres at the nuclear periphery may favour
82 ectopic recombinations (Barry et al. 2003), which can be responsible for the generation of new
83 mosaic antigens.

84 Surface antigenic variation has been previously studied on a limited set of genes in
85 *Pneumocystis carinii* infecting specifically rats. The molecular mechanism was then assumed to
86 be also active in *P. jirovecii*, as suggested by studies using PCR-based technologies. Antigen
87 diversity was believed to be generated by recombination between members of a single family of
88 ca. 80 subtelomeric genes encoding isoforms of the major surface glycoprotein (*msg*) (Keely et
89 al. 2005; Keely et al. 2009; Stringer 2007). A single of these isoforms would be expressed in
90 each cell thanks to its localization downstream of a subtelomeric expression site, the upstream
91 conserved element (UCS) present at a single copy in the genome. The UCS includes the
92 promoter of transcription, the protein start, and the leader sequence responsible for translocation
93 of the protein into the endoplasmic reticulum for final incorporation into the cell wall (Kutty et
94 al. 2001; Kutty et al. 2013). The mechanism for exchange of the expressed *msg* gene is thought
95 to be by recombination at a 33 bps long sequence which is present both at the end of UCS and
96 beginning of each *msg* (the conserved recombination junction element, CRJE). The exchange of
97 the expressed gene seems relatively frequent and would explain how different *msg* genes can be
98 expressed in each population (Kutty et al. 2001). The CRJE sequence encodes at its end a
99 potential lysine-arginine recognition site for Kexin endonuclease which might be involved in the
100 maturation of the antigen. Kutty *et al* (Kutty et al. 2008) provided evidence for frequent
101 recombinations among *msg* genes creating potentially mosaic genes. All these observations were
102 made using conventional cloning procedures and PCRs, and these mechanisms have yet to be
103 understood in a more extensive genomic context.

104 The first genome sequence of *P. jirovecii* released was obtained using technologies generating
105 short reads which prevented assembly of long repetitive sequences such as centromeres,
106 telomeres, and subtelomeres including *msg* genes (Cissé et al. 2012). A second study used a
107 mixture of techniques which generated more complete chromosomes of *P. jirovecii*, *P. carinii*,

108 and *Pneumocystis murina* (infecting specifically mice) (Ma et al. 2016). These latter authors used
109 PCRs coupled with long read sequencing to reconstruct the subtelomeres. This allowed
110 discovering new subtelomeric gene families related to *msg*. The number of these families as well
111 as the number of their members present in each *Pneumocystis* species varied. They further
112 described the specific arrangement of the members of each family within the subtelomeres, and
113 suggested, based on RNA sequencing, that all *msg* genes are expressed in a given population of
114 *P. carinii* or *P. murina*. However, they did not discuss the function of these proteins, the
115 mechanisms involved in their expression and gene variation, or the global strategy of antigenic
116 variation of these fungi.

117 The aim of the present study was to analyse in details the mechanisms of surface antigenic
118 variation in *P. jirovecii*. To that purpose, we used the PacBio sequencing technology generating
119 long DNA reads to assemble a set of subtelomeres of a single *P. jirovecii* strain from a
120 bronchoalveolar lavage fluid specimen (BALF) of a single patient. The analysis of this dataset
121 and laboratory experiments permit a new classification and the characterization of six
122 subtelomeric *msg* families, demonstrate the presence of pseudogenes, and provide important new
123 insights into the molecular mechanisms responsible for antigenic variation. Moreover, our
124 observations suggest a unique strategy of antigenic variation consisting in the continuous
125 production of new subpopulations composed of cells which are antigenically different. This
126 strategy may be associated to the particular non-sterile niche within lungs.

127

128 **Results**

129 Most if not all *P. jirovecii* infections are polyclonal (Alanio et al. 2016). In order to facilitate the
130 study of the mechanisms of antigenic variation, one patient infected with a vastly dominant strain
131 was selected by multitarget genotyping. The genome of a single *P. jirovecii* strain was assembled
132 into 219 contigs using PacBio sequencing and a dedicated bioinformatics strategy for reads
133 processing.

134

135 **Identification of subtelomeric *msg* genes and pseudogenes**







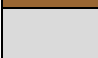
136 Automated gene prediction performed poorly in the subtelomeric regions as compared to the
137 core of the genome, due to abundant stretches of low-complexity DNA, numerous pseudogenes,
138 residual assembly errors in homopolymers, and the lack of a start codon in many *msg* genes. The
139 *msg* genes were detected by sequence homology using generalized profiles (Bucher and Bairoc
140 1994) derived from previously published sequences. A total of 113 *msg* genes with sizes ranging
141 from 331 to 3337 bps were found on 37 different contigs, only two genes being perfectly
142 identical (*msg* 52 and 61, Additional file 1: Table S1). Most of them (N=85) contained a single
143 large exon and zero to two small exons at their 5' end. The remaining 28 genes harboured many
144 stop codons in all frames and were considered as pseudogenes (Additional file 2: supplementary
145 note 1).

146

147 **Characterization of the *msg* gene families**

148 We are proposing a classification of the *msg* genes into six families (Table 1) based on of the
149 integration of four independent lines of evidence: sequence homology, gene structure, protein
150 property, and recombination events. The global picture that emerged is coherent and the details
151 on the different points are presented below.

Table 1. Characteristics of the *msg* families identified in *P. jirovecii*.

Family name	Color in figures	Gene							Protein							Corresponding Ma et al (2016) family ^a
		No. genes full-length / partial / pseudo-	Mean full-length (bps) ± st dev	Location in subtelomere relatively to telomere	Presumptive TATA box (bps to ATG, range)	CRJE	No. 5'-end introns	Average pairwise identity (%) ± st dev	Signal peptide	C-terminus			GPI-anchor signal	No. N-glycosylation site	Average pairwise identity (%) ± st dev	
										ST-rich region	ST-rich region	PE-rich region				
<i>msg-I</i>		11 / 16 / 16	3071 ± 39	proximal	- ^b	+	0	71 ± 7	- ^b	+	+	-	+	4-10	54 ± 8	<i>msg-A1</i>
<i>msg-II</i>		11 / 3 / 4	3155 ± 31 ^c	central	21-28	-	2	83 ± 13	+	+	+	-	+	2-14	73 ± 16	<i>msg-A3</i>
<i>msg-III</i>		7 / 2 / 1	3146 ± 55	central	18-24	-	2	83 ± 10	+	+	+	-	+	7-11	70 ± 13	<i>msg-A3</i>
<i>msg-IV</i>		6 / 1 / 2	2023 ± 45	central	29-36	-	1	72 ± 14	+	-	-	-	-	0-8	49 ± 17	<i>msg-B</i>
<i>msg-V</i>		8 / 6 / 1	3056 ± 126	central	30-67	-	1	66 ± 5	+	-	+	+	+	5-12	44 ± 4	<i>msg-D</i>
<i>msg-VI</i>		6 / 1 / 0	1222 ± 189	distal	33-56	-	1	45 ± 7	+	-	+	+	+	0-1	21 ± 5	<i>msg-E</i>
<i>msg outlier</i>		6 / 1 / 4	variable	central/distal	NA ^d	+/-	variable	NA	+/-	+/-	+/-	-	+/-	variable	NA	NA

^a The family *msg-C* described by Ma et al (2016) was not identified here (Supplementary note 6).

^b The promoter including the signal peptide for this family is within the UCS present at a single copy per genome.

^c The *msg3* gene was not used to calculate this value because it is ca. 900 bps shorter than the other genes of the family, although it presents all features of the family (see alignment in Figure S2).

^d Not applicable.

1 Figure 1a shows the results of the analysis of 61 *msg* genes containing an exon equal or larger
2 than 1.6 kb. Based on the multiple sequence alignments (MSAs) of the CDS and of their
3 predicted proteins, two phylogenetic trees were computed using RAxML. The different gene
4 families are clearly individualised as clades, with the exceptions of (i) *msg-II* which appears as a
5 sub-clade of *msg-I*, and (ii) *msg-I* which seems to include two sub-clades. Using an alternative
6 classification method that does not rely on a single particular MSA (JACOP, Fig. 1a), the
7 placement of *msg-II* as a sub-clade of *msg-I* was not confirmed, whereas the sub-clades of *msg-I*
8 were. Owing on the differences in the gene structures and on the recombination events reported
9 below, we believe that (i) *msg-I* and *msg-II* should be treated separately, and (ii) *msg-I* should be
10 considered as a single family including two sub-clades. Figure 1b shows the analysis of trimmed
11 CDS sequences allowing the placement of the *msg-VI* family which appeared as a clade on its
12 own, while the classification of the other families remained essentially unchanged. Figure S1
13 shows that most pseudogenes could be attributed to one *msg* family and their often longer
14 branches further account for their pseudogenic nature (Additional file 2).

15 Manual curation of the *msg* genes led to their classification in full-length, partial, and
16 pseudogenes (Table S1). Table 1 shows the characteristics of each family identified by the
17 analysis of the sequences of the full-length genes, as well as of their alignments (Fig. S2). Except
18 those of the family *msg-I*, each *msg* gene presented one or two introns at its 5' end, as well as a
19 presumptive TATA box upstream of the ATG and an initiator motif (Cap signal) at presumptive
20 sites of initiation transcription (Fig. 2a and S2). The members of the family I had only the
21 conserved recombination junction element (CRJE) at the beginning of their single exon. These
22 observations suggested that members of family I can be expressed only upon recombination of

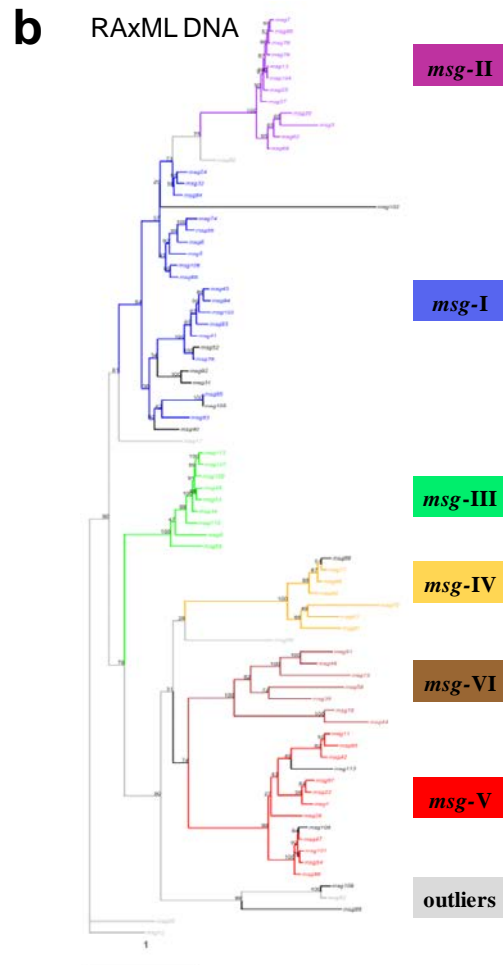
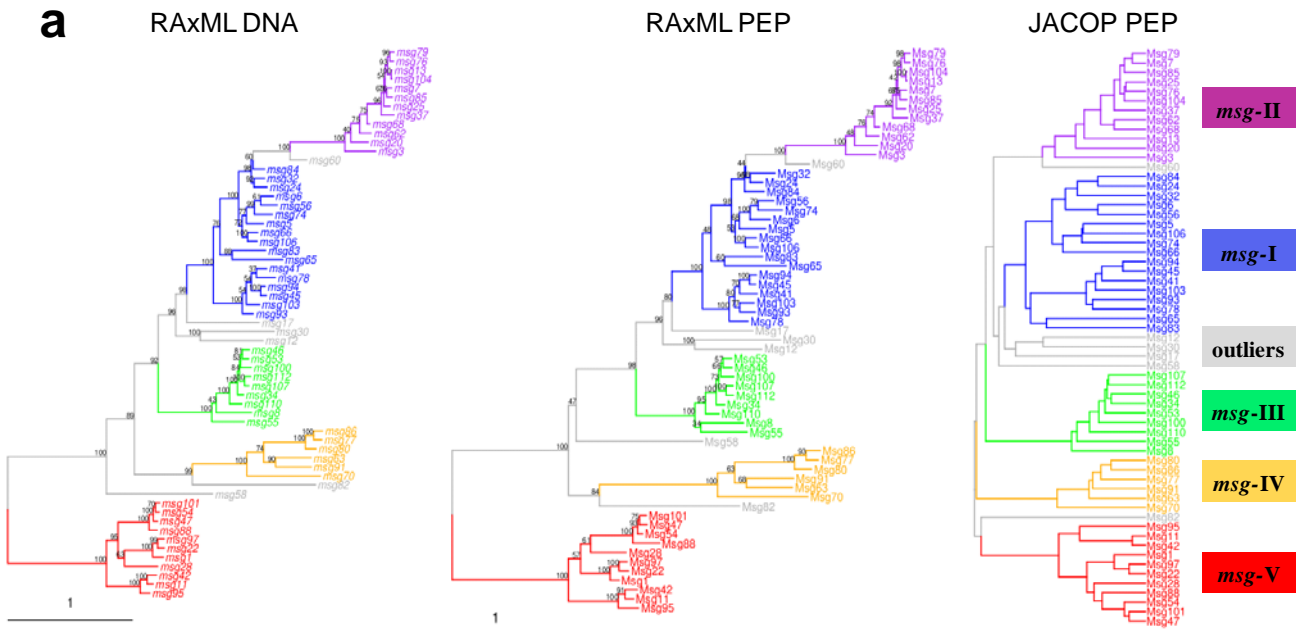


Fig. 1

24

25 **Fig. 1**

26 Classification trees of *P. jirovecii* *msg* genes and Msg proteins. The different families are
27 represented in colours and their characteristics are summarised in Table 1. A few unclassified
28 outliers are in grey. Scale, mean substitution / site. **(a)** RAxML DNA and PEP are maximum
29 likelihood trees of nucleotide and amino acid sequences of the 61 genes with an exon larger than
30 1.6 kb. Members of family V were defined as the out-group (1000 bootstraps). JACOP PEP is a
31 hierarchical classification based on local sequence similarity, a method that does not rely on a
32 particular multiple sequence alignment. **(b)** Maximum likelihood tree of the 61 genes with an
33 exon larger than 1.6 kb, plus 18 genes with an exon smaller than 1.6 kb. The sequences were
34 trimmed from position 1540 of the first alignment up to their end, and re-aligned to construct the
35 tree (1000 bootstraps). Seven of the 18 genes with an exon smaller than 1.6 kb constitute the *msg*
36 family VI shown in brown, whereas the remaining 11 shown in black belong to the other *msg*
37 families.

38

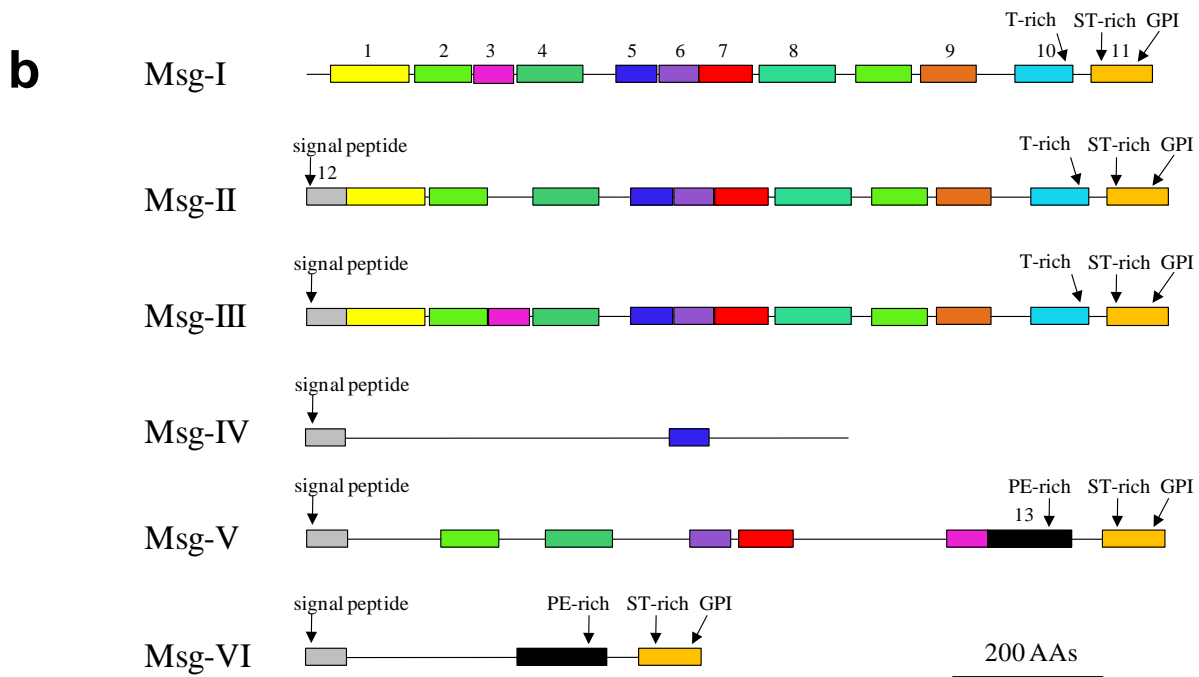
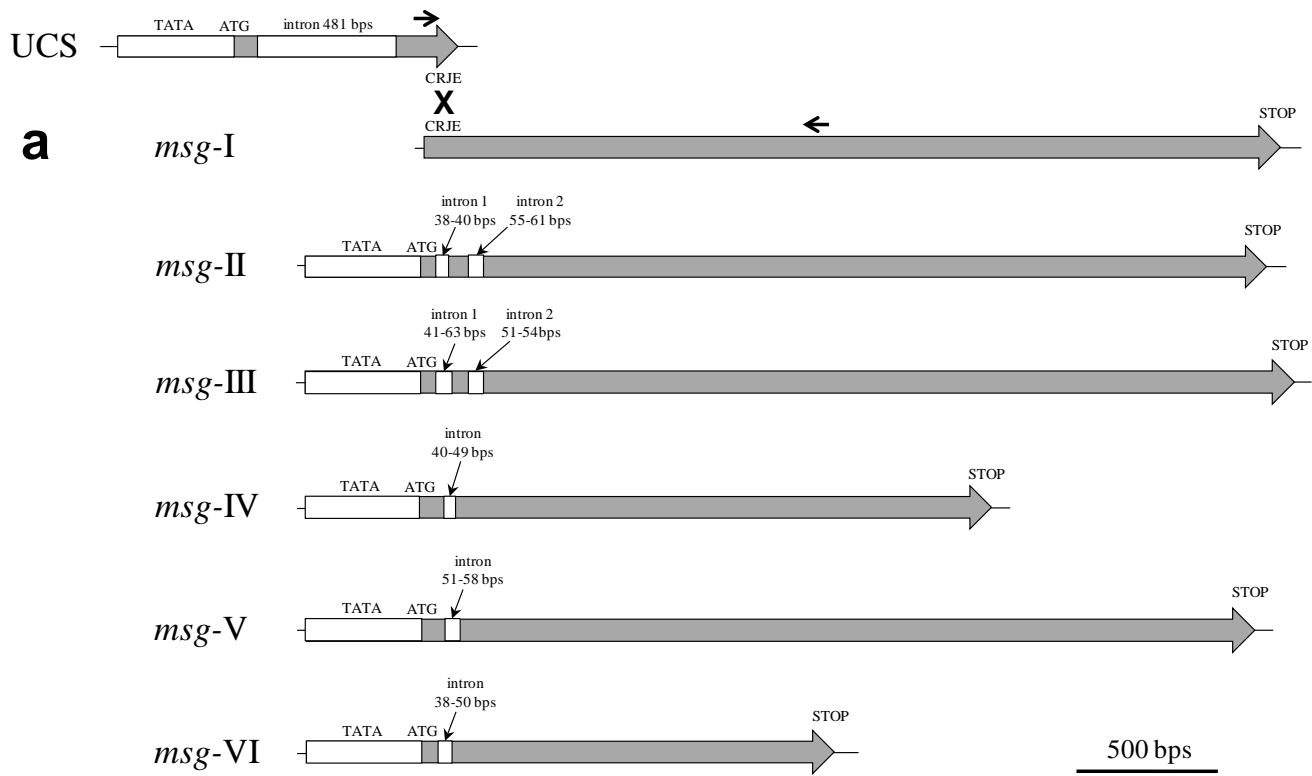


Fig. 2

40

41 **Fig. 2**

42 Diagrams of the structure of *P. jirovecii* *msg* genes and Msg proteins belonging to six families.

43 (a) Features of the *msg* genes of each family derived from the analysis of the full-length genes.

44 The UCS and recombination between CRJE sequences are figured. The approximate position of

45 PCR primers used for identification of the *msg*-I expressed genes linked to the UCS are shown

46 by arrows (Supplementary note 4). (b) Features of Msg proteins of each family derived from the

47 analyses of the full-length proteins. The 13 domains identified using MEME analysis are shown.

48 The logos of these domains are shown in Figure S4.

49

50 their CRJE with that of the single copy UCS which encompasses a promoter, whereas all
51 members of the other five families are expressed independently. Three of the six full-length
52 outlier genes seemed not expressed since they had no CRJE and missed a TATA box (Table S1).
53 Twenty-six partial genes were truncated by the end of the contig so that only three *bona fide*
54 partial genes were identified, which, however, missed TATA box, signal peptide, and/or GPI-
55 anchor signal, and thus were probably not expressed or not correctly processed (*msg* 44, 89, and
56 99).

57

58 **Characterization of the Msg protein families**

59 Analysis of the sequences and alignments (Fig. S3) of the full-length proteins of each family
60 revealed that each Msg protein, except those of family I, presented a signal peptide at its N-
61 terminus (Fig. 2b). Proteins of family I probably acquire a signal peptide upon fusion of their
62 encoding gene with the UCS. Except those of family IV, each Msg protein presented a GPI-
63 anchor signal at its C-terminus. These observations suggested that all Msg proteins are attached
64 externally to the cell wall, except those of family IV which would be secreted in the environment
65 or attached to the cell wall through another mechanism than GPI.

66 The possible conservation of motifs among the proteins of the six families was investigated
67 using Multiple Expectation-Maximization for Motif Elicitation (MEME analysis) (Bailey and
68 Elkan 1994). Thirteen conserved motifs were identified which arrangement was fairly diagnostic
69 within each family (Fig. 2b). Most motifs included several conserved cysteines and leucines,
70 which resembled to the previously identified Pfam MSG domain (Fig. S4). Interestingly,
71 conserved leucines were often separated by two to six residues. The beginning of motif 10
72 corresponded to the end of the previously identified Pfam Msg2_C domain. Accordingly, Pfam

73 predictions identified one to five MSG domains, often partial, per protein of all families, and a
74 single Msg2_C domain in each Msg-I protein (Fig. S5; Table S2). The Msg2_C domain was not
75 predicted in families II and III although they harboured the corresponding motif 10, suggesting
76 that this domain is divergent in these families. Ncoils predictor revealed three to five coiled-coil
77 motifs spread along members of families I, II, and III, whereas unstructured regions were
78 predicted at the C-terminus of Msg proteins of families I, III, V, and VI (Fig. S5).

79 Except those of family IV, each Msg protein harboured at its C-terminus two MEME motifs
80 which included a region enriched in specific residues: threonine (T-rich; motif 10), serine and
81 threonine (ST-rich; 11), or proline and glutamine (PE-rich; 13)(Fig. 2b; Table 1). The T-rich
82 region in family I included generally a stretch of nine to 15 Ts, which was not present in families
83 II and III (Fig. S3). The PE-rich region in family V was enriched in proline residues relatively to
84 that present in family VI (Fig. S3). Four to 14 potential sites of nitrogen-linked glycosylation of
85 asparagines were predicted to be present in each Msg protein, except in family VI which
86 presented no or only one such site (Table 1; Fig. S3). The localization of these glycosylation sites
87 was widespread along the protein and fairly conserved within each family.

88

89 **Arrangement of the *msg* families within the subtelomeres**

90 Consistent with subtelomeric localization, the *msg* genes were grouped at one end of their
91 contig when flanking non-*msg* genes were also present (in 20 of 37 contigs; Fig. 3 and S6a). All
92 *msg* genes identified were oriented towards one end of the contig, *i.e.* presumably towards the
93 telomere (no telomeric repeats were identified for an unknown reason; Supplementary note 2).
94 Except pseudogenes which were dispersed all over the subtelomeres, all members of family I
95 were the closest to the end of their contig, *i.e.* proximal to the telomere (Fig. 3 and S6). By

106 contrast, all members of family VI were the closest to the flanking non-*msg* genes present on
107 their contig, *i.e.* distal to the telomere. Members of the four remaining families were localized
108 centrally in the subtelomeres, between those of families I and VI. There were up to three *msg*-I
109 genes grouped at the end of 19 contigs. Members of the other five families did not show any
110 clear grouping pattern.

111

112 **Identification of the expression site of *msg*-I genes and of the genes linked to it**

113 Each infection by *P. jirovecii* is believed to involve a mixture of cells expressing different *msg*-I
114 genes under the control of the expression site, *i.e.* the UCS which is present at a single copy per
115 genome (Kutty et al. 2001). Consequently, the UCS was expected to be linked to different *msg*-I
116 genes in our DNA sample and thus cannot be unequivocally assembled, which plausibly explains
117 its absence from the PacBio assembly. A single UCS was retrieved from our DNA sample using
118 PCRs based on published sequences and it could be linked to one of the PacBio contigs
119 (Supplementary note 3). Consistently, this contig was linked to chromosome 1 which also carries
120 the UCS in the Ma et al (2016) assembly (Table S3). The UCS retrieved from our sample was
121 identical to that of Ma et al (2016), except few small changes not modifying the encoded protein
122 (Fig. S7). Interestingly, the CRJE sequence at the end of the UCS and beginning of each *msg*-I
123 gene presented an imperfect inverted repeat which was never pointed out so far (Fig. S7).

124 In order to identify the *msg*-I genes linked to the UCS in our sample, we amplified by PCR
125 the junction between these elements using one primer within the UCS and either (i) one primer
126 generic for many *msg*-I genes (Kutty et al. 2001), or (ii) one primer specific to a given *msg*-I
127 gene of the PacBio assembly (Supplementary note 4; Fig. 2a). Eighteen different *msg*-I genes
128 were found fused in frame to the UCS at the CRJE sequence, two being pseudogenes of the

129 family I with an upstream CRJE sequence, and four new *msg-I* sequences not present in the
130 PacBio assembly. The 14 *msg-I* genes found linked to the UCS which were present in the PacBio
131 assembly are identified in Figures 3 and S6 by asterisks. Three specific *msg-I* genes linked to the
132 UCS represented 74% of the subclones of the generic PCR analyzed, suggesting that sub-
133 populations of cells expressing given *msg-I* genes were of different sizes in our sample
134 (Supplementary note 4). These observations suggested that recombination between the CRJE
135 sequence of the UCS with that of different *msg-I* genes occurred at a high frequency in the single
136 *P. jirovecii* population studied here.

137

138 **Set of assembled subtelomeres**

139 The flanking non-*msg* genes allowed attributing 20 of our 37 contigs (Fig. 3 and S6a) to 15 of
140 the 20 full-length chromosomes described by Ma et al (2016) because they were also present in
141 the latter assembly (Table S3). All the remaining 17 contigs without flanking non-*msg* genes
142 (Fig. S6b) could have been assembled from the same subtelomeres as other contigs. Thus, we
143 assembled at least 20 subtelomeres out of the 40 potentially present in each cell. Given the
144 presence of a large number of subpopulations expressing different *msg-I* genes in our sample, the
145 set of subtelomeres present in each cell varied considerably. It is likely that the set we assembled
146 corresponded to a core of subtelomeres which was present in a majority of cells of the population
147 so that it could be assembled unequivocally.

148

149 **Recombination between *msg* genes**

150 Evidence of recombination events between *msg-I* genes was previously provided (Kutty et al.
151 2008). We investigated this issue among the different *msg* families using three different

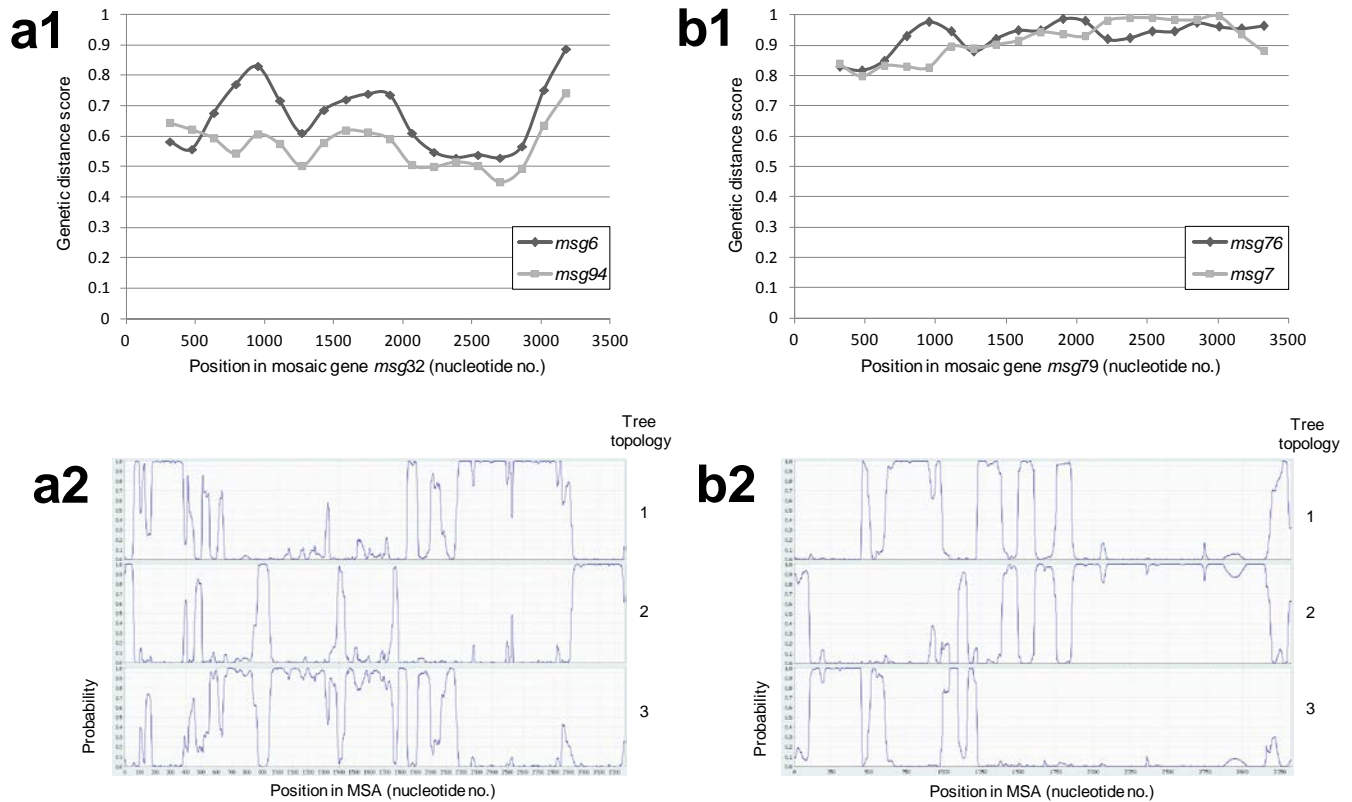
152 numerical methods: two allowing analyses of large sets of genes for screening, and one
153 analyzing only four genes at a time for more sensitive analysis. Two to 18 potential mosaic genes
154 and their putative parent genes were detected within each family I to IV, involving sometimes
155 partial or pseudogenes (Fig. 4; Table 2). On the other hand, only one potential mosaic gene was
156 identified in family V and none in family VI ($P = 0.06$). Eight of the 30 mosaic genes detected
157 shared with one parent a perfectly identical fragment of ca. 100 to 1000 bps, often close to the
158 site of the predicted recombination events (Fig. 4b and S8). These latter cases suggested very
159 recent recombination events. The putative parent genes of mosaic genes were randomly
160 distributed among the two sub-clades of family I, suggesting that this family must be considered
161 as a single entity (Supplementary note 5).

162 One to four potential recombination events per mosaic gene were generally identified using
163 the two screening methods. These events were most often confirmed by the more sensitive
164 method which, however, detected many other potential recombination events (Fig. 4 and S8).
165 Consistent with the single mosaic gene detected in families V and VI, the frequency of
166 recombination events appeared lower in these families than in the others (Fig. S9). This
167 correlated with an average pairwise identity lower within each of these two families than within
168 the others (45-66 versus 71-83%, Table 1). The predicted sites of the recombinations reported by
169 all three methods were distributed randomly along the *msg* genes for all families, and did not
170 contain any specific DNA sequence motifs (Fig. 4, S8, and S9). This suggested homologous
171 rather than site-specific recombination events.

172 In contrast, we were unable to detect recombination events between different *msg* families,
173 even using the more sensitive method (Fig. S10).

174

175



176

177

178

179 **Fig. 4**

180 Examples of detection of potential mosaic genes. (a) Mosaic gene *msg32*. (a1) The set of 11 full-
181 length *msg-I* genes was analyzed using the Recombination Analysis Tool. This method measures
182 genetic distances in windows sliding along the MSA. The genetic distance scores of the putative
183 parent genes at the middle of each window are plotted against the position in the mosaic gene.
184 The predicted recombination site is at position ca. 600, at the cross-over of the curves. The
185 second screening method Bellerophon, which is based on a similar analysis, identified a
186 recombination event at position 392. (a2) Analysis of the mosaic gene *msg32* with its putative

187 **(Legend Fig. 4 continued)**

188 parent genes together with the randomly chosen *msg84* of the same family using the more
189 sensitive method TOPALi based on the Hidden Markov Model. This method analyses only four
190 sequences at a time and calculates the probabilities of the three possible tree topologies at each
191 residue of the MSA. A recombination event is also detected at position ca. 400-600, but several
192 other recombination events are predicted. **(b)** Mosaic gene *msg79*. This gene shares an almost
193 identical fragment of 947 bps with its putative parent *msg7* (see alignment in Fig. S8c). **(b1)** The
194 set of 11 full-length *msg-II* genes was analyzed using the Recombination Analysis Tool. The
195 predicted recombination sites are at positions ca. 400, 1300, 2100, and 3100. The Bellerophon
196 method did not identify this mosaic gene. **(b2)** Analysis of the mosaic gene *msg79* with its
197 putative parent genes together with the randomly chosen *msg85* of the same family using
198 TOPALi based on the Hidden Markov Model. Recombination events are also detected at
199 positions ca. 400, 1500, and 3100, but not at 2100.

Table 2. Potential mosaic genes detected within each *msg* family ^a.

<i>msg</i> family	No. <i>msg</i> genes				Non-mosaic ^b	No. potential <i>msg</i> mosaic genes				% mosaic
	Full-length	Partial	Pseudo	Total		Full-length ^c	Partial	Pseudo ^d	Total ^b	
I	11	16	16	43	25	8	1	9	18	42
II	11	3	4	18	13	4	1	0	5	28
III	7	2	1	10	6	3	0	1	4	40
IV	6	1	2	9	7	1	0	1	2	22
V ^e	8	6	1	15	14	1	0	0	1	7
VI ^e	6	1	0	7	7	0	0	0	0	0

^a Detected using the Recombination Analysis Tool and / or Bellerophon bioinformatics screening methods among three different sets of genes of each *msg* family: full-length, full-length plus partial genes, full-length plus pseudogenes.

^b The number of potential mosaic genes among the *msg* families was almost significantly different (P = 0.06, Chi-square test).

^c Six full-length mosaic genes were detected twice but with different pairs of putative full-length parent genes according to the set of genes analysed (four, one, one of respectively family I, II, III). One mosaic gene of family I was detected twice: once with one full-length gene and one pseudogene as parents, and once with two partial genes as parents. All ten remaining were detected once with a pair of full-length parents.

^d Six mosaic pseudogenes of family I had two pseudogenes as parents. Two of family I had one full-length gene and one pseudogene as parents. The three remaining had a pair of full-length parents.

^e Several potential recombination events were detected for these two families using the more sensitive method TOPALi based on the Hidden Markov Model (Fig. S9).

1

2 **Comparison to the *msg* superfamily previously proposed**

3 The 146 *P. jirovecii* *msg* genes larger than 1.6kb reported by Ma et al (2016), out of a total of
4 179, were added into our DNA phylogenetic tree. They all clustered within our families, except
5 11 outliers (Fig. S11). The correspondence between the two sets of families is given in Table 1.

6 The comparison of the two studies is detailed in the Supplementary note 6.

7

8 **Discussion**

9 Antigenic surface variation plays a crucial role in escaping the human immune system and
10 adhering to host cells for important microbial pathogens. In the present study, we unravelled the
11 mechanisms used by the fungus *P. jirovecii* for this purpose. Our observations show that its
12 surface glycoproteins diversified during the evolution into a superfamily including six families
13 each with its own structure, function, independent mosaicism, and expression mode.

14

15 **Structure and function of Msg glycoproteins**

16 Proteins of Msg family I were previously demonstrated to adhere human epithelial cell through
17 binding to fibronectin and vitronectin (Pottratz et al. 1991; Limper et al. 1993). The ST-rich
18 regions present in *P. jirovecii* Msg glycoproteins except those of family IV are sites of oxygen-
19 linked glycosylation commonly involved in cell to cell adhesion (Dranginis et al. 2007).
20 Moreover, most of these glycoproteins were predicted to be adhesins (Supplementary note 7).
21 Consistently, their structure fits the model of modular organization of fungal adhesins with ST-
22 rich regions at the C-terminus and a ligand binding domain at the N-terminus (Dranginis et al.
23 2007; Linder and Gustafsson 2008). Linder and Gustafsson (2008) proposed that, in addition to
24 their role in adhesion, the oxygen-linked glycosylations of the ST-rich region confer rigidity to
25 the protein in order to present outward the ligand domain. Thus, the N-terminus regions of the *P.*
26 *jirovecii* adhesins may correspond to ligand binding domains. The fate and function of the
27 glycoproteins of family IV remain enigmatic since they lack the ST-rich region, are only weakly
28 predicted as adhesins (Supplementary note 7), and may not be attached to the cell wall in
29 absence of a GPI anchor signal. The conserved leucines separated by two to six residues present
30 in all *msg* families are similar to leucine zipper motifs which are often involved in protein-
31 protein non-specific binding and protein dimerisation (Hakoshima 2005). This latter function is

32 also carried out by the PE-rich region present in *msg* family V and VI (Williamson 1994). The
33 conserved coiled-coil domains discovered in *Msg* families I to III are often involved in the
34 formation heteromultimers and protein complexes (Strauss and Keller 2008; Hitchcock-
35 DeGregori and Barua 2017). The unstructured regions at the C-terminus present in four *Msg*
36 families are not informative because these regions can have several different functions (Best
37 2017). These observations suggest that the *Msg* adhesins may form homo- or hetero-oligomers at
38 the cell surface, possibly implying a further level of antigen variation which has never been
39 envisaged so far.

40

41 **Mosaicism of *msg* genes**

42 Our observations suggest that a continuous and random creation of mosaic genes by homologous
43 recombinations occurs mostly, if not exclusively, within each *msg* family. Very interestingly
44 within the scope of protein annotation, this mechanism permits by itself to define the members of
45 a protein family without having to rely upon the cutting of a phylogenetic tree at an arbitrary
46 height. The frequency of these recombinations remains to be quantified precisely, but is likely to
47 be reduced in *msg* families V and VI. The genetic mechanisms involved in the creation of mosaic
48 genes may include a single homologous recombination leading to a telomere exchange, or two
49 homologous recombinations leading to a gene fragment conversion or exchange (models are
50 shown in Fig. S12a). Such recombinations could also produce partial genes if they occur between
51 homologous regions which are not located at the same position along the recombining genes.
52 Our results suggest that this is rare because we identified only three partial *msg* genes out of 113.
53 This conclusion is also consistent with the fact that different motifs are conserved along the
54 sequence of the *Msg* proteins of each family. Our data suggest that pseudogenes might also be
55 involved in the generation of mosaic genes, and thus might constitute a reservoir of sequences

56 that can be integrated into functional antigens. The pseudogenes may result from accumulation
57 of mutations in absence of expression and thus of selective pressure. This phenomenon could be
58 enhanced by mutation and recombination rates within the subtelomeric gene families higher than
59 in the rest of the genome (Barry et al. 2003). The presence of the pseudogenes in the
60 subtelomeres might simply correspond to the state between their birth and their future decay.
61 However, they could also be maintained within the subtelomeres through indirect selective
62 pressure because of their role as reservoir of sequences for the creation of mosaic genes.

63

64 **Mutually exclusive expression of *msg-I* genes**

65 Our conclusions concerning the mutually exclusive expression of the *msg-I* genes are in
66 agreement with previous studies, but bring support for the involvement of telomere exchange
67 which has been previously hypothesized (Sunkin and Stringer 1996). The exchange of the single
68 expressed gene by recombination at the CRJE sequences might be facilitated by the localization
69 of the *msg-I* genes closest to the telomeres because this may in turn facilitate telomere exchanges
70 (a model is shown in Fig. S12b). These recombinations could be homologous in nature because
71 the full identity over 33 bps might be sufficient as it is the case in fungal cousins (Hua et al.
72 1997). However, they could also be site-specific because the imperfect inverted repeat present in
73 the CRJE is a common motif used by site-specific recombinases (Turan and Bode 2011). Up to
74 three *msg-I* genes were present at the end of the subtelomeres. There is no reasons to exclude that
75 transfer of more than one *msg-I* gene to the expression site at once also occurs, followed by
76 polycistronic expression. The polypeptide produced could be then chopped by the endoprotease
77 Kex1 at the end of each CRJE and each *Msg-I* anchored to the cell wall separately through its
78 own GPI signal. Interestingly, we detected *msg-I* pseudogenes linked to the UCS using PCR in
79 our sample. The cells expressing such truncated antigens may not be selected over time during

80 the infection because of their likely deficiency in adhesion to the host cells. They might
81 constitute a cost inherent to such system of antigenic variation based on frequent recombination
82 events.

83

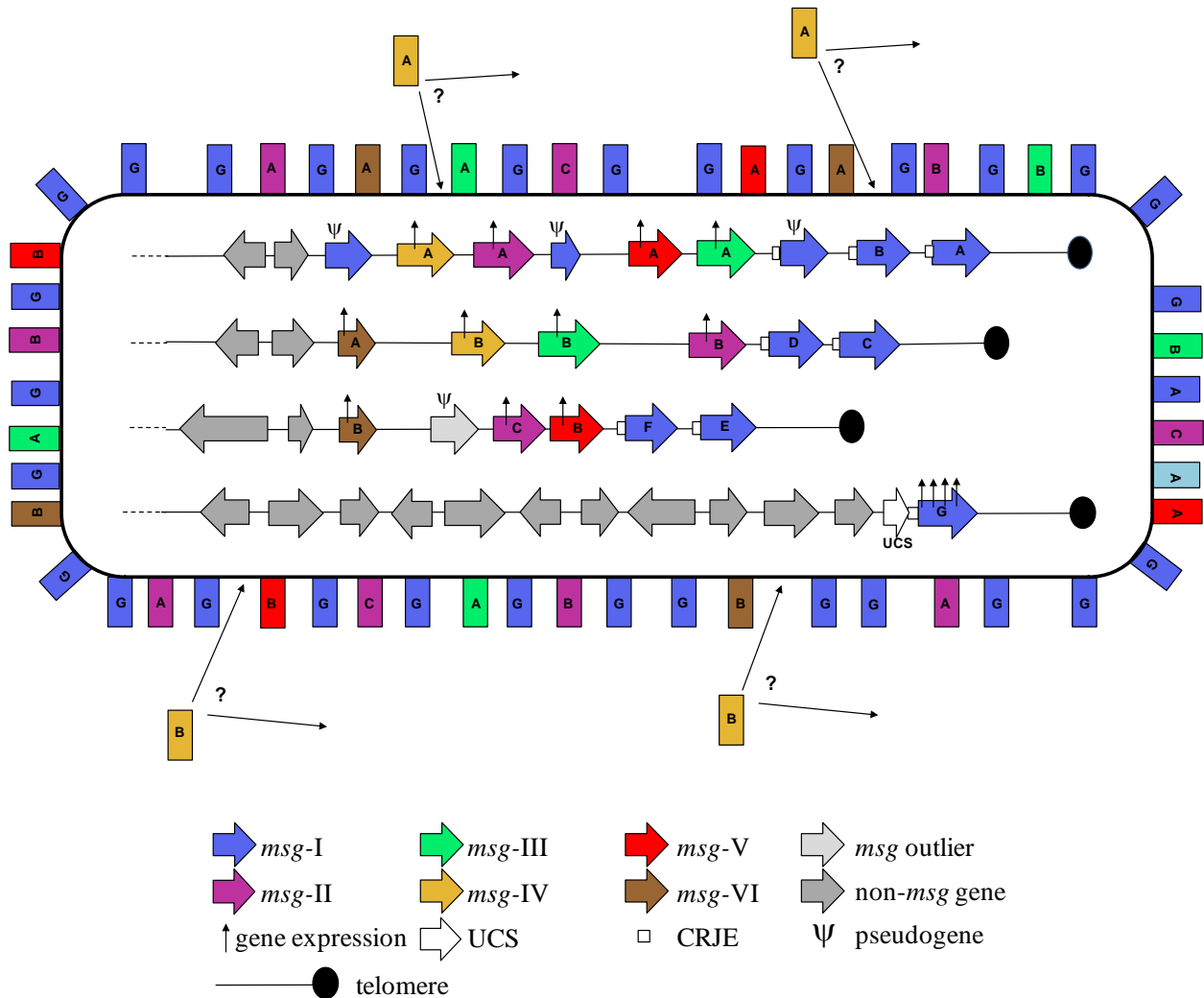
84 **Expression of *msg* families**

85 RNAseq analyses suggested that the vast majority of the *msg* genes of all families were
86 expressed in *P. carinii* and *P. murina* populations (Ma et al. 2016). As far as *P. jirovecii* is
87 concerned, alignment of our previous RNAseq data (Cissé et al. 2012) with the subtelomeres
88 assembled in the present study was compatible with the same conclusion, although the data were
89 from different clinical isolates (results not shown). Expression of most *msg*-I genes at the
90 population level is consistent with the numerous subpopulations of cells expressing different
91 *msg*-I genes that we observed. As far as *msg* families II to VI are concerned, the RNAseq data
92 are compatible with constitutive or temporally regulated expression of all genes in each cell
93 driven by the promoter present upstream of each of these genes. However, they are also
94 compatible with mutually exclusive or partially exclusive expression of these genes thanks to
95 silencing of promoters, or through another unknown mechanism.

96

97 **Model of expression of *msg* genes and cell surface structure**

98 In absence of data at the single cell level which would unravel the mode of expression of *msg*
99 genes, we propose a model of cellular expression of the *msg* families and cell surface structure
100 (Fig. 5). This model is based on the working hypothesis that all *msg* genes except those of family
101 I are expressed constitutively in both trophic forms and asci. The UCS is a strong promoter
102 (Kutty et al. 2013), probably leading to a majority of a single isoform of adhesive Msg-I antigens
103 on the cellular surface. This is consistent with the fact that Msg-I antigens are the most abundant



104
105

106

107

108 **Fig. 5**

109 Model of expression of *msg* genes and cell surface structure. The level of gene expression is
 110 figured by the number of arrows. The different isoforms of each *msg* family are differentiated by
 111 capital letters. The fate of the *Msg-IV* proteins remains to be determined (see text). The INT1
 112 and other surface proteins are not figured.

113

114 proteins in *Pneumocystis* species. All the different isoforms of adhesive glycoproteins of the
115 other families than family I present in the cell would be expressed and present on the cell surface
116 at the same time. The surface of *P. carinii* trophic cells was shown to harbour also the protein
117 INT1 participating to adhesion (Kottom et al. 2008). Recently, a transcription factor responsible
118 for expression of (a) still unidentified adhesive surface protein(s) has been reported in *P. carinii*
119 trophic cells (Kottom and Limper 2016). Genes encoding orthologs of these proteins are also
120 present in *P. jirovecii* genome (results not shown). Moreover, Kottom and Limper (2016)
121 mentioned that other uncharacterized genes which are important in binding to mammalian hosts
122 are present in *P. carinii* genome. Thus, the structure of the cell surface of both trophic forms and
123 asci is made of a complex mixture of different glycoproteins. The latter conclusion is valid
124 whatever the mode of expression of families II to VI is in reality.

125

126 **Strategy of antigenic variation**

127 The exchange of the *msg-I* isoform expressed and the generation of new mosaic genes of all *msg*
128 families would lead to a continuous segregation of subpopulations with a new mixture of
129 glycoproteins at the cell surface. Thus, the strategy of the fungus would consist in the continuous
130 generation of cells which are antigenically different. This strategy is further suggested by other
131 characteristics of *Pneumocystis* spp. First, there is a high variability of the subtelomeres between
132 *P. jirovecii* isolates⁴ which is consistent with frequent subtelomeric recombinations. The
133 subtelomeres of the isolate we studied here also differed greatly from those of the same
134 chromosomes reported by Ma et al (2016) (Supplementary note 6). Second, sexuality could be
135 obligatory in the cell cycle (Cushion and Stringer 2010; Hauser 2014) because most ectopic
136 recombinations between subtelomeres occur during meiosis, within the bouquet of telomeres
137 formed (Barry et al. 2003). The likely homothallic sexuality of *Pneumocystis* spp (Almeida et al.

138 2015) avoids the need to find a compatible partner and thus increases mating frequency, which is
139 believed to favor genetic diversity (Roach and Heitman 2014). Moreover, the genetic diversity
140 might be enhanced by mating between the numerous co-infecting strains which are generally
141 present in *P. jirovecii* infections (Alanio et al. 2016). Third, the presence of several *msg* families
142 may allow the formation of Msg hetero-oligomers that we envisage above, which would further
143 enhance the cell surface complexity.

144

145 **Strategies of antigenic variation in different human pathogens**

146 The mechanisms and hypothesised strategy of antigenic variation unravelled here appear unique
147 among human pathogens. *Candida glabrata* contain one subtelomeric family of ca. 20 adhesins
148 (Deitsch et al. 2009). *Trypanosoma brucei* presents a large reservoir of sequences used to create
149 mosaic genes of a single surface antigen family made of about a thousand of genes located in
150 subtelomeres as well as on minichromosomes (Deitsch et al. 2009). In the latter organism,
151 pseudogenes provide segments to mosaic functional antigens (Hall et al. 2013), a phenomenon
152 which might also occur in *P. jirovecii*. *Plasmodium falciparum* harbours one subtelomeric
153 antigen family of ca. 60 members (Deitsch et al. 2009). These three organisms present a single
154 gene family subject to mutually exclusive expression involving silencing in several cases. Thus,
155 their populations are homogenous antigenically but may vary over time when the expressed gene
156 is exchanged. Such strategy might be imposed by sterile niches such as blood and urinary tract.
157 This contrasts sharply with the putative strategy of antigenic variation of *P. jirovecii* consisting
158 in the continuous production of a mixture of cells antigenically different. The latter strategy may
159 be associated to the particular niche within lungs since it tolerates the presence of low abundant
160 fungi as members of the natural lung microbiota. This strategy might allow presenting most cells
161 as different organisms to the immune system and thus to be tolerated during colonisation. A

162 similar strategy might be used by *Candida albicans* living in non-sterile mucosal niches. Indeed,
163 its unique adhesin family presents a high number of serine CUG codons which are ambiguously
164 translated into serine or leucine, thus creating variability from individual genes (Rizzetto et al.
165 2015).

166 *Trypanosoma* and *Plasmodium* also differ from *Pneumocystis* spp in that they infect two
167 different hosts rather than one. This undoubtedly exerts a different selective pressure on their
168 antigenic variation system. The *Pneumocystis* spp differ considerably in their *msg* families (Ma
169 et al. 2016), as well as in the fine structure of the Msg adhesins (Mei et al. 1998). It is likely that
170 these differences are involved in the strict host species specificity of these fungi. Further work
171 aiming at understanding the relation between structure and function of the different Msg
172 glycoproteins is needed to further decipher both antigenic variation and host specificity of these
173 fungi.

174

175 **Methods**

176

177 **Bronchoalveolar lavage fluid specimens.** Fresh BALFs positive for *P. jirovecii* using
178 Methenamine-silver nitrate staining (Musto et al. 1982) were supplemented with 15% v/v
179 glycerol, frozen in liquid Nitrogen, and stored at -80°C. Only those with more than one ml
180 available and heavy fungal load were stored. Seventeen specimens were stored between 2012
181 and 2014, and used for the selection procedure described here below.

182

183 **DNA extraction and identification of an infection with a single *P. jirovecii* strain.** Genomic
184 DNA was extracted from 0.2 to 0.4 ml of BALF specimen using QIAamp® DNA Mini kit
185 (Qiagen), and resuspended in 50 µl of elution buffer. Four genomic regions were amplified by
186 PCR from genomic DNA extracted as described previously (Hauser et al. 1997). Each PCR
187 product was cloned into the plasmid pCR™4-TOPO using the TOPO TA cloning Kit for
188 Sequencing (Life Technologies). Both strands of the insert of 15 clones for each genomic region
189 were sequenced with M13 primers using the BigDye Terminator kit and the ABI Prism 3100
190 automated sequencer (both from PerkinElmer Biosystems). Among the 17 clinical specimens
191 collected, only one generated identical sequences for all clones of all genomic regions. Since ca.
192 15 clones per genomic region were analyzed, a second eventual co-infecting strain in this
193 specimen should not represent more than ca. 7% of the *P. jirovecii* population. This specimen
194 was selected for all experiments performed in the present study. It was from a HIV-infected
195 patient.

196

197 **Enrichment in *P. jirovecii* DNA and random amplification.** The DNA of the selected
198 specimen was enriched in *P. jirovecii* DNA using the NEBNext® Microbiome DNA Enrichment

199 Kit based on the absence of CpG methylation (Biolabs), purified by ethanol precipitation in
200 presence of 10 µg glycogen (Thermo Fisher Scientific), and resuspended in 50 µl of 1X TE
201 Buffer. This enrichment raised the proportion of *P. jirovecii* DNA from a few percent to ca. 55%
202 as determined *a posteriori* by high throughput sequencing. Because only small amounts of DNA
203 are recoverable from a clinical specimen and in absence of an *in vitro* culture system, sufficient
204 amount of DNA for high throughput PacBio sequencing was obtained by random amplification.
205 Five µl of DNA was randomly amplified in a 50 µl reaction using the Illustra GenomiPhi HY
206 DNA Amplification Kit (GE Healthcare). This amplification proved to create artificial molecules
207 made of inverted repeats of several kb which were revealed by PacBio sequencing. The reads
208 from these molecules were eliminated by bioinformatics (see below). DNA was then purified
209 using QIAamp® DNA blood mini kit (Qiagen) followed by ethanol precipitation in presence of
210 10 µg glycogen. Amplified DNA fragments were sized (mean 8.6 kb) and quantified using
211 Fragment AnalyzerTM (Advanced Analytical).

212

213 **High throughput PacBio sequencing.** Five µg of amplified DNA were used to prepare a
214 SMRTbell library with the PacBio SMRTbell Template Prep Kit 1 according to the
215 manufacturer's recommendations (Pacific Biosciences). The resulting library was size selected
216 on a BluePippin system (Sage Science) for molecules larger than 5 kb. The recovered library was
217 sequenced on one SMRT cell with P6/C4 chemistry and MagBeads on a PacBio RSII system
218 (Pacific Biosciences) at 240 min movie length.

219

220 **Read filtering and *P. jirovecii* genome assembly.** The flow chart of the filtering and assembly
221 procedure is shown in Figure S13a and the details for each step are described here. PacBio sub-
222 reads were extracted from the raw h5-files using DEXTRACTOR

223 (<https://github.com/thegenemyers/DEXTRACTOR/>). The average length of the extracted sub-
224 reads was 5.2 kb with a maximum of 42 kb. We removed human derived reads by mapping them
225 against the human reference genome using blasr (smrtpipe2.3, cut-off: corrected score < 55000).
226 Reverse-complementary artificial reads created by the random amplification were next filtered
227 out (cut-off:match length >=1000 bps) after mapping them onto themselves using DALIGNER
228 (<https://github.com/thegenemyers/DALIGNER/>)(V1.0, options:-A -I). The cleaned reads were
229 assembled using the tool FALCON (Chin et al. 2016)(V0.2, options: length_cutoff=8000m
230 length_cutoff_pr=1000). PacBio reads were re-mapped onto the assembly using BLASR and
231 used to evaluate and flag remaining human contigs. Human derived contigs were subsequently
232 removed. A total of 2.2 Gb of *P. jirovecii* DNA sequences corresponding to a 200-fold coverage
233 of the genome were gathered. The assembly was polished to remove residual PacBio errors using
234 Quiver (Chin et al. 2013)(smrtpipe2.3, 5 iterations). The final polished genome assembly
235 included 8.1 Mb in 219 gap-free contigs ranging from 234 bps to 386 kb with a NG50 of 108Kb,
236 and 57% of the genome in 28 contigs larger than 100 kb. The *P. jirovecii* PacBio assembly
237 obtained in the present study covered 96% of that we previously obtained using other sequencing
238 methods⁶, and contained ca. 0.5 Mbp of subtelomeric sequences. The combination of both our
239 assemblies covered 97% the assembly of Ma et al (2016). Controls consisting in PCR
240 amplification of specific subtelomeric regions from the same DNA sample confirmed the
241 accuracy of the nucleotide sequence of the polished PacBio assembly, although few errors in
242 repetitive homopolymer regions were detected (Supplementary note 8).

243

244 **Gene predictions and msg annotations.** Genes were predicted on the assembly using Augustus
245 (Stanke et al. 2006)(version 2.5.5) and a specifically trained model for *Pneumocystis* (Cissé et al.
246 2012). In order to detect novel and more distant homologous *msg* genes in the assembly, we

247 chose a generalized profile based approach (Bucher and Bairoc 1994)(Fig. S13b). A DNA profile
248 was generated based on a previously described *msg* gene in *P. carinii* (Wada and Nakamura
249 1996)(GenBank D82031.1) and a protein profile based on Msg-Rucl 21 (European Nucleotide
250 Archive ABQ51002.1) using a Smith-Waterman-Algorithm (Smith and Waterman 1981). The
251 profiles were calibrated against the scrambled genome (window approach, size=60). Using
252 pfsearchV3 (Schuepbach et al. 2013), the assembled genome was searched for homologues
253 matches with the DNA profile. Curated matches were extracted and aligned against each other
254 using MAFFT (Katoh and Standley 2013)(version 7.305). After manual curation and trimming,
255 the alignments were divided in five groups based on neighbourhood joining (% identity) using
256 Jalview (Waterhouse et al. 2009)(v2.8.1). One representative candidate per group was selected
257 and a new profile based on its sequence generated and calibrated as described here above. These
258 DNA *msg* profiles were used to find and annotate a first set of 75 *msg* genes in the assembly. A
259 combination of Blastx, genewise, in-house tools, and manual curation was applied using the
260 protein Msg profile to extend and correct these annotations to the set of 113 *msg* genes analysed
261 in the present study. The *msg* genes reported here were all manually curated with respect to their
262 start, stop and intron coordinates.

263

264 **Construction of phylogenetic trees.** For the DNA and protein based phylogenetic analysis, the
265 CDS for each annotated *msg* gene was manually corrected (up to five corrections), extracted, and
266 translated into its protein sequence. Both CDS and protein sequences were aligned against each
267 other using MAFFT (Katoh and Standley 2013)(mafft-linsi -genafpair), and the multiple
268 sequence alignment used to infer a phylogenetic tree with RAxML (Stamatakis
269 2014)(PROTGAMMAGTR for proteins and with GTRGAMMA for CDS, 1000 bootstraps). The
270 *msg* genes of family V were defined as out-group and the final tree rooted. Proteins were further

271 classified using JACOP (Sperisen and Pagni 2005)(<http://myhits.isb-sib.ch/cgi-bin/jacop/>). In
272 order to add pseudogenes and published *msg* genes from Ma et al (2016) equal or exceeding 1.6
273 kb, we injected the new sequences into the prior DNA-based multiple-alignment using MAFFT
274 (Katoh and Standley 2013)(--addfull). They were added to the original tree using the
275 evolutionary placement algorithm (EPA) from RAxML. These trees were converted into a
276 compatible format with the tool guppy from the pplacer suite (Matsen et al. 2010) (v1.1alpha14,
277 tog). Genes with an exon smaller than 1.6 kb were added to the original DNA-based multiple-
278 alignment using MAFFT (Katoh and Standley 2013)(--addfragments). The alignment was
279 trimmed and re-aligned using MAFFT (Katoh and Standley 2013). A new tree was then build
280 with RAxML (GTRGAMMA, 1000 bootstraps). All trees were analyzed and visualized using R
281 (R Core Team 2013)(3.3.2) and GGTREE (Yu et al. 2017)(v1.6.9).

282

283 **Gene and protein sequences analyses.** Alignments of full-length gene or protein sequences
284 were carried out using MAFFT (Katoh and Standley 2013). Canonical TATA box and Cap signal
285 (Bucher 1990), as well as canonical donor and acceptor sequences of *Pneumocystis* introns
286 (Thomas et al. 1999; Slaven et al. 2006), were identified by visual inspection of the alignments
287 and sequences of the genes. Signal peptide and GPI anchor signal were identified using
288 respectively Phobius (Käll et al. 2004)(<http://phobius.sbc.su.se/>) and GPI-SOM (Frankhauser and
289 Mäser 2005)(<http://gpi.unibe.ch/>) with default settings. Canonical potential sites NXS/T of
290 Nitrogen-linked glycosylation (Linder and Gustafsson 2008) were identified by visual inspection.
291 Conserved domains were searched using Multiple Expectation–Maximization for Motif
292 Elicitation (Bailey and Elkan 1994)(MEME, <http://meme-suite.org/tools/meme>). MEME analysis
293 of the 49 full-length *Msg* proteins of all families except outliers was carried out using default
294 settings, except minimum and maximum motif width of respectively 50 and 100 residues, any

295 number of sites per sequence option, and maximum of 13 motifs searched. HMMER (Finn et al.
296 2011)(biosequence analysis using profile hidden Markov models,
297 <http://www.ebi.ac.uk/Tools/hmmer/search/hmmscan>) was used with default settings on full-
298 length proteins for the following embedded predictions: Pfam, unstructured regions (Intrinsically
299 Unstructured Proteins, IUPRED), and coiled-coil motifs (Ncoils predictor). Pairwise identities
300 between full-length *msg* genes and Msg proteins were calculated using the multi-way alignment
301 type of Clone Manager 9 professional edition software.

302

303 **Search for potential mosaic genes.** Two screening methods were first used: Recombination
304 Analysis tool (Etherington et al. 2005)(RAT, <http://cbr.jic.ac.uk/dicks/software/RAT/>) and
305 Bellerophon (Huber et al. 2004)(<http://comp-bio.anu.edu.au/bellerophon/bellerophon.pl>).
306 MAFFT (Katoh and Standley 2013) alignments of various set of genes were analysed with both
307 methods. RAT was used with default settings, *i.e.* using windows of one tenth of the length of
308 the alignment and increment size equal to half of the window size. Bellerophon was used with
309 default settings, *i.e.* windows of 300 bps and Huber-Hugenholtz correction. RAT can detect
310 several recombination events whereas Bellerophon reports a single one per mosaic gene. The
311 more sensitive method TOPALi v2.5 (Milne et al. 2004)(<http://www.topali.org/>) which is based
312 on a Hidden Markov Model (HMM) was then applied on the potential mosaic genes and its
313 putative parent genes detected with the two screening methods. These three genes were aligned
314 using MAFFT (Katoh and Standley 2013) with an additional gene chosen randomly in the same
315 *msg* family since TOPALi requires four genes in input. The efficacy of the three methods to
316 detect mosaic genes was assessed by the analysis of artificial chimera produced *in silico* with
317 related genes, as well as with sets of orthologous genes from different fungal species (results not
318 shown). Only the RAT method is suitable for the search of recombination events among proteins.

319 The vast majority of the events detected at the protein level corresponded to those detected at the
320 DNA level (results not shown).

321

322 **PCR amplification and sequencing.** PCRs were performed in a final volume of 20 μ l with 0.35
323 U of High Fidelity Expand polymerase (Roche Diagnostics), using the buffer provided, each
324 dNTP at a final concentration of 200 μ M, and each primer at 0.4 μ M. PCR conditions included
325 an initial denaturation step of 3 min at 94°C, followed by 35 cycles consisting of 30 s at 94°C, 30
326 s at the annealing temperature, and 1 min per kb to be amplified at 72°C. The reaction ended
327 with 5 min of extension at 72°C. The annealing temperature and the MgCl₂ concentration were
328 optimized for each set of primers and ranged from 51 to 60°C and from 3 to 6 mM, respectively.
329 Sequencing both strands of the PCR products was performed with the two primers used for PCR
330 amplification, as well as the Big Dye Terminator DNA sequencing kit and ABI PRISM 3100
331 automated sequencer (both from Perkin-Elmer Biosystems).

332

333 **Ethics**

334 The protocol was approved by the institutional review board (Commission cantonale d'éthique de
335 la recherche sur l'être humain). All patients provided an informed written consent which was part
336 of procedure for the admittance in the hospital. This consent was documented by the fact that
337 they did not ask their samples not to be used for research. The samples were treated
338 anonymously.

339

340 **Availability of data and materials**

341 PacBio raw reads (accession SRR5533719) and PacBio assembly (accession NJFV000000000)
342 have been deposited at DDBJ/ENA/GenBank linked to BioProject PRJNA382815 and

343 BioSample SAMN06733346. The version of the PacBio assembly described in this paper is
344 version NJFV01000000.

345

346 **Acknowledgments**

347 Computations were performed at the Vital-IT Center for High-Performance Computing of the
348 Swiss Institute of Bioinformatics (<http://www.vital-it.ch>). We thank Michel Monod, Dominique
349 Sanglard, and Laurent Keller for critical reading.

350

351 **Supplementary information**

352 The table of contents of the Supplementary information is in Additional file 2.docx, together
353 with supplementary notes and figures. Supplementary Tables are in Additional file 1.docx.

354

355 **Funding**

356 This work was supported by the Swiss National Science Foundation, grant 310030_146135 to
357 P.M.H. and M.P. This Foundation had not role in any steps of the study.

358

359 **Authors' contributions**

360 P.M.H. and M.P. designed the study. S.R. and A.L. collected and prepared the samples, and
361 carried out the laboratory experiments. K.M. collected and prepared samples. E.S.-S. performed
362 genome assembly. E.S.-S. and P.M.H. performed bioinformatics analyses. E.S.-S., P.M.H., S.R.,
363 A.L., and M.P. analysed the data. P.M.H., E.S.-S., and M.P. wrote the manuscript.

364

365 **Competing interests**

366 The authors declare no competing interests.

367 **References**

368

369 Alanio A, Gits-Muselli M, Mercier-Delarue S, Dromer F, Bretagne S. 2016. Diversity of
370 *Pneumocystis jirovecii* during infection revealed by ultra-deep pyrosequencing. *Front*
371 *Microbiol* **7**: 733.

372 Almeida JMGCF, Cissé OH, Fonseca Á, Pagni M, Hauser PM. 2015. Comparative Genomics
373 suggests Primary Homothallism of *Pneumocystis* species. *Mbio* **6**: e02250-e02214.

374 Bailey TL, Elkan C. 1994. Fitting a mixture model by expectation maximization to discover
375 motifs in biopolymers. *Proc Sec Int Conf Int Syst Mol Biol* 28-36 (AAAI Press, Menlo Park,
376 California).

377 Barry JD, Ginger ML, Burton P, McCulloch R. 2003. Why are parasite contingency genes often
378 associated with telomeres? *Int J Parasit* **33**: 29-45.

379 Best RB. 2017. Computational and theoretical advances in studies of intrinsically disordered
380 proteins. *Curr Opin Struct Biol* **42**: 147-154.

381 Brown GD, Denning DW, Gow NAR, Levitz SM, Netea MG, White TC. et al. 2012. Hidden
382 killers: human fungal infections. *Sci Transl Med* **4**: 165rv13.

383 Bucher P. 1990. Weight matrix descriptions of four eukaryotic RNA polymerase II promoter
384 elements derived from 502 unrelated promoter sequences. *J Mol Biol* **212**: 563-578.

385 Bucher P, Bairoc AA. 1994. Generalized profile syntax for biomolecular sequence motifs and its
386 function in automatic sequence interpretation. *Proc Int Conf Intell Syst Mol Biol* **2**: 53-61.

387 Chin CS, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, et al. 2016. Phased
388 diploid genome assembly with single-molecule real-time sequencing. *Nat Meth* **13**: 1050-
389 1054.

- 390 Chin CS, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, et al. 2013. Nonhybrid,
391 finished microbial genome assemblies from long-read SMRT sequencing data. *Nat Meth* **10**:
392 563-569.
- 393 Cissé OH, Pagni M, Hauser PM. 2012. *De novo* assembly of the *Pneumocystis jirovecii* genome
394 from a single bronchoalveolar lavage fluid specimen from a patient. *MBio* **4**: e00428-00412.
- 395 Cushion MT, Smulian AG, Slaven BE, Sesterhenn T, Arnold J, Staben C, et al. 2007.
396 Transcriptome of *Pneumocystis carinii* during fulminate infection: carbohydrate metabolism
397 and the concept of a compatible parasite. *PLoS ONE* **2**: e423.
- 398 Cushion MT, Stringer JR. 2010. Stealth and opportunism: alternative lifestyles of species in the
399 fungal genus *Pneumocystis*. *Annu Rev Microbiol* **64**: 431-452.
- 400 Deitsch KW, Lukehart SA, Stringer JR. 2009. Common strategies for antigenic variation by
401 bacterial, fungal and protozoan pathogens. *Nat Rev* **7**: 493-503.
- 402 Dranginis AM, Rauceo JM, Coronado JE, Lipke PN. 2007. A Biochemical Guide to Yeast
403 Adhesins: Glycoproteins for Social and Antisocial Occasions. *Microbiol Mol Biol Rev* **71**:
404 282-294.
- 405 Etherington GJ, Dicks J, Roberts IN. 2005. Recombination Analysis Tool (RAT): a program for
406 the high-throughput detection of recombination. *Bioinfo* **21**: 278-281.
- 407 Fankhauser N, Mäser P. 2005. Identification of GPI anchor attachment signals by a Kohonen
408 self-organizing map. *Bioinfo* **21**: 1846-1852.
- 409 Finn RD, Clements J, Eddy SR. 2011. HMMER web server: interactive sequence similarity
410 searching. *Nucl Ac Res* **39**: W29-W37.
- 411 Hakoshima T. 2005. Leucine Zippers. eLS.
- 412 Hall JPJ, Wang H, Barry JD. 2013. Mosaic *VSGs* and the scale of *Trypanosoma brucei* antigenic
413 variation. *PLoS Pathog* **9**: e1003502.

- 414 Hauser PM. 2014. Genomic insights into the fungal pathogens of the genus *Pneumocystis*:
415 obligate biotrophs of humans and other mammals. *PLoS Pathog* **10**: e1004425.
- 416 Hauser PM, Francioli P, Bille J, Telenti A, Blanc DS. 1997. Typing of *Pneumocystis carinii* f.
417 sp. *hominis* by single-strand conformation polymorphism of four genomic regions. *J Clin*
418 *Microbiol* **35**: 3086-3091.
- 419 Hitchcock-DeGregori SE, Barua B. 2017. Tropomyosin structure, function, and interactions: a
420 dynamic regulator. *Sub Biochem* **82**: 253-264.
- 421 Hua SB, Qiu M, Chan E, Zhu L, Luo Y. 1997. Minimum length of sequence homology required
422 for *in vivo* cloning by homologous recombination in yeast. *Plasmid* **38**: 91-96.
- 423 Huber T, Faulkner G, Hugenholtz P. Bellerophon: a program to detect chimeric sequences in
424 multiple sequence alignments. *Bioinfo.* 2004;20:2317-9.
- 425 Käll L, Krogh A, Sonnhammer ELL. 2004. A Combined Transmembrane Topology and Signal
426 Peptide Prediction Method. *J Mol Biol* **338**: 1027-1036.
- 427 Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7:
428 improvements in performance and usability. *Mol Biol Evol* **30**: 772-780.
- 429 Keely SP, Renauld H, Wakefield AE, Cushion MT, Smulian AG, Fosker N, et al. 2005. Gene
430 arrays at *Pneumocystis carinii* telomeres. *Genetics* **170**:1589–1600.
- 431 Keely SP, Stringer, JR. 2009. Complexity of the MSG gene family of *Pneumocystis carinii*.
432 *BMC Gen* **10**: 367.
- 433 Kottom TJ1, Kennedy CC, Limper AH. 2008. *Pneumocystis PCINT1*, a molecule with integrin-
434 like features that mediates organism adhesion to fibronectin. *Mol Microbiol* **67**: 747-761.
- 435 Kottom TJ, Limper AH. 2016. Evidence for a *Pneumocystis carinii* Flo8-like transcription
436 factor: insights into organism adhesion. *Med Microbiol Immunol* **205**: 73-84.

- 437 Kutty G, Ma L, Kovacs JA. 2001. Characterization of the expression site of the major surface
438 glycoprotein of human-derived *Pneumocystis carinii*. *Mol Microbiol* **42**: 183-193.
- 439 Kutty G, Shroff R, Kovacs JA. 2013. Characterization of *Pneumocystis* major surface
440 glycoprotein gene (*msg*) promoter activity in *Saccharomyces cerevisiae*. *Euk Cell* **12**: 1349-
441 1355.
- 442 Kutty G, Maldarelli F, Achaz G, Kovacs JA. 2008. Variation in the major surface glycoprotein
443 genes in *Pneumocystis jirovecii*. *J Infect Dis* **198**: 741-749.
- 444 Linder T, Gustafsson CM. 2008. Molecular phylogenetics of ascomycotal adhesins—a novel
445 family of putative cell-surface adhesive proteins in fission yeasts. *Fung Gen Biol* **45**: 485-
446 497.
- 447 Limper AH, Standing JE, Hojman OA, Castro M, Neese, LW. 1993. Vitronectin binds to
448 *Pneumocystis carinii* and mediates organism attachment to cultured lung epithelial cells.
449 *Infect Immun* **61**: 4302-4309.
- 450 Ma L, Chen Z, Wei Huang D, Kutty G, Ishihara M, Wang H, et al. 2016. Genome analysis of
451 three *Pneumocystis* species reveals adaptation mechanisms to life exclusively in mammalian
452 hosts. *Nat com* **7**: 10740.
- 453 Ma L, Kutty G, Jia Q, Imamichi H, Huang L, Atzori C, et al. 2002. Analysis of variation in
454 tandem repeats in the intron of the major surface glycoprotein expression site of the human
455 form of *Pneumocystis carinii*. *J Infect Dis* **186**: 1547-1554.
- 456 Matsen FA, Kodner RB, Armbrust, EV. 2010. pplacer: linear time maximum-likelihood and
457 Bayesian phylogenetic placement of sequences onto a fixed reference tree. *BMC Bioinfo* **11**:
458 538.

- 459 Mei Q, Turner RE, Sorial V, Klivington D, Angus CW, Kovacs JA. 1998. Characterization of
460 major surface glycoprotein genes of human *Pneumocystis carinii* and high-level expression
461 of a conserved region. *Infect Immun* **66**: 4268-4273.
- 462 Milne I, Wright F, Rowe G, Marshal DF., Husmeier D, McGuire G. 2004. TOPALi: software for
463 automatic identification of recombinant sequences within DNA multiple alignments. *Bioinfo*
464 **20**: 1806-1807.
- 465 Musto L, Flanigan M, Elbadawi A. 1982. Ten-minute silver stain for *Pneumocystis carinii* and
466 fungi in tissue sections. *Arch Pathol Lab Med* **106**: 292-294.
- 467 Pottratz ST, Paulsrud J, Smith JS, Martin WJ II. 1991. *Pneumocystis carinii* attachment to
468 cultured lung cells by *Pneumocystis* gp 120, a fibronectin binding protein. *J Clin Invest* **88**:
469 403-407.
- 470 R Core Team. 2013. R: a language and environment for statistical computing. R Foundation for
471 Statistical Computing, Vienna, Austria. URL <http://www.R-project.org/>.
- 472 Ramana J, Gupta, D. FaaPred: 2010. A SVM-based prediction method for fungal adhesins and
473 adhesin-like proteins. *PLoS ONE* **5**: e9695.
- 474 Rizzetto L, Weil T, Cavalieri D. 2015. Systems level dissection of *Candida* recognition by
475 lectins: a matter of fungal morphology and site of infection. *Pathog* **4**: 639-661.
- 476 Roach KC, Heitman J. 2014. Unisexual reproduction reverses Muller's ratchet. *Genetics* **198**:
477 1059-1069.
- 478 Schuepbach T, Pagni M, Bridge A, Bouqueleret L, Xenarios I, Cerutti L. 2013. pfsearchV3: a
479 code acceleration and heuristic to search PROSITE profiles. *Bioinfo* **29**: 1215-1217.
- 480 Slaven BE, Porollo A, Sesterhenn T, Smulian AG, Cushion MT, Meller J. 2006. Large-scale
481 characterization of introns in the *Pneumocystis carinii* genome. *J Eukar Microbiol* **53**: S151-
482 153.

- 483 Smith TF, Waterman MS. 1981. Identification of common molecular subsequences. *J Mol Biol*
484 **147**: 195-197.
- 485 Sperisen P, Pagni M. 2005. JACOP: a simple and robust method for the automated classification
486 of protein sequences with modular architecture. *BMC Bioinfo* **6**: 216.
- 487 Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of
488 large phylogenies. *Bioinfo* **30**: 1312-1313.
- 489 Stanke M, Keller O, Gundunz I, Hayes A, Waack S, Morgenstern B. 2006. AUGUSTUS: *ab*
490 *initio* prediction of alternative transcripts. *Nucl Ac Res* **34**: W435-439.
- 491 Strauss HM, Keller S. 2008. Pharmacological interference with protein-protein interactions
492 mediated by coiled-coil motifs. *Hand Exp Pharmacol* **186**: 461-482.
- 493 Stringer JR. 2007. Antigenic Variation in *Pneumocystis*. *J Eukaryot Microbiol* **54**: 8-13.
- 494 Sunkin SM, Stringer JR. 1996. Translocation of surface antigen genes to a unique telomeric
495 expression site in *Pneumocystis carinii*. *Mol Microbiol* **19**: 283-295 .
- 496 Thomas CF JR., Loef EB, Limper AH. 1999. Analysis of *Pneumocystis carinii* introns. *Infect*
497 *Immun* **67**: 6157-6160.
- 498 Turan S, Bode J. 2011. Site-specific recombinases: from tag-and-target- to tag-and-exchange-
499 based genomic modifications. *FASEB J* **25**:4088-4107.
- 500 Wada M, Nakamura Y. 1996. Unique telomeric expression site of major-surface-glycoprotein
501 genes of *Pneumocystis carinii*. *DNA Res* **3**: 55-64.
- 502 Waterhouse AM, Procter JB, Martin DM, Clamp M, Barton GJ. 2009. Jalview version 2-a
503 multiple sequence alignment and analysis workbench. *Bioinfo* **25**:1189-1191.
- 504 Williamson MP. 1994. The structure and function of proline-rich regions in proteins. *Biochem J*
505 **297**: 249-260.

506 Yu G, Smith DK, Zhu H, Guan Y, Lam TT-Y. 2017. GGTREE: an R package for visualization
507 and annotation of phylogenetic trees with their covariates and other associated data. *Meth*
508 *Ecol Evol* 8: 28-36.