

1 Revised version: clean.

2 **Expression pattern of the *Pneumocystis jirovecii***
3 **major surface glycoprotein superfamily in patients**
4 **with pneumonia**

5 Emanuel Schmid-Siegert ^{1,a}, Sophie Richard ², Amanda Luraschi ^{2,b}, Konrad
6 Mühlethaler ³, Marco Pagni ^{1*}, Philippe M. Hauser ^{2*}

7

8 * These authors contributed equally to the work.

9

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

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

13 ³ Institute for Infectious Diseases, University of Bern, Bern, Switzerland

14 ^a Present address: Selexis SA, Chemin des Aulx 14, 1228 Plan-les-Ouates, Switzerland

15 ^b Present address: Resistell AG, Muttenz, Switzerland

16

17 Correspondence : Philippe.Hauser@chuv.ch

18

19 **Running title:** *Pneumocystis* surface proteins expression (40 characters)

20

21 **Abstract :** 154 words

22 **Text :** 3499 words

23 **Brief summary (40 words)**

24 The six families of surface glycoproteins of the human pathogen *Pneumocystis jirovecii* were
25 studied by RNA sequencing. A number of different isoforms of most families were expressed
26 simultaneously, suggesting that individual cells with distinct surface of properties compose
27 each population.

28

29 **Competing interests**

30 The authors declare that no competing interests exist.

31

32 **Funding information**

33 This work was supported by the Swiss National Science Foundation, grant 310030_165825.

34 This Foundation had no role in any steps of the study.

35 **Abstract**

36 **Background.** The human pathogen *Pneumocystis jirovecii* harbors six families of major
37 surface glycoproteins (MSG) encoded by a single gene superfamily. MSGs are presumably
38 responsible for antigenic variation and adhesion to host cells. The genomic organization
39 suggests that a single member of family I is expressed at a time per cell, whereas members of
40 the other families are simultaneously expressed.

41 **Methods.** We analyzed RNA sequences expressed in several clinical samples, using specific
42 weighted profiles for reads sorting and single nucleotide variants calling to estimate the
43 diversity of the expressed genes.

44 **Results.** A number of different isoforms of at least four MSG families were expressed
45 simultaneously, including of family I for which confirmation was obtained in the wet
46 laboratory.

47 **Conclusion.** These observations suggest that every single *P. jirovecii* population is made of
48 individual cells with distinct surface properties. Our results enhance our understanding of the
49 unique antigenic variation system and cell surface structure of *P. jirovecii*.

50

51

52 **Keywords:** surface antigenic variation, RNAseq, virulence factor, colonization factor

53 **Introduction**

54 The fungal genus *Pneumocystis* belongs to the subphylum Taphrinomycotina of the
55 Ascomycota [1]. It encompasses extracellular parasites that colonize the lungs of mammals [2].
56 The species infecting humans is *Pneumocystis jirovecii*, while rats and mice harbor respectively
57 *Pneumocystis carinii* and *Pneumocystis murina*. Should the human immune system decline, *P.*
58 *jirovecii* can become an opportunistic pathogen that causes severe pneumonia, which can be
59 lethal if not treated (*Pneumocystis* pneumonia, PCP). An *in vitro* method of long-term culture
60 for *Pneumocystis* species is still not available.

61 The lack of a culture method complicates the study of *P. jirovecii* pathogenicity. Its major
62 surface glycoproteins (MSG) constitute potentially a crucial factor involved in colonization
63 and/or virulence. These proteins are thought to generate surface antigenic variation allowing
64 escape from the human immune system [3-5]. Besides, the MSGs are thought to be involved
65 in adhesion to host cells [6, 7]. MSGs are encoded by six families of hypervariable genes
66 located at all subtelomeric regions of the ca. 20 chromosomes of *P. jirovecii* [8-10]. These
67 families form the largest surface protein superfamily known among fungi with approximately
68 160 *msg* genes per genome (families I to VI include respectively ca. 85, 20, 10, 20, 20, and 5
69 genes encoding isoforms; family IV is the only one potentially not anchored in the cell surface
70 because of the lack of a glycosylphosphatidylinositol signal). MSGs of family I (MSG-I) are
71 thought to be the most abundant at the *P. jirovecii* cell surface [11], although only one *msg*-I
72 gene out of ca. 80 present in the genome is probably expressed at a time in every single cell.
73 Indeed, firstly, the latter mutually exclusive expression would rely on the expression of a single
74 gene under the control of a transcription promoter that is present at a single copy per genome
75 (within a upstream conserved sequence, UCS), whereas the other genes of family I have no
76 promoter. Secondly, *Pneumocystis* cells are mostly haploid except transiently during the sexual
77 cycle [12-15]. However, at the population level, several different *msg*-I isoforms linked to the

78 UCS at the DNA level are observed, and thus are presumably expressed [9, 16]. In *P. carinii*,
79 the diversity of *msg-I* isoforms expressed has been observed also at the RNA level [17], as well
80 at the protein level that revealed a focal distribution of the epitopes within the lung [18]. The
81 exchange of the expressed *msg-I* isoform would occur upon recombination at a sequence of 33
82 bps that is present both at the end of the UCS including the promoter and at the beginning of
83 each *msg-I* gene (the conserved recombination junction element, CRJE). By contrast, each gene
84 of the other five MSG families II to VI possesses its own promoter allowing potentially
85 independent and simultaneous expression [9]. In addition to the exchange of the *msg-I* isoform
86 expressed, antigenic variation is thought to rely on recombinations between the genes of each
87 MSG family that generate gene mosaicism [3, 9, 19].

88 We propose a model for the surface antigenic variation system of *P. jirovecii* consisting
89 in the continuous segregation of new subpopulations that are antigenically different (Figure 1).
90 The aim of the present study was to challenge this model by investigating the expression of the
91 different *P. jirovecii* MSG families at the RNA level in clinical samples of patients with PCP.
92 However, the classical approach of mapping RNA sequencing (RNAseq) reads onto genomic
93 sequences was inadequate for the following reasons. Firstly, ambiguous mapping of the
94 RNAseq reads and thus erroneous signals could result from (i) the repetitive nature of these
95 genes because they are made of conserved domains [8-10], and (ii) the presence of identical
96 sequences in several *msg* genes due to the recombinations creating mosaicism by gene
97 conversion. Secondly, the set of subtelomeres present in the cells of each subpopulation
98 probably varies considerably because (i) a different *msg-I* gene may be linked to the UCS, and
99 (ii) recombinations may occur frequently between *msg* genes. Consequently, only the most
100 stable part of the subtelomeres that is present in the majority of the cells can be assembled, at
101 least until single cell sequencing will be adapted to *P. jirovecii*. Thus, a complete set of
102 subtelomeres that could be used as a reference for mapping RNAseq reads does not exist, even

103 when analyzing the genome assembly and RNAseq reads from the same *P. jirovecii* isolate. In
104 order to circumvent these limitations, we developed dedicated bioinformatics procedures in
105 order to analyze the expression of the MSG genes.

106 **Material and methods**

107 **RNA extraction and whole transcriptome amplification**

108 Total RNAs were extracted from the bronchoalveolar lavage fluid specimens (BALFs; see
109 supplementary information) using the RiboPure yeast kit (Ambion). The whole transcriptome
110 was amplified from each RNA preparation using the SeqPlex[™] RNA Amplification Kit
111 (Sigma). The procedure resulted in cDNAs with a mean size of 100 to 150 bps as revealed
112 using a 2100 Bioanalyzer system (Agilent Technologies). According to the manufacturer, such
113 size, *i.e.* below 200-400 bps, indicates degradation of input RNA. Only patient 2's BALF led
114 to a mean cDNA size of 250 bps. RNA degradation was expected because of the uncontrolled
115 period between collection of the BALFs from the patient and arrival in our laboratory, as well
116 as the complex and varying microbiota present in these samples. We previously observed
117 varying RNA degradation in BALFs by amplifying specific transcripts [20, 21]. Because the
118 SeqPlex[™] RNA Amplification Kit generated too small cDNAs for most samples (100 to 150
119 bps), the absence of genomic DNA in the RNA preparations was checked on larger cDNAs
120 (ca. 800 bps) obtained from the same RNA preparations using the REPLI-g WTA Single Cell
121 Kit (Qiagen) also involving random amplification. This check consisted in (i) the lack of
122 amplification in the absence of reverse transcription (*i.e.* directly on RNA), and (ii) the lack of
123 intron in the PCR product from the unrelated gene encoding β -tubulin, as we performed
124 previously [20, 21].

125

126 **Probes design**

127 The procedure of enrichment in *P. jirovecii* RNA using bait probes was derived from that
128 described for *Candida albicans* [22, 23]. Our SureSelect capture library included a total of
129 43,793 biotinylated bait probes of 120-nucleotide that were designed using the eArray software
130 and the 1x tiling option (Agilent Technologies, <https://earray.chem.agilent.com/earray/>). The

131 probes were head-to-tail and non-overlapping. They covered completely the length of each
132 ORF. However, no probes was placed when less than 120 nucleotides remained at the end of
133 ORF, so that a small bias was introduced at the 3' end that was more pronounced for small
134 genes. The probes covered a total of 4,135 *P. jirovecii* ORFs, *i.e.* all 3,772 of the reference
135 assembly of the *P. jirovecii* genome Pneu-jiro_RU7_V2 including 181 *msg* genes ([8], 8.4Mb
136 haploid genome; ca. 4.0Mb Orfeome), 83 genes and 25 pseudogenes of all six MSG families
137 that we described previously [9], as well as 255 *msg* genes from previous publications (Table
138 S1). The *msg* genes and pseudogenes were added in order to ensure enrichment in *msg* cDNAs
139 (see text). The UCS including the single copy promoter of MSG family I was also included
140 (locus T551_00002). The mitochondrial ORFs and the single copy rDNA genes were not
141 covered, whereas the ORFs encoding ribosomal proteins were. The probes that mapped onto
142 the human transcriptome present in UCSC resources using BLAT (<https://genome.ucsc.edu/>)
143 were discarded. There was an average of 10.0 probes for each ORF or pseudogene.

144

145 **Preparation of RNAseq libraries with enrichment in *P. jirovecii* RNA**

146 RNA libraries for RNAseq were prepared using the Agilent SureSelect^{XT} targeted cDNA
147 Enrichment Kit for multiplexed Illumina sequencing (Agilent Technologies, manufacturer's
148 reference G9611A), using the 200 ng sample preparation procedure without shearing of DNA,
149 and 16 PCR cycles for the step "Indexing and sample processing for multiplex sequencing".
150 The only change to the manufacturer's instructions has been that the samples were dried 3 to 5
151 min at room temperature rather than at 37°C after all purifications with AMPure XP beads.
152 Briefly, double-stranded cDNA was produced with adapters ligated to both ends of the cDNAs,
153 allowing subsequent amplification using primers matching the adapters. The addition of
154 primers to the cDNAs during the procedure was checked using the Agilent Bioanalyser. Each
155 library received a different index that allowed several libraries to be sequenced together

156 (multiplexing). Amplified double-stranded cDNA was incubated at 65°C for 24 h with our
157 capture library of biotinylated probes described above. The hybridized sequences were
158 captured with magnetic streptavidin beads. They were next linearly amplified using provided
159 primers and indexed in a new PCR. The non-enriched sample 1NE was sequenced directly after
160 the SeqPlex[™] RNA Amplification Kit, without applying the Agilent SureSelect^{XT} targeted
161 cDNA Enrichment Kit.

162

163 **RNA sequencing**

164 Libraries resulting from the Agilent SureSelect^{XT} targeted cDNA Enrichment Kit were
165 sequenced on an Illumina MiSeq with a Micro Reagent Kit v2 (300 cycles). Sequencing data
166 were processed using Illumina bcl2fastq2 conversion software v2.20. RNAseq paired Illumina
167 reads were merged using BMAP (v. 37.82). Merged reads mapping onto the human reference
168 genome (grch38) using HISAT2 (v.2.26.0) were discarded, and non-mapping reads were
169 extracted using a combination of samtools (v1.8, parameters: view -h -b -f 4) and bedtools
170 (v2.26.0, parameters: bamtofastq). The obtained human-filtered merged reads were de-
171 duplicated using the program cd-hit-dup from cd-hit (v 4.6.8). The proportion of *P. jirovecii*
172 sequences versus human ones in each samples of RNAseq reads was determined using a splice
173 aware mapper with standard settings (STAR 2.6.0c). The other procedures are described in the
174 supplementary information.

175 **Results**

176 **RNAseq of patients' clinical samples with enrichment in *P. jirovecii* RNA**

177 In order to study the expression of the *P. jirovecii* MSG superfamily, we analyzed total RNAs
178 extracted from the bronchoalveolar lavage fluid (BALF) specimens of six patients with PCP.
179 The proportion of *P. jirovecii* RNA in such samples being low (ca. 3%)[24], an enrichment
180 step was required. To that aim, we used the Agilent SureSelect^{XT} targeted cDNA Enrichment
181 Kit relying on hybridization to bait probes covering the whole *P. jirovecii* orfeome (complete
182 set of protein coding sequences). In order to ensure enrichment in *msg* cDNAs (complementary
183 DNA synthesized from RNA), probes were also derived from published *msg* gene sequences
184 (Table S1). The latter approach was based on the presence of conserved motifs within *msg*
185 genes and the fact that the kit allows mismatches in order to detect sequence variants. The
186 enrichment procedure increased the proportion of *P. jirovecii* RNA to 20-60% (Table 1, see
187 enriched and non-enriched samples of patient 1). After elimination of the human reads, the
188 merged Illumina RNAseq paired sequence reads (150 to 250 nucleotides) were deduplicated in
189 order to avoid biases due to the PCR amplification steps included in the sample preparation.
190 These samples of reads, which characteristics are given in Table 1, were subsequently
191 analyzed.

192

193 **Assignment of RNAseq reads and MSG expression analysis**

194 Specific weighted profiles (similar to consensus sequences) based on published *msg* gene
195 sequences were generated for each of the six MSG families. Despite that a single gene sequence
196 exists, specific profiles were also generated for eight control genes that can be considered as
197 housekeeping (except superoxide dismutase). These profiles were used to assign RNAseq reads
198 to one MSG family or control gene using a conservative best hit approach (Table 2; see
199 methods). The reproducibility of the procedure was assessed by the similarity of the results

200 obtained for two independent analyses of patient 2's BALF (samples 2Ea and 2Eb).
201 Importantly, the results of the enriched and non-enriched samples from patient 1 were also
202 similar (samples 1E and 1NE). This latter result validated the procedure of enrichment in *P.*
203 *jirovecii* RNA, including in *msg* transcripts. All eight samples of reads from the six patients
204 gave comparable results. The vast majority of the assigned reads was from MSG family I (77.1
205 to 95.0%), the second most important population was from MSG family III (2.8 to 17.7%), and
206 all other MSG families and control genes were less represented (0.1 to 4.2%). Transcripts of
207 MSG family VI, superoxide dismutase, and beta tubulin were not detected. These observations
208 demonstrated that, at the population level, genes of all MSG families are expressed, except
209 possibly those of family VI. Genes of families I and III are expressed at a higher level than the
210 other families and all control genes. The level of expression of the genes of family I was the
211 highest, *i.e.* 20 to 50 times higher than the housekeeping genes investigated.

212

213 ***In silico* estimation of the diversity of the MSG genes expressed**

214 The diversity of the genes expressed of each MSG family was estimated by calling single
215 nucleotide variants (SNVs) within a window sliding along the alignment of the RNAseq reads
216 with the specific weighted profile. The optimal size of the window to count haplotypes for all
217 MSG families was determined to be 30 bps (Figure S1A). Because the RNAseq reads were
218 deduplicated, the number of haplotypes (sequences with specific SNVs) obtained can be
219 considered as a surrogate of the number of the different *msg* isoforms expressed. This number
220 depended on the lowest proportion of the reads among those present in the window used to
221 support each haplotype, especially for family I (Figure S1B). In order to avoid detecting
222 sequencing errors while being sensitive enough, the proportion of 0.01 was chosen for all our
223 analyses because it is the usual error rate within Illumina reads [25, 26]. For all MSG families
224 in all samples of RNAseq reads, the number of haplotypes identified was most often

225 proportional to the read coverage along the profile (for example family I in sample 3E, Figure
226 2A). The only exceptions were for family I in samples 1E and 2Ea that had more reads than
227 the other samples (Figure 2B and 2C). The peaks of coverage at 3' and especially 5' regions in
228 Figure 2 are likely to result from RNA degradation, as well as to gene-specific degradation
229 pattern [27]. Interestingly, the two samples 1E and 2Ea with sufficient coverage provided
230 drastically reduced numbers of *msg-I* haplotypes at the same four locations along the profile,
231 at positions ca. 100, 2000, 2300, and 2500 bps. These positions might correspond to conserved
232 regions between protein domains where recombinations between these genes occur
233 preferentially. We calculated the median of the numbers of haplotypes obtained along each
234 MSG profile for all samples of RNAseq reads (Table 3; a number of values could not be
235 obtained because of insufficient read coverage). All these values should be considered as
236 minimal because of the conservative parameters used and the dependency on coverage. The
237 observed number of haplotypes varied from three to 21 for MSG family I, and from one to four
238 for the other families. The two samples 1E and 2Ea with a sufficient coverage for family I
239 provided both a value of ca. 20. These results suggested that (i) family I presents the highest
240 diversity of isoforms expressed during *P. jirovecii* infection, and (ii) that a number of different
241 isoforms of each MSG family are expressed, except possibly for families V (only one haplotype
242 detected) and VI (no reads detected).

243

244 ***In vitro* assessment of the diversity of *msg-I* isoforms expressed**

245 We assessed at the DNA level the diversity of *msg-I* isoforms expressed in the *P. jirovecii*
246 population infecting each of the six patients. To that aim, the repertoire of these genes was
247 amplified from each BALF's genomic DNA using primers localized in the UCS containing the
248 single copy promoter and at the end of the genes. The PCR product was subcloned and several
249 subclones were sequenced. The samples of all six patients presented a significant diversity of

250 *msg*-I isoforms expressed, *i.e.* 27 to 80% of 10 to 15 subclones sequenced were unique (Table
251 4). Only two patients shared a single sequence (patients 1 and 2).

252 Discussion

253 The human pathogenic fungus *P. jirovecii* harbors most probably a system of surface antigenic
254 variation ensuring presumably both escape from host immunity and adhesion to target cells.
255 This system involves six families of hypervariable surface glycoproteins, the MSGs, family I
256 being under mutually exclusive expression at the individual cell level. In the present study, we
257 analyzed the pattern of expression of the genes encoding these proteins at the RNA level. The
258 results were similar in six patients with PCP. The *msg* transcripts included members of at least
259 five families. The level of expression of families I and III was higher than those of the other
260 MSG families and housekeeping genes. Family I was by far expressed at the highest level. A
261 number of different isoforms of at least four of the six families were expressed. Importantly,
262 the six patients that we investigated were each co-infected with several *P. jirovecii* strains (see
263 methods), so that the results are means from several strains. Nevertheless, the similarity of the
264 results of the six patients suggests that all strains expressed similarly the MSGs.

265 We did not detect transcripts of the MSG family VI. In *P. murina*, the proteins of this
266 family are present only at the surface of the ascospores, within asci or recently released from
267 the asci [28]. This observation suggested a particular regulation of this family during the cell
268 cycle. Because asci represent generally a minority of ca. 5% in the infecting *P. jirovecii*
269 populations [29], it is possible that the transcripts encoding these proteins were present in
270 amounts too low to be detected by our procedure. Interestingly, the genes of this family are all
271 localized in the distal region of the subtelomeres, *i.e.* most distant from the telomeres and
272 closest to the genomic genes. Moreover, the recombinations between them are less frequent
273 than between the genes of the other MSG families [9]. A relationship between chromosomal
274 location, expression in ascospores, and low frequency of mosaicism is likely to exist.

275 Our observations support several aspects of our model for the *P. jirovecii* antigenic
276 variation system (Figure 1). The observed expression of several *msg*-I isoforms within each

277 infecting population is consistent with the hypothesis of a continuous segregation of
278 subpopulations expressing each a different single isoform. The value of 20 different isoforms
279 that we observed in two patients is close to the single value of 18 reported so far [9]. The latter
280 value as well as those reported in the present work being all minimal estimations, it is likely
281 that a higher diversity of family I is actually expressed. This hypothesis is also suggested by
282 the peaks up to ca. 30-35 haplotypes that we observed for this family (Figure 2B and 2C, at
283 position 227). As far as the other MSG families are concerned, the mean number of haplotypes
284 observed were always dependent on the read coverage. Nevertheless, peaks up to 10-20
285 haplotypes were observed for all families in sample 1E with the highest coverage, suggesting
286 that all families might also present an important diversity of genes expressed. The expression
287 of several isoforms of all MSG families that we observed, except possibly of family V and VI,
288 is compatible with the postulated independent expression thanks to the promoter that each gene
289 possesses. However, our analyses did not allow assessing if all genes of each family were
290 transcribed. Thus, it remains to determine if these latter genes are constitutively expressed,
291 subject to a regulation during the cell cycle, and/or silenced due to the proximity of the
292 telomeres (by the “telomere position effect”)[30].

293 The expression at a very high level of *msg-I* isoforms we observed in *P. jirovecii* is
294 consistent with our model and previous studies at the protein level [8, 11]. This high expression
295 might be due to transcription enhancement driven by the intron present in the UCS that is larger
296 than that present in the other MSG families [9, 10]. On the other hand, the high expression of
297 family III is a new feature. The latter may not be due to the presence of two introns of common
298 size for *P. jirovecii* (40 to 60 bps) close to the promoter of *msg-III* genes because a similar
299 arrangement is present for *msg-II* genes that are not over-expressed.

300 We analyzed immunocompromised patients with active PCP. However, the antigenic
301 surface variation system of *P. jirovecii* has probably evolved in immunocompetent humans

302 without PCP, and thus be above all a colonization factor. Colonized individuals include
303 potentially several categories of humans, *e.g.* infants experiencing primo-infection, transient
304 carriers such as healthcare workers in contact with PCP patients, patients with chronic lung
305 diseases, pregnant women, and elderly people [31]. In these colonized individuals, the number
306 of subpopulations expressing a different *msg-I* isoform could be reduced by the valid immune
307 system. This hypothesis deserves to be tested in order to better understand the surface antigenic
308 variation system of *P. jirovecii*.

309 In conclusion, our results enhance the understanding of the mechanisms involved in the
310 surface antigenic variation of *P. jirovecii*, as well as of its cell surface structure. The postulated
311 strategy to produce continuously subpopulations that are antigenically distinct would be unique
312 among human pathogens, and might be associated with the non-sterile niche within lungs [9,
313 32]. By contrast, pathogens that occupy sterile niches (blood, tissue), such as *Plasmodium* and
314 *Trypanosoma*, rely on cell populations that are antigenically homogeneous. Further work is
315 needed to decipher the unique surface antigenic variation system of *P. jirovecii*.

316 **Acknowledgments**

317 We thank Dr. Michael Walter of the Agilent Diagnostic and Genomics group for his help in
318 the design of the probes and RNA extraction protocol. We also thank Thierry Schuepbach for
319 help in the initial phase of the analyses of the haplotypes numbers. Sequencing was performed
320 at the Lausanne Genomic Technologies Facility, University of Lausanne, Switzerland.
321 Computations were performed at the Vital-IT Center for High-Performance Computing of the
322 Swiss Institute of Bioinformatics (<http://www.vital-it.ch>).

323 **References**

- 324 1. Redhead SA, Cushion MT, Frenkel JK, Stringer JR. *Pneumocystis* and *Trypanosoma cruzi*:
325 nomenclature and typifications. *J Eukaryot Microbiol* **2006**; 53:2–11.
- 326 2. Gigliotti F, Limper AH, Wright T. *Pneumocystis*. *Cold Spring Harb Perspect Med* **2014**;
327 4:a019828.
- 328 3. Keely SP, Renauld H, Wakefield AE, et al. Gene arrays at *Pneumocystis carinii* telomeres.
329 *Genetics* **2005**; 170, 1589-600.
- 330 4. Stringer JR. Antigenic variation in *Pneumocystis*. *J Euk Microbiol* **2007**; 54:8–13.
- 331 5. Keely SP, Stringer JR. Complexity of the MSG gene family of *Pneumocystis carinii*. *BMC*
332 *Genomics* **2009**; 10:367.
- 333 6. Pottratz ST. *Pneumocystis carinii* interactions with respiratory epithelium. *Semin Respir*
334 *Infect* **1998**; 13:323-9.
- 335 7. Hoving JC. *Pneumocystis* and interactions with host immune receptors. *PLoS Pathogens*
336 **2018**; 14:e1006807.
- 337 8. Ma L, Chen Z, Huang DW, et al. Genome analysis of three *Pneumocystis* species reveals
338 adaptation mechanisms to life exclusively in mammalian hosts. *Nat Commun* **2016**;
339 7:10740.
- 340 9. Schmid-Siegert E, Richard S, Luraschi A, Mühlethaler K, Pagni M, Hauser PM.
341 Mechanisms of surface antigenic variation in the human pathogenic fungus *Pneumocystis*
342 *jirovecii*. *mBio* **2017**; 8, e01470-17.
- 343 10. Ma L, Chen Z, Huang, et al. Diversity and complexity of the large surface protein family
344 in the compacted genomes of various *Pneumocystis* species. *mBio* **2020**; 11:e02878-19.
- 345 11. Kutty G, Shroff R, Kovacs JA. Characterization of *Pneumocystis* major surface
346 glycoprotein gene (*msg*) promoter activity in *Saccharomyces cerevisiae*. *Eukaryot Cell*
347 **2013**; 12:1349–355.

- 348 12. Stringer JR, Cushion MT. The genome of *Pneumocystis carinii*. FEMS Immunol Med
349 Microbiol **1998**; 22:15–26.
- 350 13. Wyder MA, Rasch EM, Kaneshiro ES. Quantitation of absolute *Pneumocystis carinii*
351 nuclear DNA content. Trophic and cystic forms isolated from infected rat lungs are haploid
352 organisms. J Eukaryot Microbiol **1998**; 45:233–9.
- 353 14. Martinez A, Aliouat EM, Standaert-Vitse A, Werkmeister E, Pottier M, Pincon C, Dei-Cas
354 E, Aliouat-Denis SCM. Ploidy of cell-sorted trophic and cystic forms of *Pneumocystis*
355 *carinii*. PLoS ONE **2011**; 6:e20935.
- 356 15. Hauser PM, Cushion MT. Is sex necessary for the proliferation and transmission of
357 *Pneumocystis*? PLoS Pathogens **2018**; 14:e1007409.
- 358 16. Kutty G, Ma L, Kovacs JA. Characterization of the expression site of the major surface
359 glycoprotein of human-derived *Pneumocystis carinii*. Mol Microbiol **2001**; 42:183-93.
- 360 17. Linke MJ, Smulian AG, Stringer JR, Walzer PD. Characterization of multiple unique
361 cDNAs encoding the major surface glycoprotein of rat-derived *Pneumocystis carinii*.
362 Parasitol Res **1994**; 80:478–86.
- 363 18. Angus CW, Tu A, Vogel P, Qin M, Kovacs JA. Expression of variants of the major surface
364 glycoprotein of *Pneumocystis carinii*. J Exp Med **1996**;183:1229-34.
- 365 19. Kutty G, Maldarelli F, Achaz G, Kovacs JA. Variation in the major surface glycoprotein
366 genes in *Pneumocystis jirovecii*. J Infect Dis **2008**; 198:741-9.
- 367 20. Richard S, Almeida JMGCF, Cissé OH, Luraschi A, Nielsen O, Pagni M, Hauser PM.
368 Functional and expression analyses of the *Pneumocystis* MAT genes suggest obligate
369 sexuality through primary homothallism within host lungs. mBio **2018**; 9:e02201-17.
- 370 21. Luraschi A, Richard S. Almeida JMGCF, Pagni M, Cushion MT, Hauser PM. Expression
371 and immuno-staining analyses suggest that *Pneumocystis* primary homothallism involves

- 372 trophic cells displaying both Plus and Minus pheromone receptors. *mBio* **2019**; 10:e01145-
373 19.
- 374 22. Amorim-Vaz S, Tran VDT, Pradervand S, Pagni M, Coste AT, Sanglard D. RNA
375 Enrichment method for quantitative transcriptional analysis of pathogens *in vivo* applied
376 to the fungus *Candida albicans*. *mBio* **2015**; 6:e00942-15.
- 377 23. Amorim-Vaz S, Sanglard D. Novel approaches for fungal transcriptomics from host
378 samples. *Front Microbiol* **2016**; 6:1571.
- 379 24. Cissé OH, Pagni M, Hauser PM. Genome sequence of the uncultivated fungal pathogen
380 *Pneumocystis jirovecii*. *mBio* **2012**; 4:e00428-12.
- 381 25. Minoche AE, Dohm JC, Himmelbauer H. Evaluation of genomic high-throughput
382 sequencing data generated on Illumina HiSeq and genome analyzer systems. *Genome Biol*
383 **2011**; 12:R112.
- 384 26. Heydari M, Miclotte G, Demeester P, Van de Peer Y, Fostier J. Evaluation of the impact
385 of Illumina error correction tools on *de novo* genome assembly. *BMC Bioinformatics*
386 **2017**; 18:374.
- 387 27. Xiong B, Yang Y, Fineis FR, Wang JP. DegNorm: normalization of generalized transcript
388 degradation improves accuracy in RNA-seq analysis. *Gen Biol* 2019; 20:75.
- 389 28. Bishop LR, Davis AS, Bradshaw K, Gamez M, Cissé OH, Wang H, Ma L, Kovacs JA.
390 Characterization of p57, a stage-specific antigen of *Pneumocystis murina*. *J Infect Dis*
391 **2018**; 218:282–90.
- 392 29. Aliouat-Denis CM, Martinez A, Aliouat E, Pottier M, Gantois N, Dei-Cas E. The
393 *Pneumocystis* life cycle. *Mem Inst Oswaldo Cruz* **2009**; 104 :419–26.
- 394 30. Barry JD, Ginger ML, Burton P, McCulloch R. Why are parasite contingency genes often
395 associated with telomeres? *Int J Parasitol* 2003; 33:29–45.

- 396 31. Morris A, Norris KA. Colonization by *Pneumocystis jirovecii* and its role in disease. Clin
397 Microbiol Rev **2012**; 25:297–317.
- 398 32. Hauser PM. Is the unique camouflage strategy of *Pneumocystis* associated with its
399 particular niche within host lungs? PLOS Pathogens **2019**; 15:e1007480.

400 **Figure legends:**

401 **Figure 1.** Model for the antigenic variation system of *P. jirovecii*. Only four chromosome ends
402 out of ca. 40 are shown in each cell. The fungus segregates continuously new cells expressing
403 each a new *msg-I* isoform, as well as in minority all mosaic isoforms of the other MSG families,
404 except possibly of family VI (see text). Consequently, each *P. jirovecii* population is
405 subdivided into several subpopulations that are antigenically different and that could multiply
406 or not. The single *msg-I* isoform is expressed at a high level, whereas the isoforms of the other
407 families are transcribed at a low level. The proximal location of the *msg-I* genes within the
408 subtelomeres, closest to the telomeres, suggests that the exchange of the expressed isoform
409 might be facilitated by the concomitant exchange of the telomere through a single
410 recombination between two CRJEs. The recombinations between the *msg* genes of each family
411 that generate gene mosaicism contributing to antigenic variation are not figured. CRJE,
412 conserved recombination junction element. UCS, upstream conserved sequence.

413

414 **Figure 2.** *In silico* estimation of the diversity of *msg-I* isoforms present among samples of
415 RNAseq reads. SNVs called within a 30 bps sliding window defined 30 bps haplotypes. The
416 black line shows the number of haplotype identified along the *msg* gene, whereas the red line
417 shows the number of reads analyzed in each window. (A) RNAseq reads sample 3E. (B) sample
418 1E. (C) sample 2Ea.

419 **Supplementary information**

420 **Figure S1.** Analyses of the parameters for the *in silico* estimation of the median number of
421 MSG haplotypes among RNAseq reads along a specific weighted profile. Sample 1E with the
422 highest number of reads was used (families V and VI are absent because of insufficient read
423 coverage or no reads detected, respectively). (A) sliding window size using a supporting
424 proportion of reads among those present in the window of 0.01. (B) proportion of reads among
425 those present in the window supporting each haplotype using a sliding window of 30 bps.

426

427 **Table S1.** Targeted ORFs/pseudogenes, number of probes, probe sequences.

428

429 **Supplementary material and methods**

430 **Ethics approval and consent to participate**

431 All patients provided informed written consent, which was part of the procedure for admittance
432 in both hospitals. The admittance paperwork included the possibility to ask that their samples
433 not be used for research. The samples were treated anonymously and were collected through a
434 routine procedure at the hospital. This protocol was approved by the institutional review boards
435 (www.Swissethics.ch).

436

437 **Bronchoalveolar lavage fluid specimens**

438 BALFs positive for *P. jirovecii* were collected at the University Hospitals of Lausanne (patients
439 1, 3, 4, 6) and Bern (2 and 5) between 2014 and 2017. The diagnostic was made using staining,
440 methenamine silver nitrate in Lausanne [1] and immunofluorescence in Bern (MONOFLUO
441 *Pneumocystis jirovecii* IFA Test kit, BioRad). BALFs were supplemented with 15% (vol/vol)
442 glycerol or with three volumes of RNAlater (Ambion) as quickly as possible upon reception,
443 frozen in liquid nitrogen, and stored at -80°C. This storage occurred the same day as BALF

444 collection from the patient, except for patients 4 and 6 (one and two days later, respectively).
445 All six BALFs were co-infected with several *P. jirovecii* strains as revealed by multitarget
446 genotyping as we described previously [2].

447

448 **Creation of MSG family specific profiles**

449 We created weighted profiles specific for each MSG family [3]. The profiles were based on the
450 CDS (coding DNA sequences) that we published previously [2]. For each family, CDS were
451 aligned using MAFFT (v7.310, parameters:--globalpair --maxiterate 1000). The alignments
452 were manually curated and translated into weighted profiles using pftools (v2.3.5.d,
453 parameters: pfmake -2 -m). Specific profiles were also derived for the eight control genes based
454 on the single *P. jirovecii* mRNA sequence present in GenBank (actin 1, T551_01287; alpha
455 tubulin T551_01836; beta tubulin, T551_02025; dihydrofolate reductase, T551_02883;
456 dihydropteroate synthase, U66281; elongation factor 2, T551_00854; elongation factor 3,
457 T551_02306; superoxide dismutase, T551_02671). Weighted profiles were aligned against a
458 random shuffled DNA sequence collection and scores were used to calibrate the profiles with
459 pftools (v.3, parameters: --method evd_full --heuristic-db profile --profile-sampling 25 --
460 mode 1).

461

462 **MSG expression analysis**

463 The obtained human-filtered and de-duplicated reads were aligned against all six MSG families
464 and control genes profiles with pftools in both senses (v3, parameters: -N -o 4). The results
465 were further filtered to obtain a conservative best hit by removing reads with a calibrated score
466 below 10, as well as those that had a difference between the first and second best family scores
467 below 20% of the score for the first best hit.

468

469 **Haplotype estimation**

470 Consistently, SNVs observed in enriched and non-enriched samples 1E and 1NE were almost
471 identical both in position and proportion (not shown). SNVs were characterized from the
472 obtained alignments (psa-format) in a sliding-window approach. Only reads spanning the entire
473 analyzed window were considered and reads containing identical SNVs grouped as potential
474 haplotypes. Each windows needed 20 reads or more to be considered “callable” and only
475 haplotypes supported by a minimum four reads with identical SNVs and 0.01 frequency
476 (number of haplotype supporting reads/total reads in window) were considered. The median
477 number of haplotypes for each MSG family was calculated over “callable” regions, whereas
478 regions with insufficient reads coverage were ignored. The window size was optimized based
479 on the sample 1E and calling haplotypes with 10 bps increasing windows sizes (10 to 100 bps,
480 frequency cut-off of 0.01, Figure S1A). Optimization of the frequency cut-off was similarly
481 conducted on sample 1E with a fixed window size of 30 bps and frequencies of 50, 10, 5, 1,
482 0.5 and 0.01 (Figure S1B). Curves of median number of haplotypes and coverage were plotted
483 using GraphPad Prism (v8.0.1).

484

485 **PCR amplification of the expressed *msg-I* isoforms repertoires**

486 In order to avoid contaminations, PCRs were set up and analyzed in separate rooms, and
487 negative controls were systematically performed at each experiment. PCRs were performed
488 using the High Fidelity Expand polymerase (Roche) on genomic DNA extracted from the
489 BALF using the QIAamp DNA minikit (Qiagen). The primers were GPI-rev (AAA TCA TGA
490 ACG AAA TMA YCA TTG C) and GK135 (GAC AAG GAT GTT GCT TTT GAT; Kutty et
491 al., 2001), generating a PCR product of ca. 3050 bps covering the entire *msg-I* expressed genes.
492 Primers were synthesized by Microsynth (Balgach, Switzerland). The final concentration of
493 MgCl₂ in the reaction was 3 mM. Each reaction began with denaturation for 3 min at 94 °C,

494 followed by 35 cycles of 15 sec at 94°C, 30 sec at the annealing temperature of 55°C, and 3
495 min 30 sec of elongation at 72 °C. The reactions ended with extension of 10 minutes at 72°C.
496 The realtime PCR used to determine the number of *P. jirovecii* copies in BALF given in Table
497 1 has been described previously [4].

498

499 **PCR product subcloning and sequencing**

500 Subcloning of the PCR products was carried out using the TOPO TA cloning Kit (Life
501 Technologies, Inc.). Sanger sequencing of ca. 850 bps of the 5' end of the cloned *msg-I* gene
502 was performed on minipreparation of subclone plasmid DNA [5] using primer M13 forward
503 (GTA AAA CGA CGG CCA GT). When the insert was in the inadequate direction, the
504 subclone was sequenced again using primer M13 reverse (AGC GGA TAA CAA TTT CAC
505 ACA GG). Because only one strand was sequenced, and thus that few sequencing errors were
506 possible, subclones with one or two bps different were considered identical. This concerned
507 only two to five subclones in five of the six patients studied. After alignment and trimming, the
508 identities between subclones' sequences were determined using the MAFFT identity matrix
509 tool (<https://www.ebi.ac.uk/Tools/msa/mafft/>).

510

511 **Sequence data accession number**

512 The RNA-seq data without human sequences analyzed in this work are accessible under the
513 BioProject PRJNA608830 with bioSample accession numbers SAMN14209915 to
514 SAMN14209922.

515

516 **References**

517 1. Musto L, Flanigan M, Elbadawi A. Ten-minute silver stain for *Pneumocystis carinii* and
518 fungi in tissue sections. Arch Pathol Lab Med **1982**; 106:292–4.

- 519 2. Schmid-Siegert E, Richard S, Luraschi A, Mühlethaler K, Pagni M, Hauser PM.
520 Mechanisms of surface antigenic variation in the human pathogenic fungus *Pneumocystis*
521 *jirovecii*. mBio **2017**; 8, e01470-17.
- 522 3. Bucher P, Bairoch A. A generalized profile syntax for biomolecular sequence motifs and
523 its function in automatic sequence interpretation. Proc Int Conf Intel Syst Mol Biol **1994**;
524 2:53– 61.
- 525 4. Richard S, Almeida JMGCF, Cissé OH, Luraschi A, Nielsen O, Pagni M, Hauser PM.
526 Functional and expression analyses of the *Pneumocystis* MAT genes suggest obligate
527 sexuality through primary homothallism within host lungs. mBio **2018**; 9:e02201-17.
- 528 5. Birnboim HC, Doly J. A rapid alkaline extraction procedure for screening recombinant
529 plasmid DNA. Nucl Ac Res 1979; 7:1513–23.

Table 1. Characteristics of the sets of RNAseq reads from BALFs of six patients with PCP.

| | RNAseq reads sample ^a | | | | | | | |
|---|----------------------------------|------|-----|-----|-------|------|-----|------|
| | 1E | 1NE | 2Ea | 2Eb | 3E | 4E | 5E | 6E |
| Underlying disease ^b | HIV+ | HIV+ | KT | KT | ALL | PNET | MM | HIV+ |
| <i>P. jirovecii</i> copies in BALF by realtime PCR (x10 ⁶ /ml) | 22.9 | 22.9 | 0.7 | 0.7 | 1,111 | 3.1 | 4.2 | 57.4 |
| Read pairs (x10 ⁶) | 12.6 | 8.1 | 2.7 | 1.0 | 1.0 | 11.0 | 0.5 | 2.7 |
| Merged read pairs (x10 ⁶) | 5.5 | 4.3 | 2.0 | 0.9 | 0.9 | 1.0 | 0.4 | 2.5 |
| <i>P. jirovecii</i> reads in merged read pairs (%) | 55 | 0.2 | 21 | 23 | 62 | 26 | 52 | 57 |
| Human reads in merged read pairs (%) | 23 | 89 | 50 | 57 | 7 | 39 | 13 | 5 |
| Deduplicated merged read pairs without human (x10 ³) | 1'499 | 464 | 703 | 126 | 183 | 153 | 73 | 383 |

^a Code sample name: 1 to 6, patient number. E, enriched in *P. jirovecii* cDNA using Sureselect kit. NE, non-enriched. a and b, duplicate from the same patient's BALF (patient 2).

^b KT, Kidney transplantation. ALL, Acute lymphocytic leukemia. PNET, primitive neuroectodermal tumor. MM, multiple myeloma.

Table 2. Proportion (%) of RNAseq reads assigned to one MSG family or control gene. ^a

| | RNAseq reads sample ^b | | | | | | | |
|--------------------------------|----------------------------------|-------|--------|--------|-------|-------|-------|--------|
| | 1E | 1NE | 2Ea | 2Eb | 3E | 4E | 5E | 6E |
| | 182,389 ^c | 1,213 | 61,248 | 10,641 | 7,324 | 7,106 | 8,948 | 49,497 |
| MSG family I (A1) ^d | 79.4 | 87.4 | 82.5 | 78.3 | 95.0 | 77.1 | 87.6 | 86.1 |
| MSG family II (A3) | 3.3 | 2.5 | 3.1 | 2.6 | 1.1 | 2.1 | 3.4 | 3.3 |
| MSG family III (A3) | 6.0 | 5.4 | 8.2 | 14.5 | 2.8 | 17.7 | 7.1 | 6.1 |
| MSG family IV (B) | 1.6 | 1.1 | 2.1 | 2.4 | 0.2 | 3.0 | 0.8 | 4.3 |
| MSG family V (D) | 0.3 | 0.2 | 1.4 | 0.1 | 0.2 | 0 | 0 | 0 |
| MSG family VI (E) | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Actin 1 | 0.3 | 0.3 | 0.1 | 0 | 0 | 0 | 0 | 0 |
| Alpha tubulin | 4.2 | 1.8 | 1.8 | 1.4 | 0.2 | 0 | 0 | 0 |
| Beta tubulin | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Dihydrofolate reductase | 1.8 | 0.4 | 0.5 | 0.7 | 0.5 | 0.1 | 0 | 0.1 |
| Dihydropteroate synthase | 0.3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Elongation factor 2 | 1.4 | 0.5 | 0.3 | 0.7 | 0 | 0 | 1.2 | 0 |
| Elongation factor 3 | 1.3 | 0.5 | 0.2 | 0 | 0 | 0 | 0 | 0.1 |
| Superoxide dismutase | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

^a RNAseq reads were assigned to one MSG family or gene using specific weighted profiles.

^b Code sample name: 1 to 6, patient number. E, enriched in *P. jirovecii* cDNA using Sureselect kit. NE, non-enriched. a and b, duplicates from the same patient's BALF (patient 2).

^c Total number of merged reads (without human, de-duplicated) assigned to one MSG family or control gene.

^d The nomenclature of the MSG families of Ma et al. [8] is in parentheses.

Table 3. Median number of expressed haplotypes for each MSG family.^a

| | RNAseq reads sample ^b | | | | | | | |
|---------------------|----------------------------------|-----------------|--------|--------|-------|-------|-------|--------|
| | 1E | 1NE | 2Ea | 2Eb | 3E | 4E | 5E | 6E |
| | 165,244 ^c | 1,172 | 59,594 | 10,418 | 7,273 | 7,099 | 8,850 | 49,398 |
| I (A1) ^d | 18 | nd ^e | 21 | 4 | 3 | nd | 2 | 6 |
| II (A3) | 4 | nd | 4 | nd | nd | nd | nd | nd |
| III (A3) | 1 | nd | 4 | nd | nd | nd | nd | nd |
| IV (B) | 2 | nd | 3 | nd | nd | nd | nd | nd |
| V (D) | nd | nd | nd | 1 | nd | nd | nd | nd |
| VI (E) | nd | nd | nd | nd | nd | nd | nd | nd |

^a The median number of expressed haplotypes was determined for each MSG family using SNVs calling within the RNAseq reads within a sliding window.

^b Code sample name: 1 to 6, patient number. E, enriched in *P. jirovecii* cDNA using Sureselect kit. NE, non-enriched. a and b, duplicates from the same patient's BALF (patient 2).

^c Total number of MSG RNAseq reads analyzed.

^d The nomenclature of the MSG families of Ma et al. [8] is in parentheses.

^e not determined because of insufficient read coverage (or no reads for family VI).

Table 4. Diversity of *msg-I* isoforms expressed. ^a

| | Patient no. | | | | | |
|--------------------------------------|---------------|---------------|---------------|---------------|---------------|---------------|
| | 1 | 2 | 3 | 4 | 5 | 6 |
| No. subclones sequenced | 12 | 11 | 10 | 10 | 10 | 15 |
| % unique subclones (no.) | 58 (7) | 27 (3) | 80 (8) | 70 (7) | 30 (3) | 33 (5) |
| % mean identity of subclones (range) | 69 (66-81) | 72 (68-75) | 68 (63-80) | 69 (63-78) | 98 (97-98) | 73 (65-88) |

^a The repertoire of expressed *msg-I* genes (*i.e.* linked to and thus under the control of the single copy promoter present in the UCS) was amplified by PCR from the genomic DNA of the patient's BALF, the PCR product was subcloned in a plasmid, and ca. 850 bps of the 5' end of each subclone were sequenced.

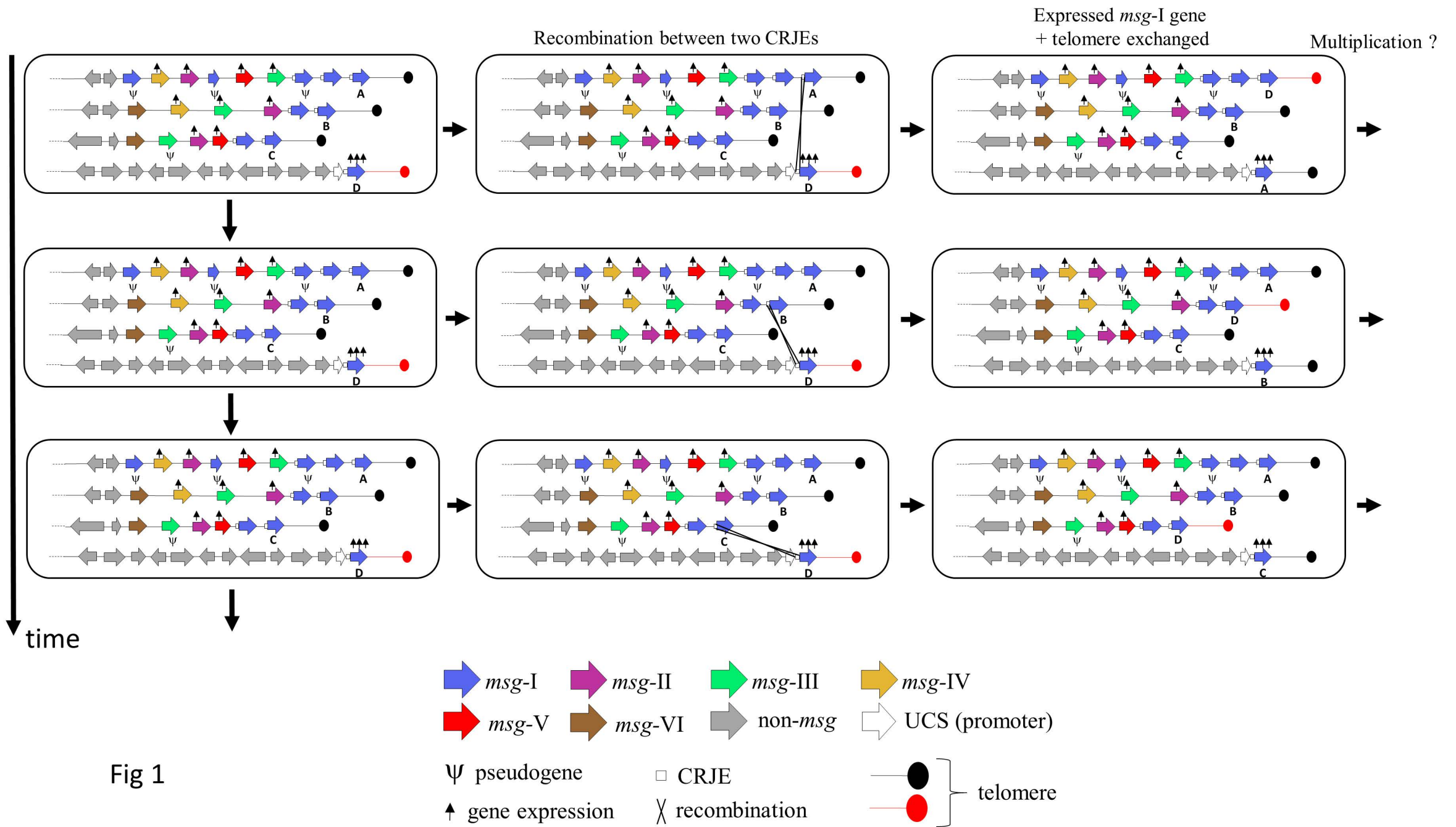


Fig 1

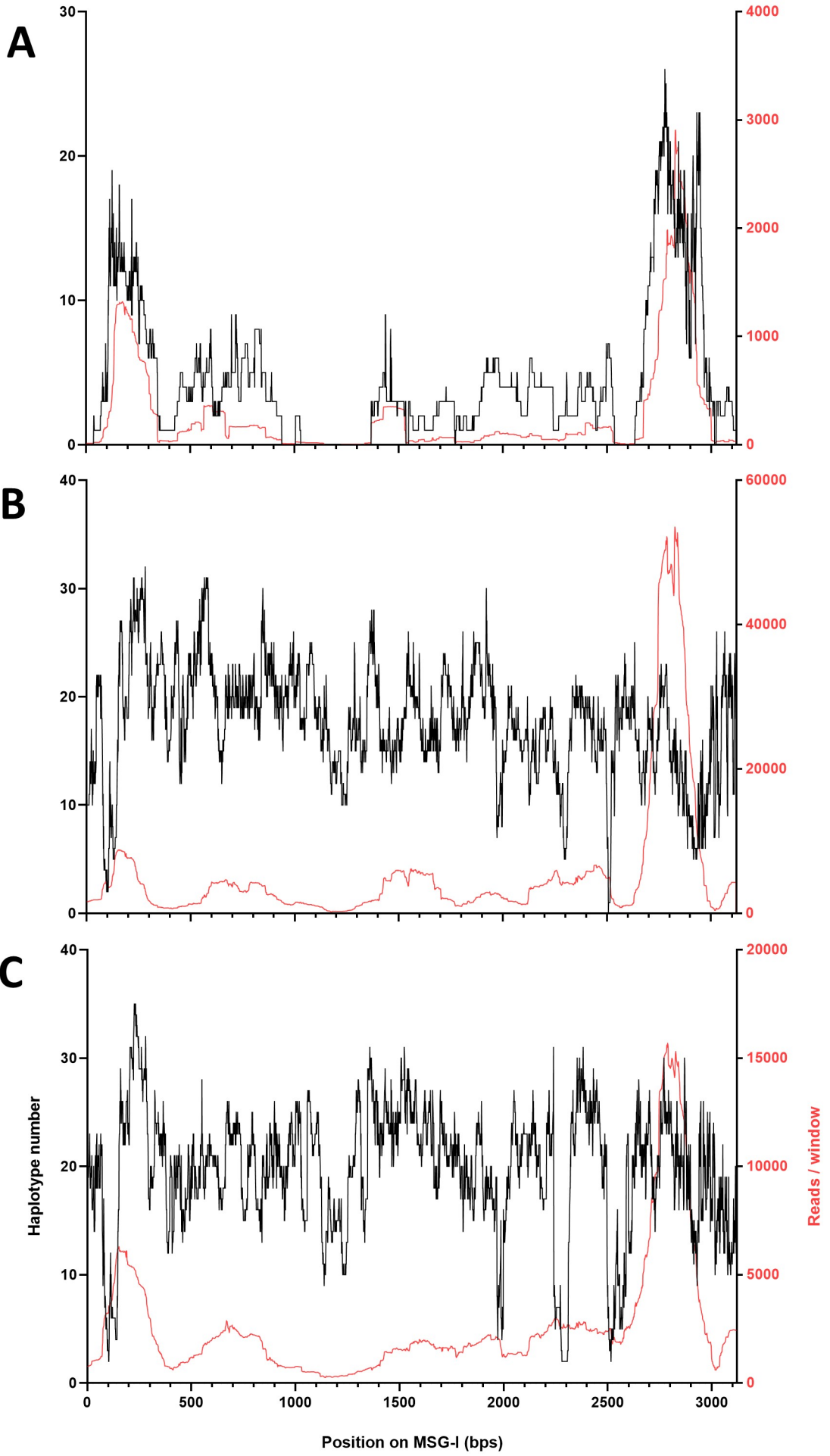


Fig 2

1 **Supplementary materials**

2 **Figure S1.** Analyses of the parameters for the *in silico* estimation of the median number of
3 MSG haplotypes among RNAseq reads along a specific weighted profile. Sample 1E with the
4 highest number of reads was used (families V and VI are absent because of insufficient read
5 coverage or no reads detected, respectively). (A) sliding window size using a supporting
6 proportion of reads among those present in the window of 0.01. (B) proportion of reads among
7 those present in the window supporting each haplotype using a sliding window of 30 bps.

8

9 **Table S1.** Targeted ORFs/pseudogenes, number of probes, probe sequences.

10

11 **Supplementary material and methods**

12 **Ethics approval and consent to participate**

13 All patients provided informed written consent, which was part of the procedure for admittance
14 in both hospitals. The admittance paperwork included the possibility to ask that their samples
15 not be used for research. The samples were treated anonymously and were collected through a
16 routine procedure at the hospital. This protocol was approved by the institutional review boards
17 (www.Swissethics.ch).

18

19 **Bronchoalveolar lavage fluid specimens**

20 BALFs positive for *P. jirovecii* were collected at the University Hospitals of Lausanne (patients
21 1, 3, 4, 6) and Bern (2 and 5) between 2014 and 2017. The diagnostic was made using staining,
22 methenamine silver nitrate in Lausanne [1] and immunofluorescence in Bern (MONOFLUO
23 *Pneumocystis jirovecii* IFA Test kit, BioRad). BALFs were supplemented with 15% (vol/vol)
24 glycerol or with three volumes of RNAlater (Ambion) as quickly as possible upon reception,
25 frozen in liquid nitrogen, and stored at -80°C. This storage occurred the same day as BALF

26 collection from the patient, except for patients 4 and 6 (one and two days later, respectively).
27 All six BALFs were co-infected with several *P. jirovecii* strains as revealed by multitarget
28 genotyping as we described previously [2].

29

30 **Creation of MSG family specific profiles**

31 We created weighted profiles specific for each MSG family [3]. The profiles were based on the
32 CDS (coding DNA sequences) that we published previously [2]. For each family, CDS were
33 aligned using MAFFT (v7.310, parameters:--globalpair --maxiterate 1000). The alignments
34 were manually curated and translated into weighted profiles using pftools (v2.3.5.d,
35 parameters: pfmake -2 -m). Specific profiles were also derived for the eight control genes based
36 on the single *P. jirovecii* mRNA sequence present in GenBank (actin 1, T551_01287; alpha
37 tubulin T551_01836; beta tubulin, T551_02025; dihydrofolate reductase, T551_02883;
38 dihydropteroate synthase, U66281; elongation factor 2, T551_00854; elongation factor 3,
39 T551_02306; superoxide dismutase, T551_02671). Weighted profiles were aligned against a
40 random shuffled DNA sequence collection and scores were used to calibrate the profiles with
41 pftools (v.3, parameters: --method evd_full --heuristic-db profile --profile-sampling 25 --
42 mode 1).

43

44 **MSG expression analysis**

45 The obtained human-filtered and de-duplicated reads were aligned against all six MSG families
46 and control genes profiles with pftools in both senses (v3, parameters: -N -o 4). The results
47 were further filtered to obtain a conservative best hit by removing reads with a calibrated score
48 below 10, as well as those that had a difference between the first and second best family scores
49 below 20% of the score for the first best hit.

50

51 **Haplotype estimation**

52 Consistently, SNVs observed in enriched and non-enriched samples 1E and 1NE were almost
53 identical both in position and proportion (not shown). SNVs were characterized from the
54 obtained alignments (psa-format) in a sliding-window approach. Only reads spanning the entire
55 analyzed window were considered and reads containing identical SNVs grouped as potential
56 haplotypes. Each windows needed 20 reads or more to be considered “callable” and only
57 haplotypes supported by a minimum four reads with identical SNVs and 0.01 frequency
58 (number of haplotype supporting reads/total reads in window) were considered. The median
59 number of haplotypes for each MSG family was calculated over “callable” regions, whereas
60 regions with insufficient reads coverage were ignored. The window size was optimized based
61 on the sample 1E and calling haplotypes with 10 bps increasing windows sizes (10 to 100 bps,
62 frequency cut-off of 0.01, Figure S1A). Optimization of the frequency cut-off was similarly
63 conducted on sample 1E with a fixed window size of 30 bps and frequencies of 50, 10, 5, 1,
64 0.5 and 0.01 (Figure S1B). Curves of median number of haplotypes and coverage were plotted
65 using GraphPad Prism (v8.0.1).

66

67 **PCR amplification of the expressed *msg-I* isoforms repertoires**

68 In order to avoid contaminations, PCRs were set up and analyzed in separate rooms, and
69 negative controls were systematically performed at each experiment. PCRs were performed
70 using the High Fidelity Expand polymerase (Roche) on genomic DNA extracted from the
71 BALF using the QIAamp DNA minikit (Qiagen). The primers were GPI-rev (AAA TCA TGA
72 ACG AAA TMA YCA TTG C) and GK135 (GAC AAG GAT GTT GCT TTT GAT; Kutty et
73 al., 2001), generating a PCR product of ca. 3050 bps covering the entire *msg-I* expressed genes.
74 Primers were synthesized by Microsynth (Balgach, Switzerland). The final concentration of
75 MgCl₂ in the reaction was 3 mM. Each reaction began with denaturation for 3 min at 94 °C,

76 followed by 35 cycles of 15 sec at 94°C, 30 sec at the annealing temperature of 55°C, and 3
77 min 30 sec of elongation at 72 °C. The reactions ended with extension of 10 minutes at 72°C.
78 The realtime PCR used to determine the number of *P. jirovecii* copies in BALF given in Table
79 1 has been described previously [4].

80

81 **PCR product subcloning and sequencing**

82 Subcloning of the PCR products was carried out using the TOPO TA cloning Kit (Life
83 Technologies, Inc.). Sanger sequencing of ca. 850 bps of the 5' end of the cloned *msg-I* gene
84 was performed on minipreparation of subclone plasmid DNA [5] using primer M13 forward
85 (GTA AAA CGA CGG CCA GT). When the insert was in the inadequate direction, the
86 subclone was sequenced again using primer M13 reverse (AGC GGA TAA CAA TTT CAC
87 ACA GG). Because only one strand was sequenced, and thus that few sequencing errors were
88 possible, subclones with one or two bps different were considered identical. This concerned
89 only two to five subclones in five of the six patients studied. After alignment and trimming, the
90 identities between subclones' sequences were determined using the MAFFT identity matrix
91 tool (<https://www.ebi.ac.uk/Tools/msa/mafft/>).

92

93 **Sequence data accession number**

94 The RNA-seq data without human sequences analyzed in this work are accessible under the
95 BioProject PRJNA608830 with bioSample accession numbers SAMN14209915 to
96 SAMN14209922.

97

98

99

100 **References**

- 101 1. Musto L, Flanigan M, Elbadawi A. Ten-minute silver stain for *Pneumocystis carinii* and
102 fungi in tissue sections. Arch Pathol Lab Med **1982**; 106:292–4.
- 103 2. Schmid-Siegert E, Richard S, Luraschi A, Mühlethaler K, Pagni M, Hauser PM.
104 Mechanisms of surface antigenic variation in the human pathogenic fungus *Pneumocystis*
105 *jirovecii*. mBio **2017**; 8, e01470-17.
- 106 3. Bucher P, Bairoch A. A generalized profile syntax for biomolecular sequence motifs and
107 its function in automatic sequence interpretation. Proc Int Conf Intel Syst Mol Biol **1994**;
108 2:53– 61.
- 109 4. Richard S, Almeida JMGCF, Cissé OH, Luraschi A, Nielsen O, Pagni M, Hauser PM.
110 Functional and expression analyses of the *Pneumocystis* MAT genes suggest obligate
111 sexuality through primary homothallism within host lungs. mBio **2018**; 9:e02201-17.
- 112 5. Birnboim HC, Doly J. A rapid alkaline extraction procedure for screening recombinant
113 plasmid DNA. Nucl Ac Res **1979**; 7:1513–23.

114

Sequence data accession number 463

The RNA-seq data without human sequences analyzed in this work are accessible under the BioProject PRJNA608830 with bioSample accession numbers SAMN14209915 to SAMN14209922 in the European Nucleotide Archive (<https://www.ebi.ac.uk/ena>).

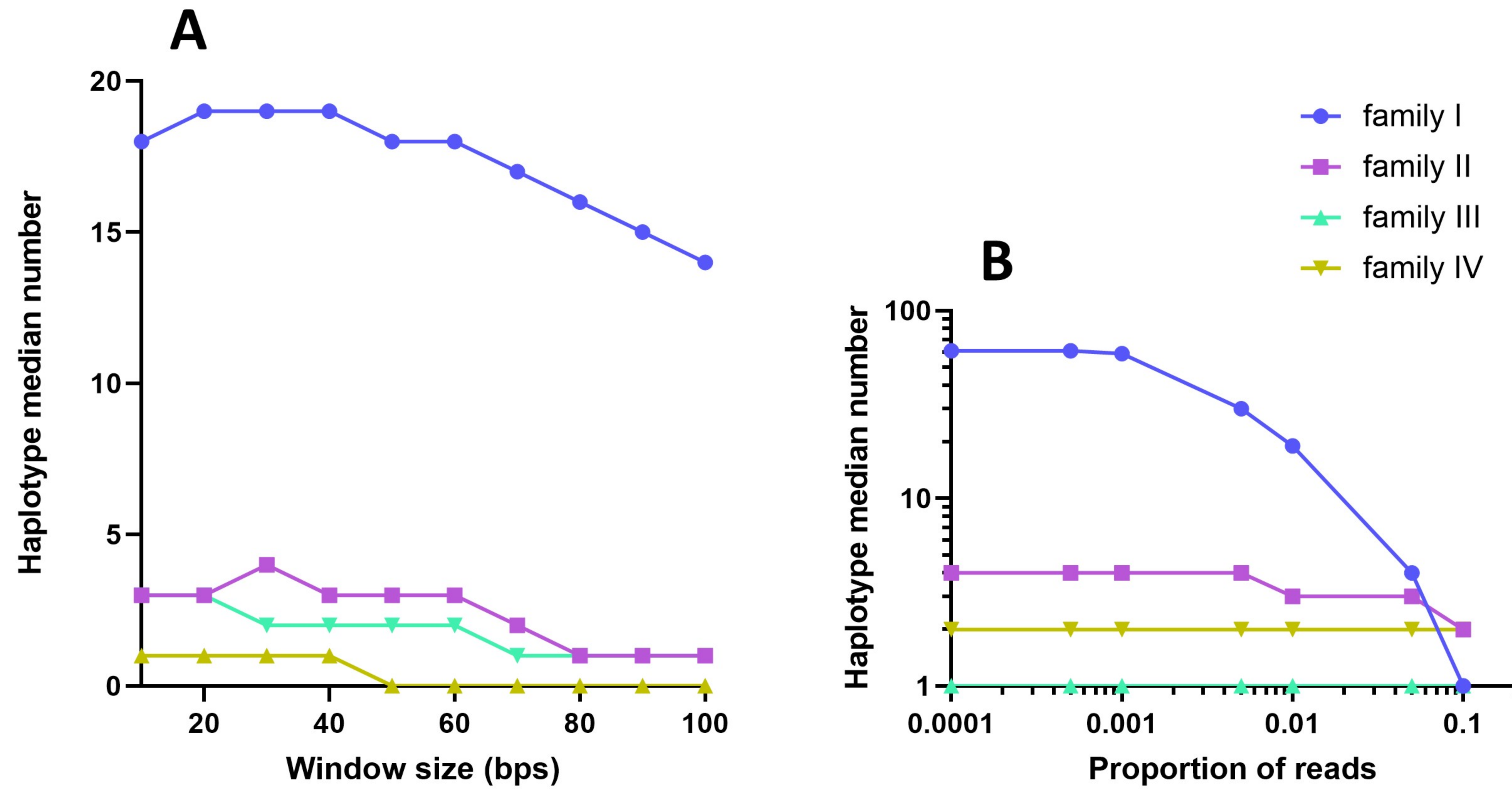


Fig S1