

1 **CRISPR system acquisition and evolution of an obligate intracellular**
2 ***Chlamydia*-related bacterium**

3
4 Claire Bertelli^{1,2}, Ousmane Cissé¹, Brigida Rusconi¹, Carole Kebbi-Beghdadi¹, Antony Croxatto¹,
5 Alexander Goesmann³, François Collyn¹, Gilbert Greub^{*1}

6 ¹ Center for Research on Intracellular Bacteria, Institute of Microbiology, University Hospital Center
7 and University of Lausanne, Lausanne, Switzerland, ² SIB Swiss Institute of Bioinformatics, Lausanne,
8 Switzerland ³ Bioinformatics and Systems Biology, Justus-Liebig-University Giessen, Gießen, Germany

9
10 *Corresponding author:

11 Gilbert Greub, MD PhD
12 Institute of Microbiology
13 University of Lausanne
14 1011 Lausanne
15 SWITZERLAND
16 Phone : +41-21-314 49 79
17 Fax : +41-21-341 40 60
18 e-mail : gilbert.greub@chuv.ch

19
20 Running Title: The Protochlamydia CRISPR system

21
22

23 **ABSTRACT**

24 Recently, a new *Chlamydia*-related organism, *Protochlamydia naegleriophila* KNic, was discovered
25 within a *Naegleria* amoeba. To decipher the mechanisms at play in the modeling of genomes from
26 the *Protochlamydia* genus, we sequenced *de novo* the full genome of *Pr. naegleriophila* combining
27 the advantages of two second-generation sequencing technologies. The assembled complete
28 genome comprises a 2,885,111 bp chromosome and a 145,285 bp megaplasmid. For the first time
29 within the *Chlamydiales* order, a CRISPR system, the immune system of bacteria, was discovered on
30 the chromosome. It is composed of a small CRISPR locus comprising eight repeats and the associated
31 *cas* and *cse* genes of the subtype I-E. A CRISPR locus was also found within *Chlamydia* sp. Diamant,
32 another *Pr. naegleriophila* strain whose genome was recently released, suggesting that the CRISPR
33 system was acquired by a common ancestor of these two members of *Pr. naegleriophila*, after the
34 divergence from *Pr. amoebophila*. The plasmid encodes an F-type conjugative system similar to that
35 found in the Pam100G genomic island of *Pr. amoebophila* suggesting an acquisition of this
36 conjugative system before the divergence of both *Protochlamydia* species and the integration of a
37 putative *Pr. amoebophila* plasmid into its main chromosome giving rise to the Pam100G genomic
38 island. Overall, this new *Pr. naegleriophila* genome sequence enables to investigate further the
39 dynamic processes shaping the genomes of *Chlamydia*-related bacteria.

40

41 INTRODUCTION

42 A large diversity prevails in the order *Chlamydiales*, as suggested by the discovery of a large number
43 of *Chlamydia* and *Chlamydia*-related bacteria belonging to nine different families (Greub 2010;
44 Everett et al. 1999; Horn 2011) and the cross-examination of metagenomics data (Lagkouvardos et al.
45 2014). The family *Parachlamydiaceae* comprises five genera that are each represented by a small
46 number of isolated strains. The genus *Protochlamydia* was lately enriched by the isolation of a
47 *Naegleria* endosymbiont that presented 97.6% identity in the 16S rRNA with *Pr. amoebophila* UWE25
48 and was thus named *Pr. naegleriophila* strain KNic (Casson et al. 2008). Since other members of the
49 *Parachlamydiaceae* family were suspected to be associated with lung infections (Greub 2009), a
50 diagnostic PCR specific for *Pr. naegleriophila* was developed and applied to bronchoalveolar lavages.
51 *Pr. naegleriophila* DNA was detected in the bronchoalveolar lavage of an immunocompromised
52 patient with pneumonia by two PCRs targeting different genomic regions and the presence of the
53 bacterium in the sample was confirmed by direct immunofluorescence (Casson et al. 2008). These
54 results indicated a potential role of *Pr. naegleriophila* in lower respiratory tract infections.

55 A recent study including *Chlamydia* genomes and other members of the *Planctomycetes-*
56 *Verrucomicrobia-Chlamydia* superphylum suggested that the branch leading to *Chlamydia* was
57 shaped mainly by genome reduction and evidenced limited occurrence of gene birth, duplication and
58 transfer (Kamneva et al. 2012), as it is the case in other strict intracellular pathogens (Darby et al.
59 2007). On the contrary, the occurrence of large families of paralogs in the genome of *Chlamydia-*
60 related bacteria suggested an evolution by extensive gene duplication (Domman et al. 2014; Eugster
61 et al. 2007). The chromosome sequence of *Pr. amoebophila* UWE25 exhibited little evidence for the
62 occurrence of lateral gene transfer (Horn et al. 2004). However, a number of probable lateral gene
63 transfers were identified between *Parachlamydia* and other amoeba-infecting bacteria such as
64 *Legionella* (Gimenez et al. 2011), a process that may take place within the amoeba itself (Bertelli &
65 Greub 2012). The *Pr. amoebophila* genome presents a genomic island (Pam100G) that encodes a
66 type IV secretion system of the F-type that might be involved in conjugative DNA transfer (Greub et
67 al. 2004). A similar system was found on the plasmid of *Simkania negevensis* (Collingro et al. 2011)
68 and a partial operon was described in *Parachlamydia acanthamoebae* (Greub et al. 2009), suggesting
69 active DNA transfer capabilities in the ancestor of the *Chlamydiales* and some of its descendants.

70 Small interspaced repetitions were initially observed in *E. coli* (Ishino et al. 1987) and they were then
71 named CRISPR, an acronym for Clustered Regularly Interspaced Short Palindromic Repeats (Jansen et
72 al. 2002). Although found in 50% of bacteria and in 90% of archaea (Weinberger et al. 2012), a
73 CRISPR system had never been previously reported in a member of the order *Chlamydiales*

74 (Makarova et al. 2011). The CRISPR locus usually consists of a variable number (up to 587) of 23-47bp
75 repeats with some dyad symmetry, but not truly palindromic, interspaced by 21-72 bp spacers
76 (Horvath & Barrangou 2010). Associated with these repeats are 2 core cas genes and additional
77 subtype-specific genes putatively providing mechanistic specificity (Koonin & Makarova 2013).
78 Similarity between spacers and extrachromosomal elements first suggested a role in immunity
79 against phage infection and more generally against conjugation or transformation by acquisition of
80 external DNA (Bolotin et al. 2005). The CRISPR-Cas system was shown to mediate an antiviral
81 response thus inducing resistance to phage infection (Deveau et al. 2010), notably in *E. coli* (Brouns
82 et al. 2008). More recently, CRISPR-Cas systems were shown to regulate stress-related response,
83 changing gene expression and virulence traits in several pathogens among which is the intracellular
84 bacteria *Francisella novicida* (Sampson & Weiss 2014; Louwen et al. 2014).

85 In this contribution, we sequenced and analyzed the complete genome of *Pr. naegleriophila* strain
86 KNic and discovered two potentially antagonistic systems, a type IV secretion system likely implicated
87 in conjugative DNA transfer and a CRISPR system that generally controls foreign DNA acquisition.
88 Furthermore, the complete genome sequence of a new species within the genus *Protochlamydia*
89 offered the possibility to look into the genome dynamics throughout evolution by comparing *Pr.*
90 *naegleriophila* KNic gene content and genome architecture to its closest relatives of the family
91 *Parachlamydiaceae*.

92

93 RESULTS

94 Chromosome features and evolution

95 *Pr. naegleriophila* KNic possesses a 2'885'090 bp circular chromosome with a mean GC content of
96 42.7%. The genome size and the GC content are surprisingly high compared to the most closely-
97 related species, *Pr. amoebophila* (Table 1), but it is consistent with its closest relative *Chlamydia* sp.
98 Diamant, another *Pr. naegleriophila* strain (hereafter referred to as *Pr. naegleriophila* Diamant). The
99 chromosome of *Pr. naegleriophila* strain KNic is predicted to encode 2,415 proteins and exhibits four
100 ribosomal operons and 43 tRNAs, more than any other *Chlamydiales* (Table 1). Two types of spacers
101 are found between the 16S and the 23S rRNA: either a simple intergenic spacer or a spacer
102 containing two tRNAs for Ala and Ile.

103 The cumulative G versus C nucleotide bias (GC skew) presents a typical pyramidal shape (Figure S1)
104 that is expected in the absence of particular large genomic islands and confirms the assembly
105 accuracy. The GC skew of *Pr. naegleriophila* is smoother than that of *Pr. amoebophila*, and does not

106 present the small inversion in the slope that is caused by the *Pr. amoebophila* genomic island (**Figure**
107 **S1**) (Greub et al. 2004). The replication origin (*ori*) and the terminus of replication (*ter*), at the
108 minimum and maximum of the curve (**Figure S1**), respectively, show an almost perfectly balanced
109 chromosome with 49.8% of the base on one arm, i.e., between *ori* and *ter*, and 50.2% on the other
110 arm, i.e., between *ter* and *ori*.

111 The two strains of *Pr. naegleriophila* are highly collinear as shown in the alignment of available
112 complete and nearly complete (<5 contigs) genomes of the family *Parachlamydiaceae* (**Figure 1**).
113 Within genus comparison of *Pr. naegleriophila* and *Pr. amoebophila* shows the occurrence of 13
114 recombination and inversion events. As expected, further distantly-related organisms from a
115 different genus exhibit a lower collinearity and an increasing number of recombination events, which
116 is correlated to the phylogenetic distance.

117 **pPNK is an F-type conjugative megaplasmid**

118 The bacterial chromosome was circularized, leaving behind several contigs with a 23-fold coverage,
119 slightly higher than the 16-fold chromosomal coverage. These contigs formed a 145,285 bp large
120 plasmid - the largest known plasmid in the order *Chlamydiales*. The plasmid pPNK presents a GC
121 content of 37.2% and includes 160 genes among which are several transposase and integrase
122 remnants, doc proteins, and systems for the maintenance of the plasmid (*parA* and PNK_p0119) that
123 are all characteristic of extra-chromosomal elements.

124 The plasmid also encodes a type IV secretion system with highest similarity to the F-type system
125 found in the genomic island of *Pr. amoebophila* UWE25 (Greub et al. 2004), in the plasmid of
126 *Chlamydia* sp. Rubis (hereafter named *Parachlamydia* sp. Rubis) and *S. negevensis* (Collingro et al.
127 2011) and to the remnants *traU*, *traN* and *traF* present in members of the family *Parachlamydiaceae*
128 (Greub et al. 2009; Collingro et al. 2011) (**Figure 2**). The type IV secretion system of *Parachlamydia*
129 sp. Rubis is located on a 30 kb long contig that should be circularized as a small plasmid to retain the
130 colinearity of the *tra* operon with other bacteria. This small plasmid would therefore contain almost
131 exclusively the *tra* operon as well as core genes for plasmid replication such as *parA*. *Parachlamydia*
132 sp. Rubis and KNic *tra* operons share a striking colinearity. The comparison of gene conservation
133 shows that *traN* has undergone different rearrangements in both *Pr. amoebophila* strains, and *traC*
134 was split in strain R18. On the other hand, *Parachlamydia* sp. Rubis, *S. negevensis* and *Pr.*
135 *naegleriophila* KNic, the three bacteria that possess the *tra* operon on a plasmid, retained intact
136 genes. Moreover, these bacteria present Ti-type *traA* and *traD* genes downstream that share
137 similarity to and other amoeba-infecting bacteria such as *Rickettsia bellii* and *Legionella* spp..

138 **A CRISPR –Cas system for the first time within *Chlamydiales***

139 In *Pr. naegleriophila*, the CRISPR locus comprises eight 28bp-long repeats separated by 33bp-long
140 spacers. The upstream operon of CRISPR-associated genes from the *E. coli* subtype I-E consists of the
141 core genes *cas1-2*, the type I gene *cas3* and subtype-specific genes *cse1-2*, *cas5*, *cas6e* and *cas7*
142 (**Figure 3**). An almost identical cas operon and a CRISPR locus were identified in *Pr. naegleriophila*
143 ‘Diamant’ (**Figure 3**). On the contrary, this system is absent from other *Parachlamydiaceae* such as
144 strains *Pr. amoebophila* UWE25, E12 and R18. Although a confirmed CRISPR locus is predicted by
145 CRISPRfinder (Grissa et al. 2007) in the recently released genomes of *Neochlamydia* sp. (Domman et
146 al. 2014; Ishida et al. 2014), no *cas* genes could be identified and the repeats were found to be due to
147 a highly repeated protein sequence.

148 The CRISPR spacers could give an interesting imprint of recent invasions by extrachromosomal
149 elements, but unfortunately no significant homology was found by BLASTN against the non-
150 redundant nucleotide database (nt) for strains KNic and Diamant (Sup. table 1 and sup. table 2).
151 Genes surrounding this locus are found in conserved order in all *Protochlamydia* species indicating
152 that this CRISPR region has most likely been acquired by horizontal gene transfer after the
153 divergence of *Pr. naegleriophila* from *Pr. amoebophila*. The gene operon structure is commonly
154 found in bacteria and two species present particular homology to *Pr. naegleriophila* KNic CRISPR
155 locus: *Anaeromyxobacter dehalogenans*, a *Deltaproteobacteria* from soil and *Rhodothermus marinus*,
156 a *Bacteroidetes* (**Figure 3**).

157

158 **DISCUSSION**

159 The genome sequence of *Pr. naegleriophila* presented in this contribution permitted to grasp initial
160 hints on mechanisms triggering the evolution of *Protochlamydia*. We could evidence the presence of
161 a CRISPR-locus in the chromosome for the first time in the order *Chlamydiales*. The sequencing data
162 also revealed the presence of a plasmid that encodes a type IV secretion system and is only partially
163 similar to the genomic island of *Pr. amoebophila*.

164 Amoebae were proposed to act as a reservoir of different amoebae-resisting bacteria where
165 horizontal gene transfer may preferentially take place (Moliner et al. 2010). The presence of a an F-
166 like conjugation plasmid putatively involved in DNA transfer in *Pr. naegleriophila* stresses the
167 likelihood of gene exchange with other bacteria or with the eukaryotic host. The maintenance of
168 intact *tra* genes in bacteria possessing the *tra* operon on a plasmid, suggest that the system has

169 retained functionality, whereas it has evolved towards pseudogenisation and deletion after being
170 integrated in the genome of *Pr. amoebophila* strains and *P. acanthamoebae* strains.

171 The presence of an F-type conjugative operon in the plasmid or in the chromosome of various strains
172 combined with the lack of conjugative operon in the plasmid or in the chromosome of the
173 *Waddliaceae*, *Criblamydiaceae* (Bertelli et al. 2015, 2014) and some *Parachlamydiaceae* challenges
174 the most parsimonious scenario proposed by Collingro *et al* (Collingro et al. 2011) that plasmids
175 evolved from a single conjugative plasmid acquired by an ancestor of the *Parachlamydiaceae*,
176 *Waddliaceae*, and *Simkaniaceae*. In favor of this hypothesis is the shared presence of the Ti-type *traA*
177 and *traD* in the paraphyletic *Parachlamydia sp. Rubis*, *Pr. naegleriophila* KNic, as well as *S.*
178 *negevensis*. If this hypothesis is correct, the plasmid and its *tra* operon were integrated within the
179 chromosome at least twice in the genus *Parachlamydia* and *Protochlamydia*. The *tra* operon was
180 completely lost several times, in the families *Waddliaceae* and *Criblamydiaceae*, in the genus
181 *Neochlamydia*, and in some strains of *Protochlamydia* and *Parachlamydia* (**Figure 2**). Furthermore, it
182 has been partially lost in the *Parachlamydia* genus, where only a few genes remain. An alternative
183 scenario of separate acquisition of *tra* operon by an ancestor of *Simkania* and members of the family
184 *Parachlamydiaceae* may be envisioned. In any case, this highlights the highly dynamic nature of the
185 *Chlamydia*-related genomes and the potential of the *tra* operon to be readily transferred, acquired,
186 and lost among these bacteria.

187 A CRISPR-Cas system has been reported in approximately 50% of bacteria with sporadic distribution
188 patterns suggesting that CRISPR loci are subject to frequent horizontal gene transfer, a hypothesis
189 supported by the presence of CRISPR loci on plasmids (Haft et al. 2005). The CRISPR locus of *Pr.*
190 *naegleriophila* and its associated genes have most probably been acquired horizontally but the
191 proteins have insufficient homology to infer a direct transfer from a given organism. This CRISPR-Cas
192 system is of a different subtype than that of another intracellular amoeba-resisting bacteria *F.*
193 *novicida* ruling out the possibility of intra-amoebal transfer between these organisms. The
194 functionality and the exact role of this CRISPR-Cas system in *Pr. naegleriophila* remains to be
195 determined by laboratory experiments but by similarity to the type IE locus present in *E. coli*, we can
196 hypothesize that it plays a role in preventing DNA acquisition or protecting against phages. Although
197 *Pr. naegleriophila* is an obligate intracellular bacteria, it may still be exposed to phages similarly to
198 other *Chlamydia* species (Śliwa-Dominiak et al. 2013). The difference in CRISPR spacers between *Pr.*
199 *naegleriophila* strains KNic and Diamant clearly highlights the dynamic and likely functional status of
200 the system, as well as the exposure of such obligate intracellular bacteria to DNA of foreign origin.
201 The absence of similarity between CRISPR spacers and sequences of the non-redundant nucleotide
202 database underlines the currently limited knowledge on phages and extrachromosomal DNA

203 circulating in amoebae-resisting bacteria, especially those growing in the ubiquitous amoeba
204 *Naegleria*.

205 Based on the complete genome sequence of *P. acanthamoebae* UV-7 and *Pr. amoebophila* UWE25 as
206 well as four draft genomes, Dommann *et al.* (Domman et al. 2014) suggested the occurrence of few
207 rearrangements within genera of the family *Parachlamydiaceae*. Beyond the difficulty to make such
208 an assessment based on the alignment of highly fragmented draft genome sequences available in
209 public databases, the genomes used in the latter study are different strains of the same species and
210 might therefore not reflect genome evolution at the genus level. Indeed, our comparison shows the
211 absence of rearrangements between the two *Pr. naegleriophila* strains KNic and Diamant, but an
212 increasing number of genome rearrangements with further distantly-related organisms. This
213 highlights the need for complete genomes to precisely unravel the bacterial evolution occurring by
214 recombination. The complete genome sequence of *Pr. naegleriophila* represents a first step toward
215 the understanding of mechanisms triggering genome evolution and evolutionary pressures at play in
216 the *Parachlamydiaceae* family.

217

218 MATERIALS AND METHODS

219 Culture and purification of *Pr. naegleriophila*

220 *Protochlamydia naegleriophila* strain KNic was grown in *Acanthamoeba castellanii* ATCC 30010 at
221 32°C using 75 cm² cell culture flasks (Becton Dickinson, Franklin Lakes, USA) with 30 ml of peptone-
222 yeast extract glucose broth. *Pr. naegleriophila* were purified from amoebae by a first centrifugation
223 step at 120 x *g* for 10 min. Then, remnants from amoebae were removed from the re-suspended
224 bacterial pellet by centrifugation at 6500 x *g* for 30 min onto 25% sucrose (Sigma Aldrich, St Louis,
225 USA) and finally at 32000 x *g* for 70 min onto a discontinuous Gastrographin (Bayer Schering Pharma,
226 Zurich, Switzerland) gradient (48%/36%/28%).

227 Genome sequencing, assembly and gap closure

228 *Pr. naegleriophila* genomic DNA was isolated with the Wizard Genomic DNA purification kit (Promega
229 Corporation, Madison, USA). Reads obtained with Genome Sequencer 20™ (Droege and Hill 2008) by
230 Roche Applied Science (Penzberg, Germany) were assembled using Newbler V1.1.02.15 yielding 93
231 large contigs with a mean 16x coverage. Scaffolding on *Pr. amoebophila* strain UWE25 and PCR-based
232 techniques were used to close the gaps between those contigs. Solexa 35 bp reads obtained from
233 sequencing with Genome Analyzer Gallx (Bennett 2004) by Fasteris (Plan les Ouates, Switzerland)
234 were then mapped to the final assembly with BWA (Li & Durbin 2009) and visualized with Consed
235 (Gordon & Green 2013). Homopolymer errors were corrected in the plasmid and the chromosome
236 sequence after manual inspection of discrepancies covered by >2 reads with a phred quality score of
237 the base >10. Sequence start was placed in an intergenic region closest to the minimum of the GC
238 skew, as determined with a sliding window of 100nt.

239 Genome annotation

240 GenDB 2.4 pipeline (Meyer et al. 2003) was used for a first automatic annotation of the genome that
241 was followed by manual curation of annotation. Coding sequence (CDS) prediction was performed
242 using Prodigal (Hyatt et al. 2010). All predicted CDS were submitted to similarity searches against nr,
243 Swissprot, InterPro, Pfam, TIGRfam and KEGG databases. Putative signal peptides, transmembrane
244 helices and nucleic acid binding domains were predicted using respectively SignalP (Petersen et al.
245 2011), TMHMM (Krogh et al. 2001) and Helix-Turn-Helix (Dodd & Egan 1990). Genome annotation
246 was manually curated with a scheme as proposed in (Bertelli et al. 2015). The complete and
247 annotated genome sequences have been deposited in the European Nucleic Archive under the
248 project PRJEB7990 with accession numbers LN879502 and LN879503.

249 **Genome analysis**

250 To identify CRISPR repeats, the genome sequences were submitted to CRISPRFinder (Grissa et al.
251 2007). The spacers within CRISPR locus of *Pr. naegleriophila* strains KNic and Diamant were
252 submitted to BLASTN (Altschul et al. 1997) homology searches against the non-redundant nucleotide
253 database. For phylogenetic reconstruction, multiple sequence alignments were performed with
254 Muscle V3.7 (Edgar 2004), and a neighbor-joining tree was reconstructed using Mega 6 (Tamura et al.
255 2013) with 1000 bootstrap, poisson distribution, gamma equal to 1. The two nearly complete
256 genomes of *Chlamydia* sp. Rubis and *Pr. naegleriophila* Diamant were reordered with Mauve (Darling
257 et al. 2004) by similarity to the closest available complete genome sequence *P. acanthamoebae* UV7
258 and *Pr. naegleriophila* KNic, respectively. These genomes and the complete genomes were aligned
259 using Mauve and the alignment was represented using GenoPlotR (Guy et al. 2010). Genomic islands
260 were predicted using IslandViewer (Dhillon et al. 2015). Home-made scripts for data analysis and
261 visualization were written in R (Cran 2010).

262

263 **ACKNOWLEDGMENTS**

264 We are grateful to Sébastien Aeby (University of Lausanne, Switzerland) for his technical help during
265 the gap closure stage. We would like to thank Burkhard Linke (Justus-Liebig-University Giessen,
266 Germany) for his assistance in maintaining this GenDB project. We acknowledge technical assistance
267 by the Bioinformatics Core Facility at JLU Giessen and access to resources financially supported by
268 the BMBF grant FKZ 031A533 within the de.NBI network. Part of the computations was performed at
269 the Vital-IT (<http://www.vital-it.ch>) Center for high-performance computing of the SIB Swiss Institute
270 of Bioinformatics.

271

272

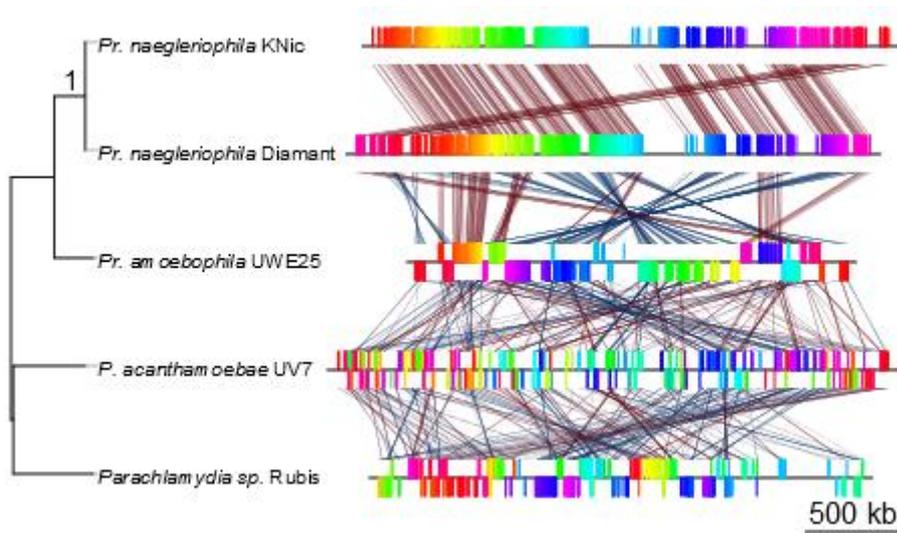
273 **REFERENCES**

- 274 Altschul SF et al. 1997. Gapped BLAST and PSI-BLAST: A new generation of protein database search
275 programs. *Nucleic Acids Res.* 25:3389–3402.
- 276 Bertelli C et al. 2015. Sequencing and characterizing the genome of *Estrella lausannensis* as an
277 undergraduate project: training students and biological insights. *Front. Microbiol.* 6.
- 278 Bertelli C, Goesmann A, Greub G. 2014. *Criblamydia sequanensis* Harbors a Megaplasmid Encoding
279 Arsenite Resistance. *Genome Announc.* 2.
- 280 Bertelli C, Greub G. 2012. Lateral gene exchanges shape the genomes of amoeba-resisting
281 microorganisms. *Front. Cell. Infect. Microbiol.* 2.
- 282 Bolotin A, Quinquis B, Sorokin A, Ehrlich SD. 2005. Clustered regularly interspaced short palindrome
283 repeats (CRISPRs) have spacers of extrachromosomal origin. *Microbiology.* 151:2551–61.
- 284 Brouns SJJ et al. 2008. Small CRISPR RNAs guide antiviral defense in prokaryotes. *Science.* 321:960–4.
- 285 Casson N, Michel R, Müller K-D, Aubert JD, Greub G. 2008. *Protochlamydia naegleriophila* as etiologic
286 agent of pneumonia. *Emerg. Infect. Dis.* 14:168–72.
- 287 Collingro A et al. 2011. Unity in Variety - the Pan-Genome of the Chlamydiae. *Mol. Biol. Evol.*
288 28:3253–3270.
- 289 Cran. 2010. The Comprehensive R Archive Network. *Wiley Interdiscip. Rev. Comput. Stat.* n/a–n/a.
- 290 Darby AC, Cho N-H, Fuxelius H-H, Westberg J, Andersson SGE. 2007. Intracellular pathogens go
291 extreme: genome evolution in the Rickettsiales. *Trends Genet.* 23:511–20.
- 292 Darling ACE, Mau B, Blattner FR, Perna NT. 2004. Mauve: multiple alignment of conserved genomic
293 sequence with rearrangements. *Genome Res.* 14:1394–1403.
- 294 Deveau H, Garneau JE, Moineau S. 2010. CRISPR/Cas system and its role in phage-bacteria
295 interactions. *Annu. Rev. Microbiol.* 64:475–93.
- 296 Dhillon BK et al. 2015. IslandViewer 3: more flexible, interactive genomic island discovery,
297 visualization and analysis. *Nucleic Acids Res.* 43:W104–8.
- 298 Dodd IB, Egan JB. 1990. Improved detection of helix-turn-helix DNA-binding motifs in protein
299 sequences. *Nucleic Acids Res.* 18:5019–26.
- 300 Domman D et al. 2014. Massive expansion of ubiquitination-related gene families within the
301 Chlamydiae. *Mol. Biol. Evol.*
- 302 Edgar RC. 2004. MUSCLE: a multiple sequence alignment method with reduced time and space
303 complexity. *BMC Bioinformatics.* 5:113.
- 304 Eugster M, Roten C-AH, Greub G. 2007. Analyses of six homologous proteins of *Protochlamydia*
305 *amoebophila* UWE25 encoded by large GC-rich genes (*lgr*): a model of evolution and concatenation
306 of leucine-rich repeats. *BMC Evol. Biol.* 7:231.

- 307 Everett KDE, Bush RM, Andersen AA. 1999. Emended description of the order Chlamydiales, proposal
308 of Parachlamydiaceae fam. nov. and Simkaniaceae fam. nov., each containing one monotypic genus,
309 revised taxonomy of the family Chlamydiaceae, including a new genus and five new species, and
310 standards. *Int. J. Syst. Bacteriol.* 49:415–440.
- 311 Gimenez G et al. 2011. Insight into cross-talk between intra-amoebal pathogens. *BMC Genomics.*
312 12:542.
- 313 Gordon D, Green P. 2013. Consed: a graphical editor for next-generation sequencing. *Bioinformatics.*
314 29:2936–7.
- 315 Greub G et al. 2009. High throughput sequencing and Proteomics to identify immunogenic proteins
316 of a new pathogen: The dirty genome approach. *PLoS One.* 4.
- 317 Greub G. 2010. International Committee on Systematics of Prokaryotes * Subcommittee on the
318 taxonomy of the Chlamydiae: Minutes of the inaugural closed meeting, 21 March 2009, Little Rock,
319 AR, USA. *Int. J. Syst. Evol. Microbiol.* 60:2691–2693.
- 320 Greub G. 2009. Parachlamydia acanthamoebae, an emerging agent of pneumonia. *Clin. Microbiol.*
321 *Infect.* 15:18–28.
- 322 Greub G, Collyn F, Guy L, Roten C-A. 2004. A genomic island present along the bacterial chromosome
323 of the Parachlamydiaceae UWE25, an obligate amoebal endosymbiont, encodes a potentially
324 functional F-like conjugative DNA transfer system. *BMC Microbiol.* 4:48.
- 325 Grissa I, Vergnaud G, Pourcel C. 2007. CRISPRFinder: a web tool to identify clustered regularly
326 interspaced short palindromic repeats. *Nucleic Acids Res.* 35:W52–7.
- 327 Guy L, Kultima JR, Andersson SGE. 2010. genoPlotR: comparative gene and genome visualization in R.
328 *Bioinformatics.* 26:2334–5.
- 329 Haft DH, Selengut J, Mongodin EF, Nelson KE. 2005. A guild of 45 CRISPR-associated (Cas) protein
330 families and multiple CRISPR/Cas subtypes exist in prokaryotic genomes. *PLoS Comput. Biol.* 1:e60.
- 331 Horn M et al. 2004. Illuminating the evolutionary history of chlamydiae. *Science.* 304:728–730.
- 332 Horn M. 2011. Phylum XXIV. Chlamydiae Garrity and Holt 2001. In: *Bergey's Manual of Systematic*
333 *Bacteriology - Volume 4: The* | Noel R. Krieg | Springer. Krieg, NR et al., editors. Springer Berlin
334 Heidelberg.
- 335 Horvath P, Barrangou R. 2010. CRISPR/Cas, the immune system of bacteria and archaea. *Science.*
336 327:167–70.
- 337 Hyatt D et al. 2010. Prodigal: prokaryotic gene recognition and translation initiation site
338 identification. *BMC Bioinformatics.* 11:119.
- 339 Ishida K et al. 2014. Amoebal endosymbiont Neochlamydia genome sequence illuminates the
340 bacterial role in the defense of the host amoebae against Legionella pneumophila. *PLoS One.* 9.

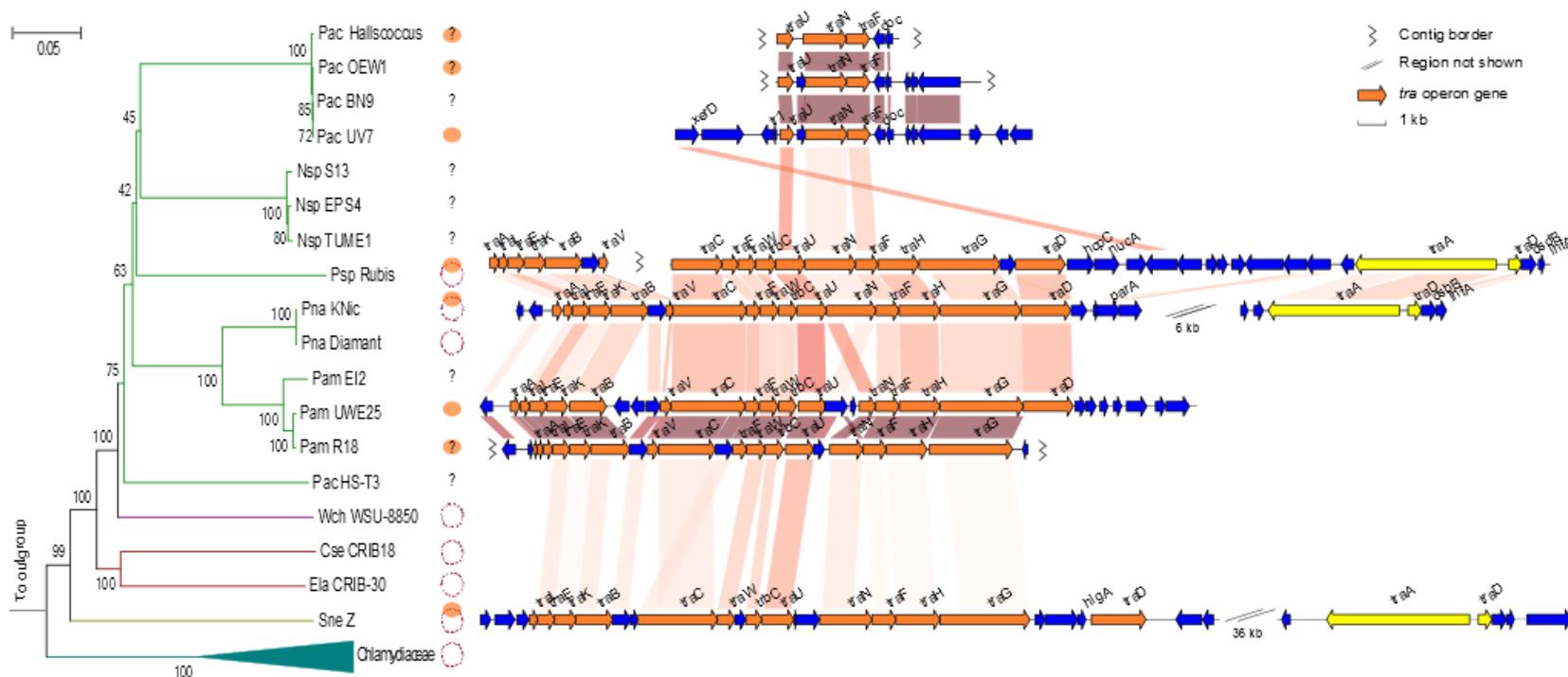
- 341 Ishino Y, Shinagawa H, Makino K, Amemura M, Nakata A. 1987. Nucleotide sequence of the *iap* gene,
342 responsible for alkaline phosphatase isozyme conversion in *Escherichia coli*, and identification of the
343 gene product. *J. Bacteriol.* 169:5429–33.
- 344 Jansen R, Embden JDA van, Gaastra W, Schouls LM. 2002. Identification of genes that are associated
345 with DNA repeats in prokaryotes. *Mol. Microbiol.* 43:1565–75.
- 346 Kamneva OK, Knight SJ, Liberles DA, Ward NL. 2012. Analysis of genome content evolution in PVC
347 bacterial super-phylum: Assessment of candidate genes associated with cellular organization and
348 lifestyle. *Genome Biol. Evol.* 4:1375–1380.
- 349 Koonin E V, Makarova KS. 2013. CRISPR-Cas: evolution of an RNA-based adaptive immunity system in
350 prokaryotes. *RNA Biol.* 10:679–86.
- 351 Krogh A, Larsson B, von Heijne G, Sonnhammer EL. 2001. Predicting transmembrane protein topology
352 with a hidden Markov model: application to complete genomes. *J. Mol. Biol.* 305:567–80.
- 353 Lagkouvardos I et al. 2014. Integrating metagenomic and amplicon databases to resolve the
354 phylogenetic and ecological diversity of the Chlamydiae. *ISME J.* 8:115–25.
- 355 Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.
356 *Bioinformatics.* 25:1754–1760.
- 357 Louwen R, Staals RHJ, Endtz HP, van Baarlen P, van der Oost J. 2014. The role of CRISPR-Cas systems
358 in virulence of pathogenic bacteria. *Microbiol. Mol. Biol. Rev.* 78:74–88.
- 359 Makarova KS et al. 2011. Evolution and classification of the CRISPR-Cas systems. *Nat. Rev. Microbiol.*
360 9:467–77.
- 361 Meyer F et al. 2003. GenDB--an open source genome annotation system for prokaryote genomes.
362 *Nucleic Acids Res.* 31:2187–2195.
- 363 Moliner C, Fournier P-E, Raoult D. 2010. Genome analysis of microorganisms living in amoebae
364 reveals a melting pot of evolution. *FEMS Microbiol. Rev.* 34:281–94.
- 365 Petersen TN, Brunak S, von Heijne G, Nielsen H. 2011. SignalP 4.0: discriminating signal peptides from
366 transmembrane regions. *Nat. Methods.* 8:785–6.
- 367 Sampson TR, Weiss DS. 2014. CRISPR-Cas systems: new players in gene regulation and bacterial
368 physiology. *Front. Cell. Infect. Microbiol.* 4:37.
- 369 Śliwa-Dominiak J, Suszyńska E, Pawlikowska M, Deptuła W. 2013. Chlamydia bacteriophages. *Arch.*
370 *Microbiol.* 195:765–771.
- 371 Tamura K, Stecher G, Peterson D, Filipski A, Kumar S. 2013. MEGA6: Molecular Evolutionary Genetics
372 Analysis version 6.0. *Mol. Biol. Evol.* 30:2725–9.
- 373 Weinberger AD, Wolf YI, Lobkovsky AE, Gilmore MS, Koonin E V. 2012. Viral diversity threshold for
374 adaptive immunity in prokaryotes. *MBio.* 3:e00456–12.
- 375

376 **FIGURES**



378 **Figure 1. Genomic rearrangements in the *Parachlamydiaceae* family**

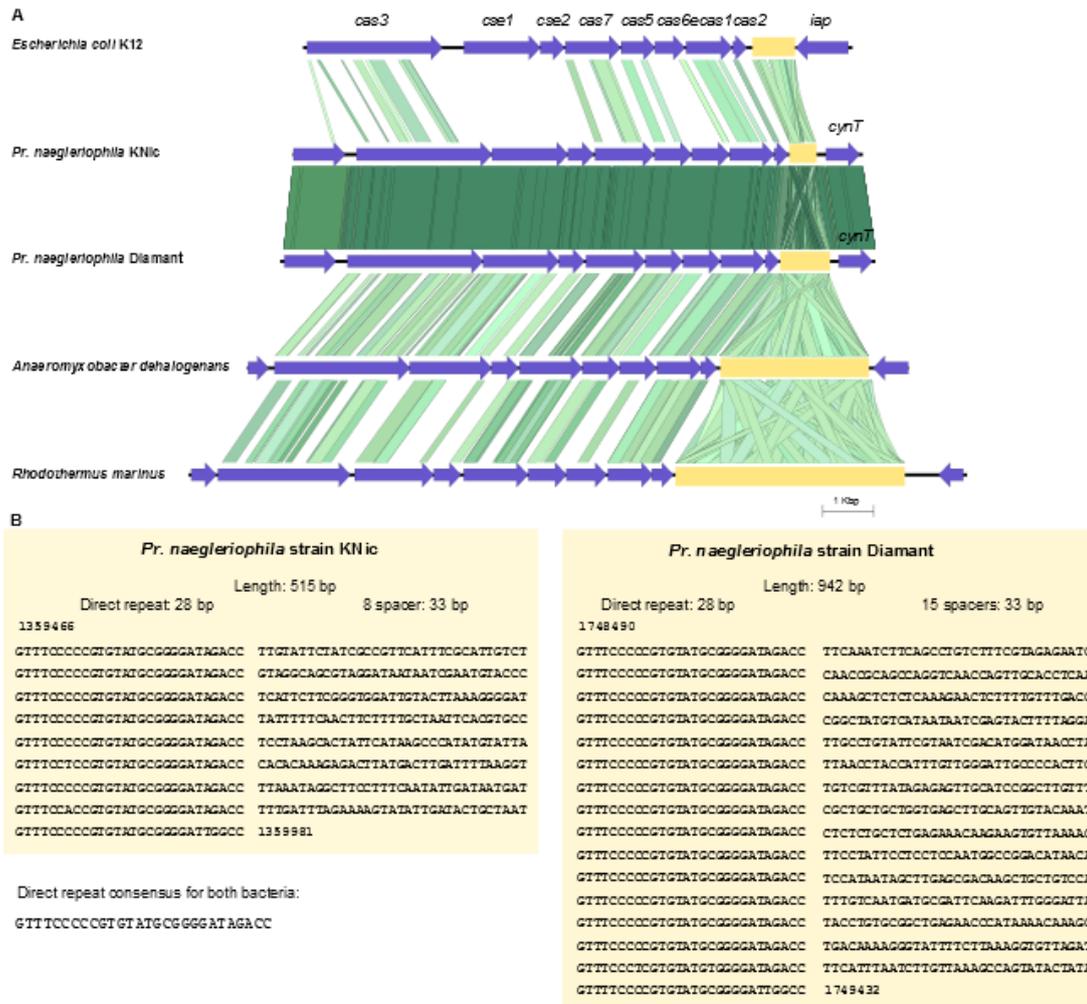
379 Left side, the phylogenetic branching of bacterial strains as inferred by a Neighbor-Joining tree
380 reconstruction based on 5 conserved proteins (DnaA, FtsK, HemL, FabI and SucA). Right side,
381 visualization of genomic rearrangements in the family *Parachlamydiaceae*. The two strains of the
382 species *Protochlamydia naegleriophila* are highly collinear, with no apparent rearrangement other
383 than due to the choice of the start of the genome sequence. With increasing distances between
384 organisms, the genomes show increasing number of rearrangements.



385

386 **Figure 2. *Chlamydiales* order, plasmids and type IV secretion system**

387 The left panel represents a neighbor joining tree of bacteria belonging the order *Chlamydiales* whose genome sequences is available based on four
 388 conserved proteins (DnaA, FtsK, HemL, FabI). The presence of a plasmid in each strain is represented by a small circular DNA molecule, and the draft
 389 genomes with no known plasmid described are indicated by a question mark as plasmids may be hidden among the numerous contigs. Orange ovals indicate
 390 the presence of a conjugative *tra* operon on the plasmid or in the bacterial chromosome. The right panel shows the conservation of the type IV secretion
 391 system *tra* operon and the surrounding genes. Pac: *P. acanthamoebae*, Nsp: *Neochlamydia* sp., Psp: *Protochlamydia* sp., Pna: *Pr. naegleriophila*, Pam: *Pr.*
 392 *amoebophila*, PacHS-T3: *Parachlamydiaceae* bacterium HS-T3.



393

394 **Figure 3. CRISPR locus and its associated genes**

395 A. CRISPR associated genes (CAS) consist of eight coding sequences, *cas3*, *cse1*, *cse2*, *cas7*, *cas5*,
 396 *cas6*, *cas1* and *cas2*, shown in blue within their genomic environment. Green lines connecting the
 397 genes in different organisms represent BLAST sequence homology with a gradient from light green to
 398 dark green for low to high percentage sequence identity, respectively. Genes neighboring the CRISPR
 399 locus present homology in *Pr. naegleriophila* genomes, but not to other genomes showing that the
 400 site of CRISPR locus insertion in *Pr. naegleriophila* genomes is different than in other bacteria. CRISPR
 401 repeats are found directly downstream of the CAS operon, as highlighted by the yellow box. B. Direct
 402 repeats and spacer sequences are detailed below.

403

404 TABLES

Species	Strain	Status	Scaffolds	Genome size	GC content	% coding	CDS	tRNAs	rRNA genes	Plasmid size	Plasmid CDS	Plasmid GC content
<i>Protochlamydia amoebophila</i>	UWE25	Complete	1	2,414,465	34.7		1855	35	7	-		
<i>Protochlamydia amoebophila</i>	EI2	Draft	178	2,397,675	34.8		1797	36	3	NA		
<i>Protochlamydia amoebophila</i>	R18	Draft	795	2,881,499	34.8		2025	41	13	NA		
<i>Protochlamydia naegleriophila</i>	KNic	Complete	1	2,885,090	42.7		2415	43	12	145,285	160	37.2
<i>Protochlamydia naegleriophila</i>	Diamant	Draft	4*	2,864,073	42.8		2424	39	7	91,928	98	40.9
<i>Parachlamydia acanthamoebae</i>	UV7	Complete	1	3,072,383	39		2531	40	10	-		
<i>Parachlamydia acanthamoebae</i>	Hall's coccus	Draft	95	2,971,261	39		2474	35	3	NA		
<i>Parachlamydia acanthamoebae</i>	OEW1	Draft	162	3,008,885	39		2321	38	4	NA		
<i>Parachlamydia acanthamoebae</i>	Bn9	Draft	72	2,999,361	38.9		2498°	NA	NA	NA		
<i>Parachlamydiaceae</i> bacterium	HS-T3	Draft	34	2,307,885	38.7		2003	39	3	NA		
<i>Parachlamydia sp.</i>	Rubis	Draft	3*	2,701,449	32.4		2446	36	5	80,697°	107	40.2
										39,075°	40	29.8
<i>Neochlamydia sp.</i>	EPS4	Draft	112	2,530,677	38.1		1843	36	4	NA		
<i>Neochlamydia sp.</i>	TUME1	Draft	254	2,546,323	38		1834	36	4	NA		
<i>Neochlamydia sp.</i>	S13	Draft	1342	3,187,074	38		2175	42	10	NA		

* Following removal of the plasmid(s) present according to our analyses. ° Not circularized, estimated size . NA : information not available

405

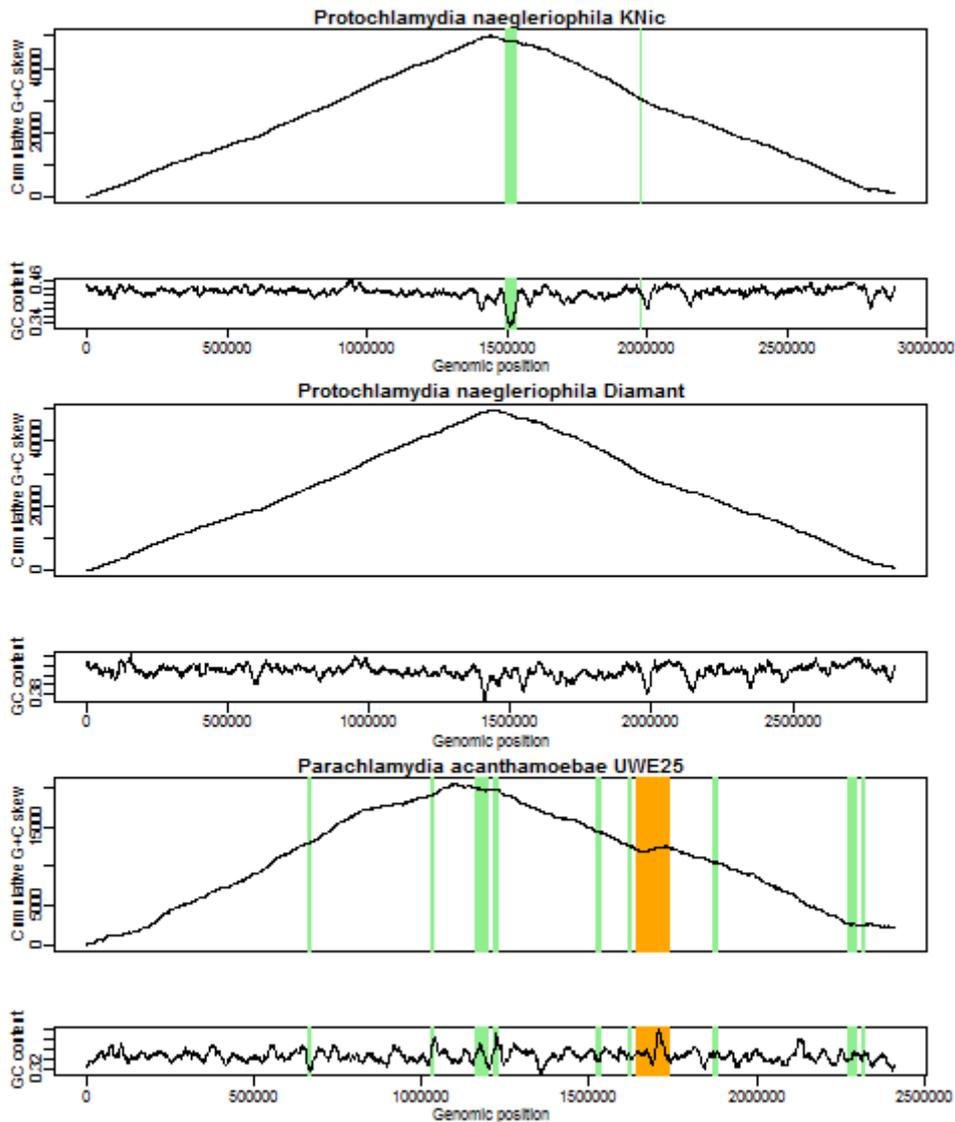
406 Table 1. Genomics characteristics of bacteria belonging to the family *Parachlamydiaceae*

407 As available on NCBI database on 22.09.2015, all genomes except KNic have been reannotated by the NCBI Prokaryotic Genome Annotation Pipeline.

408

409

410 SUPPLEMENTARY FIGURES

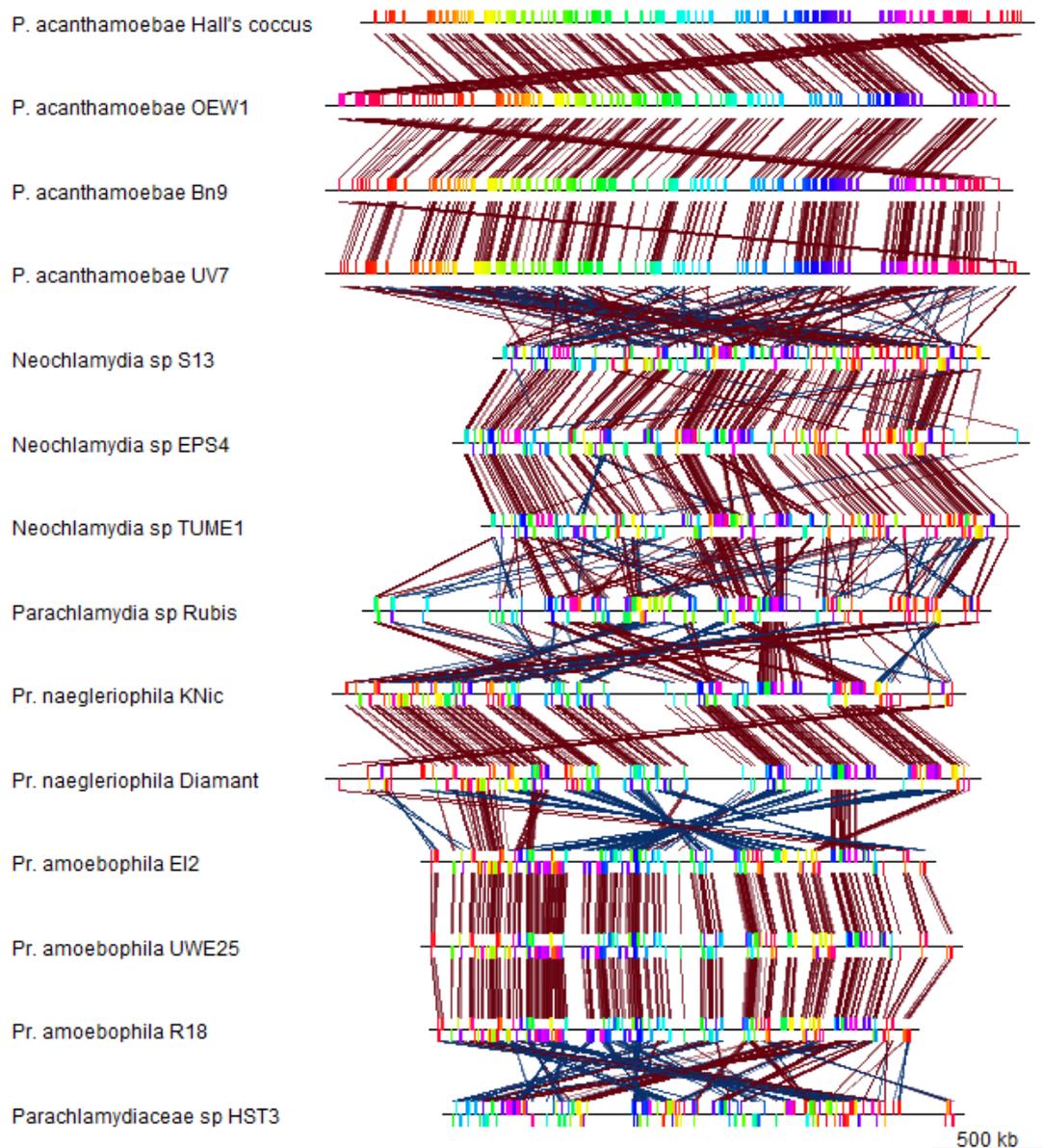


411

412 **Figure S1. Nucleotide biases and genomic islands**

413 The GC content and the bias of G versus C (GC skew) were determined in a sliding window of 1000 bp
414 along the genome sequence of *Pr. naegleriophila* KNic, *Pr. naegleriophila* Diamant and *Pr.*
415 *amoebophila* UWE25. The maximum of the GC skew indicates the putative terminus of replication.

416 The origin of replication is found at the minimum of the curve and has been used to assess the first
417 base of the genomic sequences. The contigs of *Pr. naegleriophila* Diamant have been reordered by
418 similarity to strain KNic and the putative plasmid removed. The location of predicted GIs for *Pr.*
419 *naegleriophila* KNic and *Pr amoebophila* UWE25 is indicated in green, whereas the confirmed *Pr.*
420 *amoebophila* genomic island Pam100G is indicated in orange.



421

422 **Figure S2. Genome rearrangements in the *Parachlamydiaceae***

423 This figure shows the difficulty of using draft genomes to observe genomic rearrangements. Available
424 draft genomes of members of the family *Parachlamydiaceae* were reordered using Mauve against
425 the complete genome of the most closely-related organism: *Pr. naegleriophila* Diamant and
426 *Chlamydia* sp. Rubis against *Pr. naegleriophila* KNic; *Pr. amoebophila* EI2 and R18 against *Pr.*
427 *amoebophila* UWE25; *P. acanthamoebae* BN9, Hall's coccus, and OEW1 against *P. acanthamoebae*
428 UV7; and *Parachlamydia* sp. HST3 and *Neochlamydia* sp. S13 against *Chlamydia* sp. Rubis, and
429 *Neochlamydia* sp. EPS4 and TUME1 against *Neochlamydia* sp. S13. As expected, no rearrangement
430 can be observed after the rearrangement of a fragmented draft genome against closely-related strain
431 of the same species. In the case of *Neochlamydia* sp. S13 and *Parachlamydia* sp. HST3 for which no
432 complete genome is available, numerous rearrangements can be seen, most likely as an artifact due
433 to the reduced similarity at nucleotide level between the reordered genome and the reference
434 genome.

435 **SUPPLEMENTARY TABLES**

436

Query	Percentage identity	Alignment length	Query coverage	Organism
spacer_1	95.83	24	73	Nippostrongylus brasiliensis (nematoda)
spacer_2	100	19	58	Sphingobacterium sp. PM2-P1-29 (bacteroides)
spacer_3	100	20	61	Paenibacillus polymyxa CR1 (firmicutes)
spacer_4	93.1	29	88	Tinamus guttatus (bird)
spacer_5	95.65	23	70	Sus scrofa (mammal)
spacer_6	92.86	28	85	Schistosoma margrebowiei (platyhelminthes)
spacer_7	85.71	35	97	Alteromonas macleodii str. 'Black Sea 11' (γ-prot)
spacer_8	93.1	29	88	Solanum pennellii (plant)

437

438 **Table S1 Homology of *P. naegleriophila* KNic spacers against the non-redundant nucleotide**
 439 **database by BLASTN**

440

Query	Percentage identity	Alignment length	Query coverage	Organism
spacer_1	100	19	58	uncultured archaeon
spacer_2	100	21	64	Pseudomonas aeruginosa
spacer_3	90	30	91	Populus euphratica
spacer_4	89.29	28	85	Spirometra erinaceieuropaei
spacer_5	92	25	76	Toxocara canis
spacer_6	95.65	23	70	Xenorhabdus nematophila AN6/1
spacer_7	100	19	58	Heterobasidion irregulare TC 32-1
spacer_8	87.5	32	97	Kluyveromyces marxianus
spacer_9	96	25	76	Homo sapiens
spacer_10	95.83	24	73	Phoenix dactylifera
spacer_10	100	21	64	Candidatus Nitrososphaera evergladensis SR1
spacer_11	95.83	24	73	Solanum tuberosum
spacer_12	96.3	27	82	Staphylococcus aureus subsp. aureus MSHR1132
spacer_13	100	21	64	Vibrio fischeri MJ11
spacer_14	92.31	26	79	Borrelia garinii PBi
spacer_15	92.59	27	82	Echinostoma caproni

441

442 **Table S2 Homology of *P. naegleriophila* Diamant spacers against the non-redundant nucleotide**
 443 **database by BLASTN**

444