



This is a pre- or post-print of an article published in
Pohl, S., Klockgether, J., Eckweiler, D., Khaledi, A.,
Schniederjans, M., Chouvarine, P., Tümmler, B., Häussler,
S.

The extensive set of accessory *Pseudomonas aeruginosa*
genomic components
(2014) FEMS Microbiology Letters, 356 (2), pp. 235-241.

1 **The extensive set of accessory *Pseudomonas aeruginosa***
2 **genomic components**

3 Sarah Pohl^{1,2}, Jens Klockgether³, Denitsa Eckweiler^{1,2}, Ariane Khaledi^{1,2}, Monika
4 Schniederjans^{1,2}, Philippe Chouvarine³, Burkhard Tümmler³, Susanne Häussler^{1,2*}

5 ¹ Department of Molecular Bacteriology, Helmholtz Centre for Infection Research, Braunschweig,
6 Germany

7 ² Institute for Molecular Bacteriology, TWINCORE GmbH, Centre for Clinical and Experimental
8 Infection Research, a joint venture of the Hannover Medical School and the Helmholtz Centre for
9 Infection Research, Hannover, Germany

10 ³ Clinical Research Group 'Molecular Pathology of Cystic Fibrosis and Pseudomonas Genomics',
11 Clinic for Pediatric Pneumology, Allergology and Neonatology, OE 6710, Hannover Medical School,
12 Hannover, Germany

13 *Corresponding author:

14 Susanne Häussler

15 Address: Department of Molecular Bacteriology, Helmholtz Centre for Infection Research,
16 Inhoffenstrasse 7, D-38124 Braunschweig

17 Tel: +49 (0)531 61 81 - 30 00

18 Fax: +49 (0)531 61 81 - 30 99

19 Email: susanne.haeussler@helmholtz-hzi.de

20

21 **Keywords**

22 Extended gene pool, horizontal gene transfer

23 **Running Title**

24 Extended Accessory Genome of *Pseudomonas aeruginosa*

25

26 **Abstract**

27 Up to 20% of the chromosomal *Pseudomonas aeruginosa* DNA belong to the so-
28 called accessory genome. Its elements are specific for subgroups or even single
29 strains and are likely acquired by horizontal gene transfer (HGT). Similarities of the
30 accessory genomic elements to DNA from other bacterial species, mainly the DNA of
31 γ - and β -proteobacteria, indicate a role of interspecies HGT. In this study we
32 analysed the expression of the accessory genome in 150 clinical *P. aeruginosa*
33 isolates as uncovered by transcriptome sequencing and the presence of accessory
34 genes in eleven additional isolates. Remarkably, despite the large number of *P.*
35 *aeruginosa* strains that have been sequenced to date, we found new strain-specific
36 compositions of accessory genomic elements and a high portion (10 – 20%) of genes
37 without *P. aeruginosa* homologues. Although some genes were detected to be
38 expressed/present in several isolates, individual patterns regarding the genes, their
39 functions and the possible origin of the DNA were widespread among the tested
40 strains. Our results demonstrate the unaltered potential to discover new traits within
41 the *P. aeruginosa* population and underline that the *P. aeruginosa* pangenome is
42 likely to increase with increasing sequence information.

43

44

45 **Introduction**

46 *Pseudomonas aeruginosa* bacteria inhabit various aquatic and terrestrial
47 environments and are a frequent cause of opportunistic infections in animal and
48 human hosts. Colonisation of such a broad spectrum of habitats goes along with
49 broad metabolic versatility and a high potential to adapt to new environments

50 (Ramos, 2004). In addition to many common features, *P. aeruginosa* isolates often
51 display individual phenotypic traits. This phenotypic diversity is reflected in the
52 composition of the genomes: The major part of a *P. aeruginosa* genome, the core
53 genome, is found in all strains and, with the exception of a few loci, is usually highly
54 conserved with only 0.5 – 0.7% sequence diversity (Spencer *et al.*, 2003; Cramer *et*
55 *al.*, 2011). Less conserved loci are the so-called ‘replacement islands’. Flagellin
56 glycosylation genes, O-antigen and pyoverdine biosynthesis gene clusters and the
57 major pilin gene *pilA* are present in all genomes. However, they are under
58 diversifying selection (Kung *et al.*, 2010) and can be individually shaped with regards
59 to size and gene composition. Comparative analyses revealed different types and
60 subtypes of these gene clusters (Klockgether *et al.*, 2011).

61 Besides that, up to 20% of a genome can be composed of DNA blocks specific for
62 subgroups of strains or even single strains. Such DNA blocks appear as genomic
63 islands or islets inserted between core genome parts at various chromosomal loci.
64 Their sizes range from a few hundred bp up to 200 kbp, and together they make up a
65 strain’s accessory genome. The individual composition causes variable genome
66 sizes, usually between six and seven Mbp (Klockgether *et al.*, 2011). Accessory DNA
67 blocks often display features of mobile elements such as phages, transposons or
68 integrative and conjugative elements (Kung *et al.*, 2010), indicating the importance of
69 mobile DNA in shaping the accessory genome. These mobile elements serve as
70 vehicles for associated ‘cargo’ DNA that integrates into the host genome as well.
71 Several large genomic islands and smaller insertions have been described so far
72 (Klockgether *et al.*, 2011). They were either sequenced separately, such as PAGI-1
73 (Liang *et al.*, 2001) or PAGI-5 to -11 (Battle *et al.*, 2008; Battle *et al.*, 2009), or they
74 were identified from fully sequenced strains upon comparison with existing genomes,

75 such as the large genomic islands and prophages of strain LESB58 (Winstanley *et*
76 *al.*, 2009).

77 The presence of genomic islands in the accessory genome of *P. aeruginosa* goes
78 along with the acquisition of additional functional traits. Genes typical for mobile
79 elements and ORFs encoding for hypothetical or uncharacterized proteins have been
80 found, as well as genes encoding catabolic traits or antibiotic resistance features
81 (Kung *et al.*, 2010). Contribution of genomic islands to a strain's pathogenic potential
82 has also been reported (He *et al.*, 2004). While accessory genome elements
83 generally display homology to DNA from other bacteria, especially the appearance of
84 mobile elements indicates horizontal gene transfer (HGT). These mobile elements
85 are often related in various species, mainly β - and γ -proteobacteria (Klockgether *et*
86 *al.*, 2006), implying HGT across species borders. Thus, *P. aeruginosa* has in principle
87 access to an extended pool of DNA which can contribute to a strain's accessory
88 genome. Questions on size and possible restrictions of this DNA pool and to what
89 extent it is actually exploited by single *P. aeruginosa* strains are difficult to answer,
90 but increasing insights are expected with every additional genome being sequenced.
91 Recent accessory genome analyses still revealed dozens of genes without *P.*
92 *aeruginosa* homologues, but homologues from other species – other Pseudomonads
93 as well as other genera (Bezuidt *et al.*, 2013; Klockgether *et al.*, 2013).

94 In this study we screened for the presence and the expression of accessory genes in
95 a large collection of *P. aeruginosa* isolates. Searches for high sequence identities to
96 genes from previously sequenced *P. aeruginosa* strains, other Pseudomonads and
97 other genera were performed, and the nature and distribution of acquired genes
98 within our strain pools were evaluated. Furthermore, in order to gain broader
99 knowledge on the principles of accessory genome composition, predicted functions of

100 the acquired genes were assessed. Analysis was done on full transcriptome datasets
101 for 150 individual clinical isolates and on genome sequence data for eleven *P.*
102 *aeruginosa* strains isolated from airways of cystic fibrosis patients.

103

104 **Materials & Methods**

105 ***P. aeruginosa* isolates**

106 The 150 transcriptome sequenced clinical isolates were provided by different clinics
107 or research institutions; 87 from Hannover Medical School (MHH), 40 from a strain
108 collection of the University of Freiburg, 13 from the Robert-Koch-Institute in
109 Wernigerode and 10 from the National Reference Laboratory for multidrug-resistant
110 Gram-negative bacteria in Bochum.

111 Eleven genome sequenced isolates were chosen from a collection of isolates from
112 the airways of CF patients at Hannover Medical School. Genotyping with a custom-
113 made microarray (Wiehlmann *et al.*, 2007) indicated that the selected isolates
114 represent clonal lineages that persisted in the respective patient airways for seven
115 years or longer. Analysis was performed on the sequenced genomes of the early
116 isolates.

117

118 **Sequencing**

119 Whole transcriptome-sequencing was performed using a custom-made protocol with
120 barcoded RNA-libraries to enable pooled sequencing of several samples (Dötsch *et*
121 *al.*, 2012). Briefly, bacterial cultures were grown under standard conditions (LB broth,
122 37 °C) and harvested in RNAprotect (Qiagen) at OD₆₀₀= 2. Total RNA was extracted
123 (RNeasy plus, Qiagen) and the ribosomal RNA was depleted (MICROBExpress,
124 Ambion). After mechanical fragmentation of the mRNA and ligation with RNA-
125 adapters containing 6 nt barcode sequences for multiplexing, strand-specific cDNA
126 libraries were generated. Up to 12 libraries were pooled and sequenced, after
127 additional rRNA removal, on an Illumina GenomeAnalyzer-II_x generating paired-end

128 reads of 76 or 110 bp. The raw sequence reads were sorted by their expected
129 barcodes, which were subsequently removed. Reads with more than one error in
130 their barcode were discarded; the others were trimmed using the fastq-mcf script of
131 the ea-utils package (Aronesty, 2011), removing adapters and low quality sequences.
132 Genome sequencing was done on a SOLiD 5500XL system (Life Technologies).
133 Bead-coupled fragment libraries for high-throughput sequencing were generated
134 from genomic DNA of bacteria grown at 37°C in LB-medium. DNA fragmentation and
135 processing was done using kits and standard protocols provided by Life
136 Technologies (Fragment Library Preparation: 5500 Series SOLiD™ Systems User
137 Guide (Part no. 4460960), 5500 SOLiD™ Fragment Library Core Kit (Part no.
138 4464412) and EZ Bead™ E20 System Consumables (Part no. 4453094)). The
139 resulting datasets consisted of 75 bp reads, which were trimmed by removing 3' ends
140 with bad base call qualities (< 20). Duplicated reads were removed before further
141 analysis. The reads for each isolate were aligned to the PA14 reference genome
142 (NC_008463.1) using NovoalignCS. Reads not aligned to the PA14 reference were
143 then used for accessory genome analysis.

144

145 ***De novo* Analysis**

146 The transcriptome reads were mapped against five *P. aeruginosa* reference
147 genomes (PA14, PAO1, LESB58, PACS2 and PA7) using Stampy (Lunter &
148 Goodson, 2011). Sequencing reads that did not map to any of the references were
149 used in a *de novo* transcriptome assembly approach using Velvet (Zerbino, 2008)
150 with a wide range of k-mers (27 to 37) and a minimal contig length of 100 bp.
151 202,907 contigs were blasted against all microbial genes downloaded from the
152 Microbial Genome Database (MBGD) (Uchiyama *et al.*, 2010) using a minimal hit

153 length of 100 bp and sequence similarity higher than 90% as cut-off values. An initial
154 accessory gene list containing 9242 entries was established by extracting always the
155 first coding sequence that neither belonged to nor was an ortholog of the *P.*
156 *aeruginosa* reference genomes used. Previously unmapped sequencing reads were
157 mapped against this list using Stampy, and gene coverage was calculated based on
158 the mapping results. Final analysis was performed only on genes with at least 90%
159 coverage in the respective isolate(s). All figures were produced using R for statistical
160 computing (R Core Team, 2013) and the pheatmap package (Kolde, 2013).

161 *De novo* assembly of genomic reads was done using ABySS (Birol *et al.*, 2009) with
162 k-mer size spaces optimised for each dataset to maximise the N50 value. Contigs
163 smaller than 200 bps were removed, while contigs \geq 200 bps were assembled into
164 supercontigs using the software CAP3 (Huang & Madan, 1999). Supercontigs and
165 unassembled regular contigs were then aligned to all bacterial entries of the UniProt
166 database (UniProt Consortium, 2014) using blastx to identify known proteins/genes
167 from other *P. aeruginosa* strains or other species. The results were filtered for the
168 best hit of a given region in a contig in order to remove multiple descriptions for the
169 same coding region.

170

171 **Results**

172 **Accessory Genome Analysis from Transcriptome Sequencing**

173 When the transcriptome sequencing reads of 150 clinical *Pseudomonas aeruginosa*
174 isolates were mapped against five *P. aeruginosa* reference genomes (PA14, PAO1,
175 LESB58, PACS2 and PA7), around 57% of the reads mapped to all five references,
176 and around 20% mapped to at least one of them. Sequencing reads that did not map
177 to any of the references were used in a *de novo* transcriptome assembly approach
178 resulting in an accessory gene collection containing 9242 entries. Mapping of the
179 previously unmapped sequencing reads to this accessory gene collection revealed
180 the expression of 1164 genes with at least 90% coverage. 106 of these genes have
181 previously been found in other *P. aeruginosa* strains and 1058 genes showed
182 homologues in a variety of non-*P. aeruginosa* species (Table S1).

183 The vast majority of the latter genes were found in only few isolates: 752 (71%)
184 genes were well-covered in one to five, another 129 (12%) in six to ten isolates.
185 Furthermore, the number and combination of new genes varied between the isolates
186 (Figure 1). A group of three clinical isolates contained no well-covered non-*P.*
187 *aeruginosa* homologues, and another 53 contained less than 20 new genes.
188 Nevertheless, most isolates displayed several dozens of non-*P. aeruginosa*
189 homologues (up to 162) in their accessory genomes. While subgroups of isolates
190 shared small sets of genes, no large common gene clusters were found. Apparently,
191 most of the 150 clinical isolates harbour an individual set of accessory genes.

192 Protein descriptions were obtained from the Microbial Genome Database (Table S1)
193 and grouped into the predicted protein function categories shown in Table 1. Overall,
194 40% of the 1058 transcribed genes are described to encode uncharacterised or
195 hypothetical proteins, followed by proteins related to DNA integration and

196 transposition (10%) and transcriptional regulators (6%). Other frequent functional
197 categories are metal resistance (35 proteins) and antibiotic resistance (38 proteins).
198 The metal resistance proteins, often described as copper resistance proteins,
199 originated from species like *P. mendocina* and *P. stutzeri*, but also from *Ralstonia*
200 *metallidurans*, *Achromobacter xylosoxidans* and *Stenotrophomonas maltophilia*. The
201 antibiotic resistance proteins also came from various species like *Acinetobacter*
202 *baumannii*, *E. coli*, *Klebsiella pneumoniae* and *Salmonella enterica*. Some of the
203 proteins are β -lactamases, other confer resistance to gentamicin, streptomycin or
204 tetracycline. Since the genes encoding for these proteins were found via
205 transcriptome sequencing, it cannot be differentiated whether they integrated into the
206 genomes or are localized on resistance plasmids.

207 Overall, the 1058 transcribed genes were annotated in 98 potential donor species
208 other than *P. aeruginosa*, and Table 2 presents the most common ones. Predominant
209 are other Pseudomonads like *P. stutzeri* and *P. putida*, of which 114 and 69 genes
210 were found, but also *K. pneumoniae* provided many genes to this pool. Other γ - and
211 β -proteobacteria like *R. metallidurans*, *S. enterica* and *E. coli* appeared multiple times
212 as well.

213 Table 2 also displays the number of isolates which have acquired genes annotated in
214 the respective bacterial species. Not only was a great variety of different genes from
215 other Pseudomonads found in general, but also a large number of clinical isolates
216 expressed at least one gene of *P. stutzeri* (111 isolates) or *P. putida* (52 isolates).
217 Additionally, although a more limited number of genes was annotated in the
218 Enterobacteriaceae *E. coli* (25 genes) and *S. enterica* (42 genes), homologues from
219 these species were found in as many as 58 and 38 of the 150 isolates included in
220 this approach.

221 Figure 2 illustrates that homologues from a maximum of 33 different species were
222 discovered in single clinical isolates. Interestingly, some pairs of isolates shared
223 almost identical sets of species, and these pairs were mostly phylogenetically closely
224 related. However, since not all of them share the same set of individual genes, this
225 finding indicates that rather than having a common ancestor they shared a common
226 habitat with access to a similar gene pool.

227

228 **Accessory Genome Analysis from Genome Sequencing**

229 After sequencing the genomes of eleven *P. aeruginosa* isolates from different cystic
230 fibrosis patients, a *de novo* assembly was performed with all reads that did not map
231 to the PA14 reference genome. Blastx comparison of the resulting accessory
232 genome contigs against the UniProt database led to the identification of 69 – 192 loci
233 with homologues in different species. For a majority of the accessory ORFs, the
234 closest homologues were found in other *P. aeruginosa* strains, whereas 10 – 20%
235 were from other bacterial species (Table 3). Similar to the transcriptome results,
236 these other species contained Pseudomonads as well as bacteria from other genera
237 (Table 3). In total, similarities to Pseudomonads such as *P. putida*, *P. stutzeri* or *P.*
238 *fluorescens* were more frequent (Table 4), but homologues of other proteobacteria
239 were detected for each isolate (Table S2). Most of them were γ -, few were β -
240 proteobacteria (e. g. *Ralstonia* or *Polaromonas* sp.), but also species from other
241 classes were found in individual cases (Table S2).

242 Analysis of the protein names of the identified ORFs of all analysed isolates revealed
243 that more than 50% code for uncharacterized or hypothetical proteins, which is a
244 slightly higher proportion than in the transcriptomic approach (40%). Some functional
245 categories were found for many isolates (Table S3A), like genes from replacement

246 islands and integrase/transposase as well as phage- and plasmid-related genes, all
247 reflecting the role of such elements in shaping the accessory genome. Other
248 categories appearing more often were transcriptional regulators and transporter
249 components. Some features were seen for few isolates only, such as copper
250 resistance or Type I restriction modification system genes, and some function
251 predictions appeared in single strains exclusively, e.g. type VI secretion system
252 components in isolate NA1. When focussing only on non-*P. aeruginosa* hits (Table
253 S3B), most were again uncharacterized and hypothetical proteins. Phage
254 components or other proteins typical for mobile elements were also in these shorter
255 lists. However, these lists were all highly individual: no mobile element component
256 appeared as best homologue for ORFs from many isolates.

257

258 **Discussion**

259 Data of a transcriptomic and a genomic sequencing approach were used
260 independently in this study to analyse the accessory genomes of 150 and eleven
261 clinical *Pseudomonas aeruginosa* isolates, respectively. In the transcriptome
262 analysis, between none and 162 genes without *P. aeruginosa* homologues were
263 detected to be transcribed in individual isolates (mean: 40); the genomic sequencing
264 revealed the presence of 7 – 40 (mean: 20) ‘foreign’ genes per isolate. Among the
265 non-*P. aeruginosa* species for which the accessory genes were already described, *P.*
266 *stutzeri* and *P. putida* were detected most frequently in both approaches. However,
267 genetic elements from genera other than *Pseudomonas* were found more often within
268 the transcriptome sequencing results. Most interestingly, independent of the chosen
269 approach, no large groups of isolates shared the same pattern of so far in *P.*
270 *aeruginosa* unannotated genes within their accessory genomes. In the same line,

271 many functional elements were found specifically within subgroups or even within
272 single isolates, so that, as for the new genes, no common patterns were seen in the
273 chosen collection. Contrary to this, the transcriptomic approach showed that
274 phylogenetically closely related isolates often express different genes from the same
275 species, implying that they had access to specific gene pools.

276 In summary, the extended analysis of two independent collections of clinical *P.*
277 *aeruginosa* isolates indicates that a broad gene pool from various bacterial species
278 serves as a possible source of accessory genome elements that can be accepted by
279 *P. aeruginosa*. Despite the already large number of *P. aeruginosa* genomes
280 deposited in the databases, almost all isolates analysed in this study exhibited genes
281 that were not known to belong to the accessory genome of this species so far. A
282 great variety of potential donor species as well as function predictions could be
283 detected by our approaches, demonstrating the ongoing possibility to find individual
284 genomic traits whenever the DNA or the RNA of an independent *P. aeruginosa* strain
285 is sequenced. This would imply an extended pangenome for this species, the limits of
286 which could not be unequivocally estimated.

287

288 **Acknowledgements**

289 We thank Iris F. Chaberny (Hannover Medical School, Hannover, Germany), Daniel
290 Jonas (Freiburg University Medical Centre, Freiburg, Germany), Wolfgang Witte and
291 Yvonne Pfeifer (Robert-Koch-Institute, Wernigerode, Germany) and Martin Kaase
292 and Sören Gatermann (National Reference Laboratory for multidrug-resistant Gram-
293 negative bacteria, Bochum, Germany) for providing us with clinical *P. aeruginosa*
294 isolates.

295 This work was supported by the Helmholtz International Graduate School for
296 Infection Research, the Deutsche Forschungsgemeinschaft (CRC 900 'Chronic
297 Infections: Microbial Persistence and its Control', projects A2 and A3) and the
298 European Research Council (Starting Grant 260276 'RESISTOME').

299

300 **References**

- 301 Aronesty E (2011) ea-utils: command-line tools for processing biological sequencing data
302 <http://code.google.com/p/ea-utils>.
- 303 Battle SE, Rello J & Hauser AR (2009) Genomic islands of *Pseudomonas aeruginosa*. *FEMS*
304 *Microbiol Lett* **290**: 70-8.
- 305 Battle SE, Meyer F, Rello J, Kung VL & Hauser AR (2008) Hybrid pathogenicity island PAGI-
306 5 contributes to the highly virulent phenotype of a *Pseudomonas aeruginosa* isolate in
307 mammals. *J Bacteriol* **190**: 7130-40.
- 308 Bezuidt OK, Klockgether J, Elsen S, Attree I, Davenport CF & Tümmler B (2013) Intracolonial
309 genome diversity of *Pseudomonas aeruginosa* clones CHA and TB. *BMC Genomics* **14**: 416.
- 310 Birol I, Jackman SD, Nielsen CB, *et al.* (2009) De novo transcriptome assembly with ABySS.
311 *Bioinformatics* **25**: 2872-7.
- 312 Cramer N, Klockgether J, Wrasman K, Schmidt M, Davenport CF & Tümmler B (2011)
313 Microevolution of the major common *Pseudomonas aeruginosa* clones C and PA14 in cystic
314 fibrosis lungs. *Environ Microbiol* **13**: 1690-704.
- 315 Dötsch A, Eckweiler D, Schniederjans M, Zimmermann A, Jensen V, Scharfe M, Geffers R &
316 Häussler S (2012) The *Pseudomonas aeruginosa* transcriptome in planktonic cultures and
317 static biofilms using RNA sequencing. *PLoS One* **7**: e31092.
- 318 He J, Baldini RL, Deziel E, Saucier M, Zhang Q, Liberati NT, Lee D, Urbach J, Goodman HM
319 & Rahme LG (2004) The broad host range pathogen *Pseudomonas aeruginosa* strain PA14
320 carries two pathogenicity islands harboring plant and animal virulence genes. *Proc Natl Acad*
321 *Sci U S A* **101**: 2530-5.
- 322 Huang X & Madan A (1999) CAP3: A DNA sequence assembly program. *Genome Res* **9**:
323 868-77.
- 324 Klockgether J, Reva ON & Tümmler B (2006) Spread of genomic islands between clinical
325 and environmental strains. *Prokaryotic diversity: mechanisms and significance*,(Logan NA,
326 Lapin-Scott HM & Oyston PCF, eds.), 187-200. Cambridge University Press, Cambridge,
327 United Kingdom.

328 Klockgether J, Cramer N, Wiehlmann L, Davenport CF & Tümmler B (2011) *Pseudomonas*
329 *aeruginosa* Genomic Structure and Diversity. *Front Microbiol* **2**: 150.

330 Klockgether J, Miethke N, Kubesch P, *et al.* (2013) Intraclonal diversity of the *Pseudomonas*
331 *aeruginosa* cystic fibrosis airway isolates TBCF10839 and TBCF121838: distinct signatures
332 of transcriptome, proteome, metabolome, adherence and pathogenicity despite an almost
333 identical genome sequence. *Environ Microbiol* **15**: 191-210.

334 Kolde R (2013) pheatmap: Pretty Heatmaps <http://CRAN.R-project.org/package=pheatmap>.

335 Kung VL, Ozer EA & Hauser AR (2010) The accessory genome of *Pseudomonas*
336 *aeruginosa*. *Microbiol Mol Biol Rev* **74**: 621-41.

337 Liang X, Pham XQ, Olson MV & Lory S (2001) Identification of a genomic island present in
338 the majority of pathogenic isolates of *Pseudomonas aeruginosa*. *J Bacteriol* **183**: 843-53.

339 Lunter G & Goodson M (2011) Stampy: a statistical algorithm for sensitive and fast mapping
340 of Illumina sequence reads. *Genome Res* **21**: 936-9.

341 R Core Team (2013) R: A Language and Environment for Statistical Computing
342 <http://www.R-project.org/>.

343 Ramos J-L (ed.) (2004) *Pseudomonas Volume 1: Genomics, Life Style and Molecular*
344 *Architecture*. New York: Kluwer Academics / Plenum publishers.

345 Spencer DH, Kas A, Smith EE, Raymond CK, Sims EH, Hastings M, Burns JL, Kaul R &
346 Olson MV (2003) Whole-genome sequence variation among multiple isolates of
347 *Pseudomonas aeruginosa*. *J Bacteriol* **185**: 1316-25.

348 Uchiyama I, Higuchi T & Kawai M (2010) MBGD update 2010: toward a comprehensive
349 resource for exploring microbial genome diversity. *Nucleic Acids Res* **38**: D361-5.

350 UniProt Consortium (2014) Activities at the Universal Protein Resource (UniProt). *Nucleic*
351 *Acids Res* **42**: D191-8.

352 Wiehlmann L, Wagner G, Cramer N, *et al.* (2007) Population structure of *Pseudomonas*
353 *aeruginosa*. *Proc Natl Acad Sci U S A* **104**: 8101-6.

354 Winstanley C, Langille MG, Fothergill JL *et al.* (2009) Newly introduced genomic prophage
355 islands are critical determinants of in vivo competitiveness in the Liverpool Epidemic Strain of
356 *Pseudomonas aeruginosa*. *Genome Res* **19**: 12–23.

357 Zerbino DR & Birney E (2008) Velvet: algorithms for de novo short read assembly using de
358 Bruijn graphs. *Genome Res* **18**: 821–829.

359 **Table 1.** Functional profiling of proteins encoded in the accessory genome of 150
360 clinical isolates without *P. aeruginosa* homologues in the Microbial Genome
361 Database.

Functional group	Number of proteins
uncharacterised/hypothetical	428
integrase/transposase	111
transcriptional regulator	59
antibiotic resistance	38
metal resistance	35
transporter component	26
phage protein	21
Type I restriction modification system	17
helicase	15

362

363 **Table 2.** Species distribution of non-*P. aeruginosa* homologues within the accessory
 364 genome of 150 clinical isolates.

Species	No. of genes from the species	No. of isolates with at least one gene from the species
<i>Pseudomonas stutzeri</i>	114	111
<i>Pseudomonas putida</i>	69	52
<i>Klebsiella pneumoniae</i>	64	79
<i>Acidovorax</i> sp.	57	93
<i>Ralstonia metallidurans</i>	51	41
<i>Burkholderia vietnamiensis</i>	50	42
<i>Pseudomonas mendocina</i>	48	82
<i>Alicyclophilus denitrificans</i>	44	89
<i>Salmonella enterica</i>	42	38
<i>Azotobacter vinelandii</i>	29	75
<i>Delftia acidovorans</i>	29	29
<i>Ralstonia pickettii</i>	29	56
<i>Escherichia coli</i>	25	58
<i>Stenotrophomonas maltophilia</i>	25	64
<i>Verminephrobacter eiseniae</i>	24	37
<i>Nitrosomonas europaea</i>	19	38
<i>Achromobacter xylosoxidans</i>	17	61
<i>Dechlorosoma suillum</i>	15	49
<i>Methylophaga</i> sp.	15	29

365

366 79 additional species containing less than 15 non-*P. aeruginosa* homologues are not
 367 listed.

368 **Table 3.** Numbers of proteins not known from the PA14 reference detected by
 369 genome sequencing.

Strain	No. of proteins	From other <i>P. aeruginosa</i>	From other Pseudomonads	From other genera
AA2	69	56	5	8
HK3	120	99	12	9
KB1	77	64	12	1
KK1	164	133	24	7
MF1	157	132	12	13
NA1	81	60	12	9
NM1	192	152	29	11
RP1	191	174	9	8
ST1	109	102	5	2
TR1	143	127	12	4
WU2	138	123	7	8

370

371 **Table 4.** Species distribution of non-*P. aeruginosa* homologues in eleven genome
 372 sequencing datasets.

Species	No. of genes from the species	No. of isolates with at least one gene from the species
<i>Pseudomonas</i> sp.	33	9
<i>Pseudomonas putida</i>	27	9
<i>Pseudomonas stutzeri</i>	24	8
<i>Pseudomonas syringae</i>	17	9
<i>Pseudomonas fluorescens</i>	14	8
<i>Pseudomonas chlororaphis</i>	6	2
<i>Acidovorax</i> sp.	5	2
<i>Polaromonas</i> sp.	4	2
<i>Providencia burhodogranariea</i>	4	2
<i>Pseudomonas mendocina</i>	4	3
<i>Ralstonia metallidurans</i>	4	3
<i>Stenotrophomonas maltophilia</i>	4	2
<i>Pseudomonas fulva</i>	3	3
<i>Pseudomonas pseudoalcaligenes</i>	3	3
<i>Salmonella enterica</i>	3	2

373

374 55 additional species are not listed. In nine of those two and in 46 of those one non-

375 *P. aeruginosa* homologues were detected.

376

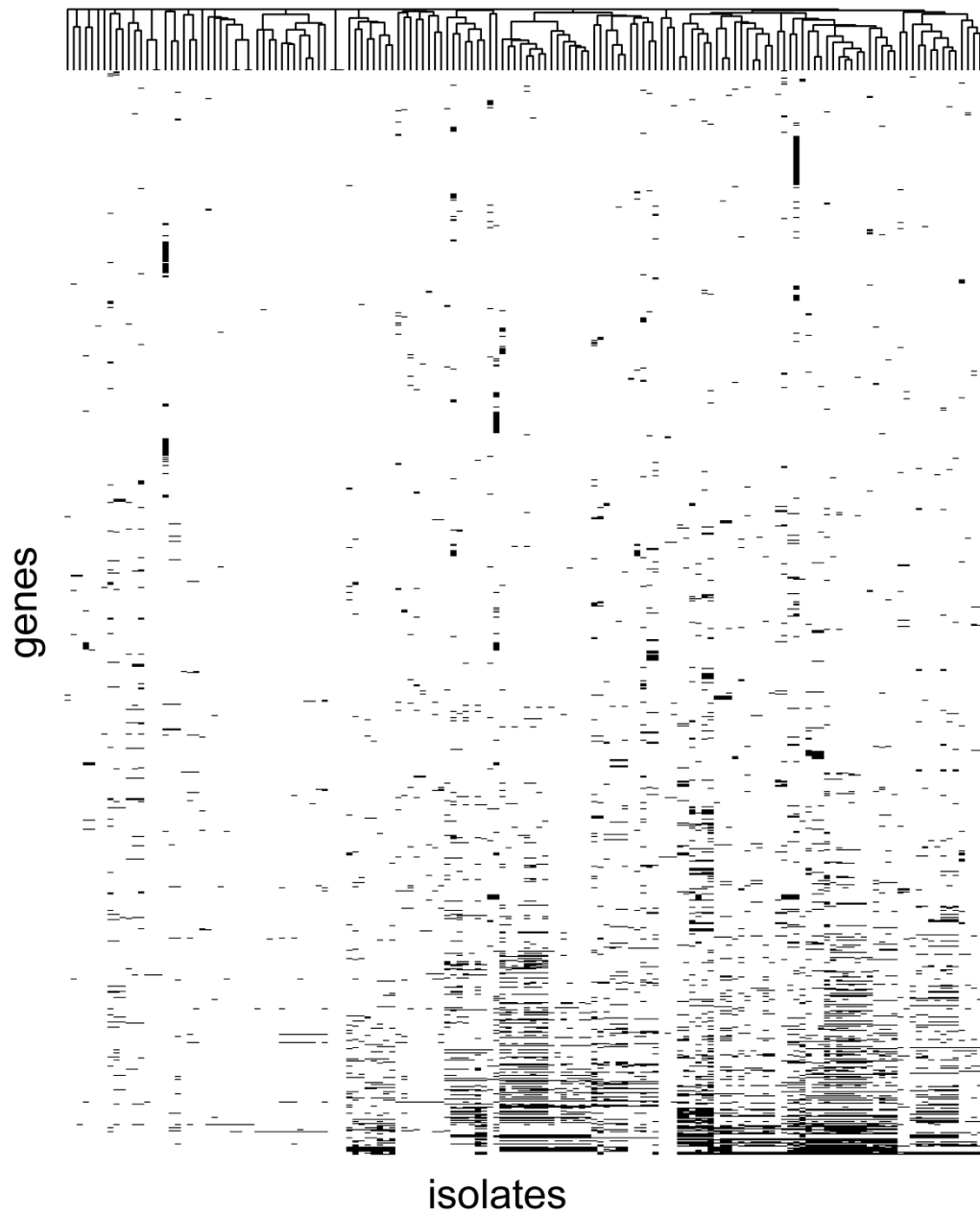
377 **Figure Legends**

378 **Figure 1.** The heat map depicts the combination of individual accessory genes that
379 are expressed by the 150 clinical isolates. The dendrogram is a visualisation of the
380 distance between the isolates based on the accessory gene composition.

381

382 **Figure 2.** The heat map depicts the probable source of the genes of the accessory
383 genomes of the 150 clinical isolates. The dendrogram is a visualisation of the
384 distance between the isolates based on the species to which their accessory genes
385 are annotated.

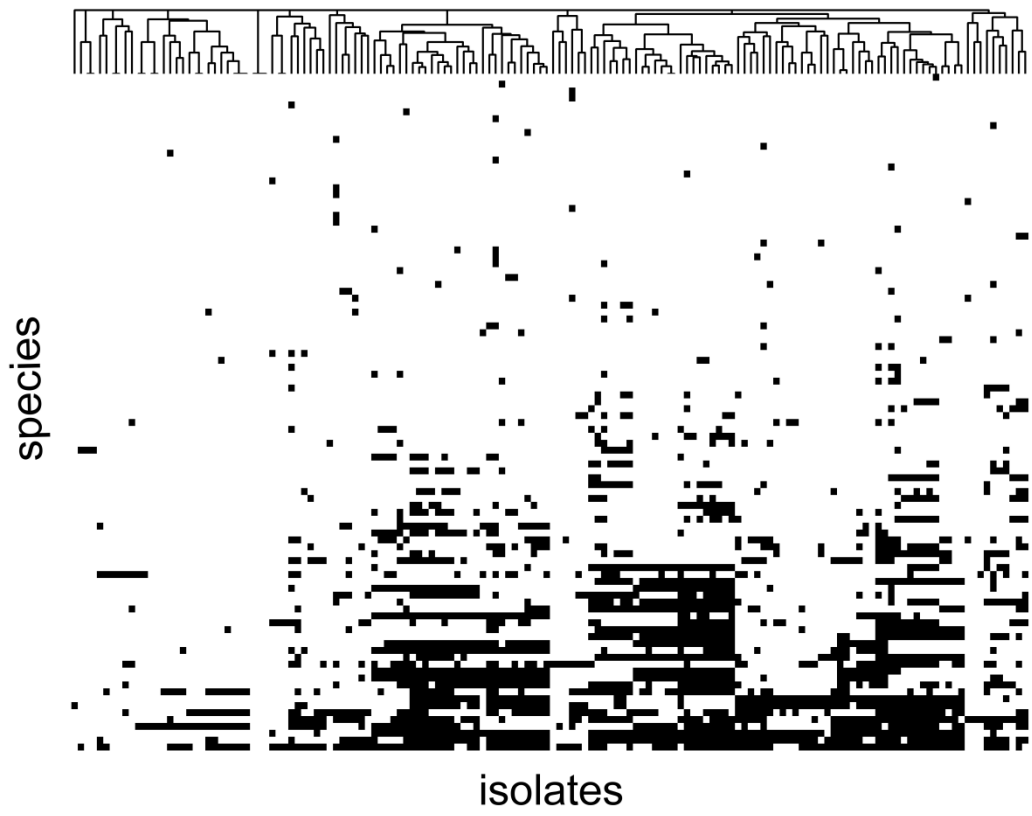
386



387

388 **Figure 1**

389



390

391 **Figure 2**