



Application of WGS data for O-specific antigen analysis and in silico serotyping of *Pseudomonas aeruginosa* isolates

Thrane, Sandra Wingaard; Taylor, Véronique L.; Lund, Ole; Lam, Joseph S.; Jelsbak, Lars

Published in:
Journal of Clinical Microbiology

Link to article, DOI:
[10.1128/JCM.00349-16](https://doi.org/10.1128/JCM.00349-16)

Publication date:
2016

Document Version
Peer reviewed version

[Link back to DTU Orbit](#)

Citation (APA):
Thrane, S. W., Taylor, V. L., Lund, O., Lam, J. S., & Jelsbak, L. (2016). Application of WGS data for O-specific antigen analysis and in silico serotyping of *Pseudomonas aeruginosa* isolates. *Journal of Clinical Microbiology*, 54(7), 1782-1788. <https://doi.org/10.1128/JCM.00349-16>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

1 Application of WGS data for O-specific antigen analysis and *in*
2 *silico* serotyping of *Pseudomonas aeruginosa* isolates

3

4 Sandra Wingaard Thrane¹, Véronique L. Taylor², Ole Lund³, Joseph S. Lam² and Lars
5 Jelsbak¹#

6

7 ¹Technical University of Denmark, Department of Systems Biology, Kgs. Lyngby,
8 Denmark

9 ²University of Guelph, Department of Molecular and Cellular Biology, Guelph,
10 Canada

11 ³Technical University of Denmark, Center for Biological Sequence Analysis,
12 Department of Systems Biology, Kgs. Lyngby, Denmark

13

14 Running title: *In silico* serotyping of *P. aeruginosa*

15

16 #Address correspondence to Lars Jelsbak, lj@bio.dtu.dk

17 **Abstract**

18

19 Accurate typing methods are required for efficient infection control. The emergence
20 of whole genome sequencing (WGS) technologies has enabled the development of
21 genomics-based methods applicable for routine typing and surveillance of bacterial
22 pathogens. In this study, we developed the *Pseudomonas aeruginosa* serotyper
23 (PAst) program, which enabled *in silico* serotyping of *P. aeruginosa* isolates using
24 WGS data. PAst has been made publically available as a web-service, and aptly
25 facilitate high-throughput serotyping analysis. The program overcomes critical issues
26 such as the loss of *in vitro* typeability often associated with *P. aeruginosa* isolates
27 from chronic infections, and quickly determines the serogroup of an isolate based on
28 the sequence of the O-specific antigen (OSA) gene cluster. Here, PAst analysis of
29 1649 genomes resulted in successful serogroup assignments in 99.27% of the cases.
30 This frequency is rarely achievable by conventional serotyping methods. The limited
31 number of non-typeable isolates found using PAst was the result of either complete
32 absence of OSA genes in the genomes or the artifact of genomic misassembly. With
33 PAst, *P. aeruginosa* serotype data can be obtained from WGS information alone.
34 PAst is a highly efficient alternative to conventional serotyping methods in relation
35 to outbreak surveillance of serotype O12 and other high-risk clones, while
36 maintaining backward compatibility to historical serotype data.

37

38 **Introduction**

39

40 *Pseudomonas aeruginosa* is a Gram-negative opportunistic pathogen and a major
41 cause of mortality and morbidity among hospitalized and compromised patients
42 including those with cystic fibrosis (CF). *P. aeruginosa* is well known for its ability to
43 cause chronic and extensively drug resistant infections (1). The outer membrane
44 lipopolysaccharide (LPS) layer is a major virulence factor of *P. aeruginosa* (2). LPS has
45 been linked to antibiotic resistance and immune evasion. Furthermore, LPS is one of
46 the receptors that determines susceptibility of the bacterium to bacteriophages and
47 pyocins (2–4). Our ability to control *P. aeruginosa* infections depends on the
48 availability of accurate typing methods. Previously, serotyping was a benchmark
49 typing method for *P. aeruginosa*. In the 1980's the International Antigenic Typing
50 Scheme (IATS) was established to classify the species *P. aeruginosa* into 20 serotypes
51 (O1-O20) (5–7). Today, serotyping is infrequently used in the clinic for typing
52 purposes, mainly because of the time consuming protocol, the need for a continuous
53 supply of serotype-specific antisera, and a high prevalence of polyagglutinating or
54 non-typeable isolates.

55

56 The loss of *P. aeruginosa* typeability has been known for decades, and has often
57 been linked to bacteria isolated from chronic infections, where typeability is lost
58 over time during the course of infection (8, 9). A study performed by Pirnay *et al* (10)
59 showed that 65% of all *P. aeruginosa* isolates examined were either non- or multi-
60 typeable and therefore assigning a particular serotype to these strains would be
61 difficult. The occurrence of these non- or multi-typeable isolates was higher when
62 evaluating isolates sampled exclusively from CF infections (10). Multi-typeability has
63 been associated with poor prognosis for CF patients, and is a trait of persistent or
64 chronic infection. This correlates with the observation that *P. aeruginosa* isolates
65 from chronic CF infections are initially resistant to human serum but evolve to
66 becoming serum sensitive over time. This is likely due to the loss of production of O-
67 antigen, which protects the bacterial cell from the human serum (8). The mechanism
68 underlying loss of typeability over time is not fully understood, but is most likely due

69 to modifications of LPS structures over extended periods of bacteria-host
70 interactions as a means to improve fitness in the host and to evade host immune
71 system, bacteriophages and antibiotic therapy.
72
73 The knowledge concerning the serotype of an isolate is important for monitoring
74 outbreaks and for understanding the structures of the LPS expressed on the surface
75 of these bacteria. O11 and O12 are more predominant than other serotypes in the
76 clinic, and intriguingly, these serotypes have been associated with multi-drug
77 resistance (MDR) (10–13). This implies that these particular LPS structures improve
78 fitness within the hosts and the hospital environments in ways that we currently do
79 not understand. Specifically for the O12 serotype, it has been shown that horizontal
80 gene transfer of LPS genes has resulted in MDR isolates and the switching of a
81 certain serotype to O12 (14). To continuously monitor LPS structure and evolution,
82 serotyping can help to improve our understanding of the isolates that successfully
83 infect patients. The continued collection of these data will also enable retrospective
84 population analysis, as serotype has been recorded for decades also prior to the
85 emergence of other DNA-based typing methods such as MLST and PCR.
86
87 *P. aeruginosa* LPS is comprised of three domains: lipid A, core oligosaccharide, and
88 O-antigen (2). Most *P. aeruginosa* isolates produce two forms of O-antigen
89 simultaneously: common polysaccharide antigen (CPA) and O-specific antigen (OSA).
90 While CPA is relatively conserved, OSA is variable and defines the serotype of an
91 isolate (2, 15). OSA is encoded in a gene cluster varying in size from just under 15 kb
92 to over 25 kb. The OSA gene cluster is flanked by the genes *ihfB/himD* and *wbpM*.
93 The 20 serotypes harbor 11 distinct OSA gene clusters, each with a high number of
94 unique genes (16). With the emergence of whole genome sequencing (WGS)
95 methods it is now possible to assign an isolate into one of 11 serogroups based on
96 the sequence and structure of the OSA gene cluster (11, 14, 17).
97
98 The present study presents a program that our group has developed for fast and
99 reliable *in silico* serotyping of *P. aeruginosa* isolates using WGS data – the
100 *Pseudomonas aeruginosa* serotyper (PAst). The program has been made publically

101 available as a web-service, and can enable high throughput serotyping analysis based
102 on analysis of the OSA gene cluster. Using PAST, issues with typeability of clinical
103 isolates can be overcome, and serotyping can be performed in a rapid and cost-
104 effective way in the clinic as whole genome sequencing of isolates become
105 accessible.
106

107 **Materials and Methods**

108

109 **PAst verification and isolates included in the study**

110 To evaluate the efficiency of the *in silico* serotyping using PAst, all available *P.*
111 *aeruginosa* genomes were acquired and analyzed. These *P. aeruginosa* genomes
112 were downloaded from NCBI and included 1120 genome assemblies (Supplementary
113 Table 1, extracted 18.08.2015). An exclusively CF-related *P. aeruginosa* dataset was
114 constructed, due mainly to the documented high level of non-typeability in
115 persistent infecting clones. The isolates described by Marvig *et al.* 2015 (475
116 genomes) (18) were used as the initial dataset. These were assembled using SPAdes
117 (19) prior to analysis. Additional CF isolates were recovered by searching for *P.*
118 *aeruginosa* genome assemblies related to CF in PATRIC (54 genomes) (20). It was
119 verified that frequently observed CF-specific strains such as DK2 and LES were part of
120 the dataset. The final dataset included 529 CF-related *P. aeruginosa* genome
121 assemblies. *In silico* serotyping of both datasets was performed using PAst in order
122 to evaluate typeability of the program. Non-typeable isolates (i.e., isolates in which
123 %coverage of reference OSA was < 95%) were manually examined for either
124 biological or technical explanations of the lack of typeability.

125

126 **PAst specifications**

127 The PAst program is developed using the programming language Perl for *in silico*
128 serotyping of *P. aeruginosa* isolates using WGS data. It is based on a BLASTn analysis
129 of the assembled input genome, against an OSA cluster database. OSA clusters with
130 > 95% coverage in the query genome represents a positive hit for a serogroup. Since
131 *P. aeruginosa* isolates have been described which either harbor multiple OSA
132 clusters or no clusters at all, the program accommodates multi-, mono- and non-
133 typeability based on analysis of the number of positive OSA hits and coverage (Figure
134 1). Compared to other studies (11, 14, 17) PAst optimizes *in silico* serotyping further
135 by distinguishing members of the O2 serogroup through identification of the
136 acquired phage-related *wzy_β* within serotypes O2 and O16 (21, 22). This enables
137 typing into 12 serogroups as opposed to the 11 described by Raymond *et al.* (16).

138 Together with a summary of the best hit(s) from the analysis and the BLAST report,
139 the user receives a multi fasta file containing the sequence(s) of the OSA cluster
140 from the analyzed isolate for use in future analysis.

141

142 **The *P. aeruginosa* OSA cluster database**

143 The database was constructed using the WGS data of the 20 *P. aeruginosa* IATS
144 serotype reference isolates (14). The genomes were assembled using SPAdes (19)
145 and the OSA clusters extracted via identification of the *ihfB/himD* gene flanking the
146 cluster upstream and the *wbpM* gene flanking the cluster downstream. The clusters
147 were aligned within their serotypes, described by Raymond *et al.* 2002 and their
148 shared structure confirmed (16). A representative cluster of each serotype was
149 selected for the database (Table 1). Also included in the database was the *wzy_β* gene
150 for distinguishing the O2 and O5 serotypes, as the two serogroups share OSA cluster
151 organization, but only the O2 and O16 serotype harbor the *wzy_β* gene present on a
152 prophage.

153

154 ***In silico* serotyping of *P. aeruginosa* isolates using PAst**

155 PAst has been implemented as a simple and user-friendly web-tool available on the
156 Center for Genomic Epidemiology (CGE) service platform
157 (<https://cge.cbs.dtu.dk/services/PAst-1.0/>). The tool accommodates raw reads, draft
158 assemblies (contigs or scaffolds) and complete genomes from all WGS platforms.
159 Raw read data are processed and assembled as previously described for other CGE
160 tools (23). Following analysis of the input data, the web-tool outputs the predicted
161 serogroup of the query genome, the %coverage of the reference OSA cluster, as well
162 as the OSA cluster sequence in multi fasta format, for the user to continue exploring
163 the OSA genes (Fig. 1). If multiple positive hits are found (multi-typeability), all the
164 identified OSA clusters are written for the user (Fig. 1). In the case of a non-typeable
165 query genome (where no OSA cluster has >95% coverage) the best hit identified is
166 written for the user together with the sequence of this hit (Fig. 1).
167 For batch analysis of larger datasets (only applicable for assembled genomes) the
168 PAst Perl program has been made available on Github:
169 <https://github.com/Sandramses/PAst>

170 **Results**

171 The PAst web server tool identifies and analyzes the nucleotide sequence of the O-
172 specific antigen (OSA) gene cluster within the provided genomes and place them into
173 one of twelve serogroups defined in Table 1. These serogroups are defined by
174 sequence similarities between the 20 IATS serotypes (16) as well as
175 absence/presence of the discriminatory *wzy*_β gene (21, 22) and are as such different
176 from previously groupings of serotypes on the basis of *in vitro* serotyping data (11,
177 14, 17). All serogroups contained three or less of the 20 IATS serotypes (Table 1).

178

179 **More than 97% of the *P. aeruginosa* dataset is typeable using PAst**

180 To evaluate the typeability efficiency of PAst all *P. aeruginosa* genome assemblies
181 available in NCBI (1120 genomes on date of extraction) were analyzed. A total of
182 97.68% (1094) of the 1120 genomes were typed unambiguously to a single
183 serogroup by PAst (Fig. 2). This means that each genome assembly had a single
184 BLAST hit of >95% OSA coverage to one sequence in our reference OSA database
185 (Fig. 2). No isolates were found to be multi-typeable and 2.32% (26 genomes) of the
186 1120 genomes were found to be non-typeable (Fig. 2). In these cases, no significant
187 BLAST hit of >95% OSA coverage to one of the sequence in the reference OSA
188 database was identified. PAst correctly determined the serogroup of the 20 IATS
189 strains as well as PAO1 (serotype O5), PA14 (serotype O10), and PAK (serotype O6).

190

191 The analysis showed that all serogroups were represented in the 1120 genomes (Fig.
192 2). Four of the 12 serogroups represented 70% of the genomes analyzed; these were
193 O3, O6, O11 and O12 (Fig. 2). The smallest serogroup was O13, which contained
194 only four genomes. We note that the same clone type could be present multiple
195 times in the dataset, and that a substantial sampling bias would therefore be
196 expected. The distribution of serotypes in our analysis thus describes what has been
197 chosen for sequencing and does not necessarily match the distribution of serotypes
198 in the actual *P. aeruginosa* population. This does not affect the high confidence of
199 PAst, as it shows that un-ambiguous typing of multiple isolates from the same
200 lineage is possible.

201

202 **PAst overcomes non-typeability issues from *in vitro* typing of CF lineages**

203 *P. aeruginosa* isolates from CF infections are often non-typeable with conventional
204 serotyping assays. To explore if our genomics-based method could enable
205 acquisition of serotype information in such isolates, we analyzed 529 genome
206 assemblies of *P. aeruginosa* isolates sampled from CF infections. This dataset
207 contained multiple examples of isolates of the same lineage that had been sampled
208 during the course of infection. This enabled us to investigate whether *in silico*
209 typeability might be lost over time as has frequently been observed for *in vitro*
210 serotyping of isolates from chronic CF infections. Interestingly, 99.81% of the
211 genomes in the CF-specific dataset could be typed to single serogroups. More
212 importantly, no multi-typeable isolates were observed and only one isolate was
213 deemed non-typeable (Fig. 3). All serogroups were represented in the dataset,
214 except for O12. The absence of O12 serotypes among CF isolates has previously been
215 reported (10). Serotypes O1, O6 and O7/O8 represented ~65% of the CF-specific
216 dataset and the smallest representation of serotypes was the O9 serogroups with
217 only two isolates from these samples (Fig. 3).

218

219 Well-known transmissible CF-specific clone types such as *P. aeruginosa* DK1 (24),
220 DK2 (25), and LES (26) are represented in the dataset due to multiple isolates being
221 sampled from various patients over several decades. Using our PAst tool, the typing
222 problems documented from *in vitro* typing of such lineages were not observed, and
223 the DK1, DK2 and LES isolates were consistently *in silico* serotyped with PAst. DK1
224 and DK2 were found to belong to the O3 serogroup, while the LES lineage belonged
225 to the O6 serogroup.

226

227 **Complete loss of O-specific antigen defining genes is a rare event**

228 Out of two WGS-based datasets (n = 1649) that were *in silico* typed with PAst, our
229 results yielded a total of 27 non-typeable isolates. The lack of typeability in these 27
230 genome assemblies was further investigated to resolve whether non-typeability in
231 these cases was due to technical or biological reasons. We found that the %OSA
232 coverage of the non-typeable isolates ranged from a minimum of 1.91% to a

233 maximum 93.96% OSA coverage (Supplementary Table 2). Of the 27 isolates
234 classified as non-typeable, thirteen were found to have OSA coverage of 0-20%,
235 whereas seven isolates had OSA coverage of 80-95% (Fig. 4). The best hit (serogroup)
236 for each of the non-typeable isolates was then examined to evaluate if certain
237 serogroups were more prone to be problematic in the PAST analysis and why. The 27
238 isolates were found to distribute across 6 serogroups (O1, O2, O6, O7, O11 and O13),
239 while 15/27 isolates showed a best hit to be typed as the O11 serogroup (Fig. 4).

240

241 The group of non-typeable isolates with a best hit to the O11 serogroup were
242 analyzed separately to identify the reason for the lack of typeability. Of the 15 O11
243 serogroup isolates, nine had an OSA coverage of 14.94-15.84% (Supplementary
244 Table 2); these corresponded to the presence of only the two flanking genes
245 *himD/ihfB* and *wbpM*. This observation shows that a best hit of a non-typeable
246 isolate to the O11 OSA cluster with a coverage of ~15% is the result of a complete
247 absence of an OSA cluster but the presence of the flanking genes. Two other isolates
248 had an OSA coverage of <2%, and corresponded to the absence of the entire OSA
249 cluster as well as the flanking genes (Supplementary table 2). In summary, a total of
250 11 of the 27 non-typeable isolates (or 11 of 1649 isolates analyzed in total) were
251 non-typeable due to a lack of the OSA cluster sequences.

252

253 **Genome mis-assembly accounts for false non-typeability**

254 Since the seven non-typeable isolates with the highest OSA coverage (80-95%) in
255 Figure 4 were all candidates for harboring complete and functioning OSA clusters,
256 we analyzed the cause of non-typeability in this group of isolates. For each of the
257 isolates, we examined whether there were mis-assembly or assembly gaps within
258 the OSA gene cluster; we also looked for the occurrence of insertion sequence (IS)
259 elements, which often cause gaps in *de novo* assembly. Indeed, five of the seven
260 isolates contained assembly gaps within their OSA cluster, which account for the
261 observed lowered OSA coverage (Table 2). The remaining two isolates had no gaps
262 within their OSA sequence (Table 2). However, both of these isolates had a best type
263 hit to the O11 serogroup, which is known to contain OSA sequences of both the O11
264 and the O17 serotypes (16) (Table 1). Interestingly, the OSA cluster in these two

265 serotypes differ only by the presence of two IS elements and a deletion in the O17
266 serotype OSA sequence (16). Alignment of the OSA sequence from the two non-
267 typable isolates to the O11 and O17 reference OSA sequences, respectively,
268 contained an O17 OSA gene cluster, which had been misassembled into
269 concatenated O11 serotype OSA clusters because of the O17 IS elements.

270 **Discussion**

271

272 The serotyping technique has been one of the standard tools for epidemiological
273 studies and infection controls for many decades. The available historical records of *P.*
274 *aeruginosa* serotypes offer a vast amount of information about *P. aeruginosa*
275 epidemiology and population structures (27–30). Although problems with non-
276 typeable isolates have been described since the implementation of the method, the
277 serotype information is still applicable today for outbreak tracking, strain typing, and
278 studies of LPS structure and evolution. The present study presents a newly
279 developed Web Server tool called PAsT, which is user friendly, reliable, and high-
280 throughput for *in silico* serotyping of *P. aeruginosa* isolates.

281

282 In contrast to conventional serology-based *in vitro* serotyping, PAsT *in silico*
283 serotyping has a very low occurrence of non-typeability. Of the 1649 analyzed
284 genomes, only 27 non-typeable isolates were detected across two separate *P.*
285 *aeruginosa* datasets. One dataset represents all available whole genome assemblies
286 of *P. aeruginosa*, while the other specifically represents genomes from CF infections,
287 which are known to contain high occurrences of non-typeability due to adaptability
288 of the bacteria into a biofilm life-style associated with chronicity of the infection (Fig.
289 1 and 2). Importantly, since the frequency of non-typeability of *in vitro* serotyped *P.*
290 *aeruginosa* isolates may amount to over 65% (10), analysis with PAsT is clearly
291 advantageous and superior compared to conventional *in vitro* serotyping.

292 Importantly, the superiority of the PAsT tool as a reliable and fast typing method is
293 consistent with other published tools for *in silico* serotyping (31–35). Similar to both
294 the SerotypeFinder (*in silico* serotyping of *E. coli* (31)), LisSero (*in silico* serotyping of
295 *Listeria monocytogenes* (34, 35)) and SeqSero (*in silico* serotyping of *Salmonella* (32))
296 PAsT resolves the OSA cluster information to the most accurate typing possible as a
297 serogroup representing 1-3 serotypes.

298

299 Interestingly, we observed a high level of conservation of the OSA gene cluster
300 within the *P. aeruginosa* genome. In contrast to certain well-documented difficulties

301 in serology-based *in vitro* serotyping, PAST identified complete OSA clusters (with
302 >95% sequence being present) in 99.27% of the analyzed genomes. As such only 12
303 of the 1649 isolates examined were found to be devoid of the OSA cluster and an
304 additional 8 isolates were found to contain only a partial OSA cluster in their
305 genomes (<80% OSA sequence compared to the reference). These findings indicate
306 that the loss of typeability of *P. aeruginosa* isolates during the course of infection is
307 either due to mutations (rather than larger deletions) or is linked to other parts of
308 the LPS biosynthesis, such as regulatory genes or transport of the structure to the
309 cell surface. A study by Bélanger *et al.* reported that mutation in any of the four *wbp*
310 genes (*wbpO*, *wbpP*, *wbpV* and *wbpM*) in the OSA gene cluster could disrupt the *P.*
311 *aeruginosa* O6 OSA biosynthesis (36). Furthermore, key genes involved in the OSA
312 assembly and translocation through the Wzx/Wzy-dependent pathway not localized
313 within the OSA cluster, for instance, *waal*, are essential for O-antigen expression
314 (37, 38). It is possible that more OSA-related genes might be present in the *P.*
315 *aeruginosa* genomes, which have not been discovered yet. Overall, our study
316 demonstrates that a complete lack of an OSA gene cluster is a rarely observed
317 phenomenon in *P. aeruginosa*.

318

319 PAST will enable further investigations of the diversity, evolution and variability of
320 the OSA clusters. For example, the sequence of the cluster is part of the output
321 material from the *in silico* serotyping which can then be readily analyzed for
322 sequence variations to provide new knowledge on the mechanisms behind loss of
323 typeability *in vitro* and *in silico*. Furthermore, PAST will enable systematic analysis of
324 serotype switching by horizontal gene transfer and genetic recombination of the
325 OSA gene cluster among different clone types. This recently described phenomenon
326 has contributed to the evolution of the multi-drug resistant *P. aeruginosa* serotype
327 O12 population that has successfully disseminated across hospitals worldwide (14).
328 It is currently unknown if there are additional cases of such serotype switching by
329 recombination.

330

331 The new PAST Web Server tool makes *in silico* serotyping of *P. aeruginosa* using WGS
332 data a fast and reliable method. The use of PAST can play an important role in future

333 surveillance of LPS evolution and possible outbreak detection. With the emergence
334 of rapidly disseminating, high-risk clones of *P. aeruginosa*, such as the O12 ST111
335 clone, new and reliable typing techniques for improved monitoring and tracking of
336 such outbreaks are becoming increasingly important (13). With the lowered cost of
337 sequencing and the increased focus on WGS of pathogens in clinics and hospital
338 settings, genomics-based tools can assist in designing future treatments and
339 containment of outbreaks.

340 **Acknowledgements**

341 Funding for this study was provided by operating grants from the Villum Foundation
342 to L.J. (VKR023113) and from the Canadian Institutes of Health Research (CIHR) to
343 J.S.L. (MOP-14687). We thank the Center for Genomic Epidemiology (CGE) at the
344 Center for Biological Sequence analysis (CBS) at DTU, especially Johanne Ahrenfeldt
345 and Rosa Allesøre, for expert assistance in setting up the PAsT web-service and
346 hosting it on their web-servers. Additional support was provided by the 'A.N.
347 Neergaard og Hustrus' Foundation to L.J., and a travel grant from Knud Højgaards
348 Foundation to S.W.T. V.L.T. was a recipient of a Cystic Fibrosis Canada Doctoral
349 Studentship, Queen Elizabeth II Graduate Scholarships in Science and Technology
350 (QEII-GSST) and J.S.L. holds a Canada Research Chair in Cystic Fibrosis and Microbial
351 Glycobiology.

352 **References**

- 353 1. **Folkesson A, Jelsbak L, Yang L, Johansen HK, Ciofu O, Høiby N, Molin S.** 2012.
354 Adaptation of *Pseudomonas aeruginosa* to the cystic fibrosis airway: an
355 evolutionary perspective. *Nat Rev Microbiol* **10**:841–51.
- 356 2. **Lam JS, Taylor VL, Islam ST, Hao Y, Kocíncová D.** 2011. Genetic and functional
357 diversity of *Pseudomonas aeruginosa* lipopolysaccharide. *Front Microbiol*
358 **2**:118.
- 359 3. **Köhler T, Donner V, van Delden C.** 2010. Lipopolysaccharide as shield and
360 receptor for R-pyocin-mediated killing in *Pseudomonas aeruginosa*. *J Bacteriol*
361 **192**:1921–1928.
- 362 4. **Nakayama K, Takashima K, Ishihara H, Shinomiya T, Kageyama M, Kanaya S,**
363 **Ohnishi M, Murata T, Mori H, Hayashi T.** 2000. The R-type pyocin of
364 *Pseudomonas aeruginosa* is related to P2 phage, and the F-type is related to
365 lambda phage. *Mol Microbiol* **38**:213–231.
- 366 5. **Liu P V, Wang S.** 1990. Three new major somatic antigens of *Pseudomonas*
367 *aeruginosa*. *J Clin Microbiol* **28**:922–925.
- 368 6. **Stanislavsky E, Lam J.** 1997. *Pseudomonas aeruginosa* antigens as potential
369 vaccines. *FEMS Microbiol Rev* **21**:243–277.
- 370 7. **Liu P V., Matsumoto H, Kusama H, Bergan T.** 1983. Survey of heat-stable,
371 major somatic antigens of *Pseudomonas aeruginosa*. *Int J Syst Bacteriol*
372 **33**:256–264.
- 373 8. **Penketh A, Pitt T, Roberts D, Hodson M, Batten J.** 1983. The relationship of
374 phenotype changes in *Pseudomonas aeruginosa* to the clinical condition of
375 patients with cystic fibrosis. *Am Rev Respir Dis* **127**:605–608.
- 376 9. **Ojeniyi B.** 1994. Polyagglutinable *Pseudomonas aeruginosa* from cystic fibrosis
377 patients - a survey. *APMS* **102**.
- 378 10. **Pirnay J-P, Bilocq F, Pot B, Cornelis P, Zizi M, Van Eldere J, Deschaght P,**
379 **Vanechoutte M, Jennes S, Pitt T, De Vos D.** 2009. *Pseudomonas aeruginosa*
380 population structure revisited. *PLoS One* **4**:e7740.
- 381 11. **Witney AA, Gould KA, Pope CF, Bolt F, Stoker NG, Cubbon MD, Bradley CR,**
382 **Fraise A, Breathnach AS, Butcher PD, Planche TD, Hinds J.** 2014. Genome
383 sequencing and characterization of an extensively drug-resistant sequence
384 type 111 serotype O12 hospital outbreak strain of *Pseudomonas aeruginosa*.
385 *Clin Microbiol Infect* **20**:O609–O618.
- 386 12. **Cholley P, Thouverez M, Hocquet D, Van Der Mee-Marquet N, Talon D,**
387 **Bertrand X.** 2011. Most multidrug-resistant *Pseudomonas aeruginosa* isolates
388 from hospitals in eastern France belong to a few clonal types. *J Clin Microbiol*
389 **49**:2578–2583.
- 390 13. **Oliver A, Mulet X, López-causapé C, Juan C.** 2015. The increasing threat of
391 *Pseudomonas aeruginosa* high-risk clones. *Drug Resist Updat* **22**:41–59.

- 392 14. **Thrane SW, Taylor VL, Freschi L, Kukavica-ibrulj I, Boyle B, Laroche J, Pirnay J-**
393 **P, Lévesque RC, Lam JS, Jelsbak L.** 2015. The widespread multidrug-resistant
394 serotype O12 *Pseudomonas aeruginosa* clone emerged through concomitant
395 horizontal transfer of serotype antigen and antibiotic resistance gene clusters.
396 *MBio* **6**:1–10.
- 397 15. **King JD, Kocíncová D, Westman EL, Lam JS.** 2009. Review: Lipopolysaccharide
398 biosynthesis in *Pseudomonas aeruginosa*. *Innate Immun* **15**:261–312.
- 399 16. **Raymond CK, Sims EH, Kas A, Spencer DH, Kutyavin T V, Ivey RG, Zhou Y,**
400 **Kaul R, Clendenning JB, Olson M V.** 2002. Genetic variation at the O-antigen
401 biosynthetic locus in *Pseudomonas aeruginosa*. *J Bacteriol* **184**:3614–3622.
- 402 17. **Kos VN, Déraspe M, McLaughlin RE, Whiteaker JD, Roy PH, Alm R a., Corbeil**
403 **J, Gardner H.** 2015. The resistome of *Pseudomonas aeruginosa* in relationship
404 to phenotypic susceptibility. *Antimicrob Agents Chemother* **59**:427–436.
- 405 18. **Marvig RL, Sommer LM, Molin S, Johansen HK.** 2015. Convergent evolution
406 and adaptation of *Pseudomonas aeruginosa* within patients with cystic
407 fibrosis. *Nat Genet* **47**:57–64.
- 408 19. **Bankevich A, Nurk S, Antipov D, Gurevich A a, Dvorkin M, Kulikov AS, Lesin**
409 **VM, Nikolenko SI, Pham S, Prjibelski AD, Pyshkin A V, Sirotkin A V, Vyahhi N,**
410 **Tesler G, Alekseyev M a, Pevzner P a.** 2012. SPAdes: a new genome assembly
411 algorithm and its applications to single-cell sequencing. *J Comput Biol* **19**:455–
412 77.
- 413 20. **Wattam AR, Abraham D, Dalay O, Disz TL, Driscoll T, Gabbard JL, Gillespie JJ,**
414 **Gough R, Hix D, Kenyon R, Machi D, Mao C, Nordberg EK, Olson R, Overbeek**
415 **R, Pusch GD, Shukla M, Schulman J, Stevens RL, Sullivan DE, Vonstein V,**
416 **Warren A, Will R, Wilson MJC, Yoo HS, Zhang C, Zhang Y, Sobral BW.** 2014.
417 PATRIC , the bacterial bioinformatics database and analysis resource. *Nucleic*
418 *Acids Res* **42**:581–591.
- 419 21. **Kaluzny K, Abeyrathne PD, Lam JS.** 2007. Coexistence of two distinct versions
420 of O-antigen polymerase, Wzy-Alpha and Wzy-Beta, in *Pseudomonas*
421 *aeruginosa* serogroup O2 and their contributions to cell surface diversity. *J*
422 *Bacteriol* **189**:4141–4152.
- 423 22. **Newton GJ, Daniels C, Burrows LL, Kropinski AM, Clarke AJ, Lam JS.** 2001.
424 Three-component-mediated serotype conversion in *Pseudomonas aeruginosa*
425 by bacteriophage D3. *Mol Microbiol* **39**:1237–1247.
- 426 23. **Larsen M V., Cosentino S, Rasmussen S, Friis C, Hasman H, Marvig RL, Jelsbak**
427 **L, Sicheritz-Pontén T, Ussery DW, Aarestrup FM, Lund O.** 2012. Multilocus
428 sequence typing of total-genome-sequenced bacteria. *J Clin Microbiol*
429 **50**:1355–1361.
- 430 24. **Markussen T, Marvig L, Gómez-lozano M, Aanæs K, Burleigh AE, Høiby N.**
431 2014. Environmental heterogeneity drives within-host diversification and
432 evolution of *Pseudomonas aeruginosa*. *MBio* **5**:1–10.
- 433 25. **Marvig RL, Johansen HK, Molin S, Jelsbak L.** 2013. Genome analysis of a

- 434 transmissible lineage of *Pseudomonas aeruginosa* reveals pathoadaptive
435 mutations and distinct evolutionary paths of hypermutators. *PLoS Genet*
436 **9**:e1003741.
- 437 26. **Jeukens J, Boyle B, Kukavica-Ibrulj I, Ouellet MM, Aaron SD, Charette SJ,**
438 **Fothergill JL, Tucker NP, Winstanley C, Levesque RC.** 2014. Comparative
439 genomics of isolates of a *Pseudomonas aeruginosa* epidemic strain associated
440 with chronic lung infections of cystic fibrosis patients. *PLoS One* **9**:e87611.
- 441 27. **Lam JS, MacDonald LA, Kropinski AM, Speert DP.** 1988. Characterization of
442 nontypable strains of *Pseudomonas aeruginosa* from cystic fibrosis patients by
443 means of monoclonal antibodies and SDS-polyacrylamide gel electrophoresis.
444 *Serodiag Immunother Infect Dis* **2**:365–374.
- 445 28. **Ojeniyi B, Wolz C, Doring G, Lam JS, Rosdahl VT, Høiby N.** 1990. Typing of
446 polyagglutinable *Pseudomonas aeruginosa* isolates from cystic fibrosis
447 patients. *Acta Pathol Microbiol Immunol Scand* **98**:423–431.
- 448 29. **Ojeniyi B, Lam JS, Høiby N, Rosdahl VT.** 1989. A comparison of the efficiency
449 in serotyping of *Pseudomonas aeruginosa* from cystic fibrosis patients using
450 monoclonal and polyclonal antibodies. *Acta Pathol Microbiol Immunol Scand*
451 **97**:631–636.
- 452 30. **Speert DP, Campbell M, Puterman ML, Govan J, Doherty C, Høiby N, Ojeniyi**
453 **B, Lam JS, Ogle JW, Johnson Z, Paranchych W, Sastry PA, Pitt TL, Lawrence L.**
454 1994. A multicenter comparison of methods for typing strains of
455 *Pseudomonas aeruginosa* predominantly from patients with cystic fibrosis. *J*
456 *Infect Dis* **169**:134–142.
- 457 31. **Joensen KG, Tetzschner AMM, Iguchi A, Aarestrup FM, Scheutz F.** 2015.
458 Rapid and easy *in silico* serotyping of *Escherichia coli* using whole genome
459 sequencing (WGS) data. *J Clin Microbiol JCM.00008–15*.
- 460 32. **Zhang S, Yin Y, Jones MB, Zhang Z, Deatherage Kaiser BL, Dinsmore B a.,**
461 **Fitzgerald C, Fields PI, Deng X.** 2015. *Salmonella* serotype determination
462 utilizing high-throughput genome sequencing data. *J Clin Microbiol*
463 **53**:JCM.00323–15.
- 464 33. **Yoshida CE, Kruczkiewicz P, Laing CR, Lingohr EJ, Victor P.** 2016. The
465 *Salmonella* In Silico Typing Resource (SISTR): An Open Web-Accessible Tool
466 for Rapidly Typing and Subtyping Draft *Salmonella* Genome Assemblies. *PLoS*
467 *One* **11**:e0147101.
- 468 34. **Kwong JC, Mercoulia K, Tomita T, Easton M, Li HY, Bulach DM, Stinear TP,**
469 **Seemann T, Howden BP.** 2015. Prospective whole genome sequencing
470 enhances national surveillance of *Listeria monocytogenes*. *J Clin Microbiol*
471 **54**:333–342.
- 472 35. **Doumith M, Buchrieser C, Glaser P, Jacquet C, Martin P.** 2004. Differentiation
473 of the Major *Listeria monocytogenes* serovars by multiplex PCR. *J Clin*
474 *Microbiol* **42**:3819–3822.
- 475 36. **Bélanger M, Burrows LL, Lam JS.** 1999. Functional analysis of genes

- 476 responsible for the synthesis of the B-band O-antigen of *Pseudomonas*
477 *aeruginosa* serotype O6 lipopolysaccharide. *Microbiology* **145**:3505–3521.
- 478 37. **Berry MC, Mcghee GC, Zhao Y, Sundin GW.** 2008. Effect of a *waaL* mutation
479 on lipopolysaccharide composition, oxidative stress survival, and virulence in
480 *Erwinia amylovora*. *FEMS Microbiol Lett* **291**:80–87.
- 481 38. **Abeyrathne PD, Daniels C, Poon KKH, Matewish MJ, Lam JS.** 2005. Functional
482 characterization of WaaL, a ligase associated with linking O-antigen
483 polysaccharide to the core of *Pseudomonas aeruginosa* lipopolysaccharide. *J*
484 *Bacteriol* **187**:3002–3012.
- 485

486 **Figure legends**

487 **FIG 1** Workflow illustrating the *in silico* serotyping of the *Pseudomonas aeruginosa*
488 serotyper (PAst).

489

490 **FIG 2** The distribution of the different serogroups (in %) identified via *in silico*
491 serotyping of the *P. aeruginosa* dataset using PAst. The analysis is based on all
492 available *P. aeruginosa* genomes assemblies (n = 1120).

493

494 **FIG 3** The distribution of the different serogroups (in %) identified via *in silico*
495 serotyping of CF specific *P. aeruginosa* isolates (n = 529) using PAst.

496

497 **FIG 4** Best-hit serotype distribution of the 27 non-typeable isolates as a function of
498 the OSA coverage.

499

500 **Tables**

TABLE 1 Serogroup definition in the PAST OSA database.

Serogroup	Reference OSA cluster	Ref. gene	Size (bp)	Serotypes within serogroup
O1	O1		18.368	O1
O2	O2	<i>wzy₆</i>	23.303	O2, O16
O3	O3		20.210	O3, O15
O5	O2		23.303	O5, O18, O20
O4	O4		15.279	O4
O6	O6		15.649	O6
O7	O7		19.617	O7, O8
O9	O9		17.263	O9
O10	O10		17.635	O10, O19
O11	O11		13.868	O11, O17
O12	O12		25.864	O12
O13	O13		14.316	O13, O14

501

TABLE 2 Non-typeable *P. aeruginosa* isolates with %OSA coverage of 80-95% with specification of assemblies.

Strain	Size (Mb)	Scaffolds	%GC	Best hit	%OSA	<i>wbpM</i>	<i>himD</i>	Gap
<i>P. aeruginosa</i> E2	635.733	196	66.4	O7	83.31	+	+	+
<i>P. aeruginosa</i> IGB83	648.065	249	66.4	O2	84.46	+	+	+
<i>P. aeruginosa</i> VRFPA04	681.803	1	66.5	O11	86.96	+	+	-
<i>P. aeruginosa</i>	627.851	176	66.1	O6	90.54	+	+	+
<i>P. aeruginosa</i> 148	664.374	128	66.1	O11	90.93	+	+	-
<i>P. aeruginosa</i> ID4365	677.663	172	66.1	O7	91.74	+	+	+
<i>P. aeruginosa</i> C2773C	671.772	200	65.9	O6	93.96	+	+	+

502







