1 **Shortomics**

2 **Comparative genomics of Czech vaccine strains of *Bordetella pertussis***

3 **Ana Dienstbier[1], Derek Pouchnik[2], Mark Wildung[2], Fabian Amman[3], Ivo L. Hofacker[3,4], Julian**

4 **Parkhill[5], Jana Holubova[6], Peter Sebo[6] and Branislav Vecerek[1*]**

5 [1]Institute of Microbiology v.v.i., Laboratory of post-transcriptional control of gene expression, 14220

6 Prague, Czech Republic, [2]Laboratory for Biotechnology and Bioanalysis, Center for Reproductive

7 Biology, Washington State University, Pullman, Washington 99164-7520, [3]University of Vienna,

8 Institute for Theoretical Chemistry, Währinger Straße 17, A-1090 Vienna, Austria, [4]University of

9 Vienna, Research group BCB, Faculty of Computer Science, Währinger Straße 24, 1090 Vienna,

10 Austria, [5]Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge,

11 UK, [6]Institute of Microbiology v.v.i, Laboratory of molecular biology of bacterial pathogens, 14220

12 Prague, Czech Republic.

13 *Corresponding author, Laboratory of post-transcriptional control of gene expression, Institute of

14 Microbiology of the CAS v.v.i., Tel: +420241062507; E-mail: vecerek@biomed.cas.cz

15 One sentence summary: This is the first report on the genomics of Czech pertussis vaccine strains

16 showing their uniqueness in terms of SNP-based phylogeny and genome organization

17 Keywords: Bordetella, pertussis, genomics, region of difference, genome rearrangement, vaccine

18 pressure

19

2

20

21 **ABSTRACT**

22 *Bordetella pertussis* is a strictly human pathogen causing the respiratory infectious disease called

23 whooping cough or pertussis. *B. pertussis* adaptation to acellular pertussis vaccine pressure ~~was~~ has

24 been ~~recently~~ repeatedly highlighted, but recent data indicate that adaptation of circulating strains

25 started already in the era of the whole cell pertussis vaccine (wP) use. We sequenced the genomes of

26 five *B. pertussis* wP vaccine strains isolated in the former Czechoslovakia in the pre-wP (1954 - 1957)

27 and early wP ~~use era~~ (1958 - 1965) eras, when only limited population travel into and out of the

28 country was possible. Four isolates exhibit a similar genome organization and form a distinct

29 phylogenetic cluster with a geographic signature. The fifth strain is rather distinct, both in genome

30 organization and SNP-based phylogeny. Surprisingly, despite isolation of this strain before 1966, its

31 closest sequenced relative appears to be a recent isolate from the US. On the genome content level,

32 the five vaccine strains contained both new and already described regions of difference. One of the

33 new regions contains duplicated genes potentially associated with transport across the membrane.

34 The prevalence of this region in recent isolates indicates that its spread might be associated with

35 selective advantage leading to increased strain fitness.

36

3

37

**INTRODUCTION**

39 *Bordetella pertussis*, the etiological agent of whooping cough (pertussis), is a strictly human Gram-
40 negative bacterium infecting the respiratory tract (Cherry, 2010). Despite massive world-wide
41 vaccination programs, pertussis remains the least-controlled vaccine-preventable infectious disease
42 and it is a major cause of infant morbidity and mortality globally (WHO, 2006). As in many other
43 countries, prior to introduction of the whole cell-based pertussis vaccine (wP), pertussis was the
44 major cause of infant mortality in the former Czechoslovakia. The morbidity due to pertussis steeply
45 declined after the compulsory vaccination was introduced in 1958 (Vysoka-Burianova *et al.*, 1976). As
46 in other countries with high vaccination coverage, the incidence of pertussis in the Czechoslovakia
47 and later on in the Czech Republic started to rise progressively since the 1990s (Raguckas *et al.*, 2007,
48 Fabianova *et al.*, 2010, Sealey *et al.*, 2016) and this trend accelerated strongly upon switch to
49 acellular pertussis subunit vaccine use in 2007 (Chlibek *et al.*, 2017).

50 While several factors are contributing to pertussis resurgence in the most developed countries, the
51 two prominent ones are the incomplete and short-lived immunity induced by current aP vaccines
52 and the genetic changes in circulating *B. pertussis* strains that lead to escape from immunity by
53 antigenic variation (Mooi *et al.*, 2014, Burdin *et al.*, 2017). However, adaptive mutations occurred
54 already in the wP era, suggesting that the major driving force of *B. pertussis* adaptation is vaccination
55 as such (Bart *et al.*, 2014). *B. pertussis* is generally considered to be a genetically monomorphic
56 pathogen (King *et al.*, 2010, Mooi, 2010) with rather limited extent of sequence variation within the
57 global population (Bart *et al.*, 2014). Nevertheless, *B. pertussis* possesses an efficient mechanism of
58 genome structure diversification due to the presence of almost 250 copies of the insertion sequence
59 IS481 in its genome. These mobile elements allow for intragenomic recombination and excision
60 and/or insertion of the flanked genome regions, leading to genome decay (Parkhill *et al.*, 2003,
61 Cummings *et al.*, 2004), genome rearrangements (Bowden *et al.*, 2016, Weigand *et al.*, 2017) and
62 gene expression alterations (Brinig *et al.*, 2006). Recently, we showed that insertion elements
63 significantly affect the global gene expression profile in *B. pertussis* (Amman *et al.*, 2018). Apparently,
64 *B. pertussis* adaptation goes beyond the changes in the genes coding for antigens that are present in
65 the acellular vaccine and involves also other virulence-associated genes and the genes coding for
66 surface-exposed proteins (Bart *et al.*, 2014). To understand how these vaccination-induced
67 adaptation changes contributed to the current re-emergence of whooping cough, it is important to
68 analyze genomes of strains from the pre-vaccine era. At present, there are over 1000 *B. pertussis*
69 genome sequences deposited in Genbank and JGI GOLD databases (Mukherjee *et al.*, 2017). cClose to

4

70    350 of them are completely assembled, of which ~~B. pertussis genomic sequences deposited at the~~

71    ~~GenBank database, of which~~ 330 are the genome sequences of the isolates from the aP vaccine era

72    (since 1990s). In this study we thus sequenced and *de novo* assembled genomes of five historical *B.*

73    *pertussis* strains that were collected in the former Czechoslovakia between 1954 and 1965. The very

74    same strains were used for the formulation of the DTwP vaccine Alditepera, produced by the

75    Institute of Sera and Vaccines in Prague (Pekarek & Rezabek, 1959, Pekarek & Rezabek, 1959). These

76    strains were not extensively passaged under laboratory conditions and represent a unique set of

77    isolates from the pre-wP vaccine (1954-1957) and early wP vaccine era (1958-1965). The content and

78    organization of the genome of these strains was compared to that of other vaccine strains and recent

79    clinical isolates.

80    **GENOME SEQUENCING AND ANNOTATION**

81    Genomes were sequenced using Illumina MiSeq (paired-end sequencing protocol) and PacBio RSII

82    platforms. Illumina data is deposited under the project PRJEB4543 in Genbank. PacBio reads were

83    assembled using HGAP SMRT Portal protocol and Illumina data was used to further polish the

84    assemblies with Pilon software (Walker *et al.*, 2014). All genomes were *de novo* assembled into single

85    contigs (Supplementary table 1), deposited in GenBank under accession numbers ERS2367611-

86    ERS2367615 and annotated using Prokka software (Seemann, 2014).

87

88    **PHYLOGENETIC ANALYSIS**

89    Genomic analysis revealed that all Czech vaccine strains belong to *ptxP1* lineage (carry pertussis toxin

90    promoter type 1). *B. pertussis* strains from *ptxP1* lineage form a phylogenetic cluster separate from

91    *ptxP3* strains which emerged in the last 25-30 years (Bart *et al.*, 2014, Weigand *et al.*, 2017).

92    Therefore, to put genomic sequences of the Czech vaccine strains into broader context with

93    previously completely sequenced *ptxP1* and *ptxP2* strains of *B. pertussis*, we performed SNP-based

94    phylogenetic analysis using the kSNP3.0 program with *k* of 23 and maximum parsimony method

95    (Gardner *et al.*, 2015). In total we analyzed 19 *ptxP1* and 4 *ptxP2* strains, which were isolated from

96    various geographic locations (USA, China, Japan, UK, Netherlands) from 1935 to 2012

97    (Supplementary Table 1). Phylogenetic analysis based on the 851 detected SNPs divided the strains

98    into six major clusters (Figure 1A). Most of the recent isolates (isolated in 2000-2012) containing

99    *ptxP1* allele clustered separately from the old *ptxP1* isolates (isolated in 1935-1965). Surprisingly, one

100   of the old Czech strains, VS67, clustered together with the strain E945, which was isolated in the USA

101   in 2005. The other four Czech strains formed a distinct cluster (cluster 4), separated from the other

102   strains by six synonymous, six non-synonymous and three intergenic SNPs (Supplementary Table 2).

103

5

104 **GENOME ORGANIZATION**

105 The sequenced genomes were aligned by progressiveMauve algorithm with default parameters

106 (Darling *et al.*, 2010). Alignment revealed that when compared to the reference strain Tohama I, all

107 genomes contain large-scale structural rearrangements (Figure 1B). According to the genome

108 organization, the strains could be classified into three groups. One group contains three strains:

109 VS377, VS401 and VS366. Strain VS393 differs from this group by a single large inversion, which,

110 among other genes, contains *fha/fim* and type III secretion system loci. VS67 differs from the other

111 four strains by two additional large-scale inversions. In order to determine whether genome

112 organization observed in the Czech strains can be also found among other already characterized

113 strains, we have extracted and compared the permutation matrices from the Mauve alignment

114 utilizing scripts published previously (Weigand *et al.*, 2017). This analysis revealed that Czech strains

115 have a unique genome organization that has not been found so far in the other *ptxP1* and *ptxP2*

116 strains (data not shown).

117

118 **REGIONS OF DIFFERENCE**

119 Mauve output was used to extract genome regions differentially present among the strains. In total

120 eight such regions of difference were identified among the studied Czech strains and Tohama I (Table

121 1, Supplementary Figure 1). All RDs are either directly or in close proximity flanked by IS481 elements

122 which indicates the mobile nature of the RDs.

123 Majority of the identified GRs have been described previously and for them we kept the previously

124 established designation (Brinig *et al.*, 2006, Bart *et al.*, 2010). Two of the newly reported RDs were

125 consecutively named RD30 and RD31 (Table 1). RD30, which is duplicated in VS67, contains genes,

126 which code for the MFS transporter and a putative membrane protein (BP2451 and BP2452). It is

127 possible that the duplicated region allows for enhanced transport of the cargo across the membrane.

128 RD31 consists of 5 genes, but the only gene with assigned function is BP0894 coding for mannose-6-

129 phosphate isomerase.

130

131 **PREVALENCE OF RDs**

132 Analysis of the association of the identified RDs with other *B. pertussis* strains showed that in many

133 cases the distinct distribution of RDs among the strains correlated with their phylogenetic

134 assignment (Supplementary Table 3). For instance, RD3 and RD5 are associated with pre-aP era *B.*

135 *pertussis* strains from phylogenetic clusters 5 and 6. RD3 is also present in the Czech strains

136 comprising cluster 4. In contrast, duplication of RD30 is prevalent in aP-era strains from the

137 phylogenetic cluster 3. This suggests that the duplication might provide the currently circulating

138 strains with a selective fitness advantage. The distribution of RD22-24 and RD26 among the strains is

5

6

139  very similar, suggesting that the functions encoded within these loci might be linked. These GRs are

140  found in almost all strains except for the cluster 6. RD29 and RD31 are missing in the Czech vaccine

141  strains VS366 and VS393, respectively, indicating that the loss of these regions might have not

142  conferred any advantage to these strains which possibly prevented their further spread.

143

144  **DISCUSSION**

145  In this study we conducted a comprehensive analysis of the genomes of five *B. pertussis* strains that

146  were collected from 1954 to 1965 in the former Czechoslovakia at times when population travel into

147  and out of the country was very limited. These representative isolates were later used for the

148  development of the local wP vaccine. In contrast to the Japanese vaccine strain Tohama I, the Czech

149  vaccine strains did not undergo massive passaging under laboratory conditions. This is the first study

150  on Czech *B. pertussis* vaccine strains and one of the very few providing complete genome assemblies

151  for the strains from the pre-wP vaccine or early wP use era. *De novo* sequencing of the genomes

152  revealed that all Czech strains contain large-scale genome structural rearrangements compared to

153  the reference strain Tohama I.

154  SNP-based phylogeny revealed that four of the strains form a separate cluster distinct from other

155  so far analyzed strains, suggesting that at the time of their isolation the geographic factors played a

156  significant role. It is tempting to speculate that following massive immunization by the wP vaccine,

157  these strains disappeared from the population and did not spread globally. On the other hand, the

158  fifth Czech strain V67 clusters together with a recent US isolate suggesting that it belongs to a

159  lineage the descendants of which may still circulate within immunized population.

160  In agreement with SNP-based phylogenetic analysis, Czech *B. pertussis* strains exhibit a genome

161  organization pattern that distinguishes them from other *ptxP1* and *ptxP2* strains, and contain

162  some new and some previously reported regions of difference (Brinig *et al.*, 2006, Bart *et al.*,

163  2010). Loss of RD3 and/or RD5 is characteristic of *B. pertussis* strains isolated from other

164  countries during the early wP use period (Kallonen *et al.*, 2011). Accordingly, RD5 was absent

165  from all Czech strains and RD3 was lost from VS67 strain. Four of the RDs (RD22-24 and RD26)

166  absent in Tohama I, are present in the Czech strains thereby supporting earlier reports which

167  demonstrate the presence of these loci in the majority of older European *B. pertussis* isolates and

168  vaccine strains (Kallonen *et al.*, 2011).

169  To conclude, our study suggests that the analysis of strains from pre-vaccine era is of high

170  importance for our understanding of *B. pertussis* evolution in the light of pertussis resurgence.

171  The impact of the various SNPs, RDs and genome re-arrangements on the physiological fitness

172  and pathogenicity of each particular *B. pertussis* isolate, however, remains to be determined.

7

181

182    **REFERENCES**

183    Amman F, D'Halluin A, Antoine R*, et al.* (2018) Primary transcriptome analysis reveals importance of
184    IS elements for the shaping of the transcriptional landscape of Bordetella pertussis. *RNA Biol,*
185    *101080/1547628620181462655*.
186    Bart MJ, van Gent M, van der Heide HG, Boekhorst J, Hermans P, Parkhill J & Mooi FR (2010)
187    Comparative genomics of prevaccination and modern Bordetella pertussis strains. *BMC genomics* **11**:
188    627.
189    Bart MJ, Harris SR, Advani A*, et al.* (2014) Global population structure and evolution of Bordetella
190    pertussis and their relationship with vaccination. *mBio* **5**: e01074.
191    Bowden KE, Weigand MR, Peng Y*, et al.* (2016) Genome Structural Diversity among 31 Bordetella
192    pertussis Isolates from Two Recent U.S. Whooping Cough Statewide Epidemics. *mSphere* **1**.
193    Brinig MM, Cummings CA, Sanden GN, Stefanelli P, Lawrence A & Relman DA (2006) Significant gene
194    order and expression differences in Bordetella pertussis despite limited gene content variation.
195    *Journal of bacteriology* **188**: 2375-2382.
196    Burdin N, Handy LK & Plotkin SA (2017) What Is Wrong with Pertussis Vaccine Immunity? The
197    Problem of Waning Effectiveness of Pertussis Vaccines. *Cold Spring Harb Perspect Biol* **9**.
198    Cherry JD (2010) The present and future control of pertussis. *Clin Infect Dis* **51**: 663-667.
199    Chlibek R, Smetana J, Sosovickova R, Fabianova K, Zavadilova J, Dite P, Gal P, Naplava P & Lzicarova D
200    (2017) Seroepidemiology of whooping cough in the Czech Republic: estimates of incidence of
201    infection in adults. *Public Health* **150**: 77-83.
202    Cummings CA, Brinig MM, Lepp PW, van de Pas S & Relman DA (2004) Bordetella species are
203    distinguished by patterns of substantial gene loss and host adaptation. *Journal of bacteriology* **186**:
204    1484-1492.
205    Darling AE, Mau B & Perna NT (2010) progressiveMauve: multiple genome alignment with gene gain,
206    loss and rearrangement. *PloS one* **5**: e11147.
207    Fabianova K, Benes C & Kriz B (2010) A steady rise in incidence of pertussis since nineties in the Czech
208    Republic. *Epidemiol Mikrobiol Imunol* **59**: 25-33.
209    Gardner SN, Slezak T & Hall BG (2015) kSNP3.0: SNP detection and phylogenetic analysis of genomes
210    without genome alignment or reference genome. *Bioinformatics* **31**: 2877-2878.
211    Kallonen T, Grondahl-Yli-Hannuksela K, Elomaa A, Lutynska A, Fry NK, Mertsola J & He Q (2011)
212    Differences in the genomic content of Bordetella pertussis isolates before and after introduction of
213    pertussis vaccines in four European countries. *Infection, genetics and evolution : journal of molecular*
214    *epidemiology and evolutionary genetics in infectious diseases* **11**: 2034-2042.

8

215 King AJ, van Gorkom T, van der Heide HG, Advani A & van der Lee S (2010) Changes in the genomic
216 content of circulating Bordetella pertussis strains isolated from the Netherlands, Sweden, Japan and
217 Australia: adaptive evolution or drift? *BMC genomics* **11**: 64.
218 Mooi FR (2010) Bordetella pertussis and vaccination: the persistence of a genetically monomorphic
219 pathogen. *Infection, genetics and evolution : journal of molecular epidemiology and evolutionary*
220 *genetics in infectious diseases* **10**: 36-49.
221 Mooi FR, Van Der Maas NA & De Melker HE (2014) Pertussis resurgence: waning immunity and
222 pathogen adaptation - two sides of the same coin. *Epidemiol Infect* **142**: 685-694.
223 Mukherjee S, Stamatis D, Bertsch J, *et al.* (2017) Genomes OnLine Database (GOLD) v.6: data updates
224 and feature enhancements. *Nucleic acids research* **45**: D446-D456.
225 Parkhill J, Sebaihia M, Preston A, *et al.* (2003) Comparative analysis of the genome sequences of
226 Bordetella pertussis, Bordetella parapertussis and Bordetella bronchiseptica. *Nature genetics* **35**: 32-
227 40.
228 Pekarek J & Rezabek K (1959) An endocrinological test for innocuity of the pertussis vaccine. *J Hyg*
229 *Epidemiol Microbiol Immunol* **3**: 79-84.
230 Pekarek J & Rezabek K (1959) The investigation of different components of pertussis vaccine
231 obtained by centrifugation. *J Hyg Epidemiol Microbiol Immunol* **3**: 67-78.
232 Raguckas SE, VandenBussche HL, Jacobs C & Klepser ME (2007) Pertussis resurgence: diagnosis,
233 treatment, prevention, and beyond. *Pharmacotherapy* **27**: 41-52.
234 Sealey KL, Belcher T & Preston A (2016) Bordetella pertussis epidemiology and evolution in the light
235 of pertussis resurgence. *Infection, genetics and evolution : journal of molecular epidemiology and*
236 *evolutionary genetics in infectious diseases* **40**: 136-143.
237 Seemann T (2014) Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**: 2068-2069.
238 Vysoka-Burianova B, Burian V, Maixnerova M, *et al.* (1976) Surveillance of pertussis in the CSSR. IV.
239 Immunological surveys of antibodies to pertussis and parapertussis in the Bohemian regions and in
240 Slovakia in 1958 - 1971. *J Hyg Epidemiol Microbiol Immunol* **21**: 229-247.
241 Walker BJ, Abeel T, Shea T, *et al.* (2014) Pilon: an integrated tool for comprehensive microbial variant
242 detection and genome assembly improvement. *PloS one* **9**: e112963.
243 Weigand MR, Peng Y, Loparev V, *et al.* (2017) The History of Bordetella pertussis Genome Evolution
244 Includes Structural Rearrangement. *Journal of bacteriology* **199**.
245 WHO (2006) Vaccine preventable deaths and the Global Immunization Vision and Strategy, 2006-
246 2015. *MMWR Morb Mortal Wkly Rep* **55**: 511-515.

247

248

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

**Dienstbier et al. Figure 1**

**A**

- B227
- B1920
- 1 ┌ J448
-   └ J447
- B226
- 2 VS67
- E945
- J445
- 3 B202
-   B201
-   VA-UT25Sm1
-   E976
-   H740
-   H375
-   I344
- 4 VS366
-   VS377
-   VS401
-   VS393
- 5 B199
-   * J446
-   BP137
-   B203
- C393
- 6 J043
-   Pelita
-   J042
-   Tohama I

Tree scale: 0.1

- ☐ pre-ACV era strains
- ☐ ACV era strains
- ✳ ptxP2 strains
- ○ Isolated in UK
- ● Isolated in Netherlands
- ○ Reference strain from Serum Institute of India
- ● Isolated in Czech Republic
- ○ Isolated in US
- ○ Isolated in Japan
- ○ Isolated in China

**B**

cya prn    fha/fim   T3SS          ptx

Tohama I

VS67

VS377, VS401, VS366

VS393

0  0.4  0.8  1.2  1.6  2  2.4  2.8  3.2  3.6  4  Mbp

**Figure 1. Phylogenetic tree and genome alignments of the analysed *B. pertussis* isolates. (A)** SNP-based maximum parsimony phylogenetic reconstruction of *ptxP1* and *ptxP2* lineages of *B. pertussis* strains. All uncollapsed nodes have bootstrap value support of >75%. The dashed line indicates the branch length shortened for better visualization. **(B)** Genome alignment of the studied *B. pertussis* Czech isolates. Homologous gene blocks are denoted by the same colour. Genome loci coding for the key virulence factors are marked above the schematic Tohama I genome representation: *cya* – adenylate cyclase; *prn* – pertactin; *fha/fim* filamentous hemagglutinin and fimbriae; T3SS – type III secretion system; *ptx* – pertussis toxin.

355x270mm (300 x 300 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

**Table 1. Distribution of regions of difference among Czech vaccine strains**

| RD | Presence/Absence in strain | Start, RD | Stop, RD | Reference |
|---|---|---|---|---|
| RD3 | Absent in VS67 | BP0911 | BP0937 | Brinig et al., 2006 |
| RD5 | Present only in Tohama I | BP1136 | BP1141 | Brinig et al., 2006 |
| RD22 | Absent only in Tohama I | BB0541 | BB0534 | Brinig et al., 2006 |
| RD23 | Absent only in Tohama I | BB0917 | BB0921 | Bart et al., 2010 |
| RD24 | Absent only in Tohama I | BB1140 | BB1158 | Brinig et al., 2006 |
| RD26 | Absent only in Tohama I | BB4888 | BB4880 | Brinig et al., 2006 |
| RD29 | Absent only in VS366 | BP2820 | BP2832 | King et al., 2008 |
| RD30 | Repeated 2x in VS67 | BP2452 | BP2451 | This study |
| RD31 | Absent only in VS393 | BP0892 | BP0896 | This study |

Dienstbier et al. Supplementary Figure 1

Graphical representation of the genetic content within the regions of difference identified in Czech vaccine strains . Regions are arranged in the VS67 strain gene order. Designation of the regions of difference is shown on the left, open reading frames are depicted as blocks and color coded according to the predicted function, their position indicates the strandedness.

287x264mm (300 x 300 DPI)

1
2
3 **Supplementary table 1. Strains used in this study**
4

| Strain | Genbank Acces. No. | Country | Year | ptxP | Reference | Additional remarks |
|---|---|---|---|---|---|---|
| B199 | CP022361 | USA: PA | 1935 | ptxp2 | Unpublished | |
| B203 | CP012128 | USA: MI | 1939 | ptxP2 | Weigand *et al.*, 2016) | Sanofi-Pasteur MSD, strain 10536 |
| B202 | CP016338 | USA: PA | 1946 | ptxP1 | Weigand *et al.*, 2016 | Lederle Laboratories, strain 134 |
| J042 | CP019869 | USA | 1947 | ptxP1 | Unpublished | |
| J043 | CP016887 | USA | 1947 | ptxP1 | Unpublished | |
| C393 | CP010963 | China | 1951 | ptxP1 | Bowden *et al.*, 2016 | CS, Chinese vaccine reference |
| E476 | CP010964 | Japan | 1954 | ptxP1 | Bowden *et al.*, 2016 | Tohama I, GlaxoSmithKline vaccine reference |
| B201 | CP013075 | USA: IN | 1955 | ptxP1 | Weigand *et al.*, 2017 | |
| B227 | CP013076 | UK | 1967 | ptxP1 | Weigand *et al.*, 2017 | |
| B226 | CP016957 | UK | 1967 | ptxP1 | Unpublished | |
| B1920 | CP009752 | Netherlands | 2000 | ptxP1 | Bart *et al.*, 2014 | |
| VA-UT25Sm1 | CP015771 | USA: ~~VA~~TX | ~~2001~~1977 | ptxP1 | Unpublished | |
| E976 | CP011175 | USA: NY | 2005 | ptxP1 | Weigand *et al.*, 2017 | |
| E945 | CP016956 | USA: CA | 2005 | ptxP1 | Unpublished | |
| H375 | CP010961 | USA: CA | 2010 | ptxP1 | Bowden *et al.*, 2016 | |
| H740 | CP011190 | USA: GA | 2011 | ptxP1 | Weigand *et al.*, 2017 | |
| I344 | CP011255 | USA: MN | 2012 | ptxP1 | Weigand *et al.*, 2017 | |
| Bp137 | CP010323 | USA | ND | ptxP2 | Akamatsu *et al.*, 2015 | Vaccine strain in Latin America |
| J448 | CP017405 | ND | ND | ptxP1 | Weigand *et al.*, 2016 | Reference strain from Serum Institute of India |
| J445 | CP017402 | ND | ND | ptxP1 | Weigand *et al.*, 2016 | Reference strain from Serum Institute of India |
| J447 | CP017404 | ND | ND | ptxP1 | Weigand *et al.*, 2016 | Reference strain from Serum Institute of India |
| J446 | CP017403 | ND | ND | ptxP2 | Weigand *et al.*, 2016 | Reference strain from Serum Institute of India |
| Pelita III | CP019957 | Japan | ND | ptxP1 | Unpublished | Indonesian reference strain, P.T. Bio Farma Indonesia |
| VS67 | ERZ500380 | Czech Republic | before 1966 | ptxP1 | This study | |
| VS393 | ERZ500382 | Czech Republic | before 1966 | ptxP1 | This study | |
| VS401 | ERZ500384 | Czech Republic | 1954 | ptxP1 | This study | |
| VS377 | ERZ500383 | Czech Republic | before 1966 | ptxP1 | This study | |
| VS366 | ERZ500381 | Czech Republic | 1957 | ptxP1 | This study | |

**REFERENCES**

Akamatsu MA, Nishiyama MY, Jr., Morone M*, et al.* (2015) Whole-Genome Sequence of a Bordetella pertussis Brazilian Vaccine Strain. *Genome Announc* **3**.
Bart MJ, Harris SR, Advani A*, et al.* (2014) Global population structure and evolution of Bordetella pertussis and their relationship with vaccination. *MBio* **5**: e01074.
Bowden KE, Weigand MR, Peng Y*, et al.* (2016) Genome Structural Diversity among 31 Bordetella pertussis Isolates from Two Recent U.S. Whooping Cough Statewide Epidemics. *mSphere* **1**.
Weigand MR, Peng Y, Loparev V, Batra D, Burroughs M, Johnson T, Juieng P, Rowe L, Tondella ML & Williams MM (2016) Complete Genome Sequences of Bordetella pertussis Vaccine Reference Strains 134 and 10536. *Genome Announc* **4**.
Weigand MR, Peng Y, Loparev V*, et al.* (2016) Complete Genome Sequences of Four Bordetella pertussis Vaccine Reference Strains from Serum Institute of India. *Genome Announc* **4**.
Weigand MR, Peng Y, Loparev V*, et al.* (2017) The History of Bordetella pertussis Genome Evolution Includes Structural Rearrangement. *J Bacteriol* **199**.

1
2
3 **Supplementary table 2. List of SNPs distinguishing the cluster 4 strains**
4
5

| SNP No. | Position | Allele | Intergenic/Synonymous/Non-synonymous | Annotation | Gene in E476 (Tohama I) |
|---|---|---|---|---|---|
| 1 | 3967296 | G/A | Synonymous | cytochrome oxidase subunit I | RD16_18905 |
| 2 | 3310990 | G/A | Synonymous | molybdate ABC transporter substrate-binding protein | RD16_15605 |
| 3 | 1684734 | G/A | Intergenic | histidine kinase | upstream of RD16_07985 |
| 4 | 1482030 | C/T | Synonymous | ATP synthase | RD16_06985 |
| 5 | 1645729 | G/A | Non-synonymous (A/T) | aminopeptidase | RD16_07800 |
| 6 | 2176522 | C/T | Intergenic | MarR family transcriptional regulator | upstream of RD16_10235 |
| 7 | 1921439 | G/A | Synonymous | hypothetical protein | RD16_09100 |
| 8 | 760046 | C/T | Non-synonymous (G/D) | DNA methylase | RD16_03690 |
| 9 | 3477892 | C/T | Non-synonymous (A/T) | ABC transporter substrate-binding protein | RD16_16365 |
| 10 | 33773 | G/A | Non-synonymous (A/V) | 2-methyl citrate dehydratase | RD16_00170 |
| 11 | 371766 | A/G | Intergenic | glutamate dehydrogenase | Upstream of RD16_01860 |
| 12 | 1336320 | C/T | Non-synonymous (E/K) | hypothetical protein | pseudogene |
| 13 | 2157160 | G/A | Synonymous | sodium transporter | pseudogene |
| 14 | 4018189 | C/T | Non-synonymous (G/D) | AraC family transcriptional regulator | RD16_19175 |
| 15 | 2401358 | C/T | Synonymous | chemotaxis protein | RD16_11360 |

**Supplementary table 3. Prevalence of identified genomic regions in strains used in this study**

| Strain | Genbank Acces. No. | RD3 | RD30 | RD29 | RD31 | RD24 | RD22/RD26 | RD23 | RD5 |
|---|---|---|---|---|---|---|---|---|---|
| B199 | CP022361 | + | | + | + | + | + | + | + |
| B203 | CP012128 | + | | + | + | + | + | + | + |
| B202 | CP016338 | | + | + | + | + | + | + | |
| J042 | CP019869 | + | | + | + | | | | + |
| J043 | CP016887 | + | | + | + | | | | + |
| C393 (CS) | CP010963 | + | | + | + | + | + | | + |
| E476 (Tohama I) | CP010964 | + | | + | + | | | | + |
| B201 | CP013075 | | + | + | + | + | + | + | |
| B227 | CP013076 | | + | + | + | + | + | + | |
| B226 | CP016957 | | + | + | + | + | + | + | |
| B1920 | CP009752 | | + | + | + | | + | + | |
| VA-UT25Sm1 | CP015771 | + | + | + | + | + | + | + | |
| E976 | CP011175 | | + | + | + | + | + | + | |
| E945 | CP016956 | | + | + | + | + | + | + | |
| H375 | CP010961 | | + | + | + | + | + | + | |
| H740 | CP011190 | | + | + | + | + | + | + | |
| I344 | CP011255 | | + | + | + | + | + | + | |
| Bp137 | CP010323 | + | | + | + | + | + | + | + |
| J448 | CP017405 | | | + | + | + | + | + | |
| J445 | CP017402 | | + | + | + | + | + | + | |
| J447 | CP017404 | | | + | + | + | + | + | |
| J446 | CP017403 | + | | + | + | + | + | + | + |
| Pelita III | CP019957 | + | | + | + | | | | + |
| VS67 | ERS2367611 | | + | + | + | + | + | + | |
| VS393 | ERS2367613 | + | | + | | + | + | + | |
| VS401 | ERS2367615 | + | | + | + | + | + | + | |
| VS377 | ERS2367614 | + | | + | + | + | + | + | |
| VS366 | ERS2367612 | + | | | + | + | + | + | |