

A peer-reviewed version of this preprint was published in PeerJ on 27 March 2018.

[View the peer-reviewed version](https://peerj.com/articles/4545) (peerj.com/articles/4545), which is the preferred citable publication unless you specifically need to cite this preprint.

Bochkareva OO, Dranenko NO, Ocheredko ES, Kanevsky GM, Lozinsky YN, Khalaycheva VA, Artamonova II, Gelfand MS. 2018. Genome rearrangements and phylogeny reconstruction in *Yersinia pestis*. PeerJ 6:e4545 <https://doi.org/10.7717/peerj.4545>

Genome rearrangements and phylogeny reconstruction in *Yersinia pestis*

Olga O Bochkareva ^{Corresp., 1}, **Natalia O Dranenko** ², **Elena S Ocheredko** ³, **German M Kanevsky** ⁴, **Yaroslav N Lozinsky** ³, **Vera A Khalaycheva** ⁵, **Irena I Artamonova** ^{1,6}, **Mikhail S Gelfand** ^{1,3,7,8}

¹ Kharkevich Institute for Information Transmission Problems, Moscow, Russia

² Department of Molecular and Chemical Physics, Moscow Institute of Physics and Technology, Moscow, Russia

³ Department of Bioengineering and Bioinformatics, Moscow State University, Moscow, Russia

⁴ Higher Chemical College at the Russian Academy of Sciences, Moscow, Russia

⁵ Stavropol State Agrarian University, Stavropol, Russia

⁶ Vavilov Institute of General Genetics, Moscow, Russia

⁷ Faculty of Computer Science, Higher School of Economics, Moscow, Russia

⁸ Skolkovo Institute of Science and Technology, Moscow, Russia

Corresponding Author: Olga O Bochkareva

Email address: olga.bochkaryova@iitp.ru

Genome rearrangements have played an important role in the evolution of *Yersinia pestis* from its progenitor *Yersinia pseudotuberculosis*. Traditional phylogenetic trees for *Y. pestis* based on sequence comparison have short internal branches and low bootstrap supports as only a small number of nucleotide substitutions have occurred. On the other hand, even a small number of genome rearrangements may resolve topological ambiguities in a phylogenetic tree.

We reconstructed the evolutionary history of genome rearrangements in *Y. pestis*. We also reconciled phylogenetic trees for each of the three CRISPR-loci to obtain an integrated scenario of the CRISPR-cassette evolution. We detected numerous parallel inversions and gain/loss events by the analysis of contradictions between the obtained evolutionary trees. We also tested the hypotheses that large within-replicore inversions tend to be balanced by subsequent reversal events and that the core genes less frequently switch the chain by inversions. Both predictions were not confirmed.

Our data indicate that an integrated analysis of sequence-based and inversion-based trees enhances the resolution of phylogenetic reconstruction. In contrast, reconstructions of strain relationships based on solely CRISPR loci may not be reliable, as the history is obscured by large deletions, obliterating the order of spacer gains. Similarly, numerous parallel gene losses preclude reconstruction of phylogeny based on gene content.

1 Genome rearrangements and phylogeny 2 reconstruction in *Yersinia pestis*

3 Olga O. Bochkareva¹, Natalia O. Dranenko², Elena S. Ocheredko³,
4 German M. Kanevsky⁴, Yaroslav N. Lozinsky³, Vera A. Khalaycheva⁵,
5 Irena I. Artamonova^{6,1}, and Mikhail S. Gelfand^{1,3,7,8}

6 ¹Kharkevich Institute for Information Transmission Problems, Moscow, Russia

7 ²Moscow Institute of Physics and Technology, Moscow, Russia

8 ³Lomonosov Moscow State University, Moscow, Russia

9 ⁴Higher Chemical College at the Russian Academy of Sciences, Moscow, Russia

10 ⁵Stavropol State Agrarian University, Stavropol, Russia

11 ⁶Vavilov Institute of General Genetics, Moscow, Russia

12 ⁷Higher School of Economics, Moscow, Russia

13 ⁸Skolkovo Institute of Science and Technology, Moscow, Russia

14 Corresponding author:

15 Olga O. Bochkareva¹

16 Email address: olga.bochkaryova@iitp.ru

17 ABSTRACT

18 Genome rearrangements have played an important role in the evolution of *Yersinia pestis* from its
19 progenitor *Yersinia pseudotuberculosis*. Traditional phylogenetic trees for *Y. pestis* based on sequence
20 comparison have short internal branches and low bootstrap supports as only a small number of nucleotide
21 substitutions have occurred. On the other hand, even a small number of genome rearrangements may
22 resolve topological ambiguities in a phylogenetic tree.

23 We reconstructed the evolutionary history of genome rearrangements in *Y. pestis*. We also reconciled
24 phylogenetic trees for each of the three CRISPR-loci to obtain an integrated scenario of the CRISPR-
25 cassette evolution. We detected numerous parallel inversions and gain/loss events by the analysis
26 of contradictions between the obtained evolutionary trees. We also tested the hypotheses that large
27 within-replicore inversions tend to be balanced by subsequent reversal events and that the core genes
28 less frequently switch the chain by inversions. Both predictions were not confirmed.

29 Our data indicate that an integrated analysis of sequence-based and inversion-based trees enhances
30 the resolution of phylogenetic reconstruction. In contrast, reconstructions of strain relationships based
31 on solely CRISPR loci may not be reliable, as the history is obscured by large deletions, obliterating
32 the order of spacer gains. Similarly, numerous parallel gene losses preclude reconstruction of phylogeny
33 based on gene content.

34 INTRODUCTION

35 *Yersinia pestis*, causing fulminant plague, has evolved clonally from an enteric pathogen, *Yersinia*
36 *pseudotuberculosis*, that, in contrast, causes a relatively benign enteric illness. Horizontal gene acquisition,
37 massive gene loss, and genome rearrangement events all have played important roles in the evolution of *Y.*
38 *pestis* from its progenitor (Achtman et al., 1999). *Y. pseudotuberculosis* and *Y. pestis* differ radically in
39 their pathogenesis despite sharing >97% identity in 75% of their chromosomal genes (Martínez-Chavarría
40 and Vadyvaloo, 2015). As only a small number of nucleotide substitutions have occurred, traditional
41 phylogenetic trees of *Y. pestis* strains based on sequence comparison have short internal branches and low
42 bootstraps. They are also significantly affected by extensive horizontal gene flow between strains due to
43 homologous recombination. Genome rearrangements are less sensitive to homologous recombination and
44 hence allow for an alternative approach to construction of phylogenetic trees, as even a small number
45 of genome rearrangements may resolve topological ambiguities in a phylogenetic tree. The *Y. pestis*
46 chromosome contains a large variety and number of insert sequences (ISs) that may have caused frequent

47 chromosome rearrangements (Liang et al., 2014).

48 Comparison of the KIM genome sequence with *Y. pestis* strain CO92 allowed to divide both genomes
49 into 27 conserved segments and the most parsimonious series of inversions for three multiple-inversion
50 regions were described (Deng et al., 2002). Further, large-scale genome rearrangements were described
51 in strains Antiqua, Nepal and Angola (Chain et al., 2006; Eppinger et al., 2010). Multiple genome
52 alignment of nine *Y. pestis* and *Y. pseudotuberculosis* genomes featured universal Locally Collinear
53 Blocks (LCBs) and yielded seven parsimonious scenarios of the inversion history (Darling et al., 2008).
54 Later, the LCB model has been used to infer the phylogenetic relationships among eight complete *Y.*
55 *pestis* genomes from the breakpoint distance matrix, yielding the conclusion that the pattern of *Y. pestis*
56 chromosome rearrangements reflects the genetic features of specific geographical areas and could be
57 applied to distinguish *Y. pestis* isolates (Liang et al., 2010). A set of gene families from 13 *Yersinia* species
58 was used to reconstruct a complete genome sequence for the ancestor, integrating information from the
59 sequences, the species tree, and the gene order (Duchemin et al., 2015).

60 Being a traditional object for the spoligotyping, a special type of genotyping based on the spacer
61 nucleotide analysis, CRISPR systems of *Y. pestis* strains often served as a model for CRISPR-based
62 evolutionary studies. All three separate genomic CRISPR loci were described in detail (Pourcel et al.,
63 2005), including numerous strains without complete genomes (Vergnaud et al., 2007; Cui et al., 2008;
64 Riehm et al., 2012; Barros et al., 2014; Riehm et al., 2015). Relationships between strains were studied
65 using the distance based on shared and differential spacers content only (Barros et al., 2014) or taking
66 into account the principles of evolutionary cassette dynamics. In particular, the evolutionary history of *Y.*
67 *pestis* based on CRISPR polymorphism was reconstructed in the form of an acyclic oriented graph (Cui
68 et al., 2008). Later, a general mathematical model of CRISPR evolution was applied to reconstruct the
69 relationships of strains for each of the three CRISPR loci (Kupczok and Bollback, 2013).

70 Here, we integrate the history at different levels of genome evolution, including gene flux, sequence
71 divergence, chromosome segmental inversions, and spacer acquisitions and deletions in CRISPR cassettes,
72 for genomes of twelve completely sequenced *Y. pestis* strains and four *Y. pseudotuberculosis* strains.

73 MATERIALS AND METHODS

74 Genomes

75 Complete genome sequences of four *Yersinia pseudotuberculosis* and twelve *Yersinia pestis*, all available
76 as of August 1st, 2013, were taken from the NCBI Genome database (Benson et al., 2013) and are listed
77 in Supplementary Table 1.

78 Construction of orthologs

79 Bidirectional best hits (BBHs) were constructed for each pair of strains using BLASTP (Zhang and
80 Madden, 1997). BLASTP hits with identity <50% or coverage of the shorter sequence <67% were
81 ignored. At the next step, if paralogs were more similar to each other than to either BBH partner, both
82 paralogs were added to the ortholog group. Then, maximal connected components were constructed. This
83 was done using ad hoc software based on the Relational Database Management System (RDBMS) Oracle
84 Database Express Edition.

85 Phylogenetic trees

86 Single-copy universal genes were used to construct a phylogenetic tree. Multiple alignments were obtained
87 using the Muscle software (Edgar, 2004). For concatenation of alignments, long insertions/deletions at
88 gene boundaries were ignored. A phylogenetic tree was constructed using the NJ algorithm with the
89 Mega7 software (Tamura et al., 2013).

90 Synteny blocks and rearrangements history

91 Synteny blocks were constructed using the Sibelia algorithm (Minkin et al., 2013) with default parameter
92 and the length of blocks more than 5000 bp. Blocks that were found in any genome more than once were
93 filtered out. The history of rearrangements was reconstructed using the MGRA algorithm (Avdeyev et al.,
94 2016).

95 **The origins and terminus of replication**

96 The origins and terminators of replication for six *Y. pestis* and two *Y. pseudotuberculosis* strain were
97 previously identified (Darling et al., 2008). We used these data to identify synteny blocks that contains
98 origins and terminators.

99 **Permutation testing for inversions**

100 For each *Y. pestis* strain, we constructed a null distribution for the percentage of chromosome length that
101 switched its location between the leading and lagging chain compared to the common ancestor. Given a
102 set of inversion lengths for each strain, we selected inversion start positions at random in the range from
103 1 to $4.5 \cdot 10^6$ (corresponding to the average genome length). Then, for each strain we obtained 10^4 of
104 random inversion sets to calculate the p-value.

105 **CRISPR analysis**

106 CRISPR-cassettes were downloaded from CRISPRdb (Grissa et al., 2007). Phylogenetic trees were
107 reconstructed manually based on the CRISPR-cassettes evolution rules. At that, two types of events
108 were allowed, addition of a new spacer at the leader end, and deletion of one or several adjacent spacers
109 from any part of a cassette. We further assumed (1) no independent addition of the same spacer to
110 two different cassettes; (2) rare, but possible independent deletions of the same cassette segments; and
111 (3) more probable single deletion of a segment including several adjacent spacers compared to several
112 subsequent deletions of the segment parts.

113 **RESULTS AND DISCUSSION**

114 **Phylogenetic trees and evolutionary events**

115 The phylogenetic tree for the analyzed *Y. pseudotuberculosis* and *Y. pestis* was constructed based on 2408
116 single-copy universal genes using a concatenation of individual nucleotide alignments (Fig. 1A). We used
117 the *Y. enterocolitica* genome to root the tree. We observed that *Y. pestis* strains form a clade within the *Y.*
118 *pseudotuberculosis* subtree, in agreement with previous genome analyses (Chain et al., 2006; Rasmussen
119 et al., 2015). As the phylogeny of *Y. pestis* was not completely resolved, we added 76 genes universal
120 for these strains and reconstructed the *Y. pestis* branch separately (Fig. 1B). There seemed to be several
121 key noise factors. A small number of nucleotide substitutions resulted in low bootstrap values in several
122 vertices, e.g. for Z176003, D182038, and D106004. On the other hand, branches of the Angola, Microtus,
123 and Antiqua strains could have been placed incorrectly due to long branch attraction.

124 To analyze the history of rearrangements we constructed whole-genome alignment using the Sibelia
125 software. We identified 166 synteny blocks, with 130 blocks found in all strains, and the remaining blocks
126 reflecting gain/loss events. The coverage of genomes by blocks exceed 87%.

127 Fixing the tree to the one based on the concatenated sequence alignment of universal genes, we
128 reconstructed 161 inversions and 62 insertions/deletions. (Fig. 2A) Two inversions occurred twice, in
129 D106004 and Z176003, and in A1122 and D182038. Parallel events could be explained by homologous
130 recombination (horizontal transfer between strains) involving a segment containing the inverted fragments.
131 Indeed, in the tree constructed using only the genes involved in the first inversion, A1122 and D182038
132 formed a separate branch (Fig. 3A), indicating a close affinity limited to this fragment. Another possibility
133 is incomplete resolution of the sequence-based tree, e.g. in the case of the D106004, Z176003, D182038
134 group with internal branches having relatively low bootstrap support (Fig. 1B). In the inversion analysis,
135 D106004 and Z176003 are sister strains with D182038 being an outgroup (Fig. 2B), whereas the sequence-
136 based tree for the genes involved in the inversion is poorly resolved and hence provides no information
137 about possible horizontal transfer.

138 To calculate the inversion rate, we performed regression analysis (Fig. 4). On average one inversion
139 occurs per 20 substitutions per genome, $R^2 = 0,87$.

140 About a dozen of parallel insertions/deletions were found. We applied optimization methods (Avdeyev
141 et al., 2016) to obtain an alternative topology with a lower number of parallel events. The optimal topology
142 based on location of the synteny blocks results in 160 inversions and 58 insertions/deletions for *Y. pestis*
143 (Fig. 2B). Here, Antiqua moves to the Microtus node, in agreement with the fact that, according to the
144 ability to ferment glycerol and to reduce nitrate, strains Antiqua, Pestoides, Microtus and Angola belong
145 to the Antiqua biovar (Chain et al., 2006). This analysis demonstrates that parallel gain and loss events

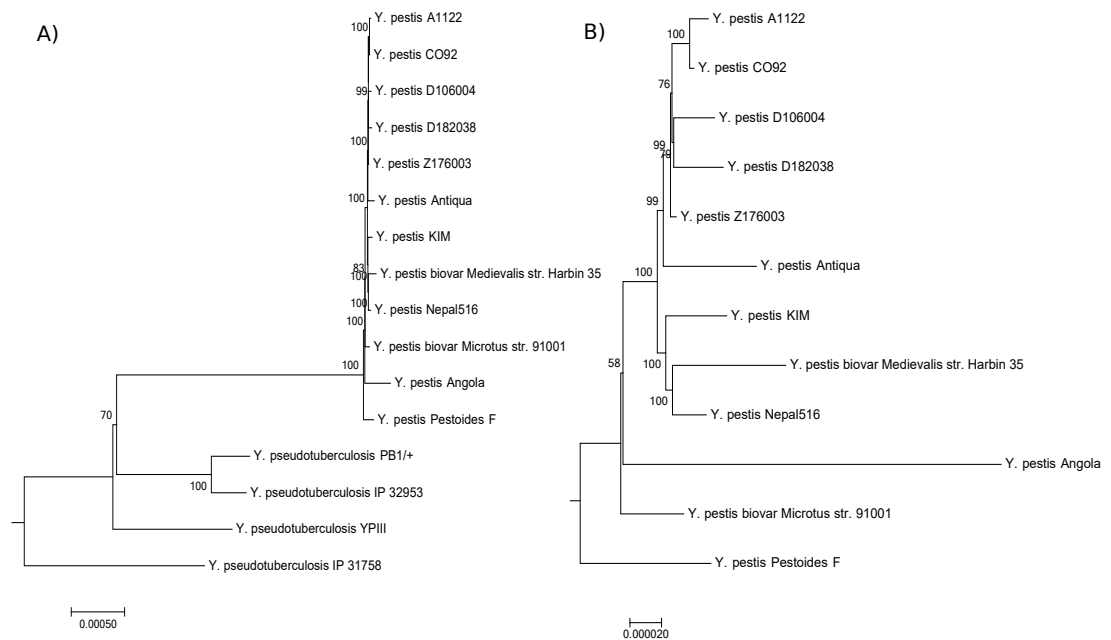


Figure 1. Phylogenetic trees (A) *Yersinia* spp., based on nucleotide alignments of 2408 single-copy universal genes. (B) *Yersinia pestis*, based on nucleotide alignments of 2484 single-copy universal genes including 76 *Y. pestis*-specific universal genes.

146 are not rare in the evolution of *Yersinia* spp., and these events should be considered only in the context of
 147 other criteria.

148 Selection acting on inversions

149 We further analyzed the selection pressure on within-replichore inversions. Inversions of common synteny
 150 blocks were separated into 41 inter-replichore and 49 within-replichore events. We constructed a null
 151 distribution for the fraction of within-replichore inversions in the *Y. pestis* history (see Methods). In the
 152 considered genomes, within-replichore inversions were over-represented with p -value $< 10^{-4}$.

153 We tested the hypothesis that large within-replichore inversions are usually balanced by subsequent
 154 (partial) reversal events. Indeed, inversions within a replichore change the leading/lagging strand A/T and
 155 G/C biases, relative gene density, and gene expression levels. Hence, they may introduce many slightly
 156 deleterious traits and be detrimental (Darling et al., 2008).

157 For strains whose evolution involved more than three within-replichore inversions, we calculated
 158 the fraction of the chromosome length that switched its chain compared to the common ancestor and
 159 calculated the p -values of the null distributions. No significant tendency for reversal could be observed
 160 (data not shown).

161 Only 560 of 2300 universal OGs have never switched the chain compared to the common ancestor.
 162 We tested whether the core genes less frequently switch the chain by inversions using the list of 139
 163 bacterial core genes (Rinke et al., 2013). No bias could be observed, as the fractions of core genes in
 164 stable and inverted synteny blocks were roughly equal (data not shown).

165 CRISPR analysis

166 CRISPR cassettes of the considered *Y. pestis* strains are shown in Fig. 5. Initially, we constructed separate
 167 phylogenetic trees for each of the three CRISPR-loci using the parsimony approach (Fig. 6, see Methods).
 168 As the number of events in each locus was small, the history of each locus could be reconstructed
 169 unambiguously.

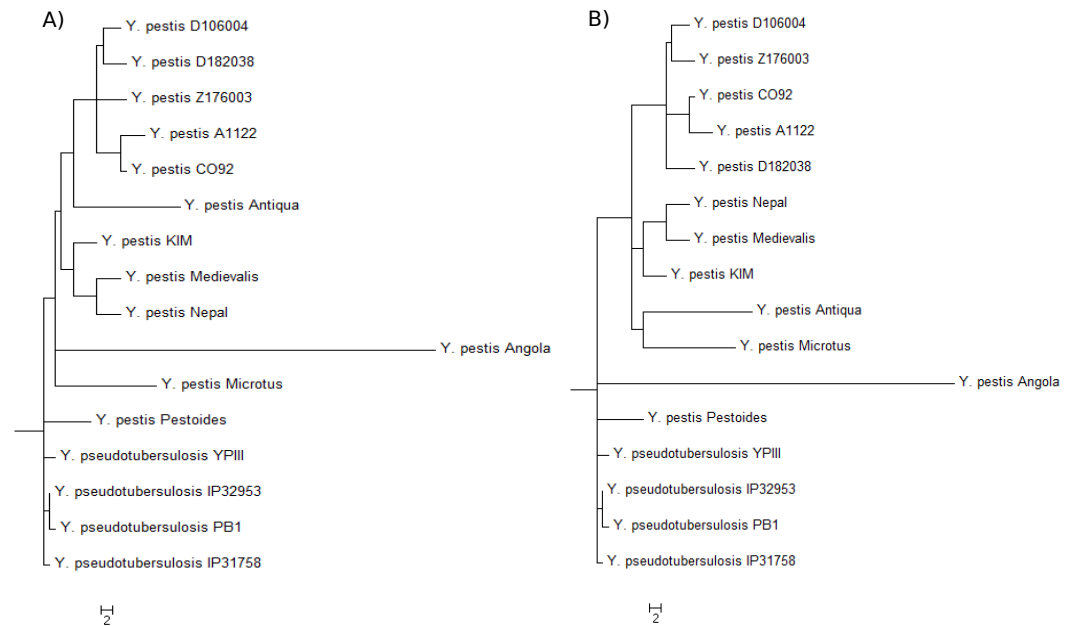


Figure 2. Phylogenetic trees *Yersinia spp.* reflecting the history of rearrangements. Branch lengths reflect the number of inversions, red top (blue bottom) numbers indicate the number of insertions (resp. deletions). (A) Optimal set of rearrangements for the sequence-based topology. (B) Optimal topology based on rearrangements.

170 However, the genome of *Y. pestis* evolves as a whole and the individual histories of the loci should
 171 be reconciled. In this case the reconstruction is ambiguous, as there are two equivalent reconstructions
 172 of the common ancestors and five equal positions of the Nepal strain on the maximum parsimony tree.
 173 Two maximum parsimony trees most compatible with the sequence tree are shown in Fig. 7. The
 174 trees constructed based on nucleotide sequences or rearrangements satisfy the rules of CRISPR-cassettes
 175 evolution (see Methods) but each of them implies two additional losses of cassette segments in comparison
 176 with the maximum parsimony tree. In particular, the sequence-based tree implies two independent parallel
 177 losses of the same segments of the main locus on the Angola and Antiqua strains branches.

178 CONCLUSIONS

179 Detailed reconstruction of evolution of bacterial strains provides a framework for epidemiological studies
 180 and analysis of acquired pathogenesis loci and drug resistance determinants.

181 Using *Y. pestis* as an example, we demonstrate that integrated analysis of sequence-based and inversion-
 182 based trees enhances the resolution of phylogenetic reconstruction. At that, inversions may resolve
 183 branches with low bootstrap support; on the other hand, sequence analysis may distinguish between
 184 parallel inversions and single inversion propagated by homologous recombination of a larger block.

185 In contrast, reconstructions of strain relationships based on solely CRISPR loci may not be reliable,
 186 as the history is greatly obscured by large deletions, obliterating the order of spacer gains. Even less
 187 reliable seem to be reconstructions based on shared spacer content. Similarly, numerous parallel gene
 188 losses preclude reconstruction of phylogeny based on gene content.

189 The hypothesis that large within-replicore inversions are usually balanced by subsequent events was
 190 not confirmed. However, it might be caused by a lack of data as inter-replicore inversions occurred rarely.
 191 The hypothesis that core genes tend not to change their chain during evolution was discarded.

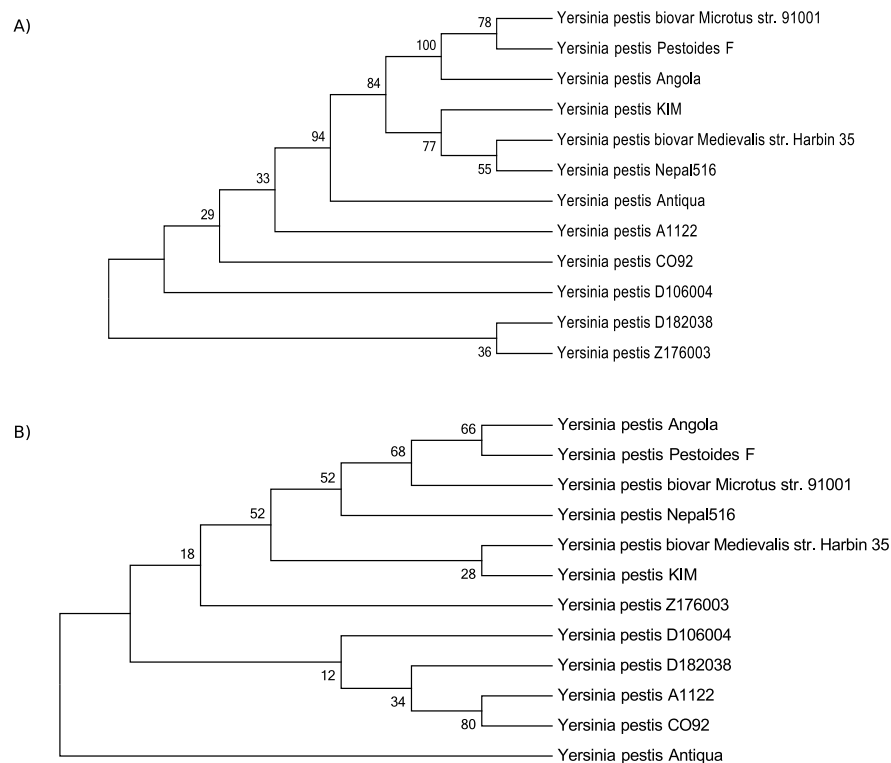


Figure 3. Phylogenetic tree for *Yersinia pestis* genes involved in the parallel inversions (A) in A1122 and D182038 and (B) in Z176003 and D106004.

192 ACKNOWLEDGMENTS

193 This study was supported by the Russian Science Foundation under grant 14-50-00150. It was initiated at
194 the Summer School of Molecular and Theoretical Biology supported by the Dynasty foundation.

195 REFERENCES

- 196 Achtman, M., Zurth, K., Morelli, G., Torrea, G., Guiyoule, A., and Carnie, E. (1999). *Yersinia pestis*, the
197 cause of plague, is a recently emerged clone of *Yersinia pseudotuberculosis*. *Proc Natl Acad Sci U S A*,
198 96(24). 14043–14048.
- 199 Avdeyev, P., Jiang, S., Aganezov, S., Hu, F., and Alekseyev, M. A. (2016). Reconstruction of ancestral
200 genomes in presence of gene gain and loss. *Journal of Computational Biology*, 23(3):150–164.
- 201 Barros, M., França, C., Lins, R., Santos, M., Silva, E., Oliveira, M., Silveira-Filho, V., Rezende, A.,
202 Balbino, V., and Leal-Balbino, T. (2014). Dynamics of CRISPR loci in microevolutionary process of
203 *Yersinia pestis* strains. *PLoS One*, 9(9):e108353.
- 204 Benson, D. A., Cavanaugh, M., Clark, K., Karsch-Mizrachi, I., Lipman, D. J., Ostell, J., and Sayers, E. W.
205 (2013). GenBank. *Nucleic Acids Res.*, 45(Database issue):D30–5.
- 206 Chain, P. S., Hu, P., Malfatti, S. A., Radnedge, L., Larimer, F., Vergez, L. M., Worsham, P., Chu, M. C., and
207 Andersen, G. L. (2006). Complete genome sequence of *Yersinia pestis* strains Antiqua and Nepal516:
208 Evidence of gene reduction in an emerging pathogen. *J Bacteriol.*, 188(12).
- 209 Cui, Y., Li, Y., Gorgé, O., Platonov, M., Yan, Y., Guo, Z., Pourcel, C., Dentovskaya, S., Balakhonov, S.,
210 Wang, X., Song, Y., Anisimov, A., Vergnaud, G., and Yang, R. (2008). Insight into microevolution of
211 *Yersinia pestis* by clustered regularly interspaced short palindromic repeats. *PLoS One*, 3(7):e2652.
- 212 Darling, A. E., Miklós, I., and Ragan, M. A. (2008). Dynamics of genome rearrangement in bacterial
213 populations. *PLoS Genet.*, 4(7):e1000128.

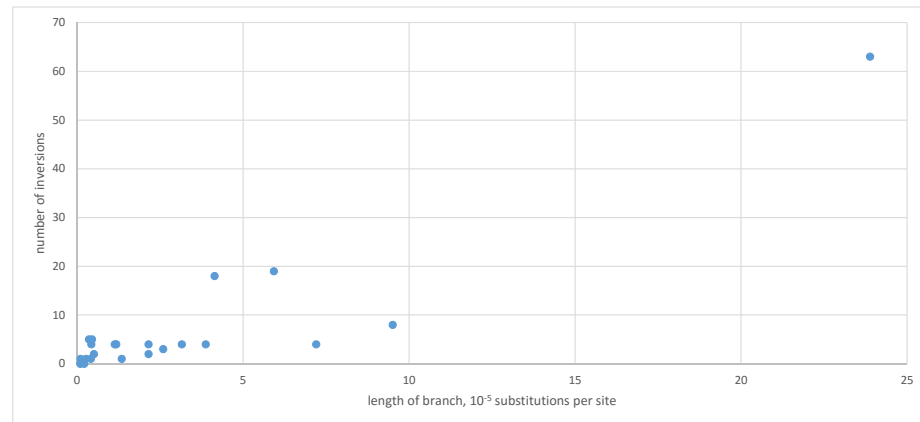


Figure 4. Correlation between inversion rates and mutation rates. Each dot corresponds to a branch in the *Y. pestis* phylogenetic tree. Horizontal axis branch length in substitution per site. Vertical axis shows the number of inversions.

- 214 Deng, W., Burland, V., Plunkett, G. I., Boutin, A., Mayhew, G., Liss, P., Perna, N., Rose, D., Mau, B.,
 215 Zhou, S., Schwartz, D., Fetherston, J., Lindler, L., Brubaker, R., Plano, G., Straley, S., McDonough, K.,
 216 Nilles, M., Matson, J., Blattner, F., and Perry, R. (2002). Genome sequence of *Yersinia pestis* KIM. *J*
 217 *Bacteriol*, 184(16):4601–4611.
- 218 Duchemin, W., Daubin, V., and Tannier, E. (2015). Reconstruction of an ancestral *Yersinia pestis* genome
 219 and comparison with an ancient sequence. *BMC Genomics*, 16.
- 220 Edgar, R. (2004). MUSCLE: multiple sequence alignment with high accuracy and high throughput.
 221 *Nucleic Acids Res*, 32(5):1792–1797.
- 222 Eppinger, M., Worsham, P. L., Nikolich, M. P., Riley, D. R., Sebastian, Y., Mou, S., Achtman, M.,
 223 Lindler, L. E., and Ravel, J. (2010). Genome sequence of the deep-rooted *Yersinia pestis* strain
 224 Angola reveals new insights into the evolution and pangenome of the plague bacterium. *J Bacteriol.*,
 225 192(6):1685–1699.
- 226 Grissa, I., Vergnaud, G., and Pourcel, C. (2007). The CRISPRdb database and tools to display CRISPRs
 227 and to generate dictionaries of spacers and repeats. *BMC Bioinformatics*, 8:172.
- 228 Kupczok, A. and Bollback, J. (2013). Probabilistic models for CRISPR spacer content evolution. *BMC*
 229 *Evol Biol.*, 13(54).
- 230 Liang, Y., Hou, X., Wang, Y., Cui, Z., Zhang, Z., Zhu, X., Xia, L., Shen, X., Cai, H., Wang, J., Xu,
 231 D., Zhang, E., Zhang, H., Wei, J., He, J., Song, Z., Yu, X., Yu, D., and Hai, R. (2010). Genome
 232 rearrangements of completely sequenced strains of *Yersinia pestis*. *J Clin Microbiol*, 48(5):1619–1623.
- 233 Liang, Y., Xie, F., Tang, X., Wang, M., Zhang, E., Zhang, Z., Cai, H., Wang, Y., Shen, X., Zhao, H., Yu,
 234 D., Xia, L., and Hai, R. (2014). Chromosomal rearrangement features of *Yersinia pestis* strains from
 235 natural plague foci in China. *Am J Trop Med Hyg.*, 91(4):722–728.
- 236 Martínez-Chavarría, L. C. and Vadyvaloo, V. (2015). *Yersinia pestis* and *Yersinia pseudotuberculosis*
 237 infection: a regulatory RNA perspective. *Front Microbiol.*, 6:956.
- 238 Minkin, I., Patel, A., Kolmogorov, M., Vyahhi, N., and Pham, S. (2013). Sibelia: A scalable and

	Main locus								Additional locus 1				Additional locus 2				
A1122 (NC017168_4)	sp1	sp2	sp3	sp4	sp5	sp6	sp7	sp8	sp5	sp4	sp3	sp2			sp3	sp2	sp1
CO92 (NC_003143_3)	sp7	sp6	sp5	sp4	sp3	sp2	sp1	sp0	sp0	sp1	sp2	sp3			sp0	sp1	sp2
D106004 (NC_017154_3)		sp7	sp6	sp5	sp4	sp3	sp2	sp1	sp4	sp3	sp2	sp1			sp1	sp2	sp3
Z176003 (NC_014029_3)		sp7	sp6	sp5	sp4	sp3	sp2	sp1	sp1	sp0	sp3				sp1	sp2	sp3
Pestoides F (NC_009381_3)			sp5	sp4	sp3	sp2	sp1	sp0	sp5	sp4	sp3	sp2	sp1	sp0	sp4	sp3	sp2
D182038 (NC_017160_3)	sp1	sp2	sp3	sp4		sp5	sp6	sp7	sp1	sp2	sp3				sp3	sp2	sp1
Antiqua (NC_008150_4)			sp5	sp4	sp3	sp2	sp1	sp0	sp2	sp1					sp0	sp1	sp2
Angola (NC_010159_2)				sp4	sp3	sp2	sp1	sp0	NO						NO		
KIM10+ (NC_004088_3)						sp0	sp1	sp2	sp0	sp1	sp2				sp2	sp1	sp0
Microtus (NC_005810_3)			sp2		sp1			sp0	sp3	sp2	sp1				NO		
Medievalis (NC_017265_2)						sp5	sp6	sp7	sp4	sp3	sp2				sp3	sp2	sp1
Nepal516 (NC_008149_1)	sp0								sp0	sp1	sp2				sp2	sp1	sp0

Figure 5. CRISPR-cassettes of completely sequenced *Y. pestis* strains. Cassette IDs and spacer numbers are given according to CRISPRdb (Grissa et al., 2007). Identical spacers are shown by the same color; unique spacers are set in frames.

- 239 comprehensive synteny block generation tool for closely related microbial genomes. *13th Workshop on*
 240 *Algorithms in Bioinformatics (WABI2013)*.
- 241 Pourcel, C., Salvignol, G., and Vergnaud, G. (2005). CRISPR elements in *Yersinia pestis* acquire new
 242 repeats by preferential uptake of bacteriophage DNA, and provide additional tools for evolutionary
 243 studies. *Microbiology*, 151(Pt 3):653–63.
- 244 Rasmussen, S., Allentoft, M. E., Nielsen, K., Orlando, L., Sikora, M., Sjögren, K.-G., Pedersen, A. G.,
 245 Schubert, M., Dam, A. V., Kapel, C. M. O., Nielsen, H. B., Brunak, S., Avetisyan, P., Epimakhov, A.,
 246 Khalyapin, M. V., Gnuni, A., Kriiska, A., Lasak, I., Metspalu, M., Moiseyev, V., Gromov, A., Pokutta,
 247 D., Saag, L., Varul, L., Yepiskoposyan, L., Sicheritz-Pontén, T., Foley, R. A., Lahr, M. M., Nielsen, R.,
 248 Kristiansen, K., and Willerslev, E. (2015). Early divergent strains of *Yersinia pestis* in Eurasia 5,000
 249 years ago. *Cell*, 163(3):571–582.
- 250 Riehm, J., Projahn, M., Vogler, A., Rajerison, M., Andersen, G., Hall, C., Zimmermann, T., Soanandrasana,
 251 R., Andrianaivoarimanana, V., Straubinger, R., Nottingham, R., Keim, P., Wagner, D., and Scholz, H.
 252 (2015). Diverse genotypes of *Yersinia pestis* caused plague in Madagascar in 2007. *PLoS Negl Trop*
 253 *Dis*, 9(6):e0003844.
- 254 Riehm, J., Vergnaud, G., Kiefer, D., Damdindorj, T., Dashdavaa, O., Khurelsukh, T., Zöller, L., Wölfel,
 255 R., Flèche, P. L., and Scholz, H. (2012). *Yersinia pestis* lineages in Mongolia. *PLoS One*, 7(2):e30624.
- 256 Rinke, C., Darling, A., Malfatti, S., Tsiamis, G., Stepanauskas, R., Schwientek, P., Sievert, S. M., Rubin,
 257 E. M., Sczyrba, A., Ivanova, N. N., Swan, B. K., Liu, W.-T., Eisen, J. A., Hugenholtz, P., Woyke,
 258 T., Gies, E. A., Dodsworth, J. A., and Hallam, S. J. (2013). Insights into the phylogeny and coding
 259 potential of microbial dark matter. *Nature*, 499(7459):431–7.
- 260 Tamura, K., Stecher, G., Peterson, D., Filipinski, A., and Kumar, S. (2013). MEGA6: Molecular Evolution-
 261 ary Genetics Analysis version 6.0. *Molecular Biology and Evolution*, 30(12):2725–2729.
- 262 Vergnaud, G., Li, Y., Gorgé, O., Cui, Y., Song, Y., Zhou, D., Grissa, I., Dentovskaya, S., Platonov, M.,
 263 Rakin, A., Balakhonov, S., Neubauer, H., Pourcel, C., Anisimov, A., and Yang, R. (2007). Analysis of

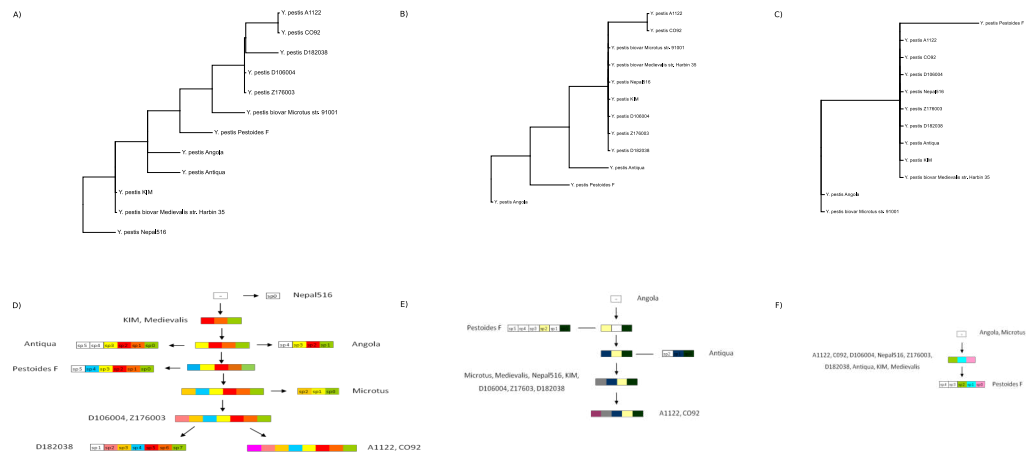


Figure 6. Cladograms (A, B, C) and schemas of evolution (D, E, F) of three CRISPR loci of *Y. pestis*. (A, D) The main, most variable, locus; (B, E) additional locus 1; (C, D) additional locus 2

264 the three *Yersinia pestis* CRISPR loci provides new tools for phylogenetic studies and possibly for the
 265 investigation of ancient DNA. *Adv Exp Med Biol*, 603:327–38.

266 Zhang, J. and Madden, T. (1997). PowerBLAST: A new network BLAST application for interactive or
 267 automated sequence analysis and annotation. *Genome Res.*, 7(6):649–656.

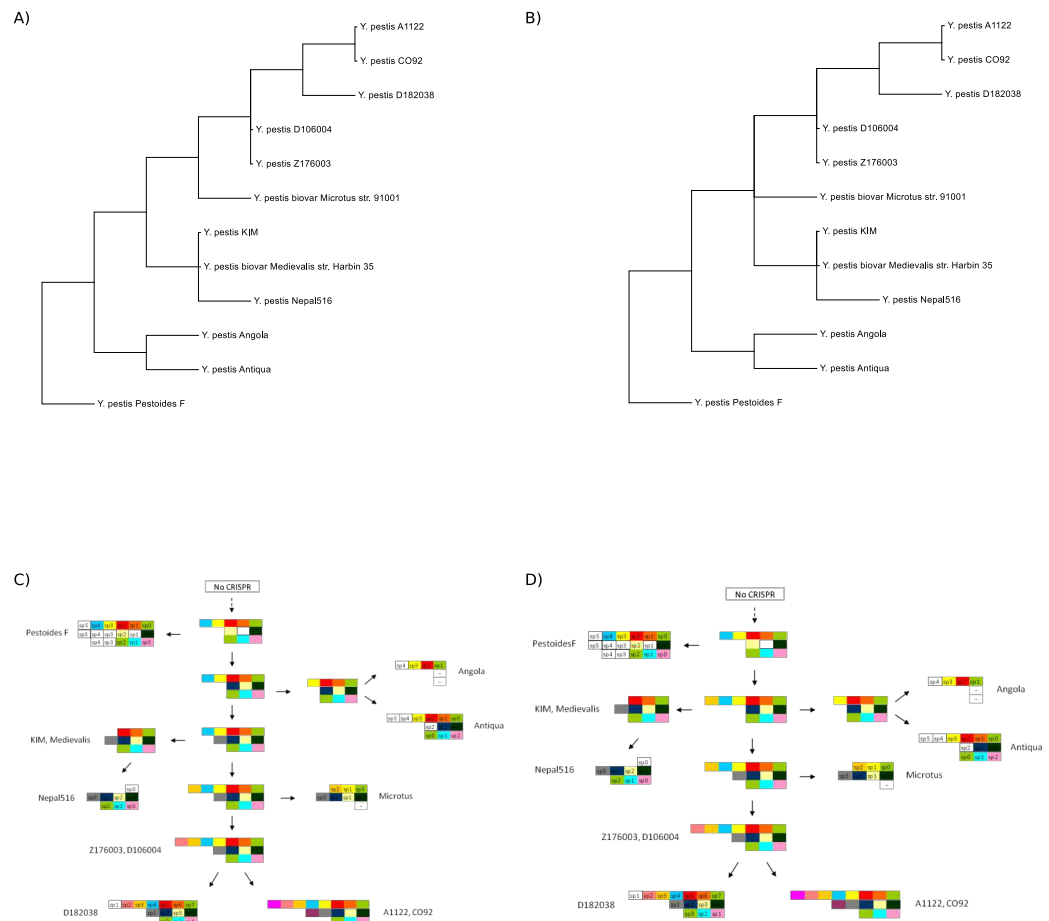


Figure 7. Cladograms (A, B) and schemas of evolution (C, D) of two integrated CRISPR-based maximum parsimony phylogenetic trees most compatible with the sequence tree