

Lawrence Berkeley National Laboratory

Lawrence Berkeley National Laboratory

Title

Complete Genome Sequence of Yersinia pestis Strains Antiqua and Nepal516: Evidence of Gene Reduction in an Emerging Pathogen

Permalink

<https://escholarship.org/uc/item/2ww0r6m3>

Authors

Chain, Patrick S.G.
Hu, Ping
Malfatti, Stephanie A.
et al.

Publication Date

2006-01-16

Peer reviewed

Title: Complete genome sequence of *Yersinia pestis* strains Antiqua and Nepal516: evidence of gene reduction in an emerging pathogen.

Running Title: *Yersinia pestis* genome sequences

Patrick S. G. Chain^{1,2§}, Ping Hu^{3§}, Stephanie A. Malfatti^{1,2}, Lyndsay Radnedge^{1¥}, Frank Larimer^{2,4}, Lisa M. Vergez^{1,2}, Patricia Worsham⁵, May C. Chu⁶, and Gary L. Andersen^{3*}

¹Biosciences Directorate, Lawrence Livermore National Laboratory, Livermore, CA 94550

²Joint Genome Institute, Walnut Creek, CA

³Center for Environmental Biotechnology, Lawrence Berkeley National Laboratory, Berkeley, CA 94720

⁴Life Sciences Division, Oak Ridge National Laboratory, Oak Ridge, TN 37831

⁵United States Army Research Institute of Infectious Diseases, Fort Detrick, MD 21702

⁶Centers for Disease Control and Prevention, Fort Collins, CO 80522

[§]These authors contributed equally to this work.

*Corresponding author.

Mailing address: Lawrence Berkeley National Laboratory, 1 Cyclotron Road, mailstop

70A3317, Berkeley, CA94720, USA

Phone: 519-495-2795

FAX: 510-486-7152

E-mail: GLAndersen@lbl.gov

[¥]Present address: Monogram Biosciences Inc., 345 Oyster Point Boulevard, South San Francisco, CA 94080

1 **ABSTRACT**

2 *Yersinia pestis*, the causative agent of bubonic and pneumonic plague, has undergone
3 detailed study at the molecular level. To further investigate the genomic diversity among this
4 group and to help characterize lineages of the plague organism that have no sequenced members,
5 we present here the genomes of two isolates of the “classical” Antiqua biovar, strains Antiqua
6 and Nepal516. The genomes of Antiqua and Nepal516 are 4.7 Mb and 4.5 Mb and encode 4,138
7 and 3,956 open reading frames respectively. Though both strains belong to one of the three
8 classical biovars, they represent separate lineages defined by recent phylogenetic studies. We
9 compare all five currently sequenced *Y. pestis* genomes and the corresponding features in *Y.*
10 *pseudotuberculosis*. There are strain-specific rearrangements, insertions, deletions, single
11 nucleotide polymorphisms and a unique distribution of insertion sequences. We found 453 single
12 nucleotide polymorphisms in protein coding regions, which were used to assess evolutionary
13 relationships of these *Y. pestis* strains. Gene reduction analysis revealed that the gene deletion
14 processes are under selective pressure and many of the inactivations are probably related to the
15 organism’s interaction with its host environment. The results presented here clearly demonstrate
16 the differences between the two Antiqua lineages and support the notion that grouping *Y. pestis*
17 strains based strictly on the classical definition of biovars (predicated upon two biochemical
18 assays) does not accurately reflect the phylogenetic relationships within this species. Comparison
19 of four virulent *Y. pestis* strains with the human-avirulent strain 91001 provides further insight
20 into the genetic basis of virulence to humans.

21

22

1 **INTRODUCTION**

2 Plague is a zoonotic disease, endemic throughout the world and is highly infectious in
3 humans. The causative agent, *Yersinia pestis*, primarily infects a wide range of rodents and is
4 transmitted via flea vectors. Throughout history, plague has ravaged human populations in three
5 major pandemic waves: Justinian’s plague (541-767 AD), which started in Africa and spread to
6 Mediterranean, the Black death of 1346 to early 19th century may have originated in central Asia
7 and spread from the Caspian Sea to Europe, and modern plague (since 1894), which began in
8 southwest China and spread globally via marine shipping routes from Hong Kong. Although
9 human disease is rare, *Y. pestis* is dangerous, highly infectious and thus has been identified as
10 having potential for use in bioterrorism or as a biological weapon.

11 It was shown that *Y. pestis* recently diverged from *Y. pseudotuberculosis* – an
12 enteropathogen, and likely comprises a clonal lineage (1, 3, 37, 40). *Y. pestis* strains have
13 historically been classified according to their ability to utilize glycerol and reduce nitrate and
14 have been grouped into three main subtypes, or biovars: Antiqua, Mediaevalis and Orientalis.
15 Isolates from the Orientalis biovar have a worldwide distribution due to spread by steamship
16 beginning 100 years ago. In contrast, isolates of Antiqua and Mediaevalis biovars are generally
17 limited to localized regions containing long-term plague foci from enzootic rodent hosts in
18 Africa and Central Asia. It has been argued that each of the biovars was associated with one of
19 the plague pandemics (14, 20, 34) and recent studies have tried to provide direct evidence
20 whether *Y. pestis* was associated with any of the historical pandemics (15, 44). DNA sequences
21 from ancient human remains dispute this assertion that different biovars were responsible for
22 each of the last three pandemics and suggest that instead, Orientalis-like *Y. pestis* may have been
23 involved in all three (15). This suggestion remains highly controversial.

1 Isolates from the biovar Antiqua have been thought to represent a more ancestral branch
2 of the plague pathogen, primarily due to their association with long-established plague foci as
3 well as sharing an additional set of genetic regions with *Y. pseudotuberculosis* and “non-
4 classical” (e.g. *Microtus* biovar) subspecies of *Y. pestis*. Our previous work using suppression
5 subtractive hybridization demonstrated a pattern of difference fragments (DFR profiles),
6 including a 15,603 bp segment of chromosomal DNA that was shared by *Y. pseudotuberculosis*
7 and a portion of both the “non-classical” subspecies of *Y. pestis* and the “classical” biovar
8 Antiqua (38). There are currently three completed genome sequences for *Y. pestis*, one each from
9 the Orientalis, the Mediaevalis, and the “non-classical” *Microtus* biovars. To get a better
10 understanding of the detailed genetic changes in a pathogen that is adapting to an intracellular
11 lifestyle, we have sequenced two isolates from the classical “Antiqua” biovar. Strain Antiqua is
12 fully virulent and possesses a DFR profile closest to *Y. pseudotuberculosis*. A *pgm*⁻ derivative of
13 the virulent strain Nepal516 was found to have a different DFR profile and is believed to
14 represent a different lineage of this biovar. Comparison with the genome sequence of the
15 previously sequenced *Y. pestis* strains as well as *Y. pseudotuberculosis* gave further insight into
16 the loci required for the adaptation to an intracellular pathogenic lifestyle. Additional insight into
17 the acquisition of virulence to humans was obtained by the comparison to the human-avirulent
18 isolate, 91001.

1 **MATERIALS AND METHODS**

2 **Bacterial strains**

3 *Y. pestis* Nepal516 was isolated from a human infection in Nepal (possibly from a 1967
4 outbreak of pneumonic plague), while strain Antiqua was isolated from a human infection in
5 Africa (1965, Republic of Congo). Both have been biochemically characterized to belong to the
6 Antiqua biovar and carry the three previously described “virulence” plasmids found in most
7 classical isolates of *Y. pestis*. Both strains have been used previously in a variety of studies (4,
8 21, 27, 36, 37, 45). The wild type Antiqua strain and a *pgm*⁻ version of the Nepal516 strain were
9 available and used in this genome sequencing project. The Nepal516 strain lacks the ~100 kb
10 *pgm* region, including the high pathogenicity island, the pesticin/yersiniabactin complex, and the
11 haemin storage locus, that are normally located between two parallel *IS100* insertion elements
12 (5, 8, 18, 22, 29, 35, 42).

13 **Construction, sequencing, and assembly**

14 Genomic DNA was isolated from *Y. pestis* strains Antiqua and Nepal516.
15 The two genomes were sequenced using the whole-genome shotgun method as previously
16 described (9). Briefly, 3kb- and 8kb-sized, randomly sheared DNA fragments were isolated and
17 cloned into pUC18 and pMCL200 respectively, for amplification in *Escherichia coli*. A larger
18 fosmid library was constructed containing approximately 40kb inserts of sheared genomic DNA
19 cloned into the pCC1Fos cloning vector. Double-ended plasmid sequencing reactions were
20 performed from all three libraries at the Department of Energy Joint Genome Institute using ABI
21 3730xl DNA Analyzers and MegaBACE 4500 Genetic Analyzers as described on the JGI
22 website <http://www.jgi.doe.gov/>.

1 Approximately 110,556 and 113,541 of sequences were assembled for Antiqua and
2 Nepal516, respectively, producing an average of 11 fold coverage across the genomes.
3 Processing of sequence traces, base calling and assessment of data quality were performed with
4 PHRED and PHRAP (P. Green, University of Washington, Seattle), respectively. Assembled
5 sequences were visualized with CONSED. The initial assemblies consisted of 154 and 113
6 contigs (≥ 20 reads per contig). Gaps in the sequence were primarily closed by resolving the
7 many repetitive regions found within the genome. The remaining gaps were closed by primer
8 walking on gap-spanning library clones or PCR products from genomic DNA. True physical
9 gaps were closed by combinatorial (multiplex) PCR. Sequence finishing and polishing added
10 roughly 300 reads and assessment of final assembly quality was completed as described (9).

11 For the genome of Nepal516, the ~70 kb pCD plasmid was underrepresented and was not
12 completed as part of the sequencing project. Nepal516 is known to contain the pCD plasmid
13 (Scott Bearden, personal communication). The existence of the pCD plasmid was verified by
14 PCR in our laboratory (data not shown). We believe that the failure of obtaining sufficient
15 quantity of pCD DNA for sequencing is due to particular laboratory conditions and has no
16 biological implications on the sequences of the chromosome, pMT and pPCP plasmids.

17 **Sequence analysis and annotation**

18 Automated gene modeling was completed by combining results from Critica, Generation,
19 and Glimmer modeling packages, and comparing the translations to GenBank's non-redundant
20 (NR) database using basic local alignment search tool for proteins (BLASTP). The protein set
21 was also searched against KEGG Genes, InterPro, TIGRFams, PROSITE, and Clusters of
22 Orthologous Groups of proteins (COGS) databases to further assess function. Manual functional
23 assignments were assessed on individual gene-by-gene basis as needed. Sequence alignment and

1 protein domain search tools (BLAST, CLUSTALW, Pfam, etc) were applied in various stages of
2 comparison. CO92 gene nomenclature is used in this work when possible, other nomenclature is
3 mentioned and used when no CO92 ortholog is available.

4 **Single nucleotide polymorphism (SNP) analysis**

5 *Yersinia* genomes are known to harbor extensive rearrangements as well as a large
6 number of insertion sequence elements (IS) and other duplicated regions (10, 13, 33, 41). These
7 repeats and insertion elements were excluded from consideration in SNP analysis. Genome-wide
8 SNP discovery was achieved by whole genome alignments using the software package
9 Mummer3 (28) and by subsequent orthologous gene alignments. For coding regions, pair-wise
10 reciprocal BLASTP analyses were performed with the five sets of *Y. pestis* proteins. An ortholog
11 pair was defined as reciprocal best top hits using a cutoff of 95% sequence identity. If an
12 ortholog was not found in any one of the five genomes, the proteins were removed from further
13 analysis. The sequences of the orthologous genes were used to find SNPs using Mummer3.
14 Whole genome comparisons were also done using Mummer3. SNPs were selected from regions
15 not covered by the ortholog alignment method described above. Synonymous and non-
16 synonymous sites were calculated as follows: for every position in the genome, it was assessed
17 whether it was located in an intergenic or a coding region; if it is in a coding region (excluding
18 coding regions from insertion elements and other repetitive elements) and the nucleotides
19 substitution results in no change in amino acid sequence, it was classified as a potential
20 synonymous SNP site, otherwise it was regarded as a potential non-synonymous site.

21 **Comparative analysis of gene deletions in *Y. pestis* genomes**

22 We analyzed the loss of function patterns in all *Y. pestis* genomes, focusing on the
23 presence and absence of protein functions. Some deletions, such as *tufB* deletion in Nepal516

1 (described below), were not included because of gene duplication. Complete datasets of proteins
2 for each *Y. pestis* genome were downloaded from published reports or from the final annotations
3 of the newly sequenced genomes. Transposases and enzymes related to insertion elements were
4 removed. The final protein datasets for deletion analysis were 3723, 3909, 3896, 3769, 3777 and
5 3867 proteins for strains CO92, KIM, Antiqua, Nepal516, 91001 and *Y. pseudotuberculosis*
6 IP32953, respectively. Pairwise alignments of proteins of all five *Y. pestis* genomes were
7 accomplished by BLASTP. A protein function was deemed absent if there was no top hit greater
8 than 95% identity or at least 75% of the query sequence. This analysis focused on the presence or
9 absence of protein and functional representation, therefore, if the protein has a closely related
10 paralog or is duplicated in the genome, it was considered present. Due to the differences in
11 annotation (particularly with smaller gene calls), we applied a cutoff criteria to remove all small
12 proteins since these were more frequently found to be differentially annotated across the *Y.*
13 *pestis* genomes. With a size filter of 75 amino acids, we are certain to have missed a small
14 number of real proteins smaller than 75 amino acids, such as the 61 amino acid carbon storage
15 regulator *crsA* (YPO3304). The final set of proteins found to be absent in at least one genome
16 was manually inspected with the aid of multiple sequence alignment tools CLUSTALW and
17 nucleotide sequence alignment tool BLASTN. If the deletion was comprised solely of repetitive
18 units, the protein was removed from this analysis, because the deletion mechanism may be
19 different in those cases and may revert frequently. Additionally, these final sets of proteins were
20 inspected in multiple genome alignments to distinguish annotation differences vs. true
21 differences in the genomes. A similar set of criteria was employed to see if homologs of these
22 proteins exist in the *Y. pseudotuberculosis* IP32953 genome.

23 **Nucleotide sequence accession number and locus tag prefixes**

1 The annotated sequences of the complete genomes of *Y. pestis* strains Antiqua and
2 Nepal516 are available at GenBank/EMBL/DDBJ (accession numbers pending). The prefixes
3 YPA and YPN are used for locus tags (gene identifier prefixes) in strains Antiqua and Nepal516,
4 respectively. When referring to specific genes throughout the text, we use the CO92 gene
5 numbers (prefixed with YPO) where possible, unless the gene does not exist in CO92 (or if it is
6 clearer to use a different prefix), then the locus tags for a different genome are used and
7 mentioned.

1 RESULTS

2 Genome overviews

3 The genomes of *Y. pestis* strains Antiqua and Nepal516 each consist of a single circular
4 chromosome and the three virulence plasmids, pMT, pCD and pPCP, which are associated with
5 most classical *Y. pestis* strains. Here we report all replicons but the pCD plasmid of Nepal516
6 (see Materials and Methods). The salient genomic features of each genome are detailed in Table
7 1, while gross chromosomal comparisons with those of the CO92, KIM and 91001 strains are
8 summarized in Figure 1. Although the global characteristics of the five genomes are quite
9 similar, a number of strain-specific insertions, deletions, rearrangements, and single nucleotide
10 polymorphisms (SNPs) were identified, along with a unique distribution of insertion sequence
11 (IS) elements (see Tables 2, 3, 4 and Supplemental data).

12 Strain-specific synonymous SNPs (sSNPs) and non-synonymous SNPs (nsSNPs)

13 The numbers of sSNPs and nsSNPs specific to one or to two genomes are shown in
14 Figure 2. There are 57 sSNPs (135 nsSNPs) specific to strain 91001 compared with all other *Y.*
15 *pestis* strains. While 27 of these sSNPs (49 nsSNPs) are shared with the ancestral *Y.*
16 *pseudotuberculosis* IP32953, the remaining 30 sSNPs (and 86 nsSNPs) differ with respect to *Y.*
17 *pseudotuberculosis* IP32953, indicating that they likely arose in 91001 since its lineage diverged
18 from the remaining *Y. pestis* sequenced isolates. Likewise, the 27 sSNPs (and 49 nsSNPs)
19 specific to 91001 (and identical to *Y. pseudotuberculosis*) are mutations predicted to have arisen
20 in the other *Y. pestis* lineage which gave rise to the remaining sequenced strains (Figure 2). No
21 SNPs (sSNPs or nsSNPs) are found to be specific to strain pairs Antiqua and KIM, Antiqua and
22 Nepal516, CO92 and KIM, or CO92 and Nepal516. However, 4 sSNPs (and 6 nsSNPs) are
23 found specifically in CO92 and Antiqua (i.e. CO92 and Antiqua share the same SNP state, while

1 all other *Y. pestis* strains have a different SNP state), and 6 sSNPs (and 11 nsSNPs) are specific
2 to KIM and Nepal516. Taken together, these data suggest a separation of these four strains into
3 two distinct branches: where Antiqua and CO92 belong to one branch and KIM and Nepal516
4 occupy the other (Figure 2). These 2 sets of sSNP and nsSNP mutations have accumulated in the
5 short period of time after the KIM/Nepal516 and CO92/Antiqua lineages diverged but before
6 each lineage further split into two (Figure 2). Thus, this analysis strongly supports the notion that
7 although strains Antiqua and Nepal516 are grouped into the same biovar (Antiqua), they
8 represent distinct lineages.

9 The genes harboring sSNPs (cumulatively for all the *Y. pestis* strains) can be distributed
10 into functional gene categories (based on Clusters of Orthologous Groups) as shown in Figure 3.
11 Non-synonymous SNPs are found distributed in 20 COG categories (Figure 3a), while
12 synonymous SNPs were distributed in 19 COGs (Figure 3b). We investigated whether SNPs
13 belonging to the various branches in Figure 2 are biased towards any functional categories, but
14 no unique or biased distribution patterns were found. All strain-specific mutations share
15 approximately 1/3 of the COG categories. Although some strain-specific SNPs are found in
16 functional categories not represented in any other strain (sSNP: three categories are unique to
17 91001 and one category is unique to Antiqua; nsSNP: two categories are unique to 91001)
18 (Figure 3a, 3b), the number of SNPs are too small to determine whether these were random
19 events. No sSNP vs nsSNP bias was readily apparent (the average nsSNP/sSNP ratio is
20 approximately 2.9) except for the large proportion of nsSNPs vs sSNPs in the cell
21 wall/membrane biogenesis COG of 91001, with a ratio of 9:2.

22 The gene sequences and gene organization in the plasmids are highly conserved. There is
23 only one synonymous SNP in the plasminogen activator protease and 5 SNPs (including small

1 deletions) in non-coding regions of all pPCP plasmids. There are 8 and 17 SNPs in the predicted
2 gene sets of all plasmids pCD and pMT, respectively. The majority of these SNPs are in
3 hypothetical proteins and are distributed more or less randomly across strains.

4 **Insertion Sequence (IS) elements and genome rearrangements:**

5 As determined by several groups already (10, 13, 33, 46), IS elements have expanded
6 tremendously in *Y. pestis* since its divergence from *Y. pseudotuberculosis*, and have served as the
7 delimiters for recombination events that have led to genomic deletions and genome
8 rearrangements. Due to its presumed continued transposition activity in the wild, *IS100* elements
9 have been successfully used for typing and grouping strains (32). We thus performed a detailed
10 analysis of the precise locations of *IS100* and the other three major IS elements (*IS1541*, *IS285*,
11 *IS1661*) as well as investigated their contribution to the observed rearrangements.

12 Each sequenced strain of *Y. pestis* has a unique set of IS elements (Table 2) and a core IS
13 set that is shared among them. It was previously observed that *Y. pestis* shares a set of 12 IS
14 elements with *Y. pseudotuberculosis* (3 of each major IS element) that appear to have been
15 acquired by *Y. pseudotuberculosis* before the evolution of *Y. pestis* based on identical locations
16 of insertion (10). In addition to this set of 12, the core set of IS elements shared among all five
17 sequenced *Y. pestis* strains is 45 *IS1541*, 15 *IS100*, 11 *IS285* and 6 *IS1661*, indicative of
18 elements present in the last common ancestor of all these sequenced strains. Several more (2
19 *IS1541*, 1 *IS100*, 2 *IS285*, 2 *IS1661*) are predicted to be shared among all 5 *Y. pestis* strains, but
20 in at least one strain, these IS elements have been subsequently involved in a deletion event
21 between two IS copies (leaving only one behind), or have been lost as part of a larger deletion.
22 Similarly, the four classical *Y. pestis* strains (Antiqua, Nepal516, KIM and CO92) also share a
23 subset of the remaining IS elements, distinct from non-classical strain 91001: 5 *IS1541*, 5 *IS100*,

1 3 IS285 with four additional IS100 elements that have likely been lost in a strain-specific manner
2 via deletion as described above. In addition, there are 4 IS elements shared between CO92 and
3 Antiqua, and 3 shared between KIM and Nepal516 (Table 2), supporting the SNP-based
4 phylogeny in Figure 2 as well as that depicted by Achtman and colleagues (2). Interestingly, a
5 small number of IS elements were shared by unexpected groups of strains that disagree with the
6 proposed phylogeny: one IS100 shared between 91001, Antiqua and CO92, one IS100 shared
7 between Antiqua and KIM; one IS100 in Antiqua, CO92 and KIM; one IS1541 in CO92 and in
8 *Y. pseudotuberculosis* IP32053; one IS285 in CO92, KIM and Nepal but not Antiqua; and one
9 IS1541 in Antiqua and in *Y. pseudotuberculosis* IP32053.

10 While IS1541 was the most active IS element between the divergence of *Y.*
11 *pseudotuberculosis* and the most recent common ancestor of the five *Y. pestis* strains, with 45
12 common IS1541 insertions among all sequenced isolates, IS100 has been the most active in more
13 recent times, though not equally among the strains (Table 2). This suggests that insertion
14 sequence activity may be punctuated and does not occur at a constant rate across different
15 strains. With the exception of Nepal516, the number of unique IS100 is significantly more than
16 for any other IS element (Table 2), though the reason for this is unclear.

17 Despite their extensive sequence similarity, the *Y. pestis* genomes appear to be in a state
18 of flux with respect to large genome rearrangements. Similar to previous observations, all breaks
19 in colinearity between the *Yersinia* genomes occurs at IS elements or other repeated sequences.
20 Differences in the GC skew patterns in *Y. pestis* genomes, including the many breaks observed in
21 Antiqua (Figure 1) are also the result of rearrangements between IS elements as previously
22 observed (10, 33, 46). Similarly, IS elements have played a large role in deletion events observed
23 in the *Y. pestis* genomes. For example, strains KIM and CO92 have undergone overlapping

1 IS1541-mediated deletions of 32 kb (YP0966-YP0994, using 91001 nomenclature) and 21.5 kb
2 (YP0966-YP0986, using 91001 nomenclature), respectively. This deletion accounted for the
3 DFR patterns observed by suppression subtractive hybridization (37). A similar strain-specific,
4 IS-mediated deletion of a phage region is described further below. Additionally, a IS1661-
5 mediated 13 kb deletion in both KIM and Nepal516 has removed a large cluster of flagellar
6 genes (YPO0738- YPO0754) including the flagellar RNA polymerase sigma factor and
7 chemotaxis membrane proteins.

8 **Functional Reduction**

9 It has been postulated that the genomes of *Y. pestis* have undergone functional reduction
10 as it made a transition from an oral-fecal pathogen causing gastroenteritis, to a vector-borne
11 pathogen causing a fatal, invasive, septicemic disease, where genes have been inactivated by
12 various mechanisms, such as deletions/insertions, frameshifts, interruptions by insertion
13 elements, and homologous or even non-homologous recombination. We determined all the genes
14 and associated functions that have been lost since *Y. pestis* emerged from *Y. pseudotuberculosis*
15 (by comparing *Y. pestis* strains and identifying functional losses) to better understand *Y. pestis*
16 diversity and evolution. We have identified a large number of strain-specific gene
17 inactivations/deletions as well as some that are specific to only two of the five genomes (Table
18 3). Other than hypothetical proteins, there is one dominant category of proteins in these strain- or
19 lineage-specific deletions: proteins contributing to the interaction of bacterium with its
20 environment or host, including membrane proteins, membrane receptors, ABC transporters,
21 flagellar proteins and chemotaxis proteins.

22 The number of strain-specific lost functions was not equal among all strains. In part, this
23 can be explained by several large deletions that effectively delete many genes (functions) in a

1 single event. While strain 91001 and Antiqua had the greatest number of strain-specific function
2 losses with 69 and 41 respectively, it is interesting to note that CO92 was found to have the
3 fewest compared to other strains. Though strains Antiqua and 91001 share several gene
4 inactivations (8 proteins), these are the result of independent mutations (homoplasmy), while those
5 lost in CO92 and Antiqua (12 proteins), as well as those lost in KIM and Nepal516 (16 proteins)
6 share lineage specific, inherited mutations (function losses) that support the phylogenetic tree.

7 Of the 69 functional losses specific to 91001, 24 are hypothetical protein, seven are
8 membrane proteins, seven are phage-related proteins, five are regulatory proteins and three are
9 transporters. Some of these 91001-specific losses of function may be related to its human-
10 avirulent phenotype. Twenty-one of the 69 belong to a single IS285-mediated deletion event
11 (YPO2108 – YPO2134).

12 Ten of the 41 Antiqua-specific lost functions were the results of two deletion events.
13 Several inactivated or deleted genes were predicted to be directly or indirectly involved in
14 interactions with environment. For example, a glutathione S transferase, YPO2367, is missing
15 from the Antiqua genome. This family of enzymes routinely responds to oxidative stress or
16 detoxification, which can be encountered during entry into phagocytic cells (12). Since bacteria
17 usually have multiple glutathione S transferases (CO92 has at least 4 based on the annotation),
18 losing one may not have a distinct phenotype. Interestingly, there were several cases where
19 similar functional inactivations/deletions were observed in two strains that affect different genes
20 that may have overlapping functions. For example, a potassium efflux pump (YPO3129) was
21 inactivated in Antiqua and a Na^+/H^+ antiporter (YPO2142) inactivated in KIM. The role of the
22 Na^+/H^+ antiporter is sodium extrusion (24) and both the Na^+/H^+ antiporter and the potassium
23 efflux are involved in pH homeostasis. Since both Antiqua and KIM are geographically limited

1 to localized regions, one possible explanation for the differential inactivations is local adaptation
2 to selectively maintain one of the two similar functions. An alternate possibility is that neither
3 gene is required, that the differences observed in these sequenced strains were random events,
4 and that both genes are not required and are on their way to being lost from the *Y. pestis* gene
5 pool. Whether the similarity of the deletion profiles mentioned above reflects adaptations to their
6 environmental niches or convergent evolution remains to be investigated.

7 Of the 20 KIM-specific inactivated genes, seven are concentrated in one deletion, and the
8 inactivation via nonsense mutation of YPO3038 (NapA, a periplasmic nitrate reductase) is
9 proposed to be one of the causes of the nitrate negative phenotype of the Medievalis biovar.
10 Nepal516 has 13 specific lost functions, including two transporters (YPO2835 and YPO1350),
11 the chromosomally-encoded type III secretion system protein YPO0266, and two classical
12 virulence factors (YPO2291 – a putative virulence factor and YPO0599 – a hemolysin/adhesin
13 mentioned further below). Some of the deletions in Nepal516 have been previously demonstrated
14 experimentally by microarray hybridizations (21). There are only four lost functions
15 (pseudogenes YPO1087 and YPO3679, along with y1377 and y2928 using KIM nomenclature)
16 that are CO92-specific losses, yet all are putative proteins without clear functional predictions.

17 No genetic regions were identified in the genomes of Antiqua or Nepal516 that were not
18 present in at least one of the previously sequenced *Y. pestis* or *Y. pseudotuberculosis* strains.
19 Though there were a number of Antiqua and Nepal516 genes or domains of genes that were
20 found to differ significantly from the other strains, many of these differences consist of varying
21 numbers of degenerate tandem repeat elements within the coding sequences of surface proteins,
22 such as in the invasin YPO3944 described further below, and were not interpreted as losses of
23 function. Surprisingly, our analysis revealed only one example of a strain-specific unique genetic

1 region with no similar DNA sequence in any of the other four sequenced strains or in the *Y.*
2 *pseudotuberculosis* genome. This is a ssDNA prophage found inserted in CO92 (YPO2271 –
3 YPO2281). This region has been found by PCR in all tested biovar Orientalis strains as well as a
4 few African strains of the biovar Antiqua (10, 19).

5 **Differences in putative virulence factors**

6 Most characterized and putative classical virulence factors are identical throughout all
7 five *Y. pestis* strains, including those found on the virulence plasmids, such as the pPCP-located
8 plasminogen activator Pla required for successful subcutaneous infection (11), and the pMT-
9 encoded murine toxin (Ymt) (23) and F1 capsular protein (16) (important for *Y. pestis* life cycle
10 and vector-mammal transmission). Similarly, loci on the chromosome are also nearly identical
11 between strains, such as the RTX-like toxin gene YPO0947, the attachment invasion locus *ail*
12 (YPO2905) and two *ail*-like genes (YPO1860, YPO2190) are virtually identical among *Y. pestis*
13 strains and with *Y. pseudotuberculosis* as well. Interestingly, a fourth *ail*-like gene (YPO2506)
14 has been deleted from the Antiqua genome.

15 Other loci, known to differ between *Y. pestis* and *Y. pseudotuberculosis*, were also
16 investigated. The *Y. pestis* invasin YPO1793 is interrupted by an IS1541 in all strains, while
17 putative adhesin YPO1562 interrupted by an IS285 in all *Y. pestis* strains except for 91001,
18 which harbors a nonsense mutation instead. Both are intact in *Y. pseudotuberculosis*. Similarly,
19 RTX transporter YPO2250 and the TccC-family insecticidal toxin YPO2312 are frameshift
20 pseudogenes in all *Y. pestis* strains but appear intact in *Y. pseudotuberculosis*.

21 A second TccC insecticidal toxin homolog, YPO2380 has been deleted only in *Y. pestis*
22 KIM. Two additional TccC toxins are found in tandem in all *Y. pestis* strains (YPO3674,
23 YPO3673), while only a single copy is found in *Y. pseudotuberculosis*. In *Y. pestis*, a second

1 family of insecticidal toxins is found upstream of these two TccC homologs and consists of a
2 complex of three genes (YPO3681, YPO3679, YPO3678). While these are present and highly
3 similar in amino acid sequence in *Y. pseudotuberculosis*, YPO3681 is inactivated by a frameshift
4 only in Antiqua, and YPO3679 harbors a frameshift in CO92 only.

5 One *Y. pestis* hemolysin/adhesin, YPO0599, located within a possible pathogenicity
6 island (YPO0641a-YPO0590) is different from that of *Y. pseudotuberculosis* at the C-terminus.
7 Several additional CDSs can be seen downstream in *Y. pseudotuberculosis* that appear to be
8 “modules”, or adhesin fragments that share high similarity to portions of the C-terminus of this
9 CDS. Different “modules” are found downstream of the *Y. pestis* gene. Interestingly Nepal516 is
10 missing a large section of the C-terminal portion of this protein, likely due to a recombination
11 between one of these modules and the corresponding section in the Nepal516 gene. The
12 remainder of this pathogenicity island is highly similar between *Y. pestis* strains and *Y.*
13 *pseudotuberculosis*. A similar module-recombination scenario is envisioned to have resulted in a
14 modified C-terminus of hemolysin/adhesin YPO2490 in 91001 compared to the other *Y. pestis*
15 strains and *Y. pseudotuberculosis*. A different mechanism, the expansion or contraction of
16 degenerate repeat units within the putative invasin YPO3944, has resulted in different sized
17 invasins in *Y. pseudotuberculosis* (5623 aa - amino acids), and the various *Y. pestis* strains
18 (91001, 3108 aa; Nepal516, 4270 aa; KIM, CO92 and Antiqua, 3013 aa). Further study is
19 required to understand any phenotypic effect these two classes of differences may have.

20 **Loss of functional TufB in Nepal**

21 The genomes of *Y. pestis* and of several other organisms have two copies of the
22 elongation factor Tu (EF-Tu). Due to its highly conserved function and ubiquitous distribution,
23 elongation factors are considered a valuable phylogenetic marker and have been used in

1 evolution studies of *Enterococci* (26), *Lactobacillus* (43) and other eubacteria (39). Interestingly,
2 in *Y. pestis*, the two copies of this gene are not as conserved as the two genes of *Escherichia coli*,
3 yet the conservation within each copy (among *Y. pestis* strains) is maintained. In *E. coli*, *tufA* and
4 *tufB* gene products only differ by a single amino acid (7) and exhibit identical physical, chemical
5 and catalytic properties (17). We found that all five *Y. pestis* genomes and *Y. pseudotuberculosis*
6 have two copies of *tuf* genes (*tufA*, *tufB*) and the general operon structure is similar to that of *E.*
7 *coli*, however, the sequence identity between the *Y. pestis* *tufA* and *tufB* is considerably lower.
8 There are 17 amino acid differences (4%) and a total of 138 nucleotide changes (11.7%) in
9 addition to a large C-terminal deletion in *tufB* of Nepal516. This result is unexpected, since
10 previous studies show that duplicate *tuf* gene within a genome differ on average by 0.7% in
11 nucleotide sequence (30). Whether the two copies of the *tuf* genes in *Y. pestis* have different
12 origins requires further investigation. All six *Yersinia* TufA proteins are 100% identical and all
13 *Y. pestis* *tufA* genes have identical nucleotide sequences while *Y. pseudotuberculosis* has a single
14 synonymous SNP (G to A). Although the *tufB* gene products differ from *tufA*, a similar
15 conservation is evident. We have found however that the TufB of Nepal516 harbors a large C-
16 terminal deletion which affects 57 aa (or 67%) of the GTP-EFTu-D3 domain, involved in
17 binding charged tRNAs and EF-Ts(6). This deletion is likely to cause loss of functionality and
18 thus we believe TufB is not functional in Nepal. It is not known whether the two copies are
19 expressed under different conditions or have slightly different functions or kinetic properties,
20 however, this deletion leads us infer that *Y. pestis* requires only one functional copy for its life
21 cycle.

22 **Comparison to the human-avirulent strain 91001**

1 A previous comparison between strain 91001 and human virulent strains CO92 and KIM
2 revealed a number of differences, including a 33-kb prophage-like sequence (YPO2096 –
3 YPO2135) that was absent in 91001 but intact in both CO92 and KIM, and suggested that this
4 difference may have contributed to this strain’s lack of virulence in human (41). This entire
5 region was also absent in *Y. pseudotuberculosis* (10) consistent with this claim. Our analysis
6 shows that this phage-like region is intact in Nepal516, but is partially deleted in Antiqua (Table
7 4). While the deletion in 91001 can be attributed to recombination between two parallel *IS100*
8 elements, the smaller deletion in Antiqua (from YPO2087 – YPO2106) is likely due to excision
9 between two *IS285* sequences (Supplemental data). Since Antiqua is a fully virulent strain, the
10 region deleted in Antiqua would not seem to contribute to human pathogenicity, however the
11 remaining portion of the prophage region deleted from 91001 may indeed contain genes that are
12 important for human virulence.

13 Previous comparisons with CO92 and KIM also revealed a list of 91001-specific
14 pseudogenes that may be related to *Y. pestis* pathogenicity and host range (41). Our gene
15 reduction analysis included two additional virulent strains and confirmed the presence or absence
16 of orthologs in *Y. pseudotuberculosis* based on our cutoffs (refer to Materials and Methods).
17 While all of these 91001-specific pseudogenes were also intact in Antiqua (Supplemental data),
18 only one (YPO2258) was also found to be inactivated in Nepal, suggesting this gene has no
19 impact on the avirulent property of 91001 in humans. Among the 91001-specific pseudogenes,
20 there are only four that are also absent in *Y. pseudotuberculosis*, YPO0733 (flagellar hook-
21 associated protein), YPO0737 (flagellin), YPO0962 (hypothetical protein) and YPO3110
22 (putative O-unit flippase). In addition, the nsSNP mutations that contribute to changes in single
23 amino acids in many proteins may affect 91001 or other *Y. pestis* strain virulence in subtle ways.

1 For example, nsSNPs were found in some genes with a possible role in virulence (e.g. the
2 Antiqua *ail* gene YPO2905 carries a nsSNP, as does the 91001 RTX toxin gene YPO0947), but
3 the significance of these substitutions is not known.

4 **DISCUSSION**

5 This work presents the complete genome sequences for the two previously unsequenced
6 *Y. pestis* major lineages (both designated Antiqua using classical nomenclature). Phylogenetic
7 relationships were elucidated clearly with the distribution of synonymous SNPs (Figure 2). Since
8 synonymous mutations do not affect protein functions (unlike nsSNPs or some IS elements),
9 their accumulation is not under selective pressure, making this the least biased method for
10 inferring evolutionary relationships. The distribution of sSNPs convincingly demonstrates that a
11 single biovar Antiqua is an inaccurate phylogenetic representation supporting previous claims
12 that categorize Antiqua strains into two groups (2, 10). Using terminology proposed previously
13 (2), lineage 1.ANT (African strain Antiqua) is closely related to Orientalis strain CO92 while
14 2.ANT (Asian strain Nepal516) is more closely related to Medievalis strain KIM. These four
15 “classical” isolates fall on a branch separate from the non-classical, human-avirulent Chinese
16 strain 91001. This analysis also revealed a relatively rapid divergence of the four distinctive
17 lineages from two ancestral lines for the “classical” *Y. pestis* strains. Although it is only possible
18 to make very crude estimations of age of descent for these four lineages, the numbers of sSNPs
19 are consistent with all of the lineages being present within the last 1,500 years of the 3 great
20 pandemics (calculation not shown).

21 Comparison of all five *Y. pestis* sequences reveals extensive DNA sequence
22 rearrangement, widespread gene reduction and strain-specific IS elements, as well as SNPs. It
23 was previously reported that *Y. pestis* strains differ greatly in genome synteny, and that repeated

1 sequences most often were found at the borders of rearrangements (10, 13). Indeed, most
2 rearrangements occur at IS elements and regardless of which genomes were chosen for two-way
3 comparisons, we identified similar numbers of rearrangements to those previously observed
4 between *Y. pestis* strains (13) and even between *Y. pestis* and *Y. pseudotuberculosis* (10) (data
5 not shown). The question remains whether these observed rearrangements have any effect on
6 transcription or whether this has an overall destabilizing influence on the genome.

7 The distribution pattern of IS elements in the sequenced strains generally supports the
8 SNP-derived phylogeny with several IS elements shared across all “classical” strains but not in
9 91001 as well as IS elements found only in the CO92/Antiqua or KIM/Nepal516 pairs of strains.
10 Only a few (5 in total) IS elements found to be shared by two or more strains did not conform to
11 the predicted phylogeny (footnotes in Table 2); similar observations have been reported
12 previously (2). Our analyses suggest that a small number of IS elements may have been precisely
13 excised from their insertion locations, that identical insertion events have occurred in 2 different
14 strains/lineages, or that there may be some limited horizontal transfer between *Y. pestis* strains
15 that have resulted in mobilizing IS elements from one strain to a different strain/lineage (or
16 alternatively, removing an IS element by introduction of wild type sequence). One example is an
17 aminotransferase (YPO3250) that is disrupted by an *IS100* in all sequenced *Y. pestis* strains
18 except Nepal516, which instead has the wild type gene and no trace of an *IS100*. These data also
19 suggest that certain IS elements may not be useful for typing or grouping strains and may explain
20 certain discrepancies in phylogenetic groupings using different methods.

21 Interestingly, the entire complement of *IS1541* (and almost all *IS1661*) elements in strain
22 91001 was acquired by the ancestor of all *Y. pestis* strains. In contrast, since 91001 diverged
23 from the other strains, it has acquired a number of strain-specific *IS100* and *IS285* elements,

1 supporting the idea of actively integrating IS elements within the genome of *Y. pestis*. With the
2 exception of Nepal516, *IS100* appears to have been more active (greater number of new
3 transposition events) than other IS elements, but the reason for this is unknown.

4 Functional reduction analysis also generally agrees with the SNP-based phylogenetic tree
5 (Figure 3) as well as with a more limited study that identified the loss of gene regions across a
6 panel of *Y. pestis* isolates using a CO92 gene-specific microarray (21). Similar to the IS and SNP
7 data, the four “classical” strains appear to share an evolutionary path distinct from strain 91001
8 based on functional reduction, and the KIM/Nepal516 and CO92/Antiqua pairs also exhibit a
9 larger number of shared function loss. The exceptions are the result of independent mutations:
10 two shared losses between KIM and Antiqua; one shared loss between Nepal516 and 91001;
11 eight shared losses between Antiqua and 91001 and 16 shared losses between CO92 and KIM
12 (Supplemental data and Table 3). The two shared function losses between KIM and Antiqua are
13 a putative siderophore biosynthetic enzyme and a putative membrane protein. The predicted
14 functions suggest that both of the proteins could be involved in interactions with the
15 environment, therefore these losses may reflect adaptations to the *Y. pestis* microenvironment.
16 Similarly, the single functional loss shared between Nepal516 and 91001 is the arabinose operon
17 regulatory protein. Although the observed shared loss of function between Antiqua and 91001
18 contained several genes, they are exclusively in the prophage region described above and it is the
19 result of independent deletion events. The shared losses between CO92 and KIM were possibly
20 from a single deletion event.

21 Strain 91001 has the highest number of strain-specific losses of function with a total of
22 69. Interestingly, all but four of the 91001-specific pseudogenes have homologs with >90%
23 identity in *Y. pseudotuberculosis*, suggesting that 91001 lost those genes while other virulent *Y.*

1 *pestis* strains retain them. It is possible that these genes may be involved in human-virulence
2 and/or fitness in the human host. Some inactivated proteins may be related to pathogenicity, such
3 as hemolysin (YPO2045), sulfatase and sulfatase modifier protein (YPO3046 and YPO3047),
4 UDP-glycosyltransferase (YPO1985) and O-unit flippase-like protein (YPO3110) (Supplemental
5 data). Hemolysin is a toxin that forms transmembrane channels and is involved in heme
6 utilization and adhesion. The precise function of the sulfatase operon (YPO3046 and YPO3047)
7 in *Y. pestis* is not known, however these enzymes belong to a family of proteins that hydrolyze
8 various sulfate esters or catalyze sulphur insertions. In mammalian cells, the oligosaccharide
9 moieties on glycoproteins, glycolipids and proteoglycans are frequently modified with sulfate.
10 Sulfatase from pathogenic bacteria have been shown to interact with mucin (47) and a previous
11 study suggested that mucin-sulphatase activity in *Burkholderia cepacia* and *Pseudomonas*
12 *aeruginosa* may contribute to their association with airway infection in cystic fibrosis patients by
13 possibly facilitating bacterial colonization (25). Thus, the deletion of the sulfatase and sulfatase
14 modifier protein in strain 91001 may have contributed to its human-avirulent phenotype. Finally,
15 the O-unit flippase is involved in translocating polysaccharide unit across the membrane while
16 UDP-glycosyltransferase (YPO1985) is typically involved in O-antigen biosynthesis. Since *Y.*
17 *pestis* is known to lack O-antigen, the actual functions of YPO3110 and YPO1985 may not
18 directly involve O-antigen but perhaps other surface polysaccharides.

19 Antiqua also had a high number of strain-specific losses, even after discounting the
20 deletion events which involved several genes (41 and 31, respectively). Interestingly, we found a
21 correlation with the observed higher *IS100* transposition activity in Antiqua, with 13 of the 31
22 inactivations due to *IS100* interruptions. The profile of Antiqua-specific loss of function contains
23 a significant amount of proteins which interact with environment, such as glutathione S-

1 reductase, chemotaxis protein, porin C protein, potassium efflux pump, insecticidal toxin,
2 flagellar motor switch protein and 6 membrane proteins without specific known functions. A
3 possible explanation for this may be that the genome has been adapting to the niche the Antiqua
4 organism occupies.

5 Discounting those genes lost in a single deletion event, the numbers of KIM-specific (14)
6 and Nepal-specific (13) functional loss are similar. Surprisingly, only three CO92-specific losses
7 of function were identified. It is possible that there was a selective advantage for the Orientalis
8 biovar to maintain a greater repertoire of genes and to maintain flexibility and be able to adapt
9 quickly to new host(s). The world-wide distribution of this group and the small number of CO92-
10 specific putative gene inactivations is consistent with this theory. A 31 amino acid deletion in
11 YPO3937 (473 amino acids) confers the glycerol negative phenotype of biovar Orientalis (33),
12 however since the deletion is below length cutoff threshold, it was not included in our study as a
13 loss of function. Unique to strain CO92 are a hypothetical protein (YPO2469), a hemolysin
14 activator protein (YPO3720) and a prophage that do not exist or have been inactivated in the
15 other sequenced *Y. pestis* strains or *Y. pseudotuberculosis*. These genes may again have been
16 retained by CO92 to maintain its ability to interact with a more variable environment.

17 Unexpectedly we found that Nepal516 has many exceptions compared with the other
18 sequenced *Y. pestis* strains, including the apparent loss of function of TufB, the much smaller
19 number of Nepal516-specific SNPs compared to that of the other strains (Figure 2), and the fact
20 that IS100 has not been as active in Nepal516 as in the other strains (Table 2). Since both
21 nsSNPs and sSNPs are equally affected, it is unlikely that this is due to selective pressure which
22 should have a neutral effect on sSNPs, but rather that the mutation rate is responsible, suggesting
23 the rate of mutation or evolution is slower in Nepal. The reason for this is not known, however a

1 possible explanation may be that this phenomenon is driven by fewer rounds of bacterial division
2 with a relatively cooler local environment and hibernation of the host(s) that fostered fewer
3 opportunities for transmission.

4 Despite the observed differences between different strains of *Y. pestis*, the sequenced
5 genomes reveal a highly conserved chromosomal backbone, reminiscent of what is observed in
6 *Bacillus anthracis* (31). Within the five genomes of *Y. pestis* compared here, a single region
7 present in strain CO92 was found to be unique (not shared with another *Y. pestis* genome),
8 though independent studies have shown that this region which encodes phage genes is present in
9 most, if not all 1.ORI strains as well as some 1.ANT strains (10, 19). We thus believe that most
10 of the genomic sequence shared among the “classical” *Y. pestis* isolates is represented within this
11 data set, though other sequences of non-classical isolates may harbor novel genomic regions not
12 revealed in these analyses.

13 **Conclusions**

14 The two completed genomes presented here, from the previously unrepresented Antiqua
15 biovar, have provided important references for SNP discovery, for the study of insertion element
16 distribution, genome rearrangement, and reductive evolution in *Y. pestis*. Comparisons of the
17 four virulent “classical” strains to the human-avirulent strain 91001 have also provided further
18 insight into *Y. pestis* human virulence. With sSNPs as the preferred method for elucidating
19 phylogenetic relationships, strains Nepal516 and Antiqua were convincingly placed in two
20 clearly separate branches, with one branch shared by strains KIM (Mediaevalis) and Nepal516,
21 and the other shared by strains CO92 (Orientalis) and Antiqua. While IS element distributions
22 and function loss across the strains generally agreed with such a phylogenetic representation,
23 certain exceptions were found and are thought to be the result of lack of selective pressure in the

1 *Y. pestis* strains inhabited niche, of possible horizontal gene exchange between *Y. pestis* strains,
2 or of homoplasy in the reductive processes. Though there is some evidence of convergent
3 evolution, whether this is the primary mechanism underlying the observed discrepancies remains
4 to be investigated. The *Y. pestis* genome is a clear example of one actively undergoing reductive
5 evolution, as its lifestyle has altered from an enteropathogen to an intracellular pathogen. The
6 genome has slowly accumulated inactivations and deletions that result in loss of function, which,
7 for the virulent strains (all strains except 91001), have little effect on pathogenicity. The
8 differences between these strains and the human-avirulent 91001 provide an ideal starting point
9 for future experiments to elucidate the mechanisms involved in *Yersinia* pathogenicity.

10

11 **ACKNOWLEDGEMENTS**

12 We thank Matt Van Ert, Ryan Easterday, Aubree Hinckley, and Maria Shin for technical
13 assistance. This work was performed under the auspices of the US Department of Energy by the
14 University of California, Lawrence Berkeley National Laboratory under Contract NO. DE-
15 AC02-05CH11231. This work was supported by an Intelligence Technology Innovation Center
16 grant.

REFERENCES

1. **Achtman, M.** 2004. Age, Descent and Genetic Diversity within *Yersinia pestis*, p. 432. In E. Carniel and B. J. Hinnebusch (ed.), *Yersinia: Molecular and Cellular Biology*. Horizon Bioscience, Norwich, UK.
2. **Achtman, M., G. Morelli, P. Zhu, T. Wirth, I. Diehl, B. Kusecek, A. J. Vogler, D. M. Wagner, C. J. Allender, W. R. Easterday, V. Chenal-Francisque, P. Worsham, N. R. Thomson, J. Parkhill, L. E. Lindler, E. Carniel, and P. Keim.** 2004. Microevolution and history of the plague bacillus, *Yersinia pestis*. *PNAS* **101**:17837-17842.
3. **Achtman, M., K. Zurth, G. Morelli, G. Torrea, A. Guiyoule, and E. Carniel.** 1999. *Yersinia pestis*, the cause of plague, is a recently emerged clone of *Yersinia pseudotuberculosis*. *PNAS* **96**:14043-14048.
4. **Adair, D. M., P. L. Worsham, K. K. Hill, A. M. Klevytska, P. J. Jackson, A. M. Friedlander, and P. Keim.** 2000. Diversity in a Variable-Number Tandem Repeat from *Yersinia pestis*. *J. Clin. Microbiol.* **38**:1516-1519.
5. **Bearden, S. W., and R. D. Perry.** 1999. The Yfe system of *Yersinia pestis* transports iron and manganese and is required for full virulence of plague. *Molecular Microbiology* **32**:403-414.
6. **Berchtold, H., L. Reshetnikova, C. O. A. Reiser, N. K. Schirmer, M. Sprinzl, and R. Hilgenfeld.** 1993. Crystal structure of active elongation factor Tu reveals major domain rearrangements. *Nature* **365**:126-132.
7. **Blattner, F. R., G. Plunkett, III, C. A. Bloch, N. T. Perna, V. Burland, M. Riley, J. Collado-Vides, J. D. Glasner, C. K. Rode, G. F. Mayhew, J. Gregor, N. W. Davis, H.**

- A. Kirkpatrick, M. A. Goeden, D. J. Rose, B. Mau, and Y. Shao.** 1997. The Complete Genome Sequence of *Escherichia coli* K-12. *Science* **277**:1453-1462.
8. **Buchrieser, C., C. Rusniok, L. Frangeul, E. Couve, A. Billault, F. Kunst, E. Carniel, and P. Glaser.** 1999. The 102-Kilobase *pgm* Locus of *Yersinia pestis*: Sequence Analysis and Comparison of Selected Regions among Different *Yersinia pestis* and *Yersinia pseudotuberculosis* Strains. *Infect. Immun.* **67**:4851-4861.
9. **Chain, P., J. Lamerdin, F. Larimer, W. Regala, V. Lao, M. Land, L. Hauser, A. Hooper, M. Klotz, J. Norton, L. Sayavedra-Soto, D. Arciero, N. Hommes, M. Whittaker, and D. Arp.** 2003. Complete Genome Sequence of the Ammonia-Oxidizing Bacterium and Obligate Chemolithoautotroph *Nitrosomonas europaea*. *J. Bacteriol.* **185**:2759-2773.
10. **Chain, P. S. G., E. Carniel, F. W. Larimer, J. Lamerdin, P. O. Stoutland, W. M. Regala, A. M. Georgescu, L. M. Vergez, M. L. Land, V. L. Motin, R. R. Brubaker, J. Fowler, J. Hinnebusch, M. Marceau, C. Medigue, M. Simonet, V. Chenal-Francisque, B. Souza, D. Dacheux, J. M. Elliott, A. Derbise, L. J. Hauser, and E. Garcia.** 2004. Insights into the evolution of *Yersinia pestis* through whole-genome comparison with *Yersinia pseudotuberculosis*. *PNAS* **101**:13826-13831.
11. **Cowan, C., H. A. Jones, Y. H. Kaya, R. D. Perry, and S. C. Straley.** 2000. Invasion of Epithelial Cells by *Yersinia pestis*: Evidence for a *Y. pestis*-Specific Invasin. *Infect. Immun.* **68**:4523-4530.
12. **De Groote, M. A., U. A. Ochsner, M. U. Shiloh, C. Nathan, J. M. McCord, M. C. Dinauer, S. J. Libby, A. Vazquez-Torres, Y. Xu, and F. C. Fang.** 1997. Periplasmic

- superoxide dismutase protects *Salmonella* from products of phagocyte NADPH-oxidase and nitric oxide synthase. PNAS **94**:13997-14001.
13. **Deng, W., V. Burland, G. Plunkett III, A. Boutin, G. F. Mayhew, P. Liss, N. T. Perna, D. J. Rose, B. Mau, S. Zhou, D. C. Schwartz, J. D. Fetherston, L. E. Lindler, R. R. Brubaker, G. V. Plano, S. C. Straley, K. A. McDonough, M. L. Nilles, J. S. Matson, F. R. Blattner, and R. D. Perry.** 2002. Genome Sequence of *Yersinia pestis* KIM. J. Bacteriol. **184**:4601-4611.
 14. **Devignat, R.** 1951. Varietes de l'espece *Pasteurella pestis*: nouvelle hypothese. Bull. W. H. O **4**:247-263.
 15. **Drancourt, M., V. Roux, L. V. Dang, L. Tran-Hung, D. Castex, V. Chenal-Francisque, H. Ogata, P. Fournier, E. Crubézy, and D. Raoult.** 2004. Genotyping, Orientalis-like *Yersinia pestis*, and Plague Pandemics. Emerg Infect Dis. **10**:1585-1592.
 16. **Du, Y., R. Rosqvist, and A. Forsberg.** 2002. Role of Fraction 1 Antigen of *Yersinia pestis* in Inhibition of Phagocytosis. Infect. Immun. **70**:1453-1460.
 17. **Furano, A. V.** 1977. The elongation factor Tu coded by the *tufA* gene of *Escherichia coli* K-12 is almost identical to that coded by the *tufB* gene. J. Biol. Chem. **252**:2154-2157.
 18. **Gehring, A. M., E. DeMoll, J. D. Fetherston, I. Mori, G. Mayhew, F. R. Blattner, C. T. Walsh, and R. D. Perry.** 1998. Iron acquisition in plague: modular logic in enzymatic biogenesis of yersiniabactin by *Yersinia pestis*. Chem Biol. **5**:573-586.
 19. **Gonzalez, M. D., C. A. Lichtensteiger, R. Caughlan, and E. R. Vimr.** 2002. Conserved Filamentous Prophage in *Escherichia coli* O18:K1:H7 and *Yersinia pestis* Biovar orientalis. J. Bacteriol. **184**:6050-6055.

20. **Guiyoule, A., F. Grimont, I. Itean, P. Grimont, M. Lefevre, and E. Carniel.** 1994. Plague pandemics investigated by ribotyping of *Yersinia pestis* strains. *J. Clin. Microbiol.* **32**:634-641.
21. **Hinchliffe, S. J., K. E. Isherwood, R. A. Stabler, M. B. Prentice, A. Rakin, R. A. Nichols, P. C. F. Oyston, J. Hinds, R. W. Titball, and B. W. Wren.** 2003. Application of DNA Microarrays to Study the Evolutionary Genomics of *Yersinia pestis* and *Yersinia pseudotuberculosis*. *Genome Res.* **13**:2018-2029.
22. **Hinnebusch, B. J., R. D. Perry, and T. G. Schwan.** 1996. Role of the *Yersinia pestis* Hemin Storage (hms) Locus in the Transmission of Plague by Fleas. *Science* **273**:367-370.
23. **Hinnebusch, B. J., A. E. Rudolph, P. Cherepanov, J. E. Dixon, T. G. Schwan, and A. Forsberg.** 2002. Role of *Yersinia* Murine Toxin in Survival of *Yersinia pestis* in the Midgut of the Flea Vector. *Science*.1069972. *Science* **296**:733-735.
24. **Ito, M., A. Guffanti, J. Zemsky, D. Ivey, and T. Krulwich.** 1997. Role of the nhaC-encoded Na⁺/H⁺ antiporter of alkaliphilic *Bacillus firmus* OF4. *J. Bacteriol.* **179**:3851-3857.
25. **Jansen, H. J., C. A. Hart, J. M. Rhodes, J. R. Saunders, and J. W. Smalley.** 1999. A novel mucin-sulphatase activity found in *Burkholderia cepacia* and *Pseudomonas aeruginosa*. *J Med Microbiol.* **48**:551-557.
26. **Ke, D., M. Boissinot, A. Huletsky, F. J. Picard, J. Frenette, M. Ouellette, P. H. Roy, and M. G. Bergeron.** 2000. Evidence for Horizontal Gene Transfer in Evolution of Elongation Factor Tu in *Enterococci*. *J. Bacteriol.* **182**:6913-6920.

27. **Klevytska, A. M., L. B. Price, J. M. Schupp, P. L. Worsham, J. Wong, and P. Keim.** 2001. Identification and Characterization of Variable-Number Tandem Repeats in the *Yersinia pestis* Genome. *J. Clin. Microbiol.* **39**:3179-3185.
28. **Kurtz, S., A. Phillippy, A. Delcher, M. Smoot, M. Shumway, C. Antonescu, and S. Salzberg.** 2004. Versatile and open software for comparing large genomes. *Genome Biology* **5**:R12.
29. **Kutyrev, V. V., A. A. Filippov, O. S. Oparina, and O. A. Protsenko.** 1992. Analysis of *Yersinia pestis* chromosomal determinants Pgm super(+) and Pst super(s) associated with virulence. *Microbial Pathogenesis* **12**:177-186.
30. **Lathe, I., Warren C., and P. Bork.** 2001. Evolution of tuf genes: ancient duplication, differential loss and gene conversion. *FEBS Letters* **502**:113-116.
31. **Medini, D., C. Donati, H. Tettelin, V. Masignani, and R. Rappuoli.** 2005. The microbial pan-genome. *Current Opinion in Genetics & Development.* **15**:589-594.
32. **Motin, V. L., A. M. Georgescu, J. M. Elliott, P. Hu, P. L. Worsham, L. L. Ott, T. R. Slezak, B. A. Sokhansanj, W. M. Regala, R. R. Brubaker, and E. Garcia.** 2002. Genetic Variability of *Yersinia pestis* Isolates as Predicted by PCR-Based IS100 Genotyping and Analysis of Structural Genes Encoding Glycerol-3-Phosphate Dehydrogenase (glpD). *J. Bacteriol.* **184**:1019-1027.
33. **Parkhill, J., B. W. Wren, N. R. Thomson, R. W. Titball, M. T. G. Holden, M. B. Prentice, M. Sebahia, K. D. James, C. Churcher, K. L. Mungall, S. Baker, D. Basham, S. D. Bentley, K. Brooks, A. M. Cerdeno-Tarraga, T. Chillingworth, A. Cronin, R. M. Davies, P. Davis, G. Dougan, T. Feltwell, N. Hamlin, S. Holroyd, K. Jagels, A. V. Karlyshev, S. Leather, S. Moule, P. C. F. Oyston, M. Quail, K.**

- Rutherford, M. Simmonds, J. Skelton, K. Stevens, S. Whitehead, and B. G. Barrell.** 2001. Genome sequence of *Yersinia pestis*, the causative agent of plague. *413*:523-527.
34. **Perry, R., and J. Fetherston.** 1997. *Yersinia pestis*--etiologic agent of plague. *Clin. Microbiol. Rev.* **10**:35-66.
35. **Perry, R. D., A. G. Bobrov, O. Kirillina, H. A. Jones, L. Pedersen, J. Abney, and J. D. Fetherston.** 2004. Temperature Regulation of the Hemin Storage (Hms+) Phenotype of *Yersinia pestis* Is Posttranscriptional. *J. Bacteriol.* **186**:1638-1647.
36. **Prentice, M. B., K. D. James, J. Parkhill, S. G. Baker, K. Stevens, M. N. Simmonds, K. L. Mungall, C. Churcher, P. C. F. Oyston, R. W. Titball, B. W. Wren, J. Wain, D. Pickard, T. T. Hien, J. J. Farrar, and G. Dougan.** 2001. *Yersinia pestis* pFra Shows Biovar-Specific Differences and Recent Common Ancestry with a *Salmonella enterica* Serovar Typhi Plasmid. *J. Bacteriol.* **183**:2586-2594.
37. **Radnedge, L., P. G. Agron, P. L. Worsham, and G. L. Andersen.** 2002. Genome plasticity in *Yersinia pestis*. *Microbiology* **148**:1687-1698.
38. **Radnedge, L., S. Gamez-Chin, P. M. McCready, P. L. Worsham, and G. L. Andersen.** 2001. Identification of Nucleotide Sequences for the Specific and Rapid Detection of *Yersinia pestis*. *Appl. Environ. Microbiol.* **67**:3759-3762.
39. **Sela, S., D. Yogev, S. Razin, and H. Bercovier.** 1989. Duplication of the *tuf* gene: a new insight into the phylogeny of eubacteria. *J. Bacteriol.* **171**:581-584.
40. **Skurnik, M.** 2003. The genus *Yersinia*: entering the functional genomic era (advances in experimental medicine and biology. Plenum US.
41. **Song, Y., Z. Tong, J. Wang, L. Wang, Z. Guo, Y. Han, J. Zhang, D. Pei, D. Zhou, H. Qin, X. Pang, Y. Han, J. Zhai, M. Li, B. Cui, Z. Qi, L. Jin, R. Dai, F. Chen, S. Li, C.**

- Ye, Z. Du, W. Lin, J. Wang, J. Yu, H. Yang, J. Wang, P. Huang, and R. Yang.** 2004. Complete genome sequence of *Yersinia pestis* strain 91001, an isolate avirulent to humans. *DNA Res.* **11**:179-197.
42. **Une, T., and R. R. Brubaker.** 1984. In vivo comparison of avirulent Vwa- and Pgm- or Pstr phenotypes of *yersiniae*. *Infect Immun.* **43**:895-900.
43. **Ventura, M., C. Canchaya, V. Meylan, T. R. Klaenhammer, and R. Zink.** 2003. Analysis, Characterization, and Loci of the *tuf* Genes in *Lactobacillus* and *Bifidobacterium* Species and Their Direct Application for Species Identification. *Appl. Environ. Microbiol.* **69**:6908-6922.
44. **Wiechmann, I., and G. Grupe.** 2005. Detection of *Yersinia pestis* DNA in two early medieval skeletal finds from Aschheim (Upper Bavaria, 6th century A.D.). *American Journal of Physical Anthropology* **126**:48-55.
45. **Wilmoth, B., M. Chu, and T. Quan.** 1996. Identification of *Yersinia pestis* by BBL Crystal Enteric/Nonfermenter Identification System. *J. Clin. Microbiol.* **34**:2829-2830.
46. **Wren, B. W.** 2003. The *yersiniae*--a model genus to study the rapid evolution of bacterial pathogens. *Nat Rev Microbiol.* **1**:55-64.
47. **Wright, D. P., C. G. Knight, S. G. Parkar, D. L. Christie, and A. M. Robertson.** 2000. Cloning of a Mucin-Desulfating Sulfatase Gene from *Prevotella* Strain RS2 and Its Expression Using a *Bacteroides* Recombinant System. *J. Bacteriol.* **182**:3002-3007.

Table 1. General genome features for *Yersinia pestis* strains Antiqua and Nepal516

	Antiqua	Nepal516
Chromosome Size (bp)	4,702,289	4,534,590
G+C content (%)	47.70	47.58
Coding sequences	4138	3956
Average gene length (bp)	953	958
Coding Density (%)	83.8	83.6
16S-23S-5S rRNAs	7	7
Transfer RNAs	68	72
pMT Size (bp)	96,471	100,918
G+C content	50.24	50.16
Coding sequences	99	104
Average gene length (bp)	832	820
Coding Density (%)	85.3	84.5
pCD Size (bp)	70,299	-*
G+C content	44.83	
Coding sequences	89	
Average gene length (bp)	601	
Coding Density (%)	76.1	
pPCP Size (bp)	10,777	10,778
G+C content	45.44	45.44
Coding sequences	9	9
Average gene length (bp)	573	573
Coding Density (%)	47.8	47.8

*pCD of Nepal516 was not completed in this study, see Materials and Methods

Table 2. Chromosome comparison between the sequenced *Y. pestis* strains

	Antiqua	Nepal516	91001	KIM	CO92
Molecular grouping (1)	1.ANT	2.ANT	0.PE4	2.MED	1.ORI
Size (bp)	4,702,289	4,534,595	4,595,065	4,600,755	4,653,728
Total IS elements					
<i>IS100</i>	75	32	30	34	44
<i>IS285</i>	24	25	23	19	21
<i>IS1541</i>	67	64	47	55	65
<i>IS1661</i>	10	8	8	8	8
Unique IS elements*					
<i>IS100</i>	39	4	15	6	13
<i>IS285</i>	2	4	11	1	1
<i>IS1541</i>	5	4	0	0	8
<i>IS1661</i>	1	0	0	0	0

*Some IS elements were shared between expected partners such as three *IS100* and one *IS285* shared between Antiqua and CO92, as well as one *IS100*, one *IS1661* and one *IS1541* shared between Nepal516 and KIM. However, five exceptions were observed: one *IS100* shared between 91001, Antiqua and CO92 but not present in Nepal or KIM; one *IS100* in 91001, Antiqua, CO92 and KIM but not in Nepal; one *IS285* in 91001, CO92, KIM and Nepal but not in Antiqua; one *IS1541* in CO92 and in *Y. pseudotuberculosis* IP32053; one *IS1541* in Antiqua and in *Y. pseudotuberculosis* IP32053

Table 3. Genome specific inactivation of genes

Deletion specific to the genome(s)	Number of proteins inactivated
CO92	4
KIM	20
Antiqua	41
Nepal516	13
91001	69
CO92, KIM	0
CO92, Antiqua	11
CO92, Nepal516	0
CO92, 91001	0
KIM, Antiqua	2
KIM, Nepal516	16
KIM, 91001	0
Antiqua, Nepal516	0
Antiqua, 91001	8
Nepal516, 91001	1

* For rows with 2 strains, the data indicate inactivations in the same CDS of both strains

Table 4. Prophage-like fragment specific to virulent *Y. pestis* strains

GENE	<i>Y. pestis</i>				<i>Y. pseudo-</i> <i>tuberculosis</i>	COG	Product
	CO92	KIM	ANTIQUA	NEPAL	IP32953		
YPO2095	+	+	-	+	-	-	hypothetical protein
YPO2096	+	*+	-	*+	-	-	hypothetical protein
YPO2097	+	*+	-	*+	-	-	hypothetical protein
YPO2098	+	+	-	+	-	R	putative phage lysozyme
YPO2099	+	+	-	+	-	-	putative prophage endopeptidase
YPO2100	+	+	-	+	-	S	phage regulatory protein
YPO2101	+	+	-	+	-	-	hypothetical protein
YPO2102	+	+	-	+	-	-	hypothetical protein
YPO2104	+	+	-	+	+	L	transposase for the IS285 insertion element
YPO2108	+	+	+	+	-	-	hypothetical protein
YPO2109	+	+	+	+	-	-	hypothetical protein
YPO2110	+	+	+	+	-	-	hypothetical protein
YPO2111	+	+	+	+	-	-	hypothetical protein
YPO2112	+	+	+	+	-	-	hypothetical protein
YPO2113	+	+	+	+	-	-	hypothetical protein
YPO2114	+	+	+	+	-	-	hypothetical protein
YPO2115	+	+	+	+	-	-	hypothetical protein
YPO2116	+	+	+	+	-	-	hypothetical protein
YPO2117	+	+	+	+	-	-	hypothetical protein
YPO2118	+	*+	*+	*+	-	-	hypothetical protein
YPO2119	+	+	+	+	-	S	putative phage tail protein
YPO2120	+	+	+	+	-	S	hypothetical protein
YPO2122	+	+	+	+	-	S	hypothetical

YPO2123	+	+	+	+	-	R	protein putative phage minor tail protein
YPO2124	+	+	+	+	-	-	hypothetical protein
YPO2125	+	+	+	+	-	-	putative phage regulatory protein
YPO2126	+	+	+	+	-	K	hypothetical protein
YPO2127	+	+	+	+	-	-	putative phage- related membrane protein
YPO2128	+	*+	*+	*+	-	-	putative phage- related lipoprotein
YPO2129	+	+	+	+	-	S	putative phage tail assembly protein
YPO2130	+	*+	*+	*+	-	-	hypothetical protein
YPO2131	+	+	+	+	-	S	putative phage host specificity protein
YPO2132	+	+	+	+	-	-	hypothetical protein
YPO2133	+	+	+	+	-	-	hypothetical protein
YPO2134	+	+	+	+	-	-	putative phage tail fiber assembly protein
YPO2135	+	*+	+	+	-	-	hypothetical protein
YPO2487	+	*+	*+	*+	*+	-	putative membrane protein
YPO2488	+	+	+	+	+	-	hypothetical protein
YPO2489	+	+	+	+	+	S	hypothetical protein

* The protein sequence was not found in the genome, however, the DNA fragment did exist in the intergenic region.

Figure 1. Circular representation of the strain Antiqua (A) and strain Nepal516 (B) chromosomes. The different rings represent (from outer to inner): 1 and 2, all genes color-coded by functional category; 3 and 4, IS elements (*IS100*, *IS285*, *IS1541*, *IS1661*); 5, deviation from average G+C content; 6, GC skew.

Figure 2. Phylogenetic ordering of *Yersinia* by SNP analysis. The number of sSNPs and the number of nsSNPs (in parentheses symbols) are illustrated at the corresponding positions.

Figure 3. Functional distribution of genes bearing SNPs. The number of genome-specific nsSNPs (Fig. 3a) and sSNPs (Fig. 3b) were grouped into COG functional classes. These were sub-categorized based on what genome(s) they were found in (light orange: 91001; red: CO92; dark red: KIM; blue: Antiqua; light blue: Nepal516; green: KIM and Nepal516; yellow: Antiqua and CO92). COG functional classes: C, energy production; D, cell division; E, amino acid metabolism; F, nucleotide metabolism; G, carbohydrate metabolism; H, coenzyme metabolism; I, lipid metabolism; J, translation; K, transcription; L, DNA replication or repair; M, cell wall/membrane biogenesis; N, cell motility; O, posttranslational modification; P, inorganic ion metabolism; Q, secondary metabolites biosynthesis, transport and catabolism ; R, general function prediction only; S, function unknown; T, signal transduction; U, intracellular trafficking and secretion; V, defense mechanism

Figure 1

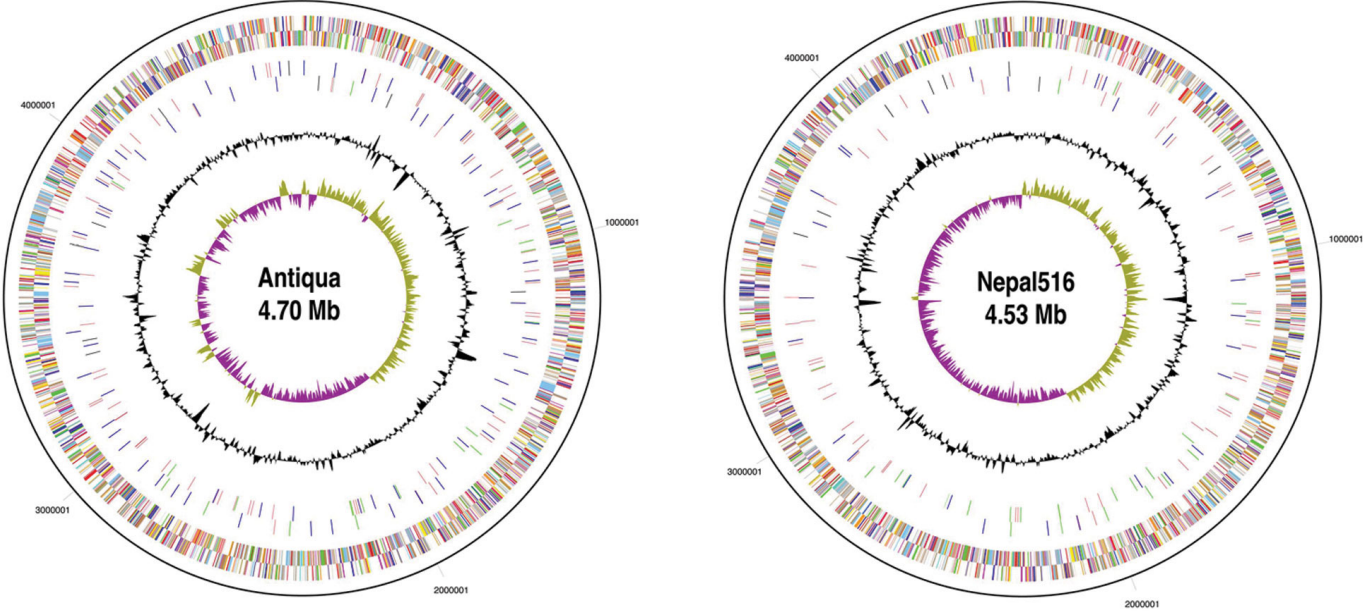


Figure 2

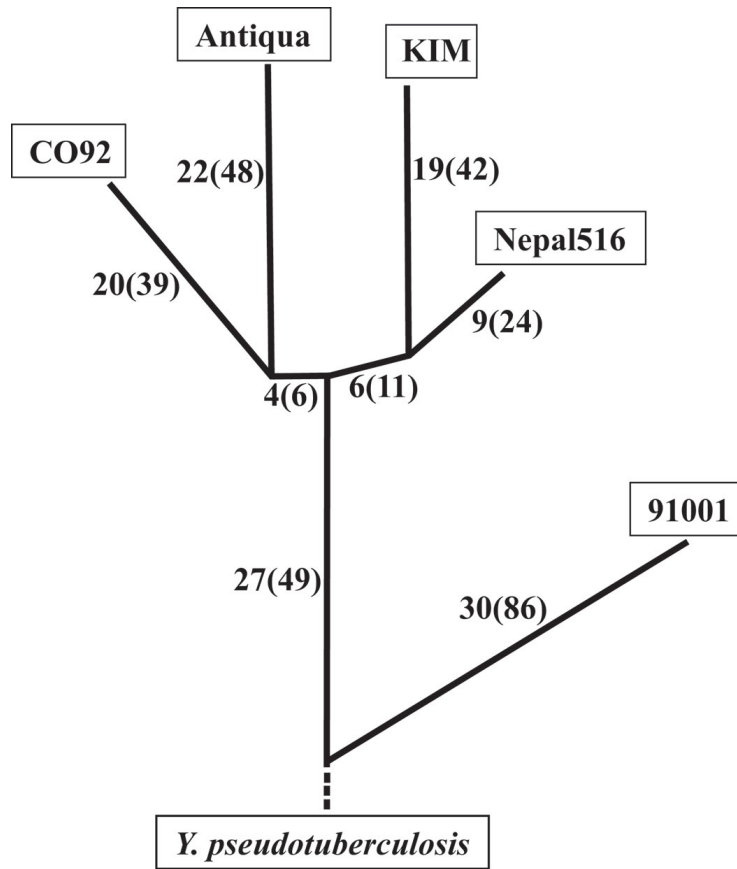


Figure 3.

(a)

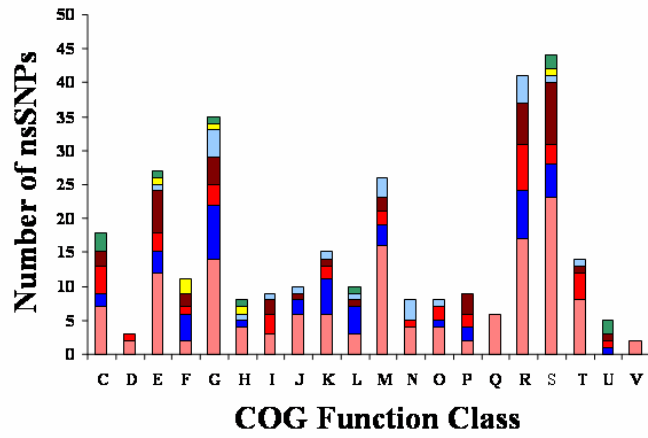


Figure 3

(b)

