


RESEARCH ARTICLE

Open Access



Genome rearrangements and selection in multi-chromosome bacteria *Burkholderia* spp.

Olga O. Bochkareva^{1,2*} , Elena V. Moroz^{1†}, Iakov I. Davydov^{3,4†} and Mikhail S. Gelfand^{1,2,5}

Abstract

Background: The genus *Burkholderia* consists of species that occupy remarkably diverse ecological niches. Its best known members are important pathogens, *B. mallei* and *B. pseudomallei*, which cause glanders and melioidosis, respectively. *Burkholderia* genomes are unusual due to their multichromosomal organization, generally comprised of 2-3 chromosomes.

Results: We performed integrated genomic analysis of 127 *Burkholderia* strains. The pan-genome is open with the saturation to be reached between 86,000 and 88,000 genes. The reconstructed rearrangements indicate a strong avoidance of intra-replicon inversions that is likely caused by selection against the transfer of large groups of genes between the leading and the lagging strands. Translocated genes also tend to retain their position in the leading or the lagging strand, and this selection is stronger for large syntenies. Integrated reconstruction of chromosome rearrangements in the context of strains phylogeny reveals parallel rearrangements that may indicate inversion-based phase variation and integration of new genomic islands. In particular, we detected parallel inversions in the second chromosomes of *B. pseudomallei* with breakpoints formed by genes encoding membrane components of multidrug resistance complex, that may be linked to a phase variation mechanism. Two genomic islands, spreading horizontally between chromosomes, were detected in the *B. cepacia* group.

Conclusions: This study demonstrates the power of integrated analysis of pan-genomes, chromosome rearrangements, and selection regimes. Non-random inversion patterns indicate selective pressure, inversions are particularly frequent in a recent pathogen *B. mallei*, and, together with periods of positive selection at other branches, may indicate adaptation to new niches. One such adaptation could be a possible phase variation mechanism in *B. pseudomallei*.

Keywords: Multi-chromosome bacteria, Genome rearrangements, Burkholderia, Pan-genome, Comparative genomics, Strain phylogeny, Positive selection

Background

The genus *Burkholderia* comprises species from diverse ecological niches [1]. In particular, *B. mallei* and *B. pseudomallei* are pathogens causing glanders and melioidosis, respectively, in human and animals [2]; *B. glumae* is a pathogen of rice [3]; *B. xenovorans* is an effective degrader of polychlorinated biphenyl, used for biodegradation of pollutants [4]; *B. phytofirmans* is a plant-beneficial endophyte that may trigger disease resistance in the host

plant [5]. *Burkholderia* genomes are unusual due to their multichromosomal organization, generally comprised of two or three chromosomes.

By definition, the pan-genome of a genus or species is the set of all genes found in at least one strain [6]. The core-genome is the set of genes shared by all strains; this gene set is usually used for accurate phylogenetic reconstruction. Genes that are not common for all considered strains but are not unique form the periphery part of a pan-genome. The pan-genome of 56 *Burkholderia* genomes was estimated to exceed 40,000 genes with no sign of saturation upon addition of more strains, and the core-genome was approximately 1000 genes [7]. A

*Correspondence: olga.bochkaryova@gmail.com

†Elena V. Moroz and Iakov I. Davydov contributed equally to this work

¹Kharkevich Institute for Information Transmission Problems, Moscow, Russia

²Center of Life Sciences Skolkovo Institute of Science and Technology, Moscow, Russia

Full list of author information is available at the end of the article



separate analysis of 37 complete *B. pseudomallei* genomes did not show saturation either [8]. The core-genome of *B. mallei* is smaller than that of *B. pseudomallei*, while the variable gene sets are larger [9].

In multi-chromosome bacterial species, the gene distribution among chromosomes is not random. The majority of genes necessary for the basic life processes usually are located in one (primary) chromosome. Other (secondary) chromosomes contain few essential genes and are mainly composed of niche-specific genes [10]. An exception is two circular chromosomes of *Rhodobacter sphaeroides* that share responsibilities for fundamental cell processes [11]. Genes from a secondary chromosome evolve faster than primary-chromosome genes and hence secondary chromosomes may serve as evolutionary test beds so that genes from secondary chromosomes provide conditional benefits in particular environments [12]. Secondary chromosomes usually evolve from plasmids [10].

Several examples of gene translocations between chromosomes in *Burkholderia* are known, e.g., the translocation between the first and the third chromosomes in *B. cenocepacia* AU 1054, affecting many essential genes [13]. Following interchromosomal translocation, genes change their expression level and substitution rate, dependent on the direction of the translocation [14].

Intra-chromosome genome rearrangements such as duplications, deletions, and inversions also play important roles in the bacterial evolution, as they strongly affect the chromosome organization and gene expression. Reconstruction of the history of genome rearrangements leads to a new class of phylogeny reconstruction algorithms [15, 16]. Chromosomal rearrangements often happen via recombination between repeated sequences, such as insertion (IS) elements [17] and rRNA operons [18]. Selection on inversion positions tends to preserve the size symmetry of the two replicores (regions of a circular chromosome between the origin and the terminus of replication), gene positions on the lagging/leading strand, and distances between genes and the origin of replication [19, 20]. Sometimes inversions are mediated by inverted paralogs. Such inversions may lead to alternating expression of these paralogs; this mechanism is known as antigenic variation by which the organism may evade host immune responses [21].

Genome rearrangement played an important role in the *B. mallei* speciation. Genomic analyses of the first sequenced *B. pseudomallei* strains and their comparison with avirulent *B. thailandensis* have shown that both chromosomes are highly syntenic between the two species, with few large-scale inversions [22, 23]. In comparison to *B. pseudomallei*, *B. mallei* genomes harbor numerous IS elements that most likely have mediated the higher rate of rearrangements [24]. In particular, IS elements of the type IS407A had undergone a significant expansion

in all sequenced *B. mallei* strains, accounting for 76% of all IS elements, and the chromosomes of these strains were dramatically and extensively rearranged by recombination across these elements [9]. The genomic reduction of *B. mallei* following its divergence from *B. pseudomallei* likely resulted in its inability to live outside the host [9, 25].

Gene gains and losses also impact the pathogenicity of species and the adaptability of an organism. The loss of a type III secretion system (T3SS)-encoding fragment in *B. mallei* ATCC 23344, compared to *B. mallei* SAVP1, is responsible for the difference in the virulence between these strains [26]. Another example is the loss of the L-arabinose assimilation operon by pathogens *B. mallei* and *B. pseudomallei* in comparison with an avirulent strain *B. thailandensis*. Introduction of the L-arabinose assimilation operon in a *B. pseudomallei* strain made it less virulent [27]. Hence, although the mechanism is not clear, there may be a link between this operon and virulence. Acquisition of the atrazine degradation and nitrotoluene degradation pathways by *B. glumae* PG1, compared to *B. glumae* LMG 2196 and *B. glumae* BGRI, likely has resulted from an adaptation since these toxic agents are used in the farming industry as a herbicide and a pesticide, respectively [28].

Here, we performed an integrated reconstruction of chromosome rearrangements for 127 complete *Burkholderia* strains, including inter-chromosome translocations, inversions, deletions/insertions, and single gene gain/loss events in the phylogenetic context. As the evolution of an obligate intracellular pathogen *B. mallei* from *B. pseudomallei* is of particular interest, we considered this branch of *Burkholderia* in additional detail, including the analysis of pan-genome statistics and identification of genes evolving under positive selection.

Methods

Available (as of 1 September 2016) complete genome sequences of 127 *Burkholderia* strains (Additional file 10: Table S1) were downloaded from the NCBI Genome database [29].

Orthologs groups

We constructed orthologous groups using Proteinortho V5.13 with the default parameters [30]. To assign GO terms to genes, we used Interproscan [31]. A GO term was assigned to an orthologous group, if it was assigned to at least 90% of genes in this group. To determine over-represented functional categories, we used the topGO v.3.6 R-package [32]. Clusters of Orthologous Groups were predicted using the eggNOG v4.5 database [33]. Protein subcellular localization was predicted using the PSORTb v3.02 web server [34]. The expression level for each orthologous group was calculated based on n data

from [35]. Using these data, we calculated RPKM, performed quantile normalization, and then calculated the average value among samples.

Pan-genome and core-genome size

To predict the number of genes in the *Burkholderia* pan-genome and core-genome, we used the binomial mixture model [36] and the Chao lower bound [37] implemented in the Micropan R-package [38]. To select the model better fitting the distribution of genes by the number of strains in which they are present, we used the Akaike information criterion with correction for a finite sample size [39, 40].

Phylogenetic trees

Trees based on nucleotide alignments

We performed codon alignment for each of the 2117 orthologous groups using Mafft v7.123b [41] and Guidance v2.01 [42]. Four orthologous groups containing sequences scored below 0.8 were excluded from further analysis. Poorly aligned residues (guidance score below 0.8) were masked. The resulting sequences were concatenated and the tree was constructed with RAxML v8.2.9 [43] using the GTR+Gamma model with 100 bootstrap runs. To ensure robustness of the tree construction we also performed calculations with 1000 replicates (see calculations at the GitHub repository).

Trees based on protein alignments

We used 1046 orthologous protein-coding genes from 127 genomes. We used Mafft v7.273 [41] in the linsi mode to align genes belonging to one orthologous group. Concatenated protein-coding sequences were used to construct the tree. We used PhyML [44] with the JTT model and discrete gamma with four categories and approximate Bayes branch supports.

Trees based on gene content

The gene content tree was constructed using the Neighbor Joining (NJ) algorithm based on the pairwise distance matrix $D_{ij} = 1 - \frac{|Strain_i \cap Strain_j|}{|Strain_i \cup Strain_j|}$, where $Strain_i$ is the set of orthologs belonging to a given strain i , ignoring paralogs.

Trees based on gene order

Trees based on gene order were built using the MLGO software (Maximum Likelihood for Gene-Order Analysis) with default parameters [16].

Trees visualization

Phylogenetic trees were visualized with FigTree v1.4.2 (<http://tree.bio.ed.ac.uk/software/figtree/>) and the Ape R-package [45].

Gene acquisition, loss, and translocations

We used GLOOME [46] for the gain/loss analysis in the evolution non-stationary model with a variable gain/loss ratio. Other parameters were set based on the character counts directly from the phyletic pattern.

To reconstruct gene translocations between chromosomes, we ordered universal single-copy orthologs and assigned a vector of ortholog presence to each strain. A component of this vector was the chromosome (1, 2, 3) harboring the ortholog in the strain. Then we subjected the obtained alignment of vectors to PAML 4.6 [47] for ancestral reconstruction with default parameters, except model = REV(GTR) and RateAncestor = 2.

Synten blocks and blocks rearrangements history

Synten blocks for closely related strains were constructed using the Sibelia software [48] with the minimal length of blocks being 5000 bp. We filtered out blocks observed in any single genome more than once. Synten blocks for distant strains were constructed using the Drimm-Synten program [49] based on locations of universal genes. The rearrangements histories for given trees topologies were constructed using the MGRA v2.2 server [50].

To distinguish between inter- and intra-replicore inversions, the origins and terminators of replication for each chromosome of each strain were determined by the analysis of peaks in GC-skew plots combined with Ori-Finder predictions [51]. Statistical significance of over-representation of inter-replicore inversions was calculated as the probability of a given number of inter-replicore inversions in the set of inversions with the given lengths. The probability of occurrence of the origin or the terminator of replication within the inversion was calculated as the ratio of the inversion length to the replicore length.

Detection of positive selection

We applied codon models for positive selection to orthologous groups common for the *B. mallei*, *B. pseudomallei*, *B. thailandensis*, *B. oklahomensis* clade. Given the low number of substitutions, it is usually not possible to reliably reconstruct the topology of a phylogenetic tree based on individual genes. On the other hand, given the high recombination rate, it is quite likely that gene evolutionary histories are different between orthologous groups. To overcome these issues we first used statistical binning [52] to group genes with similar histories, and then applied a conservative approach to detect positive selection based on multiple tree topologies.

The procedure was implemented as follows. First, we constructed a phylogenetic tree for every gene using RAxML with the GTR+Gamma model and maximum likelihood with 100 bootstrap replicates. To ensure robustness of the tree construction we also performed

calculations with 1000 replicates (see calculations in the GitHub repository). Genes with unexpectedly long branch lengths were filtered out (the maximum branch length > 0.1 or the sum of branch lengths > 0.3). Statistical binning was performed at the bootstrap incompatibility threshold of 95. For each of 25 obtained clusters we created a tree with bootstrap support using the concatenated sequence of orthologous groups belonging to the cluster.

We used two different methods to detect positive selection. The M8 vs M8a comparison allows for gene-wide identification of positive selection [53], while the branch-site model accounts for positive selection on a specific branch [54]. Each test was performed six times using different trees: the maximum likelihood tree and five random bootstrap trees. We used the minimum value of the LRT (likelihood ratio test) statistic to avoid false identification of positive selection which could be caused by an incorrect tree topology.

For the branch-site model we tested each internal branch as a foreground branch one by one; we did not test terminal branches to avoid false positives caused by sequencing errors. The results of the branch-site tests were aggregated only in the case of bipartition compatibility. We considered only bipartitions that were present in at least three tests, we also computed the minimum value of the LRT statistic. The test results were mapped back to the species tree based on the bipartition compatibility.

The strength of purifying selection is measured by $w_0 < 1$ with smaller values corresponding to stronger purifying selection. The $w_2 > 1$ parameter of the branch-site model captures positive selection, with higher values indicating stronger selection.

In both cases we used the chi-square distribution with one degree of freedom for the LRT to compute the p -value. Finally, we computed the q -value, all LRT values equal to zero were excluded from the test. We set the q -value threshold to 0.1.

Statistical methods

To estimate dependencies between various parameters such as the expression level, localization in the first/second chromosome, localization on the leading/lagging strand, we used linear models (lm function, R v3.3.2). Additional parameters such as the sum of branch lengths, alignment length, and GC-content were included as they can affect the power of the method [55]. The parameters were transformed to have a bell-shaped distribution if possible: $\log(x + 1)$ for the expression levels, $\log(x + 10^{-6})$ for the LRT statistic, and $\log(x)$ for the alignment length, sum of branch lengths, standard deviation of GC-content, and ω_0 . Continuous variables were centered at zero and scaled so that the standard deviation was equal to one. This makes the linear model coefficients directly comparable. Outliers were identified in the residual plots

and excluded from the model; the residual plots did not indicate abnormalities. For the linear models, we included potential confounding variables in the model, and kept only significant ones for the final linear model.

Results

Phylogeny and pan-genome analysis

The analysis of orthology for 127 *Burkholderia* strains yielded 757,526 orthologous groups containing two or more genes. 21,740 genes were observed in only one genome, some of them could result from mis-annotation. Alignments of 1024 single-copy common gene (hereinafter "core genes") were used for construction of the phylogenetic tree (hereinafter "the basic tree").

As the number of available *Burkholderia* genomes in GenBank is constantly increasing, we performed comprehensive pan-genome analysis. The pan-genome size for all strains is 48,000 genes with no signs of saturation, showing that the gene diversity of the *Burkholderia* species has not been captured yet (Fig. 1a). Based on these data, the binomial mixture model [36] predicts that, as more genomes are sequenced, the *Burkholderia* core-genome would reach the lower limit of 457 genes, whereas the pan-genome size would be at most 86,845. The number of new genes decreases with each new genome n at the rate $N(n) = 2557n^{-0.56}$ confirming that the pan-genome is indeed open (Additional file 1: Figure S1a). Each new genome adds about 171 genes to the pan-genome. The Chao lower bound estimate [37] of the pan-genome size is 88,080. These results are consistent with the reported pan-genome size of 56 *Burkholderia* strains [8]. The core-genome size dependence on the number of analyzed strains is shown in Fig. 1b. The number of universal genes that are present in all strains saturates at about 1050.

Additional files 2: Figure S2 and 3: Figure S3 show the core- and pan-genome size dependencies for *B. pseudomallei* and *B. mallei*, respectively. Their pan-genomes also have not reached saturation ($N(n) = 788n^{-0.53}$ for *B. pseudomallei* and $N(n) = 867n^{-0.87}$ for *B. mallei*) (Additional file 1: Figure S1b, c). These results are also consistent with the reported pan-genome size of 37 *B. pseudomallei* strains [7].

The distribution of genes by the exact number of strains in which they are present has a typical U-shape form (Fig. 2) [56, 57], with numerous unique and universal genes and fewer periphery genes. We have compared two models that are traditionally used for U-curve approximation, by the sum of three exponents (for unique genes, the periphery, and the universal genome, respectively) [58] and by the sum of two power law functions, the first term describes the genes present in a few strains (almost unique), and the second term reflects the distribution of genes present in most strains (almost universal) [56].

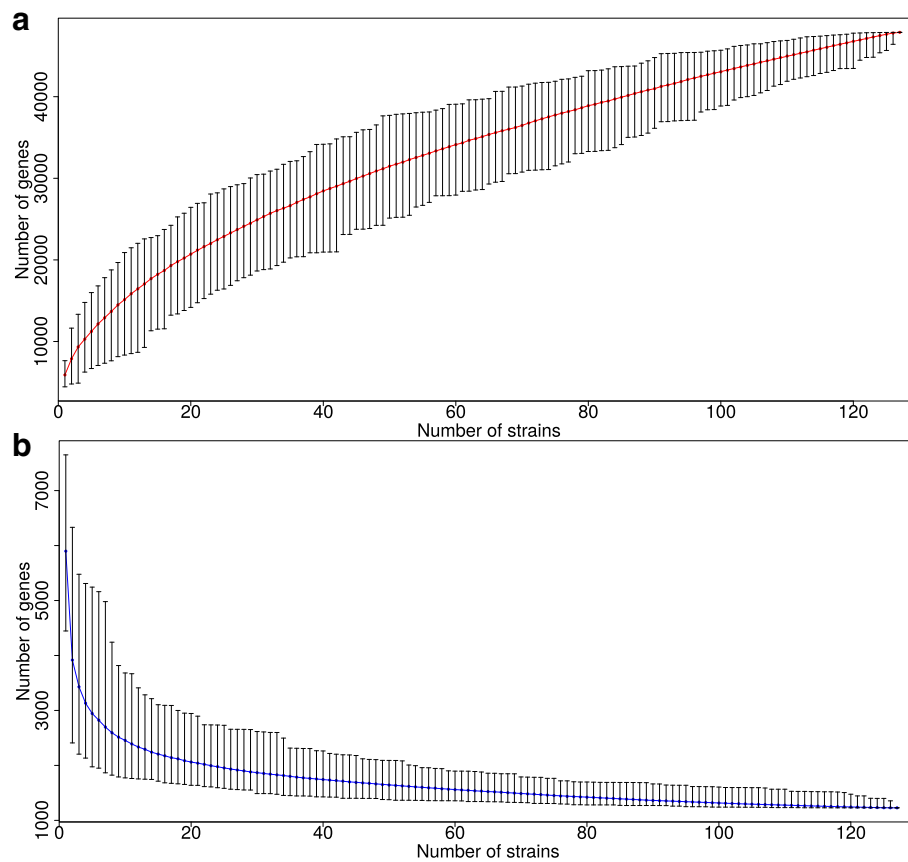


Fig. 1 Number of genes as a function of the number of sequenced *Burkholderia* genomes. **(a)** The pan-genome size, that is, the number of all genes in sequenced strains. The number of new genes decreases with each new genome n at the rate $N(n) = 2557n^{-0.56}$ confirming that the pan-genome is open. As the numbers of genomes $n \rightarrow \infty$, the pan-genome size converges to 88,080 and the core-genome size converges to 457 genes. **(b)** The core-genome size, that is, the number of common genes in sequenced strains. The core-genome for a single strain ($n = 1$) is defined as the number of genes in the strain

Application the method of the least squares with the Akaike information criterion (AIC) revealed that the approximation by the sum of three exponents recapitulates the U-shape slightly better. This is consistent with the analysis of the *Streptococcus* pan-genome [59], in which the sum of three exponents also has provided a better fit.

One possible explanation based on preliminary, unpublished observations in other bacteria could be that the power-law rule applies only to very closely related genomes, such as strains of one species, whereas more distant organisms, starting from the genus level, follow the exponential model. Indeed, for almost identical strains, gene gain and loss would be unique events, and the periphery would be vanishingly small, whereas on the other extreme, e.g. when all archaea [60] or all bacteria [61] are considered, the periphery clearly should be very large. The transition between these two modes may be a subject for additional, separate analysis and modeling.

Genes acquisition and loss

Gains and losses of genes along the phylogenetic tree were assessed, excluding plasmid genes (Additional file 4: Figure S4). *Burkholderia* species have experienced numerous gene gains and losses, that could explain their ecological diversity. In particular, a separate analysis of the *B. pseudomallei* group yielded considerable gene loss in the *B. mallei* clade. The genome reduction among the *B. mallei* strains is likely associated with the loss of genes redundant for obligate pathogens [9].

The basic tree and the gene content tree are largely consistent as the trees have the same clades with one major exception (Additional file 5: Figure S5). In the gene content tree, *B. mallei* and *B. pseudomallei* form two distinct clusters, whereas in the basic tree monophyletic *B. mallei* are nested within paraphyletic *B. pseudomallei*. The former clustering could be due to the lifestyles of *B. mallei* and *B. pseudomallei*, as both species are pathogens of animals and possess specific sets of genes. Thus even if

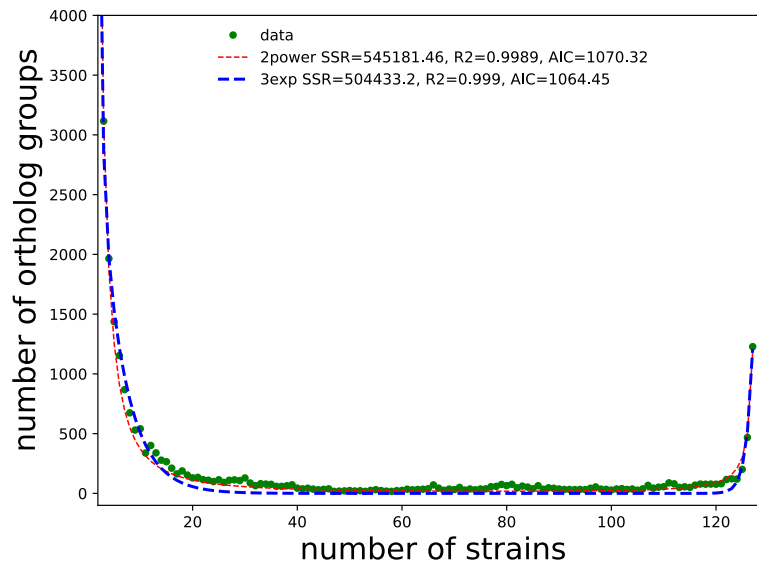


Fig. 2 Distribution of ortholog groups by the number of strains in which they are present. For each number of strains x , the number of genes y present in exactly x strains is given. The blue line corresponds to the approximation by a sum of three exponents $y = e^{-0.2x+8.4} + e^{-1.8x+11.6} + e^{0.85x-100.1}$; the red line corresponds to the approximation by a sum of two power functions $y = 21648.4x^{-1.8} + 1182.8(128 - x)^{-1.2}$. Based on the Akaike information criterion (AIC), the approximation by the sum of three exponents recapitulates the U-shape slightly better

universal genes in some *pseudomallei* strains are closer to the orthologous genes in *mallei* than to genes in other *pseudomallei* strains, these species will be distant in the gene content tree due to species-specific genes.

Although the trees are comprised of the same clades, we observed numerous contradictions in strains positions. These contradictions are likely caused by clade-specific patterns of recombination and accessory gene exchange [62].

Inter-chromosome translocations

The gene distribution among the chromosomes (Table 1) for *Burkholderia* spp. is consistent with previous observations for other multi-chromosome bacteria [10]. At that, the majority of core genes belong to the first chromosome, ten-fold less core genes are in the second chromosome, and they are almost absent in the third chromosome, the only exception resulting from a large translocation from the first to the third chromosome in *B. cenocepacia* AU 1054 [13].

Reconstruction of translocations of 1024 core genes between the chromosomes yielded 210 events (Fig. 3). The genomes of *B. cenocepacia* 895, *B. cepacia* strain LO6, and *B. contaminans* MS14 were not included in the rearrangement analysis due to likely artifacts of the genome assembly (See Additional file 6: Figure S6). Thirty-eight events were reconstructed separately for *B. mallei* and *B. pseudomallei*.

There was no statistically significant overrepresentation of GO categories in the set of translocated genes. Six genes have been translocated independently on different tree branches twice or more times, encoding aldo/keto reductase (IPR020471), HTH-type transcriptional regulator *argP* (IPR017685), gamma-glutamyltranspeptidase (IPR000101), acid phosphatase *acpA* (IPR017768), tryptophan synthase beta subunit-like PLP-dependent enzyme (IPR036052), *tonB*-dependent receptor. The reconstructed common ancestor of *Burkholderia* has 965 universal single-copy genes in the first chromosome, and 81, in the second chromosome.

We analyzed intra-chromosomal rearrangements that involve the core genes using only one representative strain from clades with closely related species. The core genes were grouped into 87 synteny blocks that contained two or more core genes in the same order in all analyzed genomes. The rearrangements history yielded no parallel events except parallel translocations between chromosomes described above. While one could expect that changes in the lifestyle and population bottlenecks could increase the mutation and recombination rates simultaneously, no correlation between the number of rearrangements and the average mutation rates of the core genes was observed (data not shown).

We then reconstructed the detailed history of rearrangements in specific clades. A large number of available genomes of closely related bacterial strains allows

Table 1 Distribution of universal orthologs among the chromosomes

Species	The first chromosome	The second chromosome	The third chromosome
<i>B. mallei</i>	930 out of 3135±103	116 out of 1842 ±142	
<i>B. pseudomallei</i>	954 out of 3498±151	92 out of 2536±178	
<i>B. thailandensis</i>	954 out of 3626±216	92 out of 2562±109	
<i>B. oklahomensis</i>	954 out of 3630±36	92 out of 2537±61	
<i>B. gladioli</i>	964 out of 3924±146	82 out of 3026±100	
<i>B. glumae</i>	964 out of 3385±120	82 out of 2524±365	
<i>B. vietnamiensis</i>	961 out of 3055±139	85 out of 2411±420	
<i>B. cepacia</i> group	962 out of 3205±146	84 out of 2528±413	0 out of 922±160
<i>B. cenocepacia</i> AU 1054	747 out of 2965	84 out of 2472	215 out of 1040
<i>B. cenocepacia</i> strain DDS 22E-1	961 out of 3296	84 out of 2831	1 out of 939
<i>B. dolosa</i> AU0158	963 out of 3084	83 out of 1861	
<i>B. ubonensis</i> MCMB22	963 out of 3216	83 out of 3035	
<i>B. pyrrocinia</i> strain DSM 10685	963 out of 3157	84 out of 2714	0 out of 838
<i>B. sp.</i> CCGE1001	968 out of 3545	78 out of 2420	
<i>B. sp.</i> CCGE1002	968 out of 3116	78 out of 2258	0 out of 1109
<i>B. sp.</i> CCGE1003	967 out of 3463	79 out of 2525	
<i>B. sp.</i> HB1	967 out of 3481	79 out of 2743	
<i>B. sp.</i> KJ006	961 out of 2917	85 out of 2132	0 out of 930
<i>B. sp.</i> OLGA172	967 out of 4023	79 out of 2998	
<i>B. sp.</i> PAMC 26561	964 out of 3034	82 out of 1437	
<i>B. sp.</i> PAMC 28687	960 out of 2991	83 out of 1367	3 out of 1509
<i>B. sp.</i> RPE64	964 out of 2907	81 out of 1422	0 out of 853
<i>B. sp.</i> RPE67	963 out of 2859	81 out of 1688	1 out of 1553
<i>B. sp.</i> TSV202	954 out of 3645	92 out of 2536	
<i>B. sp.</i> YI23	963 out of 2769	81 out of 1539	1 out of 1364

Each cell shows the number of universal genes out of the number of all genes in the chromosome. For species with more than one strain, the average number of genes and the standard deviation are shown

one to consider micro-rearrangements in the evolutionary context, revealing parallel events that may indicate the action of antigenic variation [59]. An integrated analysis of sequence-based and inversion-based trees enhances the resolution of the phylogenetic reconstruction in the case of a high rate of genome rearrangements in a population [63].

Rearrangements in the *B. cepacia* group

For 27 strains of the *cepacia* group, the average coverage of chromosomes by synteny blocks was 50% for the first, 30% for the second, and less than 10% for the third chromosome. This agrees with the preferred location of universal genes discussed above. Hereafter, the third chromosomes are not considered due to their low conservation. Fixing the tree to the basic one, we reconstructed 17 inversions and 574 insertion/deletion events. The topology of the phylogenetic tree based on the order of synteny blocks (Additional file 7: Figure S7c) is not consistent with the

basic tree, and a majority of deep nodes have low bootstrap support that may be explained by numerous parallel gain/loss events.

Only one parallel inversion of length 530 kb was found in the first chromosome of *B. cenocepacia* AU 1054 and *B. cenocepacia* J2315, the inversion breakpoints formed by the 16S-23S rRNA loci. In order to distinguish between truly parallel events and homologous recombination between these strains, we constructed a tree based on proteins encoded by genes from the inverted fragment. *B. cenocepacia* AU 1054 and *B. cenocepacia* J2315 did not change their position in the tree, and, in particular, did not cluster together (data not shown). Hence, this block was not subject to homologous recombination between these strains.

Two non-universal synteny blocks were found in different chromosomes in different strains. One block with length 8.5 kb is located in the first chromosome of *B. cenocepacia* MC0-3 and in the second chromosome of

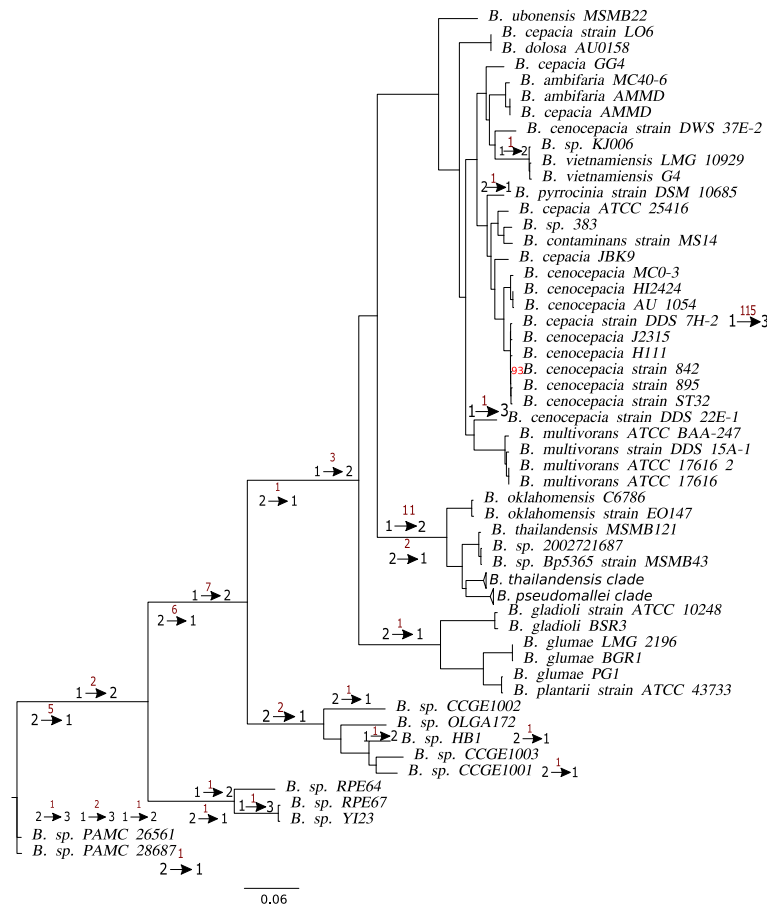


Fig. 3 Translocations in *Burkholderia* spp. The phylogenetic tree of *Burkholderia* is constructed based on the protein sequence similarity of single-copy universal genes. The bootstrap support is shown for branches where it is < 100. The red numbers above the arrows show the number of genes translocated between chromosomes on the tree branches; the black numbers mark chromosomes that have been involved in the transfer, the arrows show the direction of the transfer

B. cepacia ATCC 25416. This genomic island contains five genes that belong to the iron uptake pathway and an AraC family protein. Some genes of this cassette were also found in other *Burkholderia* species (Fig. 4a).

Another block with length 5.5 kb was found only in 17 of 30 strains belonging to the *cepacia* group (Fig. 4b). This island contains four genes encoding the acetyl-CoA carboxylase complex, glycoside hydrolase (GO:0005975 carbohydrate metabolic process), and a LysR family protein. The island is found in all *B. mallei*, *B. pseudomallei*, *B. oklahomensis*, *B. glumae*, *B. gladioli* and is absent in *B. thailandensis* and other strains. Its presence in different chromosomes and differences between the tree of this cassette (Additional file 8: Figure S8) and the basic tree indicate that this genomic island is spreading horizontally.

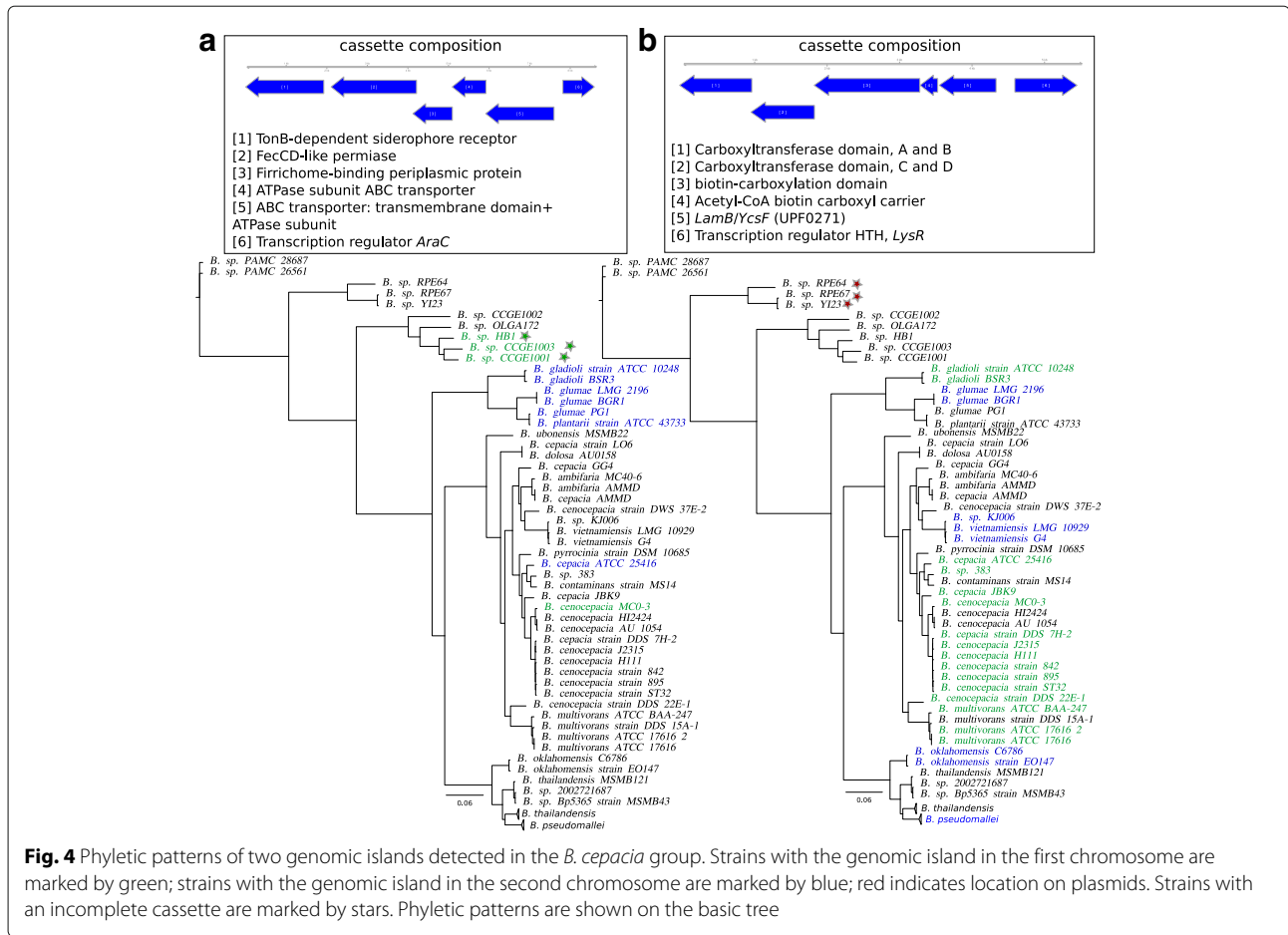
Rearrangements in the *B. mallei* clade

For fifteen *B. mallei* strains and two *B. pseudomallei* used as outgroups, we constructed 104 common synteny blocks

in both chromosomes. Only one block with length 40 kb, comprised of 24 universal genes, was translocated in the *B. mallei* clade. This block is bounded by IS elements and rRNA genes that may indicate that this translocation resulted from recombination between chromosomes.

This indicates that in these strains translocations between chromosomes are rare in comparison to within-chromosome rearrangements. Fixing the tree to the basic one, we reconstructed 88 inversions in the first chromosomes and 27 inversions in the second ones (Fig. 5). The reconstruction yields nine parallel events in the first chromosomes and three, in the second ones. The boundaries of the inversions are formed by repeated sequences (transposases).

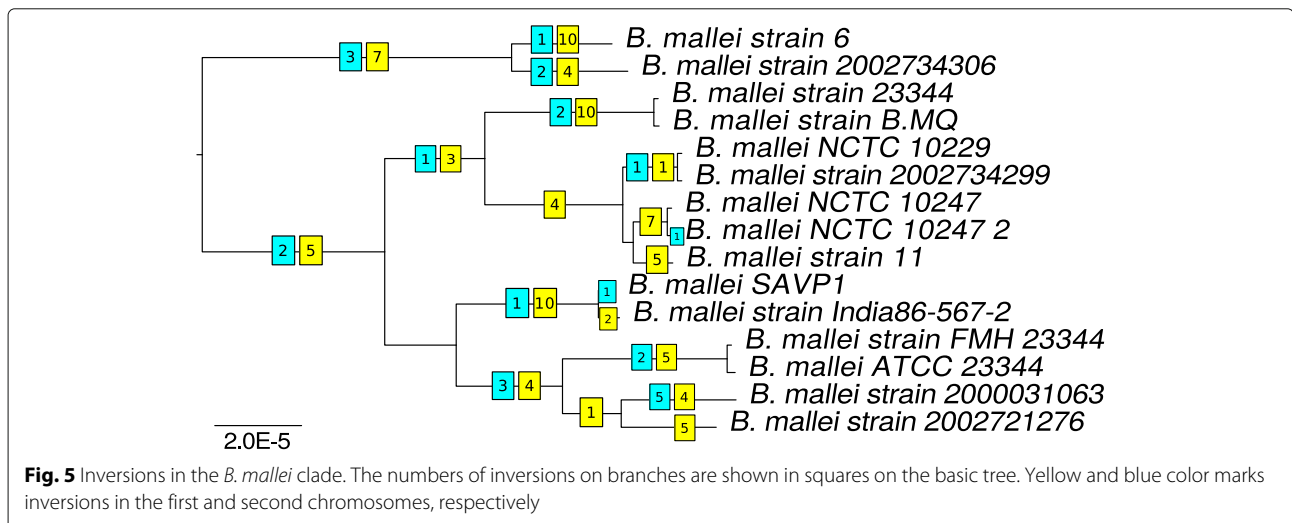
To test the possibility that the contradictions between the tree topology and the inversion history had been caused by homologous recombination, we constructed trees based on genes involved in these events. For all inverted sequences, strains do not change their positions



in the tree (data not shown). Therefore, we suppose that parallel events were caused by active intragenome recombination linked to a limited number of repeated elements.

We applied maximum likelihood optimization methods to obtain a topology based on the universal gene order. The optimized topology (Additional file 7: Figure

S7a) yielded a comparable number of parallel inversions, demonstrating that the latter were not an artifact arising from an incorrect phylogeny. We observed positive correlation between the inversion rate and the mutation rates in the core genes (Spearman test, $\rho = 0.8, p\text{-value} = 10^{-7}$) (Fig. 6).



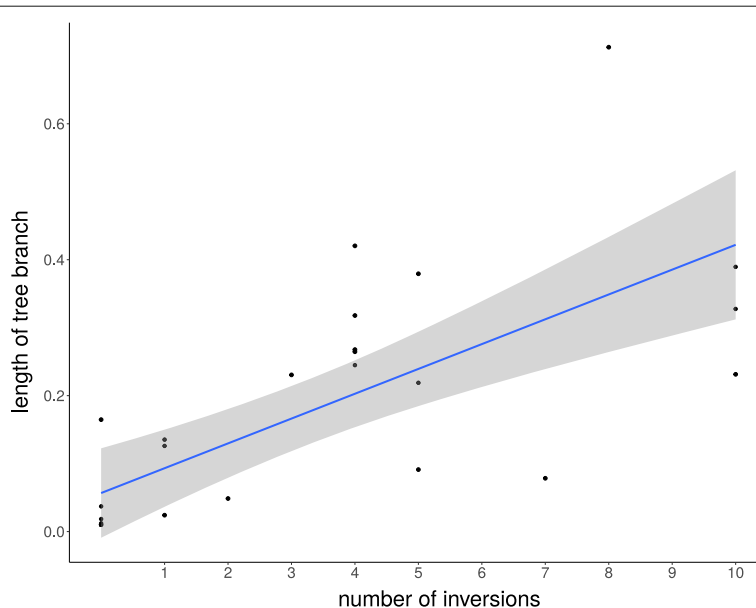


Fig. 6 The rearrangements rate as a function of the mutation rate for *B. mallei*. Each dot corresponds to a branch in the basic phylogenetic tree (Fig. 5)

Rearrangements in *B. pseudomallei*

The gene order in 51 strains of *B. pseudomallei* turned out to be significantly more stable than that in *B. mallei*, as only three inversions were reconstructed in the first chromosomes, and five, in the second chromosomes (Fig. 7a). Moreover, the average coverage of chromosomes by synteny blocks was more than 90% for the first, and 80% for the second chromosomes, revealing a stable order and gene content. Two blocks with length about 20–25 kb are swapped in *B. pseudomallei* K42 that is likely to be an assembly artifact.

Inversions in the second chromosomes with length about 1.3 Mb have the same boundaries for all seven strains despite the fact that they are located at distant branches of the phylogenetic tree (Fig. 7b). Breakpoints of these inversions are formed by six genes encoding (1,2) rhamnosyltransferase type 1 A,B; (3) drug resistance transporter (*mrB/QacA* subfamily); (4) rhamnosyltransferase type II; (5,6) components of a RND efflux system, outer membrane lipoproteins *nodT* and *emrA*.

Rearrangements in the *B. thailandensis* clade

For 15 strains *B. thailandensis*, we constructed 56 synteny blocks in both chromosomes. Two strains of *B. oklahomensis* and one *B. pseudomallei* were used as outgroups. The average coverage by blocks was 75% for the first, and 50% for the second chromosomes. Fixing the tree topology to the basic tree, we reconstructed 18 inversions and 265 insertion/deletion events (Fig. 8). *B. thailandensis* has a higher rate of inversions and deletions than *B. oklahomensis* and *B. pseudomallei*. The reconstruction yields two parallel events in the first chromosomes and one,

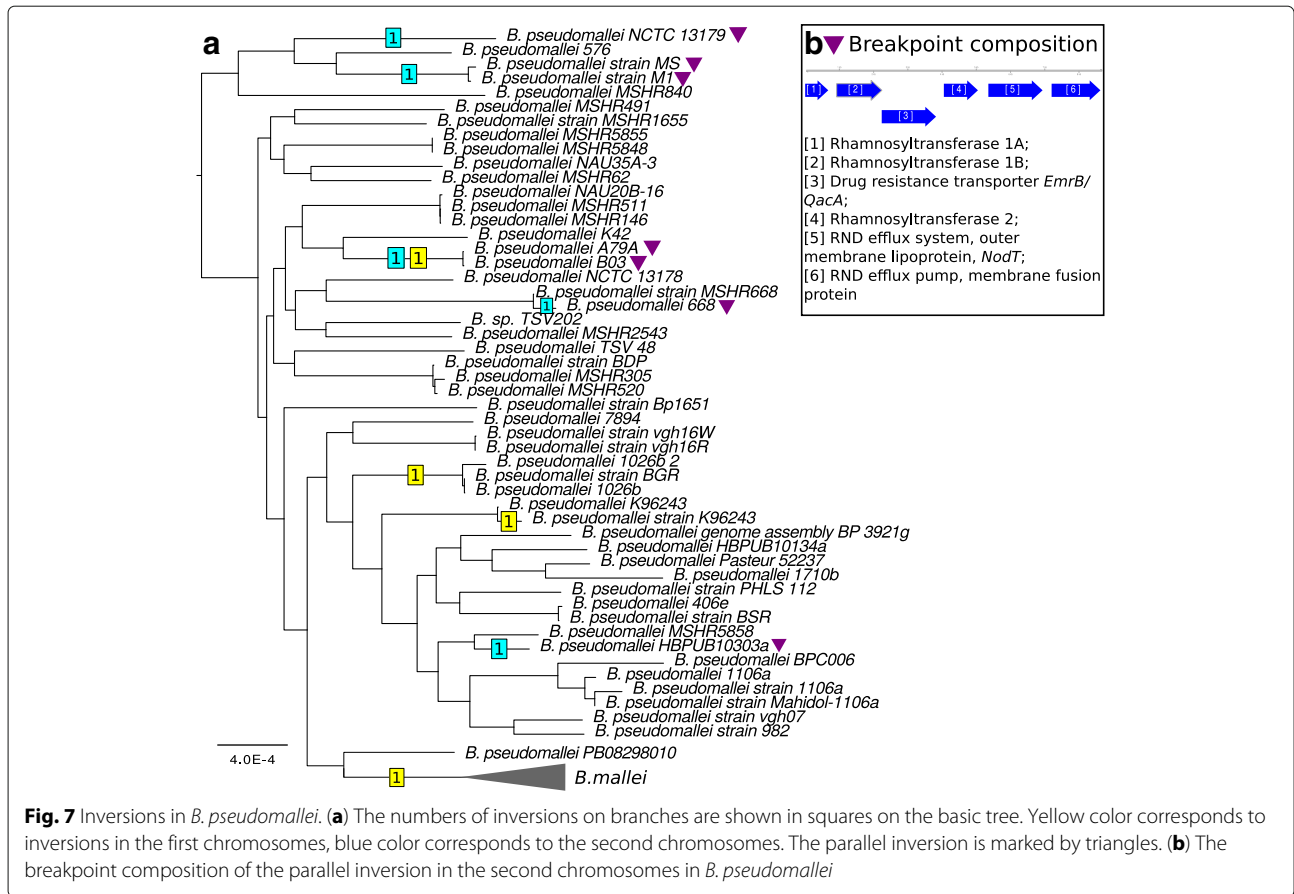
in the second ones. The boundaries of these inversions are formed by repeated sequences (transposases). For all inverted sequences, strains do not change their position in the trees based on sequences similarities of genes involved in these events (data not shown).

The topology of the phylogenetic tree based on the order of synteny blocks (Additional file 7: Figure S7b) is largely consistent with the basic tree, the only exception being a changed position of *B. thailandensis* E254 caused by parallel inversions.

Two non-universal, non-trivial translocated synteny blocks were found. One is a block with length 38 kb in the first chromosome in *B. pseudomallei*, the second chromosome in *B. oklahomensis*, and absent in the *B. thailandensis* genomes. This block is comprised of genes linked with amino acids metabolism. The second block is a parallel phage insertion with length 9 kb in the first chromosome of *B. oklahomensis* strain EO147 and in the second chromosome of *B. thailandensis* 2003015869.

Selection regimes

As the evolution of species from the *B. thailandensis*, *B. pseudomallei*, and *B. mallei* clade is of particular interest due to dramatic changes in their lifestyle, including an adaptation to intra-cellular one, for these strains we identified genes evolving under positive selection. 1842 single-copy genes common for the *B. oklahomensis*, *B. thailandensis*, *B. pseudomallei*, *B. mallei* clade were tested. We detected 197 genes evolving under positive selection using the M8 model (Additional file 11: Table S2). No GO categories were significantly overrepresented but we observed overrepresentation of outer membrane



proteins (permutation test, p -value=0.03) consistent with observations in other bacterial species [64, 65].

To identify branch-specific positive selection, we used the branch-site test. In total, we identified seventeen events (Table 2), twelve of which we successfully mapped

to the basic tree (Fig. 9). In the remaining five cases (flagellar hook protein FlgE, porin related exported protein, penicillin-binding protein, phosphoenolpyruvate-protein kinase and cytidylate kinase), the detected branches (bipartitions) of the gene trees were incompatible

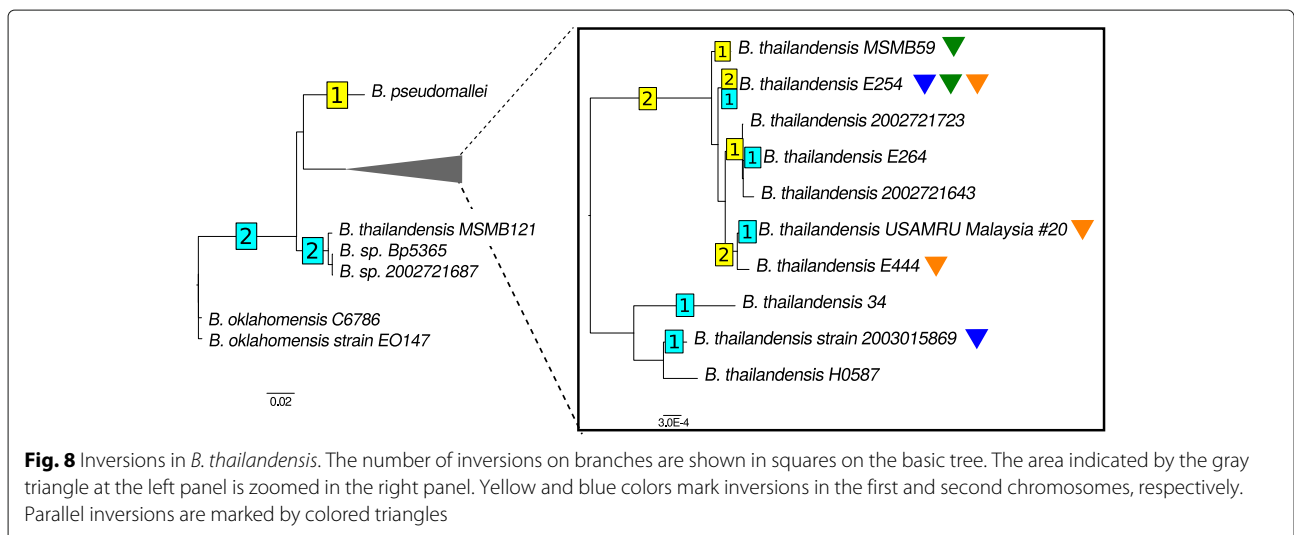


Table 2 Genes evolving under branch-specific positive selection

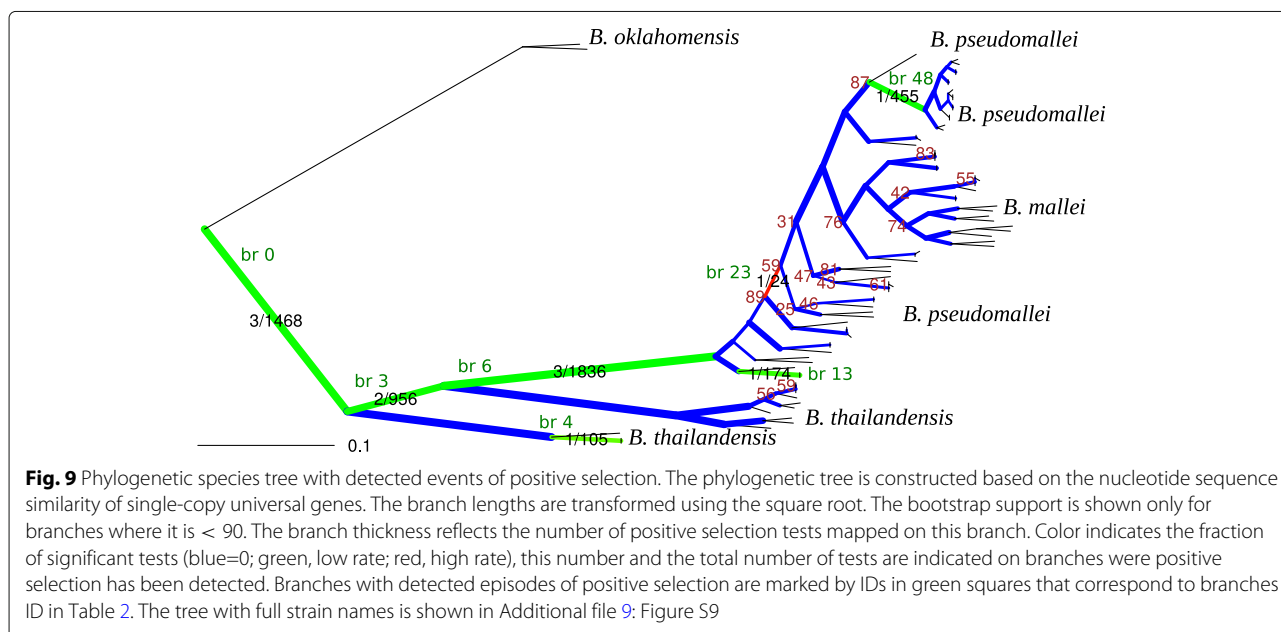
Branch ID	Function of the gene product	ω_2	p -value	COG	Localization
0	Transketolase (<i>tktA</i>)	12.7	$8 \cdot 10^{-5}$	G	CP
0	Putative aminotransferase protein	8.8	$2 \cdot 10^{-5}$	E	CP
0	Glycine cleavage system T protein (<i>gcvT</i>)	26.1	$5 \cdot 10^{-5}$	E	CP
3	Error-prone DNA polymerase (<i>dnaE2</i>)	8.8	$2 \cdot 10^{-5}$	F	CP
3	Metallo-dependent hydrolases	104	$2 \cdot 10^{-6}$	Q	CP
4	<i>LysR</i> -family transcriptional regulator	507	$6 \cdot 10^{-22}$	K	CP
6	Dyp-type peroxidase	18	$2 \cdot 10^{-7}$	P	CP
6	<i>Kipl</i> family	17	$1 \cdot 10^{-5}$	E	CP
6	<i>OmpA</i> family transmembrane protein	35	$1 \cdot 10^{-5}$	M	OM
13	Alpha/beta hydrolase fold	40	$8 \cdot 10^{-8}$	S	CP
23	Inner membrane protein <i>YajD/ElaB</i>	1000	$7 \cdot 10^{-11}$	S	NA
48	Glutamate synthase large subunit-like protein	117	$8 \cdot 10^{-8}$	E	NA
N/A	Cytidylate kinase	1000	$1 \cdot 10^{-5}$	F	CP
N/A	Flagellar hook protein (<i>FlgE</i>)	4	$1 \cdot 10^{-6}$	N	EC
N/A	Phosphoenolpyruvate-protein kinase	576	$5 \cdot 10^{-5}$	G	CP
N/A	Porin related exported protein	7.5	$1 \cdot 10^{-5}$	M	OM
N/A	Penicillin-binding protein	86	$3 \cdot 10^{-5}$	M	CM

The COG categories are coded as follows: K, transcription; M, cell wall/membrane biogenesis; N, Cell motility; G, carbohydrate transport and metabolism; E, amino acid transport and metabolism; F, nucleotide transport and metabolism; P, inorganic ion transport and metabolism; Q, secondary metabolites biosynthesis, transport and catabolism; R General function prediction only. The localization is coded as follows: CP, Cytoplasmic; OM, outer membrane; EC, Extracellular; CM, Cytoplasmic Membrane; NA, unknown (these proteins may have multiple localization sites)

with the basic tree, and thus could not be mapped to it.

Outer membrane proteins such as the flagellar hook protein FlgE, porin-related exported protein, OmpA family protein can serve as targets for the immune response.

Moreover, OmpA is known to be associated with virulence, being involved in the adhesion and invasion of host cells, induction of cell death, serum and antimicrobial resistance, and immune evasion [66]. Positive selection on the error-prone DNA polymerase, having a lower repli-



cation fidelity, might result from a need for a higher mutation rate facilitating adaptation to a new life style. Similarly, bacterial transcription factors are known to enable rapid adaptation to environmental conditions, that might explain strong positive selection of the LysR-family transcriptional regulator.

The majority of genes evolving under positive selection have been identified in the longest branches; accordingly, the fraction of events is higher in these branches. This might indicate rapid adaptation to new ecological niches during species formation. However, the branch-site test for positive selection is more powerful on longer branches, and the position of a branch in the tree might affect the power [67]. Hence, overrepresentation of positive selection events may be related to the power of the method, and does not necessary indicate the higher number of genes affected by positive selection on these branches.

We used linear modeling to identify determinants affecting purifying selection (Table 3). The strongest observed correlation is that highly expressed genes tend to evolve under stronger purifying selection, which is also consistent with previous observations [12]. The expression levels in our dataset are higher for the first chromosome (Additional file 12: Table S3), which is consistent with observations for other multi-chromosome bacterial species [68]. While the effect of correlation is not particularly high, the correlation is strongly statistically significant. This indicates that despite high stochasticity of mRNA expression, there is a statistically strong association between the expression level and gene localization and average GC content.

Longer genes tend to experience stronger purifying selection that is consistent with previously shown negative correlation between the d_N/d_S value and the median length of protein-coding genes in a variety of species [69]. However, this observation also could be explained by the greater power in detecting strong negative selection in longer genes, similarly to the increase in the power when detecting positive selection for longer genes [67].

Table 3 The linear model of average ω (negative selection, estimated using M8), non-significant variables removed from the model

	Estimate	Std. Error	t-value	p-value
Alignment length	-0.085	0.025	-3.397	0.000704
Average expression level	-0.079	0.026	-3.023	0.002561
Sum of branch lengths	0.518	0.026	19.638	$< 2 \cdot 10^{-16}$

For the full model see Additional file 12: Table S3. The model p -value is $< 2.2 \cdot 10^{-16}$; the adjusted R^2 is 0.2882

Selection on recombination events

In many bacteria, within-replichore inversions, that is, inversions with endpoints in the same replichore, have been shown to be relatively rare and significantly shorter than inter-replichore inversions [70, 71]. The pattern of inversions reconstructed for both chromosomes in *B. mallei* is consistent with both of these observations.

Inter-replichore inversions are overrepresented in the first (p -value $< 10^{-33}$) and the second (p -value $< 10^{-30}$) chromosomes. The lengths of inter-replichore inversions have a wide distribution up to the full replichore size (Fig. 10a), whereas the observed within-replichore inversions mainly do not exceed 15% of the replichore length. We observed only two longer inversions, both in *B. mallei* FMH23344. These inversions overlap with each other and may be explained by a single translocation event. This strong avoidance of inter-replichore inversions is probably caused by selection against gene movement between the leading and the lagging strands [72].

The reconstruction of translocations also revealed that genes tend to retain their position on the leading or lagging strand (two-sided binomial test, p -value=0.03, Fig. 10b). Moreover, all blocks of more than three genes retain their positions. We have not observed any difference in the level of purifying selection between genes translocated from the leading and lagging strands.

Discussion

The pan-genome of most bacterial species is open and driven by horizontal gene transfer that is known to be one of the major forces of bacterial genome evolution [73–75]. Generally, pan-genomes with a large periphery are characteristic of organisms with large long-term effective population sizes and an ability to fill a variety of new niches [75]. On the other hand, pan-genomes of obligate intracellular bacteria species such as *Chlamydia* are characterized by a large pool of universally conserved genes, a small periphery, and relatively few strain-specific genes [76].

Despite the fact that *B. mallei* is an intracellular pathogen with a relatively small population size, its pan-genome is open and characterized by a large number of accessory genes. The reason for that is likely to be numerous deletions in the second chromosomes that followed the recent change in the lifestyle and mode of pathogenicity. Strain-specific deletions yielded a large number of genes retained by only few *B. mallei* genomes hence contributing to the large periphery fraction of the pan-genome. If this explanation is correct, the size of the pan-genome periphery will gradually decrease with time, as continuing gene losses are unlikely to be compensated by acquisitions due to the limited horizontal gene transfer in intracellular parasites.

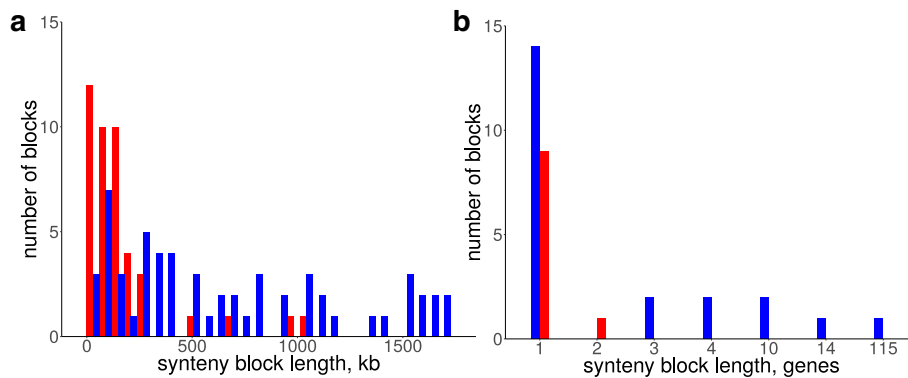


Fig. 10 Histograms of lengths of rearranged syntenic blocks. **(a)** Blocks inverted in the first chromosomes in *B. mallei*; **(b)** blocks translocated between chromosomes in *Burkholderia* spp. Blue color corresponds to syntenic blocks that have retained their position with respect to the leading/lagging strand; red color corresponds to syntenic blocks that changed the strand

An important mode of genome evolution is rearrangements of chromosome fragments. In prokaryotes with single chromosome, the prevalent type of rearrangements are the symmetrical inversions around the origin of replication. This pattern is traditionally explained by the selection against gene movement between the leading and the lagging strands caused, in turn, by overrepresentation of highly expressed genes on the leading strand [19, 70, 77]. While a number of rearrangements in some *Burkholderia* strains have been described [13, 23, 24], the increased phylogenetic coverage allowed us to actually map the events to the phylogenetic tree. Reconstruction of inversions on both chromosomes is consistent with the inversion pattern in single-chromosome bacterial species.

The pattern of inter-chromosome translocations also revealed selection against gene movement between the leading and the lagging strands. Nevertheless, this tendency is not statistically strong, in particular due to insufficient sample size. Future analysis of multi-chromosome genera such as *Vibrio* or *Brucella* [78] should allow one to validate and extend this observation.

Young pathogens such as *Yersinia pestis*, *Shigella* spp., *B. mallei* are known to have a particularly high rates and a variety of mobile elements that may be explained by fast evolution under changed selection pressure in new conditions, bottlenecks in the population history, and weaker selection against repetitive elements due to the decreased effective population size [79]. Accumulation of IS elements is most likely responsible for frequent genome rearrangement and strong genome reduction in *B. mallei* [24].

The observed parallel inversion between paralogous genes in *B. pseudomallei* encoding surface antigen protein might indicate the action of an antigen variation mechanisms leading to phenotype diversification. Strong similarity between the repeats flanking this inversion did

not allow us to map the point of recombination inside the repeats as done in [59]. The ability of clonal bacterial populations to generate genomic and phenotypic heterogeneity is thought to be of great importance for many commensal and pathogenic bacteria. While direct confirmation of this mechanism requires analysis of transcripts, parallel, independent inversions are a good lead for subsequent experimental validation [80].

Conclusions

The rearrangement rates differ dramatically in the *Burkholderia* species; from a couple of inversions in *B. pseudomallei* strains to dozens of events in *B. mallei* strains. The tree based on the alignment of universal genes and the gene content tree also show some differences, caused by excessive gene gains and losses at some branches, most notably, gene loss in *B. mallei* following a drastic change of the lifestyle.

Integrated reconstruction of chromosome rearrangements in the context of strains phylogeny reveals parallel rearrangements. In particular, we detected parallel inversions in the second chromosomes of *B. pseudomallei* with breakpoints formed by genes encoding membrane components of multidrug resistance complex, that may be linked to a phase variation mechanism. Two genomic islands, spreading horizontally between chromosomes, were detected in the *B. cepacia* group. Hence, evolutionary and functional analysis of parallel rearrangements identifies possible cases of phase variation by inversions and integration of new genomic islands that is especially important for the micro-evolution of pathogens.

The observed strong avoidance of large intra-replichore inversions is likely caused by selection against transfer of large groups of genes between the leading and the lagging strands. At that, translocated genes also

tend to retain their position in the leading or the lagging strand and this selection is stronger for large syntenies.

Overall, this study demonstrates the strength of integration of diverse approaches to the analysis of bacterial genomic evolution.

Additional file

Additional file 1: Figure S1. The number of new genes added to the pangenome upon addition of new strains. (a) *Burkholderia* spp., (b) *B. pseudomallei*, and (c) *B. mallei*. The number of new genes is plotted as a function of the number (n) of strains sequentially added (see the model in [81]). For each n, points are the values obtained for different strain combinations; red symbols are the averages of these values. The superimposed line is a fit with a decaying power law $y = A * n^B$. (PDF 1224 kb)

Additional file 2: Figure S2. Pan-genome (a) and core-genome (b) size of *B. pseudomallei* strains. (PDF 21 kb)

Additional file 3: Figure S3. Pan-genome (a) and core-genome (b) size of *B. mallei* strains. (PDF 31 kb)

Additional file 4: Figure S4. Gene flow during *Burkholderia* evolution. Red and blue numbers are, respectively, the numbers of gained and lost genes on a given branch. (PDF 414 kb)

Additional file 5: Figure S5. Comparison the topologies of phylogenetic trees based on the protein sequence similarity of single-copy universal genes and the gene content. (PDF 26 kb)

Additional file 6: Figure S6. Whole-genome alignments of *cepacia* strains that were not included in the rearrangement analysis due to likely artifacts of the genome assembly. (a) *Burkholderia* sp. 383 and *B. cepacia* strain LO6 (b) *Burkholderia* sp. 383 and *B. contaminans* strain MS14, (c) *Burkholderia* sp. 383 and *B. cenocepacia* strain 895, (d) *B. cepacia* strain LO6 and *B. cenocepacia* strain 895. (PDF 25,604 kb)

Additional file 7: Figure S7. Tanglegrams showing the differences between the tree topologies based on the protein sequence similarity of single-copy universal genes and the tree topologies based on the syntenic blocks arrangements. (a) *B. mallei* clade; (b) *B. thailandensis* clade; (c) *B. cepacia* group. (PDF 524 kb)

Additional file 8: Figure S8. Phylogenetic tree constructed based on the concatenation of alignments of genes forming the genomic island. (PDF 6 kb)

Additional file 9: Figure S9. Phylogenetic tree showing detected events of positive selection. (PDF 8 kb)

Additional file 10: Table S1. Analyzed *Burkholderia* strains with genomes characteristics. (TXT 12 kb)

Additional file 11: Table S2. Genes evolving under positive selection. (CSV 17 kb)

Additional file 12: Table S3. Linear models (a) of average ω (negative selection, estimated using M8); (b) expression level [35]. (PDF 92 kb)

Acknowledgements

We thank Pavel Shelyakin for valuable comments. Analysis of parallel inversions was performed by Alisa Rodionova at the Summer School of Molecular and Theoretical Biology (Barcelona, 2016), supported by the Zimin Foundation.

Funding

The study was supported by the Russian Science Foundation under grant 18-14-00358. Analysis of gene selection was supported by the Russian Foundation of Basic Research (grant 16-54-21004) and Swiss National Science Foundation (grant number IZLRZ3_163872) and performed in part at the Vital-IT center for high-performance computing of the Swiss Institute of Bioinformatics. The funding body had no role in the design of the study and collection, analysis, and interpretation of data and in writing the manuscript.

Availability of data and materials

The datasets supporting the conclusions of this article and used ad hoc scripts are available via the link <https://github.com/OlgaBochkaryova/burkholderia-genomics>.

Authors' information

MSG mikhail.gelfand@iitp.ru, OOB olga.bochkaryova@gmail.com, EVM elena.v.lopatina@gmail.com, IID iakov.davydov@gmail.com.

Authors' contributions

MSG conceived the study, OOB, EVM and MSG designed the study; EVM, OOB and IID developed the methods, analyzed the data; EVM, OOB and IID wrote the manuscript, MSG reviewed the paper. All authors read and approved the final version of the manuscript.

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Author details

¹Kharkevich Institute for Information Transmission Problems, Moscow, Russia. ²Center of Life Sciences Skolkovo Institute of Science and Technology, Moscow, Russia. ³Department of Ecology and Evolution & Department of Computational Biology, University of Lausanne, Lausanne, Switzerland. ⁴Swiss Institute of Bioinformatics, Lausanne, Switzerland. ⁵Faculty of Computer Science, Higher School of Economics, Moscow, Russia.

Received: 29 April 2018 Accepted: 14 November 2018

Published online: 27 December 2018

References

- Coenye T, Vandamme P. Diversity and significance of *Burkholderia* species occupying diverse ecological niches. *Environ Microbiol.* 2003;5:719–29.
- Howe C, Sampath A, Spotnitz M. The *pseudomallei* group: a review. *J Infect Dis.* 1971;124:598–606.
- Ham JH, Melanson RA, Rush MC. *Burkholderia glumae*: next major pathogen of rice? *Mol Plant Pathol.* 2011;12:329–39.
- Groris J, De Vos P, Caballero-Mellado J, Park J, Falsen E, Quensen Jr, Tiedje J, Vandamme P. Classification of the biphenyl- and polychlorinated biphenyl-degrading strain LB400T and relatives as *Burkholderia xenovorans* sp. nov. *Int J Syst Evol Microbiol.* 2004;54:1677–81.
- Rommel M. I., Nowak J., Lazarovits G. Growth enhancement and developmental modifications of in vitro grown potato (*Solanum tuberosum* spp. *tuberosum*) as affected by a nonfluorescent *Pseudomonas* sp. *Plant Physiol.* 1991;96:928–36.
- Tettelin H, Massignani V, Cieslewicz M, Donati C, Medini D, Ward N, Angiuoli S, Crabtree J, Jones A, Durkin A, Deboy R, Davidsen T, Mora M, Scarselli M, Margarit y Ros I, Peterson J, Hauser C, Sundaram J, Nelson W, Madupu R, Brinkac L, Dodson R, Rosovitz M, Sullivan S, Daugherty S, Haft D, Selengut J, Gwinn M, Zhou L, Zafar N, Khouri H, Radune D, Dimitrov G, Watkins K, O'Connor K, Smith S, Utterback T, White O, Rubens C, Grandi G, Madoff L, Kasper D, Telford J, Wessels M, Rappuoli R, Fraser C. Genome analysis of multiple pathogenic isolates of *Streptococcus agalactiae*: implications for the microbial "pan-genome". *Proc Natl Acad Sci U S A.* 2005;102:13950–5.
- Ussery DW, Kiel K, Lagesen K, Sicheritz-Pontén T, Bohlin J, Wassenaar T. The genus *Burkholderia*: analysis of 56 genomic sequences. *Genome Dyn.* 2009;6:140–57.
- Spring-Pearson S, Stone J, Doyle A, Allender C, Okinaka R, Mayo M, Broomall S, Hill J, Karavis M, Hubbard K, Insalaco J, McNew L, Rosenzweig C, Gibbons H, Currie B, Wagner D, Keim P, Tuanyok A. Pangenome analysis of *Burkholderia pseudomallei*: Genome evolution

- preserves gene order despite high recombination rates. *PLoS ONE*. 2015;10(10):0140274.
9. Losada L, Ronning C, DeShazer D, Woods D, Fedorova N, Kim H, Shabalina S, Pearson T, Brinkac L, Tan P, Nandi T, Crabtree J, Badger J, Beckstrom-Sternberg S, Saqib M, Schutzer S, Keim P, Nierman W. Continuing evolution of *Burkholderia mallei* through genome reduction and large-scale rearrangements. *Genome Biol Evol*. 2010;2:102–16.
 10. Egan ES, Fogel MA, Waldor MK. Divided genomes: negotiating the cell cycle in prokaryotes with multiple chromosomes. *Mol Microbiol*. 2005;56:1129–38.
 11. Mackenzie C, Choudhary M, Larimer F, Predki P, Stilwagen S, Armitage J, Barber R, Donohue T, Hosler J, Newman J, Shapleigh J, Sockett R, Zeilstra-Ryalls J, Kaplan S. The home stretch, a first analysis of the nearly completed genome of *Rhodobacter sphaeroides* 2.4.1. *Photosynth Res*. 2001;70:19–41.
 12. Cooper VS, Vohr S, Wrocklage S, Hatcher P. Why genes evolve faster on secondary chromosomes in bacteria. *PLoS Comput Biol*. 2010;6:1000732.
 13. Guo FB, Ning L, Huang J, Lin H, Zhang H. Chromosome translocation and its consequence in the genome of *Burkholderia cenocepacia* AU-1054. *Biochem Biophys Res Commun*. 2010;403:375–9.
 14. Morrow JD, Cooper VS. Evolutionary effects of translocations in bacterial genomes. *Genome Biol Evol*. 2012;4:1256–62.
 15. Alekseyev MA, Pevzner PA. Breakpoint graphs and ancestral genome reconstructions. *Genome Res*. 2009;19:943–57.
 16. Hu F, Lin Y, Tang J. MLGO: phylogeny reconstruction and ancestral inference from gene-order data. *BMC Bioinformatics*. 2014;15:354.
 17. Raeside C, Gaffé J, Deatherage D, Tenaillon O, Briska A, Ptashkin R, Cruveiller S, Médigue C, Lenski R, Barrick J, Schneider D. Large chromosomal rearrangements during a long-term evolution experiment with *Escherichia coli*. *MBio*. 2014;5:01377–14.
 18. Huang W-C, Chen Y, Teng L, Lien H, Chen J, Chia J. Chromosomal inversion between *rrn* operons among *Streptococcus mutans* serotype c oral and blood isolates. *J Med Microbiol*. 2008;57:198–206.
 19. Eisen JA, Heidelberg J, White O, Salzberg S. Evidence for symmetric chromosomal inversions around the replication origin in bacteria. *Genome Biol*. 2000;1:0011.
 20. Kowalczyk M, Mackiewicz P, Mackiewicz D, Nowicka A, Dudkiewicz M, Dudek M, Cebrat S. DNA asymmetry and the replicational mutational pressure. *J Appl Genet*. 2001;42(4):553–77.
 21. García-Pastor L, Puerta-Fernández E, Casadesús J. Bistability and phase variation in *Salmonella enterica*. *Biochim Biophys Acta*. 2018;S1874-9399(17):30286–9.
 22. Challacombe J, Stubben C, Klimko C, Welkos S, Kern S, Bozue J, Worsham P, Cote C, Wolfe D. Interrogation of the *Burkholderia pseudomallei* genome to address differential virulence among isolates. *PLoS ONE*. 2014;9(12):115951.
 23. Yu Y, Kim H, Chua H, Lin C, Sim S, Lin D, Derr A, Engels R, DeShazer D, Birren B, Nierman W, Tan P. Genomic patterns of pathogen evolution revealed by comparison of *Burkholderia pseudomallei*, the causative agent of melioidosis, to avirulent *Burkholderia thailandensis*. *BMC Microbiol*. 2006;26(6):46.
 24. Nierman W, DeShazer D, Kim H, Tettelin H, Nelson K, Feldblyum T, Ulrich R, Ronning C, Brinkac L, Daugherty S, Davidsen T, Deboy R, Dimitrov G, Dodson R, Durkin A, Gwinn M, Haft D, Khouri H, Kolonay J, Madupu R, Mohammud Y, Nelson W, Radune D, Romero C, Sarria S, Selengut J, Shamblyn C, Sullivan S, White O, Yu Y, Zafar N, Zhou L, Fraser C. Structural flexibility in the *Burkholderia mallei* genome. *Proc Natl Acad Sci U S A*. 2004;101(39):14246–51.
 25. Godoy D, Randle G, Simpson A, Aanensen D, Pitt T, Kinoshita R, Spratt B. Multilocus sequence typing and evolutionary relationships among the causative agents of melioidosis and glanders, *Burkholderia pseudomallei* and *Burkholderia mallei*. *J Clin Microbiol*. 2003;41:2068–79.
 26. Schutzer SE, Schlater L, Ronning C, DeShazer D, Luft B, Dunn J, Ravel J, Fraser-Liggett C, Nierman W. Characterization of clinically-attenuated *Burkholderia mallei* by whole genome sequencing: candidate strain for exclusion from Select Agent lists. *PLoS ONE*. 2008;3:2058.
 27. Moore RA, Reckseidler-Zenteno S, Kim H, Nierman W, Yu Y, Tuanyok A, Warawa J, DeShazer D, Woods D. Contribution of gene loss to the pathogenic evolution of *Burkholderia pseudomallei* and *Burkholderia mallei*. *Infect Immun*. 2004;72:4172–87.
 28. Lee HH, Park J, Kim J, Park I, Seo YS. Understanding the direction of evolution in *Burkholderia glumae* through comparative genomics. *Curr Genet*. 2016;62:115–23.
 29. Benson DA, Cavanaugh M, Clark K, Karsch-Mizrachi I, Lipman DJ, Ostell J, Sayers EW. GenBank. *Nucleic Acids Res*. 2013;45(Database issue):30–5.
 30. Lechner M, Findeiss S, Steiner L, Marz M, Stadler P, Prohaska S. Proteinortho: Detection of (co-)orthologs in large-scale analysis. *BMC Bioinformatics*. 2011;12(1):124.
 31. Jones P, David Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J, Alex M, Nuka G, Pesseat S, Quinn AF, Sangrador-Vegas A, Scheremetjew M, Yong S-Y, Lopez R, Hunter S. InterProScan 5: genome-scale protein function classification. *Bioinformatics*. 2014;30(9):1236–40.
 32. Alexa A, Rahnenfuhrer J. TopGO: Enrichment analysis for gene ontology. R package version 2.30.0. 2016.
 33. Huerta-Cepas J, Szklarczyk D, Forslund K, Cook H, Heller D, Walter MC, Rattei T, Mende DR, Sunagawa S, Kuhn M, Jensen LJ, von Mering C, Bork P. EggNOG 4.5: a hierarchical orthology framework with improved functional annotations for eukaryotic, prokaryotic and viral sequences. *Nucleic Acids Res*. 2016;44(4):286–93.
 34. Yu N, Wagner J, Laird M, Melli G, Rey S, Lo R, Dao P, Sahinalp S, Ester M, Foster L, Brinkman F. PSORTb 3.0: improved protein subcellular localization prediction with refined localization subcategories and predictive capabilities for all prokaryotes. *Bioinformatics*. 2017;26:1608–15.
 35. Lazar Adler NR, Allwood EM, Deveson Lucas D, Harrison P, Watts S, Dimitropoulos A, Treerat P, Alwis P, Devenish RJ, Prescott M, Govan B, Adler B, Harper M, Boyce JD. Perturbation of the two-component signal transduction system, BprRS, results in attenuated virulence and motility defects in *Burkholderia pseudomallei*. *BMC Genomics*. 2016;17:331.
 36. Snipen L, Almoy T, Ussery DW. Microbial comparative pan-genomics using binomial mixture models. *BMC Genomics*. 2009;10:385.
 37. Chao A. Estimating the population size for capture-recapture data with unequal catchability. *Biometrics*. 1987;43:783–91.
 38. Snipen L, Liland KH. Micropan: an R-package for microbial pan-genomics. *BMC Bioinformatics*. 2015;16:79.
 39. Akaike H. A new look at the statistical model identification. *IEEE Trans Autom Control*. 1974;19:716–23.
 40. Hurvich CM, Tsai C-L. Regression and time series model selection in small samples. *Biometrika*. 1989;76:297–307.
 41. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7: Improvements in performance and usability. *Mol Biol Evol*. 2013;30(4):772–80.
 42. Penn O, Privman E, Ashkenazy H, Landan G, Graur D, Pupko T. GUIDANCE: a web server for assessing alignment confidence scores. *Nucleic Acids Res*. 2010;38(Web Server issue):W23–8.
 43. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics*. 2014;30(9):1312–3.
 44. Guindon S, Dufayard J-F, Lefort V, Anisimova M, Hordijk W, Gascuel O. New algorithms and methods to estimate maximum-likelihood phylogenies: Assessing the performance of PhyML. *Syst Biol*. 2010;59(3):307–21.
 45. Paradis E, Claude J, Strimmer K. Ape: analyses of phylogenetics and evolution in R language. *Bioinformatics*. 2004;20:289–90.
 46. Cohen O, Ashkenazy H, Belinky F, Huchon D, Pupko T. GLOOME: gain loss mapping engine. *Bioinformatics*. 2010;26:2914–5.
 47. Yang Z. PAML: a program package for phylogenetic analysis by maximum likelihood. *Comput Appl Biosci*. 1997;13:555–6.
 48. Minkin I, Patel A, Kolmogorov M, et al. Sibelia: a scalable and comprehensive synteny block generation tool for closely related microbial genomes. In: Darling A, Stoye J, editors. *Algorithms in Bioinformatics*, number 8126 in Lecture Notes in Computer Science. Springer-Verlag: Berlin; 2013. p. 215–29.
 49. Pham SK, Pevzner PA. DRIMM-Synteny: decomposing genomes into evolutionary conserved segments. *Bioinformatics*. 2010;26:2509–16.
 50. Avdeyev P, Jiang S, Aganezov S, Hu F, Alekseyev MA. Reconstruction of ancestral genomes in presence of gene gain and loss. *J Comput Biol*. 2016;23(3):150–64.
 51. Gao F, Zhang C-T. Ori-Finder: a web-based system for finding *oriC*s in unannotated bacterial genomes. *BMC Bioinformatics*. 2008;9(1):79.
 52. Mirarab S, Bayzid MS, Boussau B, Warnow T. Statistical binning enables an accurate coalescent-based estimation of the avian tree. *Science*. 2014;346(6215):1250463.

53. Yang Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol.* 2007;24:1586–91.
54. Zhang J, Nielsen R, Yang Z. Evaluation of an improved branch-site likelihood method for detecting positive selection at the molecular level. *Mol Biol Evol.* 2005;22:2472–9.
55. Drummond DA, Bloom JD, Adami C, Wilke CO, Arnold FH. Why highly expressed proteins evolve slowly. *PNAS.* 2005;102(40):14338–43.
56. Gordienko EN, Kazanov MD, Gelfand MS. Evolution of pan-genomes of *Escherichia coli*, *Shigella* spp., and *Salmonella enterica*. *J Bacteriol.* 2013;195:2786–92.
57. Moldovan M, Gelfand M. Pangenomic definition of prokaryotic species and the phylogenetic structure of *Prochlorococcus* spp. *Front Microbiol.* 2018;9:428.
58. Makarova KS, Sorokin A, Novichkov P, Wolf Y, Koonin E. Clusters of orthologous genes for 41 archaeal genomes and implications for evolutionary genomics of archaea. *Biol Direct.* 2007;2:33.
59. Shelyakin PV, Bochkareva O, Karan AA, Gelfand MS. Comparative analysis of *Streptococcus* genomes. *bioRxiv.* 2018;447938.
60. Wolf YI, Makarova KS, Yutin N, Koonin EV. Updated clusters of orthologous genes for archaea: a complex ancestor of the archaea and the byways of horizontal gene transfer. *Biol Direct.* 2012;7(1):46.
61. Koonin EV, Wolf YI. Genomics of bacteria and archaea: the emerging dynamic view of the prokaryotic world. *Nucleic Acids Res.* 2008;36(21):6688–719.
62. Nandi T, Holden M, Didelot X, Mehershahi K, Boddey J, Beacham I, Peak I, Harting J, Baybayan P, Guo Y, Wang S, How L, Sim B, Essex-Lopresti A, Sarkar-Tyson M, Nelson M, Smither S, Ong C, Aw L, Hoon C, Michell S, Studholme D, Titball R, Chen S, Parkhill J, Tan P. *Burkholderia pseudomallei* sequencing identifies genomic clades with distinct recombination, accessory, and epigenetic profiles. *Genome Res.* 2015;25:129–41.
63. Bochkareva OO, Dranenko NO, Ocheredko ES, Kanevsky GM, Lozinsky YN, Khalaycheva VA, Artamonova II, Gelfand MS. Genome rearrangements and phylogeny reconstruction in *Yersinia pestis*. *PeerJ.* 2018;6:4545.
64. Cao P, Guo D, Liu J, Jiang Q, Xu Z, Qu L. Genome-wide analyses reveal genes subject to positive selection in *Pasteurella multocida*. *Front Microbiol.* 2017;8:961.
65. Xu Z, Chen H, Zhou R. Genome-wide evidence for positive selection and recombination in *Actinobacillus pleuropneumoniae*. *BMC Evol Biol.* 2011;11:203.
66. Sousa SA, Morad M, Feliciano JR, Pita T, Nady S, El-Hennamy RE, Abdel-Rahman M, Cavaco J, Pereira L, Barreto C, Leitão JH. Outer membrane protein A and OprF – versatile roles in gram-negative bacterial infections. *FEBS J.* 2012;279(6):919–31.
67. Yang Z, dos Reis M. Statistical properties of the branch-site test of positive selection. *Mol Biol Evol.* 2011;28:1217–28.
68. Dryselius R, Izutsu K, Honda T, Iida T. Differential replication dynamics for large and small vibrio chromosomes affect gene dosage, expression and location. *BMC Genomics.* 2008;9:559.
69. Novichkov PS, Wolf YI, Dubchak I, Koonin EV. Trends in prokaryotic evolution revealed by comparison of closely related bacterial and archaeal genomes. *J Bacteriol.* 2009;191(1):65–73.
70. Darling AE, Miklós I, Ragan MA. Dynamics of genome rearrangement in bacterial populations. *PLoS Genet.* 2008;4(7):1000128.
71. Repar J, Warnecke T. Non-random inversion landscapes in prokaryotic genomes are shaped by heterogeneous selection pressures. *Mol Biol Evol.* 2017;34(8):1902–22.
72. Zhang G, Gao F. Quantitative analysis of correlation between AT and GC biases among bacterial genomes. *PLoS ONE.* 2017;12(2):0171408.
73. Oliveira P, Touchon M, Cury J, Rocha E. The chromosomal organization of horizontal gene transfer in bacteria. *Nat Commun.* 2017;8:841.
74. Koonin E. Horizontal gene transfer: essentiality and evolvability in prokaryotes, and roles in evolutionary transitions. *F1000Research.* 2016;5(F1000 Faculty Rev):1805.
75. McInerney JO, McNally A, O'Connell MJ. Why prokaryotes have pangenomes. *Nat Microbiol.* 2017;2:17404.
76. Collingro A, Tischler P, Weinmaier T, Penz T, Heinz E, Brunham R, Read T, Bavoil P, Sachse K, Kahane S, Friedman M, Rattei T, Myers G, Horn M. Unity in variety—the pan-genome of the *Chlamydiae*. *Mol Biol Evol.* 2011;28(12):3253–70.
77. Price MN, Alm EJ, Arkin AP. Interruptions in gene expression drive highly expressed operons to the leading strand of DNA replication. *Nucleic Acids Res.* 2005;33(10):3224–34.
78. diCenzo GC, Finan TM. The divided bacterial genome: structure, function, and evolution. *Microbiol Mol Biol Rev.* 2017;81(3):00019–17.
79. Mira A, Pushker R, Rodriguez-Valera F. The neolithic revolution of bacterial genomes. *Trends Microbiol.* 2006;14(5):200–6.
80. Sekulovic O, Garrett EM, Bourgeois J, Tamayo R, Shen A, Camilli A. Genome-wide detection of conservative sitespecific recombination in bacteria. *PLoS Genet.* 2018;4:1007332.
81. Donati C, Hiller N, Tettelin H, Muzzi A, Croucher N, Angiuoli S, Oggioni M, Dunning Hotopp J, Hu F, Riley D, Covacci A, Mitchell T, Bentley S, Kilian M, Ehrlich G, Rappuoli R, Moxon E, Massignani V. Structure and dynamics of the pan-genome of *Streptococcus pneumoniae* and closely related species. *Genome Biol.* 2010;11(10):107.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

