





The Effects of Signal Erosion and Core Genome Reduction on the Identification of Diagnostic Markers

Jason W. Sahl,^{a,b} Adam J. Vazquez,^a Carina M. Hall,^a Joseph D. Busch,^a Apichai Tuanyok,^c Mark Mayo,^d James M. Schupp,^b Madeline Lummis,^a Talima Pearson,^a Kenzie Shippy,^a  Rebecca E. Colman,^b Christopher J. Allender,^a Vanessa Theobald,^d  Derek S. Sarovich,^d  Erin P. Price,^d Alex Hutcheson,^e Jonas Korchach,^e John J. LiPuma,^f Jason Ladner,^g Sean Lovett,^g Galina Koroleva,^g Gustavo Palacios,^g Direk Limmathurotsakul,^{h,i} Vanaporn Wuthiekanun,^h Gumphol Wongsuwan,^h Bart J. Currie,^d  Paul Keim,^{a,b} David M. Wagner^a

Center for Microbial Genetics and Genomics, Northern Arizona University, Flagstaff, Arizona, USA^a; Translational Genomics Research Institute, Flagstaff, Arizona, USA^b; Emerging Pathogens Institute, University of Florida, Gainesville, Florida, USA^c; Global and Tropical Health Division, Menzies School of Health Research, Darwin, Northern Territory, Australia^d; Pacific Biosciences, University of Michigan, Ann Arbor, Michigan, USA^e; Division of Pediatric Infectious Diseases, University of Michigan, Ann Arbor, Michigan, USA^f; Center for Genome Sciences, USAMRIID, Fort Detrick, Maryland, USA^g; Mahidol-Oxford Tropical Medicine Research Unit, Mahidol University, Bangkok, Thailand^h; Department of Tropical Hygiene, Faculty of Tropical Medicine, Mahidol University, Bangkok, Thailandⁱ

ABSTRACT Whole-genome sequence (WGS) data are commonly used to design diagnostic targets for the identification of bacterial pathogens. To do this effectively, genomics databases must be comprehensive to identify the strict core genome that is specific to the target pathogen. As additional genomes are analyzed, the core genome size is reduced and there is erosion of the target-specific regions due to commonality with related species, potentially resulting in the identification of false positives and/or false negatives.

IMPORTANCE A comparative analysis of 1,130 *Burkholderia* genomes identified unique markers for many named species, including the human pathogens *B. pseudomallei* and *B. mallei*. Due to core genome reduction and signature erosion, only 38 targets specific to *B. pseudomallei/mallei* were identified. By using only public genomes, a larger number of markers were identified, due to undersampling, and this larger number represents the potential for false positives. This analysis has implications for the design of diagnostics for other species where the genomic space of the target and/or closely related species is not well defined.

Received 11 May 2016 Accepted 10 August 2016 Published 20 September 2016

Citation Sahl JW, Vazquez AJ, Hall CM, Busch JD, Tuanyok A, Mayo M, Schupp JM, Lummis M, Pearson T, Shippy K, Colman RE, Allender CJ, Theobald V, Sarovich DS, Price EP, Hutcheson A, Korchach J, LiPuma JJ, Ladner J, Lovett S, Koroleva G, Palacios G, Limmathurotsakul D, Wuthiekanun V, Wongsuwan G, Currie BJ, Keim P, Wagner DM. 2016. The effects of signal erosion and core genome reduction on the identification of diagnostic markers. *mBio* 7(5):e00846-16. doi:10.1128/mBio.00846-16.

Editor Julian Parkhill, Sanger Institute

Copyright © 2016 Sahl et al. This is an open-access article distributed under the terms of the [Creative Commons Attribution 4.0 International license](https://creativecommons.org/licenses/by/4.0/).

Address correspondence to David M. Wagner, dave.wagner@nau.edu.

Whole-genome sequence (WGS) data are routinely used to develop DNA-based diagnostics for rapid and accurate identification of clinical pathogens (1, 2). Validating the specificity of diagnostic targets ensures that assays do not produce false positives (identifying a nonpathogen as a pathogen) or false negatives (not identifying a pathogen that is actually present). To avoid false positives and negatives, DNA-based diagnostics must be conserved across the target species and absent from nontarget species.

Two critical issues arise during the process of identifying specific diagnostics from bacterial genomes. First, the number of genes in the core genome (i.e., genes present in every individual of a species) tends to become smaller as the number of sequenced genomes increases (3, 4). Certain pathogens (e.g., *Yersinia pestis*) propagate clonally, are highly homogeneous, and show little variation in core genome size with additional sampling (5). In this case, the core genome size is not expected to become drastically smaller as more genomes are analyzed. In contrast, the core genome size of *Burkholderia pseudomallei* becomes significantly smaller with each new genome added (6). A second issue arises from genomes of related species, or “near neighbors,” that share

core genes with the target species. In a process of signature erosion, this genomic overlap often increases as near-neighbor genomes are added to the analysis, thus eroding the number of potential diagnostic targets. Unfortunately, near neighbors are often undersampled (or not sampled at all) during the search for diagnostic targets, which hinders efforts to identify species-specific targets.

Burkholderia represents a model genus for the demonstration of core genome reduction and signal erosion. The *Burkholderia* genus contains a diverse set of species, including plant pathogens (7) and human pathogens, such as *B. pseudomallei*, the causative agent of melioidosis (8), and *B. mallei*, the causative agent of glanders (9). The *pseudomallei* group includes *B. pseudomallei*, *B. mallei*, *B. oklahomensis*, *B. thailandensis*, and the newly described *B. humptydooensis* (10). The *B. cepacia* complex (Bcc) is a diverse group within *Burkholderia* that is associated with opportunistic infections and is comprised of at least 20 genomic species (11, 12). Most of the relationships between these species have been determined through gene marker analyses, such as the *recA* gene (13, 14) or multilocus sequence typing (15).

From a genomics perspective, *Burkholderia* whole-genome se-

TABLE 1 Summary of new genomes sequenced as part of this study

Clade	No. of genomes
<i>B. anthina</i>	8
<i>B. cenocepacia</i> 1	1
<i>B. cenocepacia</i> 2	4
<i>B. cepacia</i>	78
<i>B. diffusa</i>	12
<i>B. gladioli</i>	1
<i>B. humptydoensis</i>	5
<i>B. lata</i>	2
<i>B. latens</i>	2
<i>B. metallica</i>	1
<i>B. multivorans</i>	14
<i>B. oklahomensis</i>	2
Putative species 1	3
Putative species 2	4
Putative species 3	10
Putative species 4	7
Putative species 5	8
<i>B. pseudomallei</i>	256
<i>B. pseudomultivorans</i>	9
<i>B. pyrrocinia</i>	1
<i>B. seminalis</i>	2
<i>B. stagnalis</i>	67
<i>B. thailandensis</i>	8
<i>B. territorii</i>	33
<i>B. ubonensis</i>	254
<i>B. vietnamiensis</i>	37
Total	829

quencing efforts have focused on *B. pseudomallei* (16) and *B. mallei* (17). Recent studies have begun to sequence other *Burkholderia* spp., including members of the Bcc (18). However, large-scale, whole-genome, phylogenetics-based studies that define the overall phylogenetic structure among *Burkholderia* species using high-resolution methods are currently lacking.

In this study, we extensively surveyed the environment in Australia, the United States, and Southeast Asia for *Burkholderia* spp. We sequenced a large collection of genomes to (i) explore the genomic diversity of *Burkholderia* spp. that grow on Ashdown's agar, (ii) identify specific diagnostic markers for *B. pseudomallei* and *B. mallei*, and (iii) understand the sampling effects of core genome size reduction and signal erosion on the selection of highly specific diagnostic targets.

RESULTS

Whole-genome sequencing of *Burkholderia* spp. In this study, we analyzed the whole-genome sequences of 829 *Burkholderia* spp. that grow on Ashdown's agar (Table 1), a selective medium containing the aminoglycoside gentamicin. These isolates were collected from diverse geographic locations in the United States, Thailand, and Australia (see Table S1 in the supplemental material). To understand the effects of core genome reduction and signature erosion on the identification of highly specific diagnostic targets, the genomes of 256 diverse *B. pseudomallei/mallei* strains were sequenced, assembled, and deposited in public databases (see Table S1); these genomes were combined with 160 *B. pseudomallei/mallei* genome assemblies already in public databases. Most of the genomes ($n = 779$) in this study were sequenced on the Illumina platform, with 50 genomes also sequenced on the PacBio platform, which generated highly contiguous and often finished assemblies (see Table S1).

Core genome SNP phylogeny. To understand the phylogenetic structure of the *Burkholderia* genus, genomes sequenced in this study, as well as GenBank reference genomes (see Table S2 in the supplemental material), were aligned against *B. pseudomallei* strain K96243 (19) with NUCmer (20) and single-nucleotide polymorphisms (SNPs) were identified with NASP. The maximum-likelihood phylogeny inferred from core, orthologous SNPs ($n = 105,877$) demonstrated that all genomes sequenced in this study, with the exception of 1 *B. gladioli* genome, grouped in either the *Burkholderia cepacia* complex (Bcc) or the *B. pseudomallei* group (Fig. 1). Based on the monophyletic nature and complexity of the latter clade, we propose to name it the *B. pseudomallei* complex (Bpc). Multiple additional *Burkholderia* genomes from GenBank were analyzed and were found to be more distantly related to these two groups. As such, they were not examined in detail in this study but were included for marker-screening purposes (see "Paraburkholderia" genomes in Table S2). Our clade naming scheme is consistent with a recently published taxonomic scheme for *Burkholderia* (21).

The work performed in this study greatly expands the known genomic diversity of *Burkholderia*. For example, at the outset of this study, only two *B. ubonensis* genome assemblies were available in GenBank. This is likely due to the fact that most genome sequencing has focused on clinically relevant organisms, whereas we sampled both environmental and clinical isolates. This study adds the genomes of 254 *B. ubonensis* isolates, including three finished genomes (3 contigs) and three nearly finished genomes (4 to 5 contigs) (see Table S1 in the supplemental material). All of these genomes are publicly available and will help provide phylogenetic context for additional *Burkholderia* genomes that are sequenced, including from clinical isolates. We have also generated the first whole-genome sequences for other recently described species, such as *B. stagnalis* and *B. territorii* (12), including completed genomes, which will provide data for additional comparative studies.

Comparative genomics. Based on the topology of the core genome phylogeny (Fig. 1), pangenome statistics were calculated for each major clade (see Table 2 in the supplemental material) using the large-scale BLAST score ratio (LS-BSR) pipeline (22). The core genome of each primary clade was aligned against all surveyed genomes ($n = 1,130$) to identify species- or clade-specific markers. A marker was determined to be clade specific if it had a BSR of >0.8 in all target genomes and <0.4 in all nontarget genomes; although this definition is very conservative, it was used to identify discriminatory markers, regardless of genome assembly quality. The results demonstrate that species-specific markers were identified for most of the major clades (Table 2); a multi-FASTA file of all species-specific coding regions is publically available (<https://gist.github.com/jasonsahl/3e4132ca1d09b717fcc2>). A screen of these species-specific markers against all genomes was visualized to demonstrate their specificity to each targeted clade (Fig. 2). The stability of markers from clades with a limited number of representatives is unknown and will need to be validated with additional sequencing. Markers also were identified for the *B. cepacia* complex (Bcc) and the *B. pseudomallei* complex (Bpc), which can help to verify results obtained through diagnostic sequencing efforts.

Putative new species. Based on the phylogeny (Fig. 1), five divergent clades were identified that may represent novel species

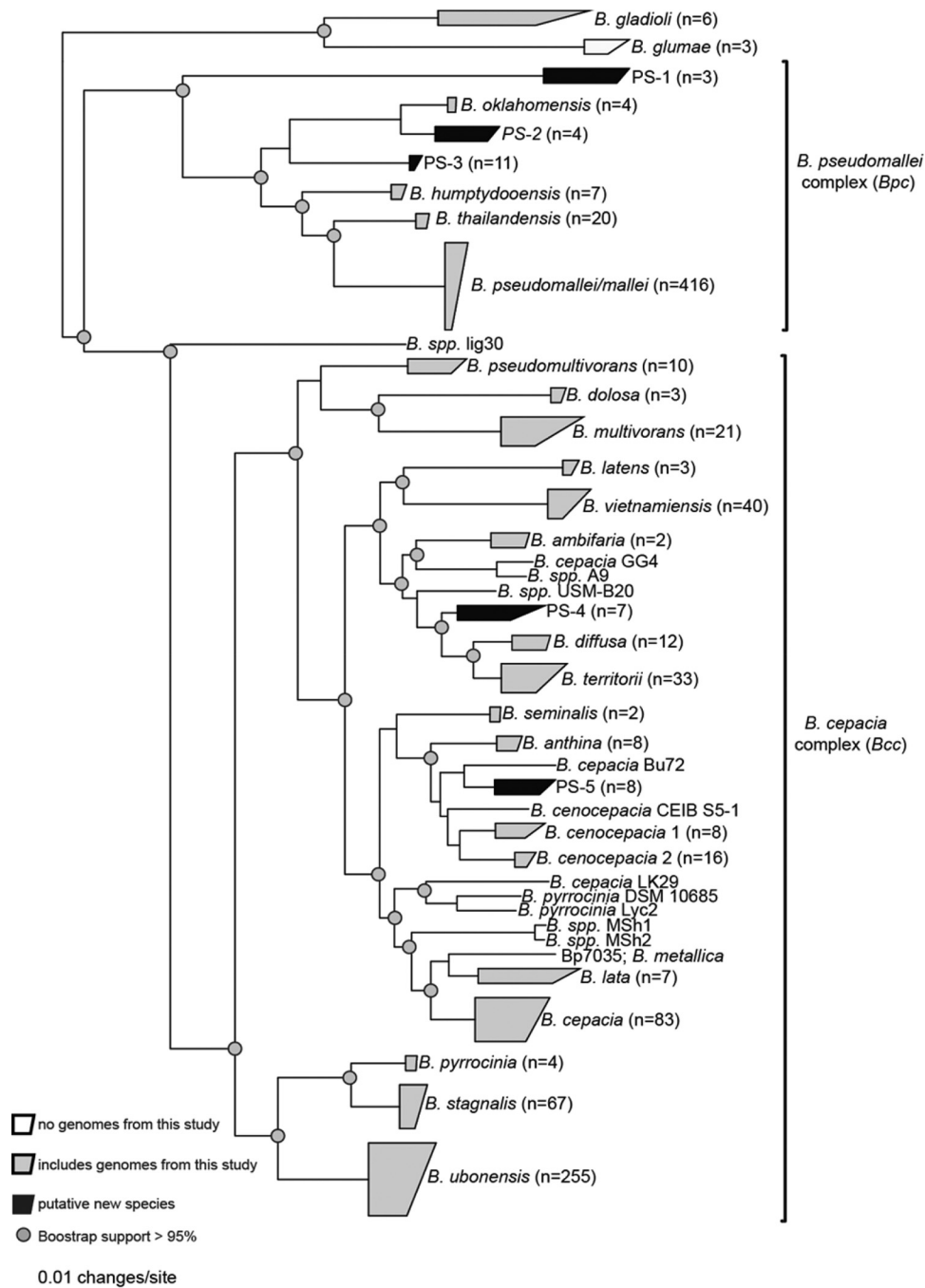


FIG 1 A core genome single-nucleotide-polymorphism (SNP) phylogeny of *Burkholderia* genomes. All SNPs were identified by aligning genome assemblies against the finished genome of *B. pseudomallei* K96243 (19) with NUCmer (20) and processed with the Northern Arizona SNP Pipeline (<http://tgennorth.github.io/NASPtgennorth.github.io/NASP>) (30). A maximum-likelihood phylogeny was inferred on the concatenated SNP alignment with RAxML version 8 (31) with 100 bootstrap replicates. Clades were collapsed with ARB (41). Putative novel species are named with PS (putative species) and the clade number.

(PS-1 through PS-5). We have generated completed or nearly completed genomes for at least one isolate from each of these clades (Table 1). A BLASTN alignment of the extracted *recA* sequences against the GenBank nucleotide database failed to identify a close match to a named species for any of these clades. To demonstrate the differences between genomes in these putative

species, one representative was compared against a genome of the nearest species, based on the closest patristic distance, or tree path distance, to the nearest monophyletic clade in the global phylogeny (Fig. 1). For each pairwise comparison, the average nucleotide identity (ANI) and DNA-DNA hybridization (DDH) values were calculated and tabulated (Table 3). The results demonstrate that

TABLE 2 Core genome statistics

Species/clade	Core genome size (CDSs)	No. of:	
		Genomes	Species-/clade-specific markers
<i>ambifaria</i>	5,408	2	71
<i>anthina</i>	5,507	8	13
<i>cenocapacia</i> 1	3,823	8	8
<i>cenocapacia</i> 2	5,076	16	22
<i>capacia</i>	4,415	83	7
<i>diffusa</i>	4,566	12	7
<i>dolosa</i>	5,451	3	436
<i>gladioli</i>	4,898	6	833
<i>glumae</i>	3,253	3	264
<i>humptydooensis</i>	5,115	7	157
<i>lata</i>	4,214	7	0
<i>latens</i>	5,348	3	105
<i>multivorans</i>	4,001	21	53
<i>oklahomensis</i>	5,681	4	141
PS-1	3,693	3	504
PS-2	4,231	4	23
PS-3	5,047	11	195
PS-4	4,366	7	0
PS-5	4,978	8	0
<i>pseudomallei</i>	2,339	392	22
<i>pseudomallei/mallei</i>	1,690	416	38
<i>pseudomultivorans</i>	4,549	10	62
<i>pyrrocinia</i>	6,397	4	153
<i>seminalis</i>	6,533	2	90
<i>stagnalis</i>	4,835	67	54
<i>thailandensis</i>	4,447	20	116
<i>territorii</i>	4,399	33	0
<i>ubonensis</i>	3,128	255	40
<i>vietnamiensis</i>	3,803	40	71

many of the clades have ANI values of <95% compared to the nearest reference genome based on its position in the phylogeny. Putative species 2 (PS-2), which is most closely related to *B. oklahomensis*, demonstrated ANI values on the border of the species threshold compared to *B. oklahomensis* genomes. All of the genomes from PS-2 have been isolated from Australia, whereas all *B. oklahomensis* genomes, including the two publicly available genomes, have been isolated from the United States (see Table S1 in the supplemental material). This physical separation, combined with the borderline ANI values, may argue for separate species, but biochemical testing is required to bolster this separation and is currently ongoing.

Core genome size reduction with additional sampling and signal erosion with the inclusion of near-neighbor genomes.

In bacteria with highly plastic genomes, the inclusion of additional isolates can cause the core genome size to decrease (3). To demonstrate this effect in *Burkholderia*, we calculated pangenome statistics on 416 *B. pseudomallei/mallei* genomes. The results demonstrated that as additional genomes were added to the analysis, the core genome size reduced to 1,684 coding DNA sequences (CDSs); annotation of these CDSs is provided in Table S3 in the supplemental material. This analysis included genomes from *B. mallei*, which has undergone significant evolutionary decay (9), and isolates from a chronic *B. pseudomallei* infection that have also undergone substantial genome reduction over time due to long-term host adaptation (23). By inclusion of a diverse set of genomes, the minimum set of genes required by all *B. pseudomallei/mallei* bacteria could be identified. From randomly subsampling

the 416 genomes at different genome levels, the sampling effect on the core genome size was visualized (Fig. 3A).

In addition to core genome size reduction, the effect of including additional near-neighbor genomes on accurate diagnostics was also investigated. The core genome from all *B. pseudomallei/mallei* genomes was aligned using LS-BSR against a randomly selected subset of near-neighbor genomes ranging from 10 to 300 genomes, with each iteration performed 100 times. By the time 300 near-neighbor genomes were randomly selected, the number of *B. pseudomallei/mallei* markers converged on the same number that was obtained using the entire set of 714 near-neighbor genomes (Fig. 3B). This result demonstrates that a significant number of near-neighbor genomes must be sequenced in order to identify a set of molecular markers that are highly discriminatory for a given clade.

When we considered the 416 *B. pseudomallei/mallei* genomes, a surprisingly small number of unique markers ($n = 38$) were identified (see Table S3 in the supplemental material). Of these markers, one in particular (BPSS0060; encodes a hypothetical protein) only contained 3 polymorphisms across all of the diverse *B. pseudomallei/mallei* isolates in our study. This gene represents a highly specific diagnostic target that is under low selection for mutation. If only *B. pseudomallei* was considered, 22 conserved markers were identified (see Table S4) in *B. pseudomallei* that were missing from *B. mallei* and all other *Burkholderia* genomes considered (Fig. 2).

If we only consider the publicly available genomes used in this study ($n = 298$), the core genome size for *B. pseudomallei/mallei* is

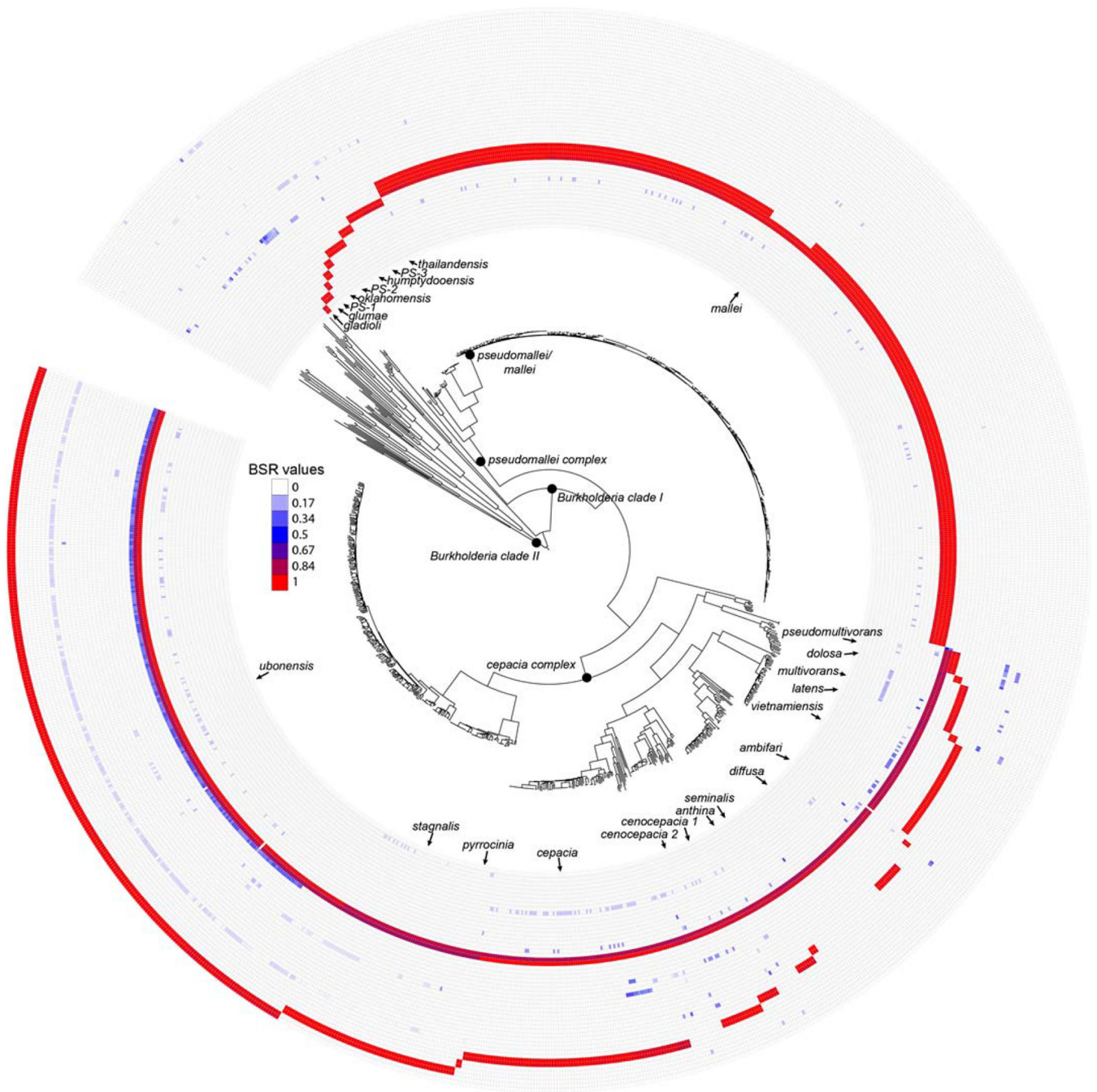


FIG 2 A core genome single-nucleotide-polymorphism (SNP) phylogeny associated with a heat map of markers unique to specific clades. The core genome phylogeny was inferred with RAXML (31) on a concatenated SNP alignment produced by aligning 1,130 genomes against the finished genome of *B. pseudomallei* K96243 (19), using NUCmer (20) in conjunction with NASP (<http://tgennorth.github.io/NASP/tgennorth.github.io/NASP>). Coding regions unique to specific clades were aligned against all genomes with LS-BSR (22), and the heat map was visualized with the Interactive Tree of Life (42). The heat map demonstrates the distribution of identified markers against all genomes screened in this study.

2,570 CDSs. When this core genome was screened against other *Burkholderia* near-neighbor genomes available in GenBank ($n = 141$), 63 markers were identified that were unique to *B. pseudomallei/mallei*. In contrast, if near-neighbor genomes sequenced in the study were also included ($n = 573$), 51 markers were identified. By not including additional nontarget reference genomes, 13 of these markers would represent false positives in screening

studies. If only target genomes in GenBank were considered, 25 false positives would be identified, demonstrating the need to include large numbers of target and nontarget genomes.

DISCUSSION

Accurate design of highly specific diagnostics is important for the detection of dangerous human pathogens in both environmental

TABLE 3 Average nucleotide identity and DNA-DNA hybridization values between representatives of putative novel species and representatives of established clades^a

Genome	Clade	Nearest genome	ANIm (%)	ANiB (%)	DDH range (%)
MSMB175	Putative species 1	<i>B. gladioli</i> BSR3	85.5	79.8	18.7–23.7
BDU8	Putative species 2	<i>B. oklahomensis</i> C6786	94.9	94.8	59.3–75.8
MSMB0852	Putative species 3	<i>Burkholderia</i> sp. strain MSMB43	92.4	91.1	44.5–52.7
MSMB0856	Putative species 4	<i>B. pyrrocinia</i> lyc 2	91.2	89.8	44.9–60.8
NRF60-BP8	Putative species 5	<i>B. cenocepacia</i> KC-01	94.1	93.5	54.5–56.9

^a ANI, average nucleotide identity; ANIm, uses NUCmer alignments; ANiB, uses BLASTN alignments; DDH, DNA-DNA hybridization.

and clinical settings. Timely pathogen identification directly from clinical specimens could inform the early treatment of potentially deadly infections. However, as our study demonstrates, the genomic targets of molecular assays need to first be thoroughly validated to avoid false positives and false negatives, which can potentially confound diagnostic tests and delay appropriate patient treatment. In this study, we highlight the importance of exploring the strict core genome size and signature erosion before designing diagnostic PCR targets. Our genome-based approach is applicable for other researchers who wish to develop diagnostic assays for other pathogens.

The effects of signature erosion and core genome size reduction were highlighted in the genus *Burkholderia*. To characterize the genomic space within *B. pseudomallei/mallei*, as well as in closely related genomes, we sequenced 829 *Burkholderia* genomes from diverse locations. A large-scale comparative genomics analysis of these genomes demonstrated that specific molecular markers were identified for many of the major *Burkholderia* clades identified from the core genome SNP phylogeny (Fig. 2). These unique coding regions were likely acquired horizontally, based on the lack of homology of these regions to other lineages within the genus (Fig. 2). To demonstrate the need to sequence a large collection of genomes to identify specific diagnostic targets, a core genome reduction analysis was performed (Fig. 3A). This analysis demon-

strated that sequencing additional genomes causes the core genome size to decline. This analysis was performed by including a large number of draft genome assemblies, which may cause genomic elements to be truncated, based on unresolvable repeats, or missing altogether, based on either insufficient coverage or assembly algorithms that remove either short contigs or regions of anomalous coverage. Based on the genome panel used in this analysis, including a large and diverse set of isolates is important to avoid selecting potential diagnostic targets that are susceptible to false-negative results when screening either clinical or environmental samples. If only genomes in public databases were selected, multiple markers would be identified that represent potential false negatives.

The other important factor to consider when designing diagnostic markers is the effect of signature erosion that can be introduced due to the inclusion of close relatives to the clade of interest. If only genomes available in GenBank were included in the analysis, 63 markers were identified that appeared to be specific to *B. pseudomallei/mallei*. However, if all non-*pseudomallei/mallei* genomes were included in the analysis, only 51 *B. pseudomallei/mallei* markers were identified, demonstrating the impact of including a comprehensive set of genomes outside the targeted species or clade. If all genomes from our study were included, only 38 *B. pseudomallei/mallei*-specific markers were identified, which

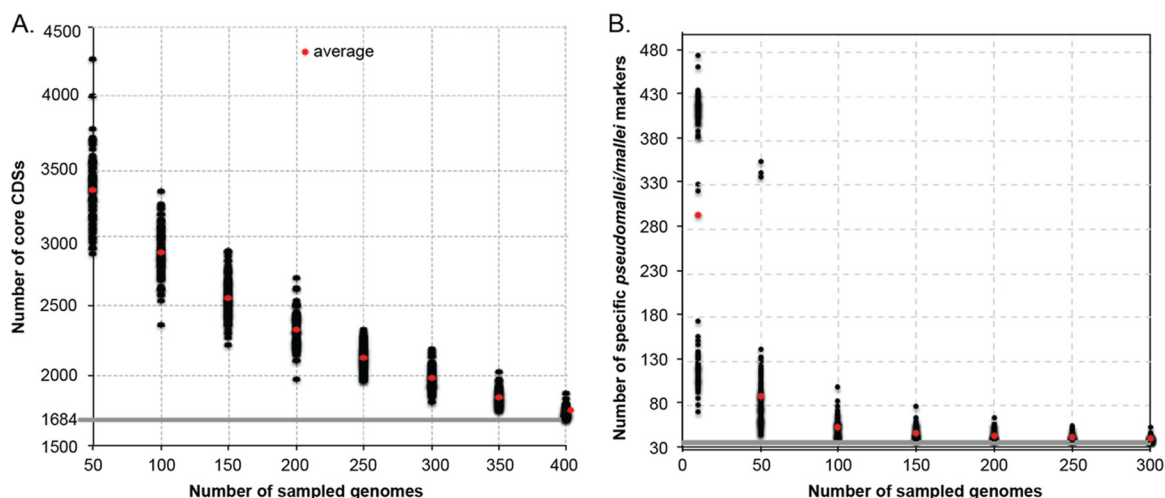


FIG 3 (A) Core genome reduction in *Burkholderia pseudomallei/mallei*. The core genome was calculated with the LS-BSR pipeline (22) on 416 genomes. For subsampling, genomes were randomly selected at different depths and the number of coding regions (CDSs) with a BLAST score ratio (BSR) of >0.8 in all genomes was calculated and plotted. For each subsampling level, 100 iterations were performed. The mean value at each level is shown in red, and each replicate is shown in black. (B) The effect of signature erosion on the design of *B. pseudomallei/mallei* diagnostic markers. Genomes outside the *B. pseudomallei/mallei* clade ($n = 714$) were randomly selected at different depths. The core genome of 416 *B. pseudomallei/mallei* genomes was screened against non-*pseudomallei/mallei* genomes with LS-BSR (22), and the number of markers with a BSR of <0.4 in non-*pseudomallei/mallei* genomes was calculated and plotted. One hundred independent replicates were processed at each sampling depth. The mean value at each level is shown in red, and each replicate is shown in black.

demonstrates the need to include diverse genomes both from the targeted clade and from clades closely related to the targeted clade.

This study both expands the known genomic diversity of the *Burkholderia* genus and provides a framework for using genomic data to design highly specific diagnostic targets. For some species, near-neighbor genomes are not available or are difficult to isolate, which complicates the identification of these targets and highlights the need for continued genome sequencing. The reported sampling effects on strict core genome size and signature erosion must be considered when interpreting surveillance results for human pathogens.

MATERIALS AND METHODS

Isolate collection, DNA extraction, genome sequencing, assembly. *Burkholderia* isolates were collected from diverse global locations, with a focus on regions where *B. pseudomallei* is highly endemic, including northern Australia and northeastern Thailand (see Table S1 in the supplemental material). Isolates were collected by the Menzies School of Health and Research, Northern Arizona University, the University of Michigan, the James Cook University, Mahidol University, and the U.S. Army Medical Research Unit (USAMRU). All final culture and DNA extraction procedures were performed at Northern Arizona University, and whole-genome sequencing (WGS) was performed at the Translational Genomics Research Institute (TGen), using Illumina, and the U.S. Army Medical Research Institute of Infectious Diseases (USAMRIID), using PacBio.

Isolates initially grown on Ashdown's agar were streaked from a single purified colony to form a lawn and then stored at -80°C in Luria Bertani (LB) broth with 20% glycerol. Cultures were grown on LB agar plates and incubated at 37°C for 24 to 48 h. High-molecular-weight DNA was extracted using the Qiagen DNeasy blood and tissue kit (catalog no. 69504; Valencia, CA) for whole-genome sequencing on the Illumina (Illumina, Inc., San Diego, CA) and Pacific Biosciences (Menlo Park, CA) platforms. Using approximately $2.7\ \mu\text{g}$ of genomic DNA (gDNA), libraries were prepared for Illumina whole-genome sequencing as previously described (24).

DNA was sequenced on multiple platforms, including Illumina HiSeq 2000, Illumina MiSeq, and PacBio. Raw Illumina reads were assembled with SPAdes version 3.5.0 (25) in conjunction with a pipeline developed to identify sequence contamination between multiplexed samples (<https://github.com/jasonsahl/UGAP>). Contigs that either showed an anomalously low depth of coverage compared to other contigs from the same assembly or aligned to other organisms multiplexed in the same lane were manually removed. Genome assembly information is shown in Table S1 in the supplemental material.

For PacBio assemblies, genomic DNA was sheared to an average size of 20 kb using g-TUBEs (Covaris, Inc.). After DNA damage repair and end repair, hairpin adapters were ligated to form a SMRTbell template. ExoIII and ExoVII treatment was used to remove failed ligation products. Size selection was performed on the Blue Pippin system (Sage Sciences) using a 0.75% dye-free agarose gel cassette, marker S1, and Hi-Pass protocol; the low cut was set on 4,000 bp. The final library assessment was obtained with the Qubit double-stranded DNA (dsDNA) BR assay and Agilent 2100 Bioanalyzer DNA 12000 chip analyses. Annealing of the sequencing primer and binding polymerase P4 to the SMRTbell template was performed according to the PacBio calculator. The polymerase-template complexes were bound to MagBeads, loaded onto SMRTcells at a final concentration of 180 pM, and sequenced with 180-min movies on the PacBio RS II instrument.

PacBio sequences were assembled *de novo* using the Hierarchical Genome Assembly Process (HGAP) (26). Draft assemblies were checked for overlapping ends using Gepard (27) and BLAST (28). Overlapping ends are typical of long-read assemblies of circular chromosomes. Redundant end sequences were trimmed to one copy, and the genome was rotated to create a new breakpoint. Reads were then realigned to the trimmed and

shifted draft assembly for correction using the Quiver algorithm. Contigs that did not have identifiable homologous ends were corrected using Quiver without further processing.

Species identification using core genome SNP phylogeny, ANI, and DDH calculation. To model the evolutionary relationships between *Burkholderia* spp., a set of reference genomes (see Table S2 in the supplemental material) was downloaded from GenBank (29) and combined with the genomes sequenced in this study. For a number of these genomes, only raw reads were available, which were assembled for use in the comparative analyses described below. All genomes were aligned against the reference genome of *B. pseudomallei* K96243 (19) using NUCmer (20). Regions that aligned more than once by a reference self-alignment (i.e., duplicated regions) were removed from downstream analyses. All SNP-based methods were wrapped by the Northern Arizona SNP pipeline (NASP) (<http://tgennorth.github.io/NASPTgennorth.github.io/NASP/>) (30). Orthologous SNPs conserved in all genomes were concatenated, and a maximum-likelihood phylogeny was inferred with RAxML version 8 (31) using the ASC_GTRGAMMA substitution model and Lewis correction (32).

For determining species differences, the average nucleotide identity (ANI) was calculated with default values in JSpecies (33). JSpecies calculates ANIb, which uses BLASTN alignments (28), or ANIm, which uses NUCmer alignments. The average values were reported over the entire length of all alignments. To find the nearest neighbor to which to query target genomes, the closest patristic distances were chosen, as calculated by DendroPy (34). Digital DNA-DNA hybridization (DDH) values were calculated with a Web service (<http://ggdc.dsmz.de>) (35), and the range of reported values is presented.

Identifying *B. pseudomallei* and *B. mallei* markers for diagnostics using comparative genomics and pangenome analysis. Coding DNA sequences (CDSs) were identified for each species with Prodigal (36) and were dereplicated with USEARCH (37). Each representative CDS was then aligned against each genome with BLAT (38), and the BLAST score ratio (BSR) (39) was calculated; these methods were all wrapped by the large-scale BLAST score ratio (LS-BSR) pipeline (22). LS-BSR was performed for each species, and the number of core CDSs (BSR of >0.8 in all genomes) in each group was calculated; a BSR of 0.8 is roughly equivalent to 80% protein identity over 100% of the length of the protein (3). These core CDSs from a given species or clade were then screened against all other genomes, and those genes with a BSR of <0.4 in all other species were identified as suitable species diagnostic markers.

The pangenome was calculated for each clade using LS-BSR in conjunction with BLAT. A CDS was determined to belong to the core genome if it had a BSR of >0.8 in all genomes queried for a given species or clade of interest. Each core CDS was then screened against all genomes in the analysis with LS-BSR. A CDS was determined to be species specific if it was in the core genome of the species or clade of interest and missing or highly divergent (BSR of <0.4) in all other *Burkholderia* genomes.

Core genome size reduction and signal erosion. To understand the sampling effect on the core genome size in *B. pseudomallei/mallei*, a set of 416 *B. pseudomallei/mallei* genomes was sampled without replacement from 1 to 400, with 100 iterations at each level. From each subsampling, a set number of genomes were randomly selected with a Python script (<https://gist.github.com/990d2c56c23bb5c2909d.git>), and the core genome (CDSs with a BSR of >0.8 in all genomes) was calculated and plotted. *B. pseudomallei* and *B. mallei* were treated as a single species for this and many of the subsequent analyses, as *B. mallei* is recognized as an equine-adapted clone within *B. pseudomallei* (40).

To understand the erosion of *B. pseudomallei/mallei*-specific targets with the inclusion of sequences from other *Burkholderia* spp., the core genome ($n = 1,684$ CDSs) from a set of 416 *B. pseudomallei/mallei* genomes was used. All *Burkholderia* near-neighbor genomes ($n = 714$) were then randomly sampled without replacement at different levels from 1 to 300. The *B. pseudomallei/mallei* core genome was then aligned against these near-neighbor genomes to identify core regions present in other

Burkholderia species, and the number of CDSs with a BSR of <0.4 in all near-neighbor genomes, indicating missing genes, was calculated and plotted.

Accession numbers. Sequence data were submitted to the Sequence Read Archive for each isolate. Furthermore, genome assemblies for all isolates were submitted to NCBI. Individual accession numbers are shown in Table S1 in the supplemental material, and all data are deposited under BioProject accession numbers PRJNA285704 and PRJNA279182.

SUPPLEMENTAL MATERIAL

Supplemental material for this article may be found at <http://mbio.asm.org/lookup/suppl/doi:10.1128/mBio.00846-16/-/DCSupplemental>.

Table S1, XLSX file, 0.1 MB.

Table S2, XLSX file, 0.1 MB.

Table S3, XLSX file, 0.1 MB.

Table S4, XLSX file, 0.1 MB.

ACKNOWLEDGMENTS

Opinions, interpretations, conclusions, and recommendations are those of the authors and do not necessarily reflect the official policy or position of the US Army, US Department of Defense, nor the US Government.

FUNDING INFORMATION

This work was funded by DOD | Defense Threat Reduction Agency (DTRA) (HDTRA1-12-C-0066 and CB10246).

REFERENCES

- Driebe EM, Sahl JW, Roe C, Bowers JR, Schupp JM, Gillece JD, Kelley E, Price LB, Pearson TR, Hepp CM, Brzoska PM, Cummings CA, Furtado MR, Andersen PS, Stegger M, Engelthaler DM, Keim PS. 2015. Using whole genome analysis to examine recombination across diverse sequence types of *Staphylococcus aureus*. *PLoS One* 10:e0130955. <http://dx.doi.org/10.1371/journal.pone.0130955>.
- Sahl JW, Morris CR, Emberger J, Fraser CM, Ochieng JB, Juma J, Fields B, Breiman RF, Gilmour M, Nataro JP, Rasko DA. 2015. Defining the phylogenomics of *Shigella* species: a pathway to diagnostics. *J Clin Microbiol* 53:951–960. <http://dx.doi.org/10.1128/JCM.03527-14>.
- Rasko DA, Rosovitz MJ, Myers GS, Mongodin EF, Fricke WF, Gajer P, Crabtree J, Sebahia M, Thomson NR, Chaudhuri R, Henderson IR, Sperandio V, Ravel J. 2008. The pangenome structure of *Escherichia coli*: comparative genomic analysis of *E. coli* commensal and pathogenic isolates. *J Bacteriol* 190:6881–6893. <http://dx.doi.org/10.1128/JB.00619-08>.
- Califf KJ, Keim P, Wagner DM, Sahl JW. 2015. Redefining the differences in gene content between *Yersinia pestis* and *Yersinia pseudotuberculosis* using large-scale comparative genomics. *Microb Genom* 1(4): 000028. <http://dx.doi.org/10.1099/mgen.0.000028>.
- Morelli G, Song Y, Mazzoni CJ, Eppinger M, Roumagnac P, Wagner DM, Feldkamp M, Kusecek B, Vogler AJ, Li Y, Cui Y, Thomson NR, Jombart T, Leblois R, Lichtner P, Rahalison L, Petersen JM, Balloux F, Keim P, Wirth T, Ravel J, Yang R, Carniel E, Achtman M. 2010. *Yersinia pestis* genome sequencing identifies patterns of global phylogenetic diversity. *Nat Genet* 42:1140–1143. <http://dx.doi.org/10.1038/ng.705>.
- Spring-Pearson SM, Stone JK, Doyle A, Allender CJ, Okinaka RT, Mayo M, Broomall SM, Hill JM, Karavis MA, Hubbard KS, Insalaco JM, McNew LA, Rosenzweig CN, Gibbons HS, Currie BJ, Wagner DM, Keim P, Tuanyok A. 2015. Pangenome analysis of *Burkholderia pseudomallei*: genome evolution preserves gene order despite high recombination rates. *PLoS One* 10:e0140274. <http://dx.doi.org/10.1371/journal.pone.0140274>.
- Coenye T, Vandamme P. 2003. Diversity and significance of *Burkholderia* species occupying diverse ecological niches. *Environ Microbiol* 5:719–729. <http://dx.doi.org/10.1046/j.1462-2920.2003.00471.x>.
- Wiersinga WJ, Currie BJ, Peacock SJ. 2012. Melioidosis. *N Engl J Med* 367:1035–1044. <http://dx.doi.org/10.1056/NEJMra1204699>.
- Losada L, Ronning CM, DeShazer D, Woods D, Fedorova N, Kim HS, Shabalina SA, Pearson TR, Brinkac L, Tan P, Nandi T, Crabtree J, Badger J, Beckstrom-Sternberg S, Saqib M, Schutzer SE, Keim P, Nierman WC. 2010. Continuing evolution of *Burkholderia mallei* through genome reduction and large-scale rearrangements. *Genome Biol Evol* 2:102–116. <http://dx.doi.org/10.1093/gbe/evq003>.
- Gee JE, Allender CJ, Tuanyok A, Elrod MG, Hoffmaster AR. 2014. *Burkholderia pseudomallei* type G in Western Hemisphere. *Emerg Infect Dis* 20:682–684. <http://dx.doi.org/10.3201/eid2004.130960>.
- Ho CC, Lau CC, Martelli P, Chan SY, Tse CW, Wu AK, Yuen KY, Lau SK, Woo PC. 2011. Novel pan-genomic analysis approach in target selection for multiplex PCR identification and detection of *Burkholderia pseudomallei*, *Burkholderia thailandensis*, and *Burkholderia cepacia* complex species: a proof-of-concept study. *J Clin Microbiol* 49:814–821. <http://dx.doi.org/10.1128/JCM.01702-10>.
- De Smet B, Mayo M, Peeters C, Zlosnik JE, Spilker T, Hird TJ, LiPuma JJ, Kidd TJ, Kaestli M, Ginther JL, Wagner DM, Keim P, Bell SC, Jacobs JA, Currie BJ, Vandamme P. 2015. *Burkholderia stagnalis* sp. nov. and *Burkholderia territorii* sp. nov., two novel *Burkholderia cepacia* complex species from environmental and human sources. *Int J Syst Evol Microbiol* 65:2265–2271. <http://dx.doi.org/10.1099/ijs.0.000251>.
- Payne GW, Vandamme P, Morgan SH, Lipuma JJ, Coenye T, Weightman AJ, Jones TH, Mahenthalingam E. 2005. Development of a recA gene-based identification approach for the entire *Burkholderia* genus. *Appl Environ Microbiol* 71:3917–3927. <http://dx.doi.org/10.1128/AEM.71.7.3917-3927.2005>.
- Ginther JL, Mayo M, Warrington SD, Kaestli M, Mullins T, Wagner DM, Currie BJ, Tuanyok A, Keim P. 2015. Identification of *Burkholderia pseudomallei* near-neighbor species in the Northern Territory of Australia. *PLoS Negl Trop Dis* 9:e0003892. <http://dx.doi.org/10.1371/journal.pntd.0003892>.
- Baldwin A, Mahenthalingam E, Thickett KM, Honeybourne D, Maiden MC, Govan JR, Speert DP, Lipuma JJ, Vandamme P, Dowson CG. 2005. Multilocus sequence typing scheme that provides both species and strain differentiation for the *Burkholderia cepacia* complex. *J Clin Microbiol* 43:4665–4673. <http://dx.doi.org/10.1128/JCM.43.9.4665-4673.2005>.
- Nandi T, Holden MT, Didelot X, Mehershahi K, Boddey JA, Beacham I, Peak I, Harting J, Baybayan P, Guo Y, Wang S, How LC, Sim B, Essex-Lopresti A, Sarkar-Tyson M, Nelson M, Smither S, Ong C, Aw LT, Hoon CH, Michell S, Studholme DJ, Titball R, Chen SL, Parkhill J, Tan P. 2015. *Burkholderia pseudomallei* sequencing identifies genomic clades with distinct recombination, accessory, and epigenetic profiles. *Genome Res* 25:129–141. <http://dx.doi.org/10.1101/gr.177543.114>.
- Nierman WC, DeShazer D, Kim HS, Tettelin H, Nelson KE, Feldblyum T, Ulrich RL, Ronning CM, Brinkac LM, Daugherty SC, Davidsen TD, Deboy RT, Dimitrov G, Dodson RJ, Durkin AS, Gwinn ML, Haft DH, Khouri H, Kolonay JF, Madupu R, Mohammad Y, Nelson WC, Radune D, Romero CM, Sarria S, Selengut J, Shamblyn C, Sullivan SA, White O, Yu Y, Zafar N, Zhou L, Fraser CM. 2004. Structural flexibility in the *Burkholderia mallei* genome. *Proc Natl Acad Sci U S A* 101: 14246–14251. <http://dx.doi.org/10.1073/pnas.0403306101>.
- Johnson SL, Baker AL, Chain PS, Currie BJ, Daligault HE, Davenport KW, Davis CB, Inglis TJ, Kaestli M, Koren S, Mayo M, Merritt AJ, Price EP, Sarovich DS, Warner J, Rosovitz MJ. 2015. Whole-genome sequences of 80 environmental and clinical isolates of *Burkholderia pseudomallei*. *Genome Announc* 3:e01282-14. <http://dx.doi.org/10.1128/genomeA.01282-14>.
- Holden MT, Titball RW, Peacock SJ, Cerdeño-Tárraga AM, Atkins T, Crossman LC, Pitt T, Churcher C, Mungall K, Bentley SD, Sebahia M, Thomson NR, Bason N, Beacham IR, Brooks K, Brown KA, Brown NF, Challis GL, Cherevach I, Chillingworth T, Cronin A, Crossett B, Davis P, DeShazer D, Feltwell T, Fraser A, Hance Z, Hauser H, Holroyd S, Jagels K, Keith KE, Maddison M, Moule S, Price C, Quail MA, Rabinovitch E, Rutherford K, Sanders M, Simmonds M, Songsivilai S, Stevens K, Tumapa S, Vesaratchavest M, Whitehead S, Yeats C, Barrell BG, Oyston PC, Parkhill J. 2004. Genomic plasticity of the causative agent of melioidosis, *Burkholderia pseudomallei*. *Proc Natl Acad Sci U S A* 101:14240–14245. <http://dx.doi.org/10.1073/pnas.0403302101>.
- Delcher AL, Salzberg SL, Phillippy AM. 2003. Using MUMmer to identify similar regions in large sequence sets. *Curr Protoc Bioinformatics* Chapter 10:Unit 10.3. <http://dx.doi.org/10.1002/0471250953.bi1003s00>.
- Depoorter E, Bull MJ, Peeters C, Coenye T, Vandamme P, Mahenthalingam E. 2016. *Burkholderia*: an update on taxonomy and biotechnological potential as antibiotic producers. *Appl Microbiol Biotechnol* 100: 5215–5229. <http://dx.doi.org/10.1007/s00253-016-7520-x>.
- Sahl JW, Caporaso JG, Rasko DA, Keim P. 2014. The large-scale Blast score ratio (LS-BSR) pipeline: a method to rapidly compare genetic con-

- tent between bacterial genomes. PeerJ PrePrints 2:e332. <http://dx.doi.org/10.7717/peerj.332>.
23. Price EP, Sarovich DS, Mayo M, Tuanyok A, Drees KP, Kaestli M, Beckstrom-Sternberg SM, Babic-Sternberg JS, Kidd TJ, Bell SC, Keim P, Pearson T, Currie BJ. 2013. Within-host evolution of *Burkholderia pseudomallei* over a twelve-year chronic carriage infection. *mBio* 4:e00388-13. <http://dx.doi.org/10.1128/mBio.00388-13>.
 24. Engelthaler DM, Hicks ND, Gillece JD, Roe CC, Schupp JM, Driebe EM, Gilgado F, Carriconde F, Trilles L, Firacative C, Ngamskulrungraj P, Castañeda E, Lazera Mdos S, Melhem MS, Pérez-Bercoff A, Huttley G, Sorrell TC, Voelz K, May RC, Fisher MC, Thompson GR, III, Lockhart SR, Keim P, Meyer W. 2014. *Cryptococcus gattii* in North American Pacific Northwest: whole-population genome analysis provides insights into species evolution and dispersal. *mBio* 5:e01464-14. <http://dx.doi.org/10.1128/mBio.01464-14>.
 25. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM, Nikolenko SI, Pham S, Pribelski AD, Pyshkin AV, Sirotkin AV, Vyahhi N, Tesler G, Alekseyev MA, Pevzner PA. 2012. SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *J Comput Biol* 19:455–477. <http://dx.doi.org/10.1089/cmb.2012.0021>.
 26. Chin CS, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, Clum A, Copeland A, Huddleston J, Eichler EE, Turner SW, Korlach J. 2013. Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing data. *Nat Methods* 10:563–569. <http://dx.doi.org/10.1038/nmeth.2474>.
 27. Krumsiek J, Arnold R, Rattei T. 2007. Gepard: a rapid and sensitive tool for creating dotplots on genome scale. *Bioinformatics* 23:1026–1028. <http://dx.doi.org/10.1093/bioinformatics/btm039>.
 28. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. *J Mol Biol* 215:403–410. [http://dx.doi.org/10.1016/S0022-2836\(05\)80360-2](http://dx.doi.org/10.1016/S0022-2836(05)80360-2).
 29. Benson DA, Karsch-Mizrachi I, Clark K, Lipman DJ, Ostell J, Sayers EW. 2012. GenBank. *Nucleic Acids Res* 40:D48–D53. <http://dx.doi.org/10.1093/nar/gkr1202>.
 30. Sahl JW, Lemmer D, Travis J, Schupp J, Gillece J, Aziz M, Driebe E, Drees K, Hicks ND, Williamson C, Hepp C, Smith DE, Roe C, Engelthaler DM, Wagner DM, Keim P. 2016. The Northern Arizona SNP pipeline (NASP): accurate, flexible, and rapid identification of SNPs in WGS datasets. *bioRxiv*<http://dx.doi.org/10.1101/037267>.
 31. Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* 30:1312–1313 <http://dx.doi.org/10.1093/bioinformatics/btu033>.
 32. Leaché AD, Banbury BL, Felsenstein J, de Oca AN, Stamatakis A. 2015. Short tree, long tree, right tree, wrong tree: new acquisition bias corrections for inferring SNP phylogenies. *Syst Biol* 64:1032–1047 <http://dx.doi.org/10.1093/sysbio/syv053>.
 33. Richter M, Rosselló-Móra R. 2009. Shifting the genomic gold standard for the prokaryotic species definition. *Proc Natl Acad Sci U S A* 106:19126–19131. <http://dx.doi.org/10.1073/pnas.0906412106>.
 34. Sukumaran J, Holder MT. 2010. DendroPy: a Python library for phylogenetic computing. *Bioinformatics* 26:1569–1571. <http://dx.doi.org/10.1093/bioinformatics/btq228>.
 35. Meier-Kolthoff JP, Auch AF, Klenk HP, Göker M. 2013. Genome sequence-based species delimitation with confidence intervals and improved distance functions. *BMC Bioinformatics* 14:60. <http://dx.doi.org/10.1186/1471-2105-14-60>.
 36. Hyatt D, Chen GL, Locascio PF, Land ML, Larimer FW, Hauser LJ. 2010. Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics* 11:119. <http://dx.doi.org/10.1186/1471-2105-11-119>.
 37. Edgar RC. 2010. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics* 26:2460–2461. <http://dx.doi.org/10.1093/bioinformatics/btq461>.
 38. Kent WJ. 2002. BLAT—the BLAST-like alignment tool. *Genome Res* 12:656–664. <http://dx.doi.org/10.1101/gr.229202>.
 39. Rasko DA, Myers GS, Ravel J. 2005. Visualization of comparative genomic analyses by BLAST score ratio. *BMC Bioinformatics* 6:2. <http://dx.doi.org/10.1186/1471-2105-6-2>.
 40. Pearson T, Giffard P, Beckstrom-Sternberg S, Auerbach R, Hornstra H, Tuanyok A, Price EP, Glass MB, Leadem B, Beckstrom-Sternberg JS, Allan GJ, Foster JT, Wagner DM, Okinaka RT, Sim SH, Pearson O, Wu Z, Chang J, Kaul R, Hoffmaster AR, Brettin TS, Robison RA, Mayo M, Gee JE, Tan P, Currie BJ, Keim P. 2009. Phylogeographic reconstruction of a bacterial species with high levels of lateral gene transfer. *BMC Biol* 7:78. <http://dx.doi.org/10.1186/1741-7007-7-78>.
 41. Ludwig W, Strunk O, Westram R, Richter L, Meier H, Yadhukumar, Buchner A, Lai T, Steppi S, Jobb G, Förster W, Brettske I, Gerber S, Ginhart AW, Gross O, Grumann S, Hermann S, Jost R, König A, Liss T, Lüßmann R, May M, Nonhoff B, Reichel B, Strehlow R, Stamatakis A, Stuckmann N, Vilbig A, Lenke M, Ludwig T, Bode A, Schleifer K-H. 2004. ARB: a software environment for sequence data. *Nucleic Acids Res* 32:1363–1371. <http://dx.doi.org/10.1093/nar/gkh293>.
 42. Letunic I, Bork P. 2007. Interactive Tree Of Life (iTOL): an online tool for phylogenetic tree display and annotation. *Bioinformatics* 23:127–128. <http://dx.doi.org/10.1093/bioinformatics/btl529>.