# Prediction of *Burkholderia pseudomallei* DsbA substrates identifies potential virulence factors and vaccine targets

Vezina, Ben[1]*, Petit, Guillaume A.[1]*, Martin, Jennifer L.[1,2], Halili, Maria A.[1]+

[1]Griffith institute for Drug Discovery, Griffith University, Building N75, 46 Don Young Rd, Nathan, QLD 4111, Australia

[2]Vice-Chancellor's Unit, University of Wollongong, Northfields Avenue, Wollongong, NSW 2500 Australia

+Corresponding author Email: m.greenup@griffith.edu.au

*These authors contributed equally

# **Abstract**

Identification of bacterial virulence factors is critical for understanding disease pathogenesis, drug discovery and vaccine development. In this study we used two approaches to predict virulence factors of *Burkholderia pseudomallei*, the Gram-negative bacterium that causes melioidosis. *B. pseudomallei* is naturally antibiotic resistant and there are no melioidosis vaccines. To identify *B. pseudomallei* protein targets for drug discovery and vaccine development, we chose to search for substrates of the *B. pseudomallei* periplasmic disulfide bond forming protein A (DsbA). DsbA introduces disulfide bonds into extra-cytoplasmic proteins and is essential for virulence in many Gram-negative organism, including *B. pseudomallei*. The first approach to identify *B. pseudomallei* DsbA virulence factor substrates was a large-scale genomic analysis of 511 unique *B. pseudomallei* disease-associated strains. This yielded 4,496 core gene products, of which we hypothesise 263 are DsbA substrates. Manual curation of the 263 mature proteins yielded 73 associated with disease pathogenesis or virulence. These were screened for structural homologues to predict potential B-cell epitopes. In the second approach, we searched the *B. pseudomallei* genome for homologues of the more than 90 known DsbA substrates in other bacteria. Using this approach, we identified 15 potential *B. pseudomallei* DsbA virulence factor substrates. Two putative *B. pseudomallei* virulence factors were identified by both methods: homologues of PenI family β-lactamase and of succinate dehydrogenase flavoprotein subunit. These two proteins could serve as high priority targets for future *B. pseudomallei* virulence factor characterization.

# Introduction

*Burkholderia pseudomallei* is a Gram-negative soil dwelling saprophyte, and an opportunistic pathogen responsible for the severe tropical disease melioidosis [1]. *B. pseudomallei* infections are difficult to treat [2-4] and are intrinsically resistant to almost all available antibiotics [5-8]. Predominant resistance factors utilised by *B. pseudomallei* include a thick, impermeable cell wall combined with efficient efflux pumps that interfere with drug activity [9]. Furthermore, *B. pseudomallei* infections are difficult to diagnose as melioidosis symptoms vary significantly, ranging from fever, pneumonia, urinary tract infections, and on rare occasions encephalomyelitis [4]. Standard treatment consists of a combination of intravenous antibiotic for two weeks to stop septicaemia, followed by a second eradication phase that can last for up to six months, with no guarantee of success [10].

More generally, antibiotic resistance is increasing at an accelerating rate among pathogenic bacteria [11]. New approaches and treatment strategies are needed including vaccination [12], novel antimicrobial compounds [13] and antivirulence strategies [14]. There is currently no successful, persistent vaccine against *B. pseudomallei* [15]. However Outer Membrane Protein A (OmpA) has been used as a subunit vaccination against melioidosis in mice [16].

Identification of *B. pseudomallei* virulence factors would contribute towards understanding pathogenesis and could aid in drug discovery and vaccine development [17]. Targeting virulence rather than viability is an approach that is hypothesized to have a number of benefits including an increased range of possible anti-virulence mechanisms compared to antimicrobial compounds, as well as the possibility of reducing selection pressure [18, 19]. Both vaccine development and anti-virulence approaches could reduce selection pressure and potentially reduce resistance development [14, 18, 19].

3

64

65    The formation of correct disulfide bonds is critical for the proper folding and function of

66    proteins [20]. In bacteria, the introduction of disulfide bonds is mediated by the DiSufide Bond-

67    forming proteins (DSB). The DSB proteins are of particular interest as an antivirulence

68    strategy, because many virulence factors contain disulfide bonds [19, 21-23]. The Disulfide

69    bond forming protein A (DsbA) is a periplasmic protein found in most Gram-negative bacteria

70    and incorporates a thioredoxin fold with two cysteines which introduce disulfide bonds into

71    substrate proteins via a redox transfer reaction [24].

72

73    Mice infected with *B. pseudomallei* DsbA knockouts (or of its redox partner DsbB) have an

74    increased rate of survival compared with mice infected with wild type *B. pseudomallei* [25,

75    26]. These findings suggest that many *B. pseudomallei* virulence factors are substrates of

76    DsbA*, as is also observed in *Escherichia coli* [27, 28]*, Klebsiella pneumoniae* [29]*, Salmonella*

77    *enterica* [30]*, Francisella tularensis* [31] and many more [22, 23, 32]. However, the full extent

78    of *B. pseudomallei* DsbA substrates has not been investigated. Identification of *B. pseudomallei*

79    DsbA substrates would help identification of infection mechanisms, and could lead to the

80    discovery of key virulence factors and potential drug and vaccine targets. Finding potential

81    DsbA substrates is assisted by the observation that: (i) DsbA is located in the periplasm, and

82    thus its substrates are likely to have a secretion signal sequence; and (ii) proteins containing

83    disulfide bonds may have an even rather than an odd number of cysteines in their sequence.

84    This last point is thought to have evolved to limit formation of mis-matched disulfide bonds

85    and therefore misfolded proteins [33, 34].

86

87    In the present study, we used two approaches to identify potential *B. pseudomallei* DsbA

88    substrates for further study as virulence factors. In one approach, we used computational

89    methods to generate a curated list of 263 putatively extra-cytoplasmic proteins from the core

90    genome of 511 disease-associated isolates of *B. pseudomallei*, 73 of which were predicted to

91    be virulence-associated. In the second approach, 15 candidate DsbA virulence factor substrates

92    were identified by sequence homology to known DsbA virulence factor substrates in other

93    bacteria.

94

# Results

## Genomic analysis to predict *B. pseudomallei* DsbA virulence factor substrates

In this approach, our strategy was to cast a wide net initially, by determining the pangenome of disease-associated isolates of *B. pseudomallei*, and then filtering from that the core genome (i.e. the highly conserved genes). The disease-associated *B. pseudomallei* core genome should then be enriched in conserved virulence factors. At the time of this analysis the NCBI database [35] contained 1577 *B. pseudomallei* isolates. Metadata notation allowed selection of 512 isolates associated with disease (i.e. isolates from swabs/clinical isolates: accession numbers of these are given in S1 Fig); other genomes were discarded. We note that only 355 of the 512 isolates were tagged 'pathogen' in the NCBI database indicating a discrepancy between NCBI assignment and user-uploaded metadata. Analysis of the pangenome, that is the core, accessory and unique genes of these 512 *B. pseudomallei* isolates (see Table 1), revealed two identical strains. Therefore for the remainder of this analysis, only the 511 unique strains were used.

**Table 1: Pangenome results of 511 disease-associated *B. pseudomallei* strains.**

| Pangenome breakdown | Classification | Number of genes | Percent of pangenome (%) |
|---|---|---|---|
| Core genes | (99% <= strains <= 100%) | 4,496 | 22.49 |
| Soft core genes | (95% <= strains < 99%) | 517 | 2.59 |
| Shell genes | (15% <= strains < 95%) | 965 | 4.83 |
| Cloud genes | (0% <= strains < 15%) | 14,013 | 70.10 |
| Total pangenome | (0% <= strains <= 100%) | 19,991 | 100 |

The pangenome is subdivided into the core (found in every strain), soft shell core (found in 95 – 99% of strains), shell (found in 15 – 95% of strains), and cloud (found in 0 – 15% of strains) genes. The total number of genes is shown, along with the percentage of total pangenome.

We found that the core genome consisted of 4,496 genes (see S2 Fig) or 22.49% of the total 19,991 pangenome. This analysis largely agrees with a previous pangenomic analysis which

117  extrapolated a modelled core genome of 4,568±16 from a much smaller set of 37 isolate

118  genomes [36]. In that approach, modelling was used to predict the core genome if the number

119  of isolates was expanded. Our approach gives an exact number because all 4,496 genes were

120  found in all 511 genomes. Notably, the dithiol oxidase redox enzyme pair DsbA and DsbB and

121  the disulfide isomerase redox relay enzymes DsbC and DsbD were all identified as core genes.

122

123  We then used the *B. pseudomallei* core genome for further analysis, because it encodes highly

124  conserved proteins - a key criteria for selecting vaccine or anti-virulence targets.

125

126  From these 4,496 core genes, 726 were predicted to encode proteins with a signal sequence

127  and which are therefore likely to be exported out of the cytoplasm and into the periplasm where

128  DsbA is localised. Of these 726 proteins, 263 have an even number of cysteines, indicating the

129  likelihood that the proteins form intramolecular disulfide bonds (see S3 Fig). We predict that

130  these 263 proteins are substrates of *B. pseudomallei* DsbA. The workflow for this analysis is

131  shown in Fig 1.

132

133  **Fig 1: Bioinformatic workflow.** From the 1,577 *B. pseudomallei* genomes found on NCBI,
134  511 were unique and associated with disease and these were used for further analysis. The
135  pangenome of these 511 genomes comprised 19,991 unique genes. 4,496 of these were
136  classified as core genes. Predicted translation of these genes gave 726 predicted extra-
137  cytoplasmic proteins. Of these extra-cytoplasmic proteins, 263 were predicted to contain an
138  even number of cysteines. We predict that these 263 proteins are substrates of *B.
139  pseudomallei* DsbA.

140

## Distribution of cysteines in the core genome of disease-related *B. pseudomallei*

143  Many bacterial extra-cytoplasmic (periplasmic and extracellular) proteins have a strong

144  preference for an even number of cysteines, which is thought to reduce the chances of non-

7

145   native disulfide bond formation [33]. We examined the cysteine distribution of encoded

146   proteins in the *B. pseudomallei* pangenome to investigate whether the previously demonstrated

147   enrichment of an even number of cysteines in extra-cytoplasmic proteins in other Gram-

148   negative bacteria [33] was also true for *B. pseudomallei.*

149

150   The distribution of cysteines in *B. pseudomallei* cytoplasmic and extra-cytoplasmic proteins

151   was calculated for the pangenome (total of 19,991 genes) and the core genome (4,496 genes)

152   (refer to Table 1). In cytoplasmic *B. pseudomallei* proteins, cysteine distribution followed a

153   Poisson law peaking at zero for the pangenome and at one for the core genome (denoted by the

154   orange lines in the histograms on Figs 2A and 2B). This distribution changed for extra-

155   cytoplasmic *B. pseudomallei* proteins. For the core genome (blue bars Fig 2B), *B. pseudomallei*

156   proteins with an even number of cysteines were over-represented compared to a typical Poisson

157   distribution. As extra-cytoplasmic proteins represent a small fraction of the total number of the

158   translated core genome and pangenome (16% and 11.5% of all proteins, respectively), we also

159   analysed the normalised frequency (Figs 2C and 2D). The core genome normalised cysteine

160   distribution reveals a sawtooth pattern with a preference for even number of cysteines with

161   peaks for two, four, six and eight cysteines (Fig 2D). In contrast, the pangenomic normalised

162   cysteine distribution for extra-cytoplasmic *B. pseudomallei* proteins does not indicate a strong

163   preference for even number of cysteines (Fig 2C). Overall, the saw-tooth pattern observed in

164   Figs 2B and 2D is similar to that described for *E. coli* exported proteins [33] although not as

165   pronounced.

166

167   **Fig 2: Cysteine distribution in the translated genome of *B. pseudomallei*.** Panel **A** shows
168   the distribution of cysteines in the pangenome (19,991 proteins). Panel **B** represents the same
169   analysis for the core genome, comprising 4,496 translated genes. Predicted number of extra-
170   cytoplasmic proteins for each number of cysteines are represented as blue bars. Similarly,
171   predicted cytoplasmic proteins are represented as orange lines. Panels **C** and **D** represent the
172   normalised frequency of cysteine-containing extra-cytoplasmic proteins. The blue line in panel

173   **D** peaks for proteins with 2, 4, 6 and 8 cysteines suggesting a preference for an even number
174   of cysteines. This trend is not observed as strongly in panel **C**, where a clear peak can only be
175   seen for two and eight cysteines. The normalised frequency was calculated by dividing the
176   number of extra-cytoplasmic proteins (having *N* number of cysteines) by the total number of
177   proteins with *N* cysteines (*N* being a number between 0 - 20 as per the data points in **C** and **D**
178   above).

179

## Functional assignment of core, extra-cytoplasmic, putative DsbA substrates

182   The next step in the genomic analysis was to predict which of the 263 putative DsbA substrates

183   are associated with virulence. Of the 263 selected proteins, 44 were annotated as

184   hypothetical/uncharacterised. The remaining 219 proteins include ABC transporter-related

185   proteins, housekeeping proteins like cytochrome C, proteins required for motility such as

186   flagellar and fimbrial proteins, enzymes such as collagenase, peptidases and proteases, as well

187   as antibiotic resistance enzymes, β-lactamases. Many oxidoreductases were also present

188   including DsbA, DsbD and others such as Gfo/Idh/MocA family, glycerol-3-phosphate

189   dehydrogenase GpsA and thioredoxin-like TlpA oxidoreductases. Redox enzymes such as

190   DsbB and DsbC are core genes with signal sequences, and they have catalytic rather than

191   structural disulfides. These two enzymes are not identified as DsbA substrates in our filter as

192   they have an odd number of cysteines.

193

194   Gene Ontology (GO) classification of the gene and gene-product function of the 263 proteins

195   reveals a variety of functions, totalling 223 GO descriptions (Fig 3) (see S4 File for a full list).

196   The highest frequency are integral components of the membrane (66 proteins), followed by

197   proteins involved in redox processes (25 proteins). Of particular interest due to their putative

198   involvement in virulence, are proteins associated with: proteolysis (20), heme binding (15),

199   hydrolase activity (9), carbohydrate metabolism (8), serine-type endopeptidase activity (7), cell

200    adhesion (6), metallo-endopeptidase activity (6), pilus formation and organisation (6), copper

201    binding (5), lipid catabolism (4), choline binding (3), triglyceride lipase activity (3),

202    aminopeptidase activity (2), porin activity (OmpA family proteins) (2), chitin catabolism (1),

203    *N*-carbamoylputrescine amidase activity (1) and toxin activity (Tat pathway signal protein) (1).

204

205    **Fig 3: Gene Ontology (GO) descriptions of predicted extra-cytoplasmic proteins with an**
206    **even number of cysteines**. The highest frequency of proteins with an even number of
207    cysteines are integral components of membranes (66 proteins), followed by proteins involved
208    in redox (oxidation-reduction) processes (25 proteins) and proteolysis (20 proteins). For ease
209    of representation and clarity, GO descriptors with less than three counts were excluded from
210    this graph. A complete graph, along with raw values can be found in S4 File.

211

212    By further inspection of the 263 core, putatively extra-cytoplasmic DsbA substrates, and by

213    using the GO descriptions to aid in predicting protein functions,73 sequences were identified

214    which were virulence-associated (Table 2). These include serine-type endopeptidases [37]

215    associated with adherence, choline binding proteins N-carbamoylputrescine amidase, essential

216    for production of putrescine, a component of Gram-negative cell walls of pathogens and key

217    virulence [39-42], many proteases and peptidases.

218

219    **Table 2: Predicted virulence-associated core, extra-cytoplasmic proteins.**

| Virulence-associated GO description | Accession numbers |
|---|---|
| Aminopeptidase activity | ABA50277.1; WP_053292838.1 |
| Bacterial-type flagellum assembly | WP_004525898.1 |
| Beta-lactamase activity | KGV04506.1 |
| Carbohydrate metabolic processes | ABA52198.1; EDO83218.1; EEH25224.1; WP_004526045.1; WP_004526830.1; WP_004553625.1; WP_053293009.1 |
| Cell adhesion/lipid metabolic/catabolic process/chitinase | WP_004193933.1 |
| Cell adhesion/pillus | EDU07436.1; WP_004193385.1; WP_038760383.1; WP_038765499.1; WP_063597677.1 |
| Chitin catabolic process | WP_076802983.1 |
| Choline binding and transport | ABA51731.1; ABN86005.1; ABN92885.1 |
| Copper ion binding | WP_004529973.1; WP_004546221.1 |

| | |
|---|---|
| Heme binding | WP_004194773.1; WP_004535805.1; WP_004536717.1; WP_004538457.1; WP_004538458.1; WP_038730764.1; WP_041189005.1; WP_043304483.1; WP_076903047.1; WP_139900217.1; WP_151277731.1 |
| Heme binding/copper ion binding | WP_029671417.1; WP_122827599.1 |
| Heme binding/proteolysis | WP_009981622.1 |
| Heme bindingcopper ion binding | WP_080248664.1 |
| Hydrolase activity | CFL10512.1; EEC34719.1; WP_004525656.1; WP_024428578.1; WP_024429096.1; WP_080300428.1 |
| Lipid metabolic/catabolic process | WP_009956690.1; WP_080248725.1 |
| Metallopeptidase/metalloendopeptidase activity | AFR18870.1; WP_004548157.1; WP_011204325.1; WP_038708181.1; WP_038730428.1; WP_076887541.1 |
| N-carbamoylputrescine amidase activity | WP_045597613.1 |
| Penicillin binding/beta-lactamase activity | EDO89205.1 |
| Pillus and pillus organisation | WP_151269450.1 |
| Porin activity | WP_004189892.1; WP_011205039.1 |
| Proteolysis/hydrolase activity | WP_011204795.1; WP_076852667.1 |
| Serine-type endopeptidase/carboxypeptidase activity | ABA50268.1; ACQ98979.1; AFR20596.1; WP_004528537.1; WP_004529035.1; WP_004553586.1; WP_011852052.1; WP_024428782.1; WP_038778478.1 |
| Toxin activity | WP_038707916.1 |
| Triglyceride lipase activity | EEH28759.1; WP_038741497.1; WP_038775093.1 |
| Xenobiotic transmembrane transporter activity | WP_004534049.1 |

220 Analysis of the 263 putative DsbA substrates revealed 73 proteins associated with virulence,
221 based on GO descriptions. Accession numbers from *B. pseudomallei* are shown, separated by
222 a semicolon.

223

## 224 **Sequence homology prediction of *B. pseudomallei* DsbA virulence**

## 225 **factor substrates**

226 To complement the genomic analysis described above we used a second approach to identify

227 DsbA substrates, by screening all *B. pseudomallei* genomes uploaded on NCBI [43] (taxid

228 28450) for homologues of known DsbA substrates. We implemented this approach because

229 some DsbA substrates might be filtered out using the genomic approach described above if the

230 substrates are not encoded by core genes, or if the gene product has an odd number of cysteines.

231

232    Over 90 DsbA substrates have been reported in the literature. We searched for *B. pseudomallei*

233    homologues of these DsbA substrates using the following criteria: (i) presence of secretion

234    signal, (ii) at least two cysteines in the mature sequence, (iii) at least 20% identity and (iv) 50%

235    coverage to a known DsbA substrate sequence. After removing duplicates, our analysis found

236    that *B. pseudomallei* encodes homologues of 15 DsbA substrates (Table 3). Two of these 15

237    are DsbA substrates in other *Burkholderia* species *B. cepacia* and *B. cenocepacia* [44-47]: a

238    metalloproteases, ZmpA and a sulfatase-like hydrolase transferase. In *B. cenocepacia,* ZmpA

239    is a wide spectrum metalloprotease, thought to cause tissue damage during infection [48].

240

241    **Table 3**: **List of *B. pseudomallei* proteins homologous to previously reported DsbA**
242    **substrates.**

| Accession Number (DsbA substrate) | Organism | Reference | *B. pseudomallei* homologue | Identity / coverage (%) | Protein function | Cys # |
|---|---|---|---|---|---|---|
| WP_059237834 | *B. cepacia* | [44] | WP_076835606.1 | 89 /100 | Sulfatase like hydrolase /transferase | 3 |
| WP_006481898 | *B. cenocepacia* | [45, 46] | WP_139900467 | 87/100 | M4 family metallopeptidase | 4 |
| gi\|89255876 | *F. tularensis* | [34] | WP_050859308 | 24/92 | lytic transglycosylase | 3 |
| gi\|89255615 | *F. tularensis* | [34] | WP_080367462 | 40/51 | Pilin | 2 |
| gi\|89255615 | *F. tularensis* | [34] | WP_076953316 | 27/92 | Pilin | 2 |
| gi\|89256194 | *F. tularensis* | [34] | WP_041862011 | 30/83 | Molybdopterin synthase adenyl transferase (MoeB) | 13 |
| gi\|89256236 | *F. tularensis* | [34] | WP_064459078 | 34/53 | DNA/RNA endonuclease | 2 |
| gi\|89256237 | *F. tularensis* | [34] | WP_050772403 | 31/90 | PenI family Beta-lactamase | 4 |
| gi\|89256856 | *F. tularensis* | [34] | WP_044360358 | 21/80 | hypothetical protein | 4 |
| gi\|89256859 | *F. tularensis* | [34] | WP_058035453 | 39/80 | Polyamine ABC transporter substrate binding protein | 3 |
| gi\|89257049 | *F. tularensis* | [34] | WP_009915682 | 54/99 | Succinate dehydrogenase | 6 |
| WP_001363619 | *E. coli* | [22] | WP_102811167 | 38/88 | Molecular chaperone | 3 |
| AAC38377 | *E. coli* | [22] | WP_082252625 | 44/93 | T3SS outer membrane ring protein | 4 |
| AAA24962 | *Heamophilus Influenza* | [22] | WP_053293022 | 47/92 | ABC transporter substrate binding protein | 4 |
| CAA43967 | *Yersinia pestis* | [22] | WP_085538626 | 32/83 | Pilus assembly protein PapD | 2 |

12

243 The accession number of the known DsbA substrate (in an organism other than *B.*
244 *pseudomallei)*, the organism and the publication reference are given in the first three columns.
245 The corresponding *B. pseudomallei* homologue is given in the fourth column. The identity and
246 coverage (number of residues in the result sequence that overlap with the search sequence) is
247 given in percent in the column "identity/coverage". The final two columns provide the protein
248 function and the number of cysteines in the predicted mature sequence. All proteins in this
249 table are known or predicted to be secreted or periplasmic.

250

251 Over 50 DsbA substrates in *Francisella tularensis* were identified by trapping and co-purifying

252 substrates bound to a DsbA variant [34]. Of these 50, we found nine homologues encoded in

253 *B. pseudomallei* (see Table 3). These include homologues of the lytic transglycosylase domain

254 containing protein (implicated in peptidoglycan rearrangement) and homologues of two pilin

255 proteins involved in the formation of pilus and flagella. Also present is an MoeB homologue;

256 MoeB is a molybdopterin synthase adenyl transferase (cytoplasmic in *E. coli* but likely

257 periplasmic in *B. pseudomallei* due to the twin-arginine translocation (TAT) signal sequence).

258 A PenI family β-lactamase homologue is also found in *B. pseudomallei*; this is a class A β-

259 lactamase that confers resistance to β-lactams including, in rare cases, ceftazidime (commonly

260 used to treat melioidosis) [49]. A succinate dehydrogenase flavoprotein subunit homologue,

261 found in the bacterial inner membrane and part of the electron transport chain, is also encoded

262 in *B. pseudomallei*. This protein is cytoplasmically oriented in *E. coli*, though again the *B.*

263 *pseudomallei* version has a TAT signal sequence suggesting a possible periplasmic

264 localisation.

265

266 A number of DsbA substrates identified in *E. coli* (reviewed in [22]) have *B. pseudomallei*

267 homologues including a molecular chaperone homologous to PapD and EscC, involved in the

268 formation of the Type III secretion system (T3SS). The T3SS assembly requires DsbA activity

269 in many Gram-negative bacteria, including *E. coli* and *S. typhimurium*. [50, 51]. Finally, a *B.*

13

270   *pseudomallei* protein homologous to the *Y. pestis* pilus assembly protein Caf1M (a molecular

271   chaperone involved with assembly of the surface capsule of the bacterium) was also identified.

272

273   Of the 15 putative *B. pseudomallei* DsbA substrates identified using this substrate homology

274   method, two were also identified in the genomic pipeline method. These are the PenI and

275   succinate dehydrogenase flavoprotein subunit homologues.

276

277   We then aligned the sequences of the Table 3 *B. pseudomallei* proteins to identify any possible

278   sequence conservation around the cysteine residues, but no pattern was identified. This lack of

279   peptide sequence motif in DsbA substrates has also been observed in *E.coli*, demonstrating the

280   difficulty of DsbA substrate prediction [52].

281

## Epitope prediction of virulence-associated proteins

283   To determine whether the DsbA substrates identified in the two methods above could

284   contribute to vaccination efforts against *B. pseudomallei*, we also predicted B-cell epitopes,

285   using a structure-informed approach. The sequences of the 73 putative, extra-cytoplasmic

286   DsbA substrates (predicted virulence factors, Table 2) along with the 15 homologous DsbA

287   substrates (Table 3) were screened against the Protein Data Bank (PDB) [53], to identify

288   structurally characterised homologues (see S6 File). Six of the 73 proteins were found to have

289   at least 80% similarity to a structurally characterised protein. Three of these six protein

290   structures were from *Pseudomonas* species, while the other three were from *Burkholderia*

291   species. Similarity was used rather than identity to account for mutations of functionally similar

292   residues. The six protein structures were then used as models to predict structurally-informed

293   B-cell epitopes of length 10-20 residues (Table 4 and Fig 4) using the SEPPA3 server.

294 **Table 4: B-cell epitope prediction.**

| Gene name | Predicted epitopes | Homologue PDB code | Accession number |
|---|---|---|---|
| beta-lactamase Toho-1 | RREPELNTALPGDER; TTMRNPNAQARDDVIA | 3W4O | KGV04506.1 |
| type 1 fimbrial protein | SSKAYTIAEGDNTF | 5N2B | WP_063597677.1 |
| triacylglycerol lipase | SSTNNTNQDALA; AYVQQVLAATGASK | 1HQD | WP_038741497.1 |
| class D beta-lactamase | VSGDPGQNNGLDR | 6NI0 | EDO89205.1 |
| triacylglycerol lipase | QQVLAVTGAQK; SHTHNTNQDAIA | 1HQD | WP_038775093.1 |
| S8 family serine peptidase | SGDEGVYECNNRGYPDGSNYTV; SNETVWNEGLDGNGKLW; YECNNRGYPDGSNYTV; MADLDASGNTGLTQ; QTNGSGGNYSDDQEG; GYSGYGYKASTGWDY | 1GA1/1NLU | WP_004553586.1 |

295 The virulence-associated putative DsbA substrates (Table 2) were screened for ≥80% similarity
296 to proteins within the PDB to account for substitution of functionally similar residues. The
297 structures were then screened for epitopes using SEPPA 3.0. Fourteen B-cell epitopes of 10 to
298 20 residues were predicted.
299

300 **Fig 4: Predicted B-cell epitopes.** Graphical representation of B-cell epitopes found in Table
301 4. Proteins are shown as white surfaces and their respective PDB ID is given in the bottom
302 left corner of each box. The epitope region is highlighted in red and the corresponding
303 homologous sequences found in *B. pseudomallei* are given in one letter code under each
304 respective structure and separated by semicolon when more than one sequence pointed to the
305 same epitope.

306

307 These epitopes provide an interesting list for further evaluation. For example, epitopes from

308 beta-lactamase Toho-1 and class D beta-lactamase could provide a useful vaccination approach

309 for *B. pseudomallei* because these directly target antibiotic resistance proteins. Similar

310 approaches have conferred protection against other bacteria in animal models [54-57].

311

312 Vaccination targeting adhesion proteins and essential virulence factors such as FimA [58, 59]

313 and type 1 fimbrial protein is a commonly used approach due to the external localisation of

314 these proteins and their exposure to host immune systems. Anti-fimbrial antibodies have been

315 shown to interfere with function and reduce disease [60, 61] and a FimA vaccine provided

316    protection against *Streptococcus parasanguis, Streptococcus mitis, Streptococcus mutans and*

317    *Streptococcus salivarius* in rats [62-64].

318

319    Vaccination against conserved, secreted enzymes such as the triacylglycerol lipase (EstA) and

320    S8 family serine peptidase enzymes may also be a useful strategy. Secreted peptidases are

321    known virulence factors in many pathogenic bacteria [37, 65] and vaccines targeting them have

322    attenuated disease in animal models [66, 67]. Two triacylglycerol lipases (WP_038741497.1

323    and WP_038775093.1) were identified as having a structural homologue in the PDB. These

324    two lipases are both core genes and share 78% similarity (72% identity, 87% query cover). and

325    their sequences were both aligned to the same PDB code, resulting in epitope variants of similar

326    sequences.

327

## Discussion

328

329   In the present study, we analysed genomes from 512 *B. pseudomallei* isolates specifically

330   associated with disease to identify core putative DsbA substrates and virulence factors.

331   Pangenomic analysis of *B. pseudomallei* has previously been performed utilising 37 isolates

332   from a variety of isolation sources [36] and concluded the pangenome to be 'open', indicating

333   that new isolates will continually increase the number of total genes, which we found to be the

334   case, based on a pangenome of 19,991 genes from 512 isolates. Previous studies comparing

335   the *B. pseudomallei* genome with the obligate pathogen *Burkholderia mallei* (responsible for

336   glanders) and the generally non-pathogenic *Burkholderia thailandensis* [68-71], identified

337   several loci likely to be involved in *B. pseudomallei* virulence. These include the capsular

338   polysaccharide gene cluster and Type III secretion needle complex [71], which were not

339   considered core genes, demonstrating the importance of large-scale analysis.

340

341   In the present study, we used two orthogonal approaches to identify a total of 278 putative

342   DsbA substrates, with 86 predicted to be virulence factors (S5 File). Of these, 73 were

343   identified by the genome analysis approach and 15 were identified by the DsbA substrate

344   homology approach. Two of the putative 86 DsbA virulence factor substrates were identified

345   in both approaches. These two are the experimentally validated bacterial virulence factors and

346   DsbA substrates succinate dehydrogenase flavoprotein subunit, and a PenI family β-lactamase

347   (both reported to be *F. tularensis* DsbA substrates) [34].

348

349   Delving deeper into the results presents some curious outcomes. For example, the well-

350   characterised *E.coli* DsbA substrate and virulence factor FlgI [27, 72] was not picked up as a

351   potential *B. pseudomallei* DsbA substrate by either method, though *B. pseudomallei* encodes

352   FlgI. The *B. pseudomallei* FlgI sequence has 4 cysteines in the translated gene product but the

353    predicted mature sequence after cleavage of the signal sequence has just one cysteine.

354    Generally, DsbA does not interact with proteins having just one cysteine. If *B. pseudomallei*

355    FlgI is a DsbA substrate (that is yet to be tested), then the most likely reasons that it was not

356    identified as a substrate by either of the two methods we used are that (i) the predicted signal

357    peptide is incorrect and/or (ii) the single cysteine of *B. pseudomallei* FlgI forms an inter-

358    molecular disulfide bond.

359

360    The finding that the two orthogonal approaches identified the same two target proteins suggests

361    that there is merit in using different theoretical approaches to select high priority targets for

362    further evaluation (in this case, the PenI family Beta-lactamase and succinate dehydrogenase

363    flavoprotein subunit). On the other hand, the fact that there were so few overlaps in the

364    predicted substrates from the two methods raises questions about the filters we applied.

365    Specifically, we found that of the 15 potential substrates identified by the substrate homology

366    method, 5 had an odd numbers of cysteines, whereas the genomic analysis filtered these

367    proteins out of consideration. We applied the even cysteine filter because previous reports

368    showed that *E. coli* exported proteins have a strong preference for an even number of cysteines.

369    This even number of cysteine preference is present in *B. pseudomallei* exported proteins (Fig

370    2) though is not as pronounced as in *E. coli*. By restricting our genomic analysis to core, extra-

371    cytoplasmic *B. pseudomallei* proteins with an even number of cysteines, some DsbA substrates

372    may therefore have been missed. There is considerable evidence that many virulence factors

373    such as adhesion and motility proteins, toxins and enzymes are extra-cytoplasmic proteins in

374    both Gram-positive and Gram-negative bacteria [21, 22, 73]. Given that extra-cytoplasmic

375    proteins in the translated core genome *of B. pseudomallei* have a slight preference for even

376    number of cysteines (Fig 2) and the identification of many virulence-associated proteins within

377    the 263 proteins in the list, the approach taken in this analysis (Fig 1) to identify DsbA

378 substrates was justified. Further, the genomic analysis focused on highly conserved proteins

379 from the core genome; accessory proteins associated with virulence would not be identified

380 using this approach. Nevertheless, the genomic analysis identified homologues of known DsbA

381 substrates in other bacteria, such as the OmpA porin, supporting the use of this approach.

382 However, attempting to identify epitopes from proteins which are not found in every disease-

383 causing isolate may present challenges for anti-virulence and vaccination attempts.

384

385 In addition, the genomic analysis identified several proteins of unknown function which could

386 represent novel virulence factors for future studies. Importantly, our theoretical approach was

387 extended to predict structurally-informed surface epitopes for several core gene DsbA

388 substrates for potential vaccine or antibody development (Table 4).

389

390 In summary, our *in silico* analysis combined a substrate homology approach and a genomic

391 analysis approach to identify more than 80 potential *B. pseudomallei* DsbA virulence factor

392 substrates, two of which we mark as high priority for experimental validation. Future

393 characterization of these proteins will aid our understanding of *B. pseudomallei* virulence and

394 could provide new targets for antivirulence drug discovery and vaccine development. The

395 approaches we report here could also be applied to identify potential DsbA virulence factor

396 substrates in other pathogenic bacteria.

397

# Methods

## Data acquisition and filtering of core, extra-cytoplasmic, putative DsbA substrates

1577 *B. pseudomallei* genomes were obtained from the genome information table from NCBI (https://www.ncbi.nlm.nih.gov/genome/genomes/476) (date accessed: 1/2/20).The biosample accession numbers were batch downloaded using Entrez. A list of assembly accession numbers can be found in S1 Fig. Metadata was then scraped for disease association using grep with the following command:

```
grep -A 1 "disease"
```

The assemblies were then downloaded using Entrez and annotated using a prokka (version 1.14.5) [74] for loop with the following command:

```
for file in *.fna; do tag=${file%.fna}; prokka --prefix "$tag" --locustag "$tag" --genus Burkholderia --
species pseudomallei --strain "$tag" --outdir "$tag"_prokka --force --addgenes "$file"; done
```

The .gff files were used as input for roary (version 3.11.2) [75] without splitting paralogues via the following command:

```
roary -e --mafft -i 90 -v -p 72 -z -s -o output -f *.gff
```

The roary output file was altered from interleaved fasta to one line per sequence

```
awk '{if(NR==1) {print $0} else {if($0 ~ /^>/) {print "\n"$0} else {printf $0}}}' input.fa > output.fa
```

The core genome was then used in the remaining analysis and core DNA sequences were translated into protein sequences using transeq [76] with the following command:

```
transeq -sequence input.fasta -outseq output.fasta -table 11 -frame 1
```

424    The core genome was then filtered based on signal sequence and then the sequence of the

425    mature exported protein, as predicted utilising SignalP 5.0 [77, 78]

426        signalp -fasta prot_core_genome_complete.fasta -format short -mature -org gram- -verbose

427

428    These sequences were then filtered for genes containing even numbers of cysteines

429        awk -F \C 'NF % 2' < input.fasta | awk "/C.*C/" | sed '/>/{$!N;/\n.*>/!P;D}' > output.fasta

430

431    This list was then annotated via screening sequences against NCBI and Gene Ontology [79]

432    using the PANNZER2 server [80].

433

## Identification of DsbA substrate homologues in *B. pseudomallei*

435    DsbA substrates were also predicted using a substrate homology search. This approach may

436    identify proteins not encoded in the core genome. The *B. pseudomallei* genome was screened

437    for homologues of known DsbA substrates using BLASTP. A starting list of confirmed DsbA

438    substrates was extracted from the literature [22, 34, 45-48, 81], and their amino acid sequences

439    used in BLAST searches [82] against the NCBI protein database [43] for homologues in *B.*

440    *pseudomallei* using default search parameters. In some cases two search proteins identified the

441    same homologue in *B. pseudomallei*. In these cases only the search protein most similar to the

442    *B. pseudomallei* homologue is given in Table 3. The results were filtered to select proteins with

443    at least 20% sequence identity and a sequence coverage of at least 50%. Protein sequences with

444    fewer than two cysteines were removed. Exported proteins were selected on the basis of

445    predicted signal sequence (SignalP 5.0 [77]) or experimental evidence of extra-cytoplasmic

446    localisation for the reported DsbA substrate in another *Burkholderia* species.

447

21

## Cysteine distribution analysis

449    Fasta files containing either the 19,991 pan genes or the 4,496 core gene of *B. pseudomallei*

450    with their corresponding amino acid sequences and descriptors were utilised to calculate the

451    distribution of cysteines with a custom Python 3.0 script (available on Github :

452    (https://github.com/gpetit99/cysteineCount_bPseudomallei/blob/master/CysCountFrequency.

453    py"). Briefly, lists of the extra-cytoplasmic  protein sequences with signal peptides removed

454    were compared to lists of the protein sequences from the whole genome to create dataframes

455    with either cytoplasmic or extra-cytoplasmic proteins. Proteins were grouped based on the

456    presence or absence of SP, and based on the number of cysteines in the mature protein. To

457    calculate the normalised frequency of cysteines for extra-cytoplasmic proteins, we divided

458    the number of extra-cytoplasmic proteins having N cysteines by the total number of proteins

459    having N cysteines (N being an integer from 0 to 73 – No protein has more than 73 cysteines

460    in the *B. pseudomallei* translated genome). This analysis was run for the core genome and

461    pangenome independently. Other statistics (e.g. number of proteins in each group) were

462    extracted from the dataframes.

463

## Epitope prediction

465    The metadata for each of the 263 proteins in the annotated list was manually inspected to select

466    for further analysis a total of 73 proteins likely related to virulence. The sequences of these 73

467    selected proteins were combined with the 15 selected proteins from the homology analysis (to

468    give 86 unique protein sequences). These were screened against the protein data bank using

469    BLAST (criteria: ≥80% positive substitutions/similarity used as a threshold) to find structurally

470    characterised homologues. These structural homologues were then used to predict B-cell

471    epitopes using SEPPA 3.0 (http://www.badd-cao.net/seppa3/index.html) with a threshold of

472    0.1 [83]. Similarity was used rather than identity to account for mutations of functionally

22

473    similar residues. Predicted B-cell epitopes were accepted if they were $10 - 20$ residues in

474    length, as described in [84].

475

476

477

# Acknowledgments

# Conflict of interest

483   The authors declare that there are no conflicts of interest.

484

485

486   .

# References

1.    White N. Melioidosis. *The Lancet* 2003;361(9370):1715-22. DOI: 10.1016/s0140-6736(03)13374-0

2.    Chakravorty A, Heath C. Melioidosis: An updated review. *Aus J Gen Pract* 2019;48:327-32. DOI: 10.31128/AJGP-04-18-4558

3.    Willcocks SJ, Denman CC, Atkins HS, Wren BW. Intracellular replication of the well-armed pathogen *Burkholderia pseudomallei*. *Curr Opin Microbiol* 2016;29:94-103. DOI: 10.1016/j.mib.2015.11.007

4.    Wiersinga WJ, Virk HS, Torres AG, Currie BJ, Peacock SJ, Dance DAB, et al. Melioidosis. *Nat Rev Dis Primers* 2018;4:17107. DOI: 10.1038/nrdp.2017.107

5.    Rhodes KA, Schweizer HP. Antibiotic resistance in *Burkholderia* species. *Drug Resist Updat* 2016;28:82-90. DOI: 10.1016/j.drup.2016.07.003

6.    Podnecky NL, Rhodes KA, Mima T, Drew HR, Chirakul S, Wuthiekanun V, *et al*. Mechanisms of resistance to folate pathway inhibitors in *Burkholderia pseudomallei*: deviation from the norm. *mBio* 2017;8(5):e01357-17. DOI: 10.1128/mBio.01357-17

7.    Held K, Gasper J, Morgan S, Siehnel R, Singh P, Manoil C. Determinants of extreme ß-lactam tolerance in the *Burkholderia pseudomallei* complex. *Antimicrob Agents Chemother* 2018;62(4)e00068-18. DOI: 10.1128/AAC.00068-18

8.    Podnecky NL, Rhodes KA, Schweizer HP. Efflux pump-mediated drug resistance in *Burkholderia*. *Front Microbiol* 2015;6:305. DOI: 10.3389/fmicb.2015.00305

9.    Schweizer HP. Mechanisms of antibiotic resistance in *Burkholderia pseudomallei*: implications for treatment of melioidosis. *Future microbiol* 2012;7(12):1389-99. DOI: 10.2217/fmb.12.116

10.   Dance D. Treatment and prophylaxis of melioidosis. *Int J Antimicrob Agent* 2014;43(4):310-8. DOI: 10.1016/j.ijantimicag.2014.01.005

512    11.    Antimicrobial resistance: global report on surveillance: World Health Organization;

513    2014.

514    12.    Kennedy DA, Read AF. Why the evolution of vaccine resistance is less of a concern

515    than the evolution of drug resistance. *Proc Natl Acad Sci U S A*. 2018;115(51):12878-86.

516    DOI: 10.1073/pnas.1717159115

517    13.    Thabit AK, Crandon JL, Nicolau DP. Antimicrobial resistance: impact on clinical and

518    economic outcomes and the need for new antimicrobials. *Expert Opin Pharmacother*

519    2015;16(2):159-77. DOI: 10.1517/14656566.2015.993381

520    14.    Rasko DA, Sperandio V. Anti-virulence strategies to combat bacteria-mediated

521    disease. *Nat Rev Drug Discov* 2010;9(2):117-28. DOI: 10.1038/nrd3013

522    15.    Johnson MM, Ainslie KM. Vaccines for the Prevention of Melioidosis and Glanders.

523    *Curr Trop Med Rep* 2017;4(3):136-45. DOI: 10.1007/s40475-017-0121-7

524    16.    Hara Y, Mohamed R, Nathan S. Immunogenic *Burkholderia pseudomallei* outer

525    membrane proteins as potential candidate vaccine targets. *PLoS One* 2009;4(8):e6496. DOI:

526    10.1371/journal.pone.0006496

527    17.    Nagpal G, Usmani SS, Raghava GPS. A Web Resource for Designing Subunit

528    Vaccine Against Major Pathogenic Species of Bacteria. *Front Immunol* 2018;9(2280). DOI:

529    10.3389/fimmu.2018.02280

530    18.    Mühlen S, Dersch P. Anti-virulence Strategies to Target Bacterial Infections. In:

531    Stadler M, Dersch P, editors. How to Overcome the Antibiotic Crisis : Facts, Challenges,

532    Technologies and Future Perspectives. Cham: Springer International Publishing; 2016. p.

533    147-83.

534    19.    Heras B, Scanlon MJ, Martin JL. Targeting virulence not viability in the search for

535    future antibacterials. *Brit J Clin Pharmacol* 2015;79(2):208-15. DOI: 10.1111/bcp.12356

536   20.    Anfinsen CB. Principles that govern the folding of protein chains. *Science*

537   1973;181(4096):223-30. DOI: 10.1126/science.181.4096.223

538   21.    Smith RP, Paxman JJ, Scanlon MJ, Heras B. Targeting Bacterial Dsb Proteins for the

539   Development of Anti-Virulence Agents. *Molecules* 2016;21(7). DOI:

540   10.3390/molecules21070811

541   22.    Heras B, Shouldice SR, Totsika M, Scanlon MJ, Schembri MA, Martin JL. DSB

542   proteins and bacterial pathogenicity. *Nat Rev Microbiol* 2009;7(3):215. DOI:

543   10.1038/nrmicro2087

544   23.    Bocian-Ostrzycka KM, Grzeszczuk MJ, Banaś AM, Jagusztyn-Krynicka EK.

545   Bacterial thiol oxidoreductases—from basic research to new antibacterial strategies. *Appl*

546   *Microbiol Biotechnol* 2017;101(10):3977-89. DOI: 10.1007/s00253-017-8291-8

547   24.    Shouldice SR, Heras B, Walden PM, Totsika M, Schembri MA, Martin JL. Structure

548   and function of DsbA, a key bacterial oxidative folding catalyst. *Antioxid Redox Signal*

549   2011;14(9):1729-60. DOI: 10.1089/ars.2010.3344

550   25.    Ireland PM, McMahon RM, Marshall LE, Halili M, Furlong E, Tay S, et al.

551   Disarming Burkholderia pseudomallei: structural and functional characterization of a

552   disulfide oxidoreductase (DsbA) required for virulence in vivo. *Antioxid Redox Signal*

553   2014;20(4):606-17. DOI: 10.1089/ars.2013.5375

554   26.    McMahon RM, Ireland PM, Sarovich DS, Petit G, Jenkins CH, Sarkar-Tyson M, et al.

555   Virulence of the Melioidosis Pathogen Burkholderia pseudomallei Requires the

556   Oxidoreductase Membrane Protein DsbB. *Infect Immun* 2018;86(5) DOI: 10.1128/IAI.00938-

557   17

558   27.    Dailey FE, Berg HC. Mutants in disulfide bond formation that disrupt flagellar

559   assembly in *Escherichia coli. Proc Natl Acad Sci USA* 1993;90(3):1043-7. DOI:

560   10.1073/pnas.90.3.1043

561    28.    Totsika M, Heras B, Wurpel DJ, Schembri MA. Characterization of two homologous

562    disulfide bond systems involved in virulence factor biogenesis in uropathogenic *Escherichia*

563    *coli* CFT073. *J Bacteriol* 2009;191(12):3901-8 DOI: 10.1128/JB.00143-09

564    29.    Kurth F, Rimmer K, Premkumar L, Mohanty B, Duprez W, Halili MA, et al.

565    Comparative sequence, structure and redox analyses of *Klebsiella pneumoniae* DsbA show

566    that anti-virulence target DsbA enzymes fall into distinct classes. *PLoS One*

567    2013;8(11):e80210. DOI: 10.1371/journal.pone.0080210

568    30.    Heras B, Totsika M, Jarrott R, Shouldice SR, Guncar G, Achard MES, et al.

569    Structural and functional characterization of three DsbA paralogues from *Salmonella enterica*

570    serovar Typhimurium. *J Biol Chem* 2010;285(24):18423-32. DOI: 10.1074/jbc.M110.101360

571    31.    Straskova A, Pavkova I, Link M, Forslund A-L, Kuoppa K, Noppa L, et al. Proteome

572    analysis of an attenuated *Francisella tularensis* dsbA mutant: Identification of potential

573    DsbA substrate proteins. *J Proteome Res* 2009;8(11):5336-46. DOI: 10.1021/pr900570b

574    32.    Hatahet F, Boyd D, Beckwith J. Disulfide bond formation in prokaryotes: history,

575    diversity and design. *Biochim Biophys Acta* 2014;1844(8):1402-14. DOI:

576    10.1016/j.bbapap.2014.02.014

577    33.    Dutton RJ, Boyd D, Berkmen M, Beckwith J. Bacterial species exhibit diversity in

578    their mechanisms and capacity for protein disulfide bond formation. *Proc Natl Acad Sci USA*

579    2008;105(33):11933-8. DOI: 10.1073/pnas.0804621105

580    34.    Ren G, Champion MM, Huntley JF. Identification of disulfide bond isomerase

581    substrates reveals bacterial virulence factors. *Mol Microbiol* 2014;94(4):926-44. DOI:

582    10.1111/mmi.12808

583    35.    Benson DA, Clark K, Karsch-Mizrachi I, Lipman DJ, Ostell J, Sayers EW. GenBank.

584    Nucleic acids research. 2015;43(Database issue):D30.

585   36.   Spring-Pearson SM, Stone JK, Doyle A, Allender CJ, Okinaka RT, Mayo M, *et al.*

586   Pangenome analysis of *Burkholderia pseudomallei*: Genome evolution preserves gene order

587   despite high recombination rates. *PLoS One* 2015;10(10):e0140274. DOI:

588   10.1371/journal.pone.0140274

589   37.   Backert S, Bernegger S, Skórko-Glonek J, Wessler S. Extracellular HtrA serine

590   proteases: An emerging new strategy in bacterial pathogenesis. *Cell Microbiol*

591   2018;20(6):e12845. DOI: 10.1111/cmi.12845

592   38.   Gosink KK, Mann ER, Guglielmo C, Tuomanen EI, Masure HR. Role of novel

593   choline binding proteins in virulence of *Streptococcus pneumoniae*. *Infect Immun*

594   2000;68(10):5690-5. DOI: 10.1128/iai.68.10.5690-5695.2000

595   39.   Nakamya MF, Ayoola MB, Park S, Shack LA, Swiatlo E, Nanduri B. The role of

596   cadaverine synthesis on *Pneumococcal* capsule and protein expression. *Med Sci (Basel)*

597   2018;6(1):8. DOI: 10.3390/medsci6010008

598   40.   Koski P, Vaara M. Polyamines as constituents of the outer membranes of *Escherichia*

599   *coli* and *Salmonella typhimurium*. *J Bacteriol* 1991;173(12):3695-9. DOI:

600   10.1128/jb.173.12.3695-3699.1991

601   41.   Yethon JA, Vinogradov E, Perry MB, Whitfield C. Mutation of the

602   lipopolysaccharide core glycosyltransferase encoded by waaG destabilizes the outer

603   membrane of *Escherichia coli* by interfering with core phosphorylation. *J Bacteriol*

604   2000;182(19):5620-3. DOI: 10.1128/jb.182.19.5620-5623.2000

605   42.   Wortham BW, Oliveira MA, Fetherston JD, Perry RD. Polyamines are required for

606   the expression of key Hms proteins important for *Yersinia pestis* biofilm formation. *Environ*

607   *Microbiol* 2010;12(7):2034-47. DOI: 10.1111/j.1462-2920.2010.02219.x

608   43.   Database resources of the National Center for Biotechnology Information. Nucleic

609   Acids Res. 2016;44(D1):D7-19. DOI: 10.1093/nar/gkv1290

610 44. Hayashi S, Abe M, Kimoto M, Furukawa S, Nakazawa T. The dsbA-dsbB disulfide

611 bond formation system of *Burkholderia cepacia* is involved in the production of protease and

612 alkaline phosphatase, motility, metal resistance, and multi-drug resistance. *Microbiol*

613 *Immunol* 2000;44(1):41-50. DOI: 10.1111/j.1348-0421.2000.tb01244.x

614 45. Corbett C, Burtnick M, Kooi C, Woods D, Sokol P. An extracellular zinc

615 metalloprotease gene of *Burkholderia cepacia*. *Microbiology*. 2003;149(8):2263-71. DOI:

616 10.1099/mic.0.26243-0

617 46. Abe M, Nakazawa T. The dsbB gene product is required for protease production by

618 *Burkholderia cepacia*. *Infect Immun* 1996;64(10):4378-80.

619 47. Kooi C, Subsin B, Chen R, Pohorelic B, Sokol P. *Burkholderia cenocepacia* ZmpB is

620 a broad-specificity zinc metalloprotease involved in virulence. *Infect Immun*

621 2006;74(7):4083-93. DOI: 10.1128/IAI.00297-06

622 48. Kooi C, Corbett C, Sokol P. Functional analysis of the *Burkholderia cenocepacia*

623 ZmpA metalloprotease. *J Bacteriol* 2005;187(13):4421-9. DOI: 10.1128/JB.187.13.4421-

624 4429.2005

625 49. Papp-Wallace KM, Becka SA, Taracila MA, Winkler ML, Gatta JA, Rholl DA, *et al.*

626 Exposing a β-Lactamase "Twist": the mechanistic basis for the high level of ceftazidime

627 resistance in the C69F variant of the *Burkholderia pseudomallei* PenI β-Lactamase.

628 *Antimicrob Agents Chemother* 2016;60(2):777-88. DOI: 10.1128/aac.02073-15

629 50. Miki T, Okada N, Danbara H. Two periplasmic disulfide oxidoreductases, DsbA and

630 SrgA, target outer membrane protein SpiA, a component of the *Salmonella* pathogenicity

631 island 2 Type III secretion system. *J Biol Chem* 2004;279(33):34631-42. DOI:

632 10.1074/jbc.M402760200

633   51.     Miki T, Okada N, Kim Y, Abe A, Danbara H. DsbA directs efficient expression of

634   outer membrane secretin EscC of the enteropathogenic *Escherichia coli* Type III secretion

635   apparatus. *Microb Pathog* 2008;44(2):151-8. DOI: 10.1016/j.micpath.2007.09.001

636   52.     Paxman JJ, Borg NA, Horne J, Thompson PE, Chin Y, Sharma P, et al. The structure

637   of the bacterial oxidoreductase enzyme DsbA in complex with a peptide reveals a basis for

638   substrate specificity in the catalytic cycle of DsbA enzymes. *J Biol Chem*

639   2009;284(26):17835-45 DOI: 10.1074/jbc.M109.011502

640   53.     H.M. Berman, J. Westbrook, Z. Feng, G. Gilliland, T.N. Bhat, H. Weissig, I.N.

641   Shindyalov, P.E. Bourne The Protein Data Bank *Nucleic Acids Research* 2000; 28: 235-7.

642   DOI: 10.1093/nar/28.1.235

643   54.     Lipsitch M, Siber GR. How Can Vaccines Contribute to Solving the Antimicrobial

644   Resistance Problem? *mBio* 2016;7(3):e00428-16. DOI: 10.1128/mBio.00428-16

645   55.     Senna JP, Roth DM, Oliveira JS, Machado DC, Santos DS. Protective immune

646   response against methicillin resistant *Staphylococcus aureus* in a murine model using a DNA

647   vaccine approach. *Vaccine* 2003;21(19-20):2661-6. DOI: 10.1016/s0264-410x(02)00738-7

648   56.     Zarantonelli ML, Antignac A, Lancellotti M, Guiyoule A, Alonso J-M, Taha M-K.

649   Immunogenicity of meningococcal PBP2 during natural infection and protective activity of

650   anti-PBP2 antibodies against meningococcal bacteraemia in mice. *J Antimicrob Chemother*

651   2006;57(5):924-30. DOI: 10.1093/jac/dkl066

652   57.     Ciofu O, Bagge N, Høiby N. Antibodies against β-lactamase can improve ceftazidime

653   treatment of lung infection with β-lactam-resistant *Pseudomonas aeruginosa* in a rat model of

654   chronic lung infection. *APMIS* 2002;110(12):881-91. DOI: 10.1034/j.1600-

655   0463.2002.1101207.x

656  58.     Fenno JC, Shaikh A, Spatafora G, Fives-Taylor P. The fimA locus of *Streptococcus*

657  *parasanguis* encodes an ATP-binding membrane transport system. *Mol Microbiol*

658  1995;15(5):849-63. DOI: 10.1111/j.1365-2958.1995.tb02355.x

659  59.     Liu C-C, Ou S-C, Tan D-H, Hsieh M-K, Shien J-H, Chang P-C. The fimbrial protein

660  is a virulence factor and potential vaccine antigen of *Avibacterium paragallinarum*. *Avian*

661  *Dis* 2016;60(3):649-55. DOI: 10.1637/11410-031316-Reg.1

662  60.     Holmgren J, Svennerholm A-M. Vaccines against mucosal infections. *Curr Opin*

663  *Immunol* 2012;24(3):343-53. DOI: 10.1016/j.coi.2012.03.014

664  61.     Singh B, Mortezaei N, Savarino SJ, Uhlin BE, Bullitt E, Andersson M. Antibodies

665  damage the resilience of fimbriae, causing them to be stiff and tangled. *J Bacteriol*

666  2016;199(1):e00665-16. DOI: 10.1128/JB.00665-16

667  62.     Viscount HB, Munro CL, Burnette-Curley D, Peterson DL, Macrina FL.

668  Immunization with FimA protects against *Streptococcus parasanguis* endocarditis in rats.

669  *Infect Immun* 1997;65(3):994-1002.

670  63.     Kitten T, Munro CL, Wang A, Macrina FL. Vaccination with FimA from

671  *Streptococcus parasanguis* protects rats from endocarditis caused by other viridans

672  streptococci. *Infect Immun* 2002;70(1):422-5. DOI: 10.1128/iai.70.1.422-425.2002

673  64.     Vandemaele F, Ververken C, Bleyen N, Geys J, D'Hulst C, Addwebi T, *et al.*

674  Immunization with the binding domain of FimH, the adhesin of type 1 fimbriae, does not

675  protect chickens against avian pathogenic *Escherichia coli*. *Avian Pathol* 2005;34(3):264-72.

676  DOI: 10.1080/03079450500112682

677  65.     Hritonenko V, Stathopoulos C. Omptin proteins: an expanding family of outer

678  membrane proteases in Gram-negative *Enterobacteriaceae*. *Mol Membr Biol* 2007;24(5-

679  6):395-406. DOI: 10.1080/09687680701443822

680    66.    Santillan DA, Andracki ME, Hunter SK. Protective immunization in mice against

681    group B streptococci using encapsulated C5a peptidase. *Am J Obstet Gynecol*

682    2008;198(1):114. e1-e6. DOI: 10.1016/j.ajog.2007.06.003

683    67.    Marana MH, Jørgensen LvG, Skov J, Chettri JK, Holm Mattsson A, Dalsgaard I, *et*

684    *al.* Subunit vaccine candidates against *Aeromonas salmonicida* in rainbow trout

685    *Oncorhynchus mykiss*. *PLoS One* 2017;12(2):e0171944. DOI: 10.1371/journal.pone.0171944

686    68.    Ong C, Ooi CH, Wang D, Chong H, Ng KC, Rodrigues F, et al. Patterns of large-

687    scale genomic variation in virulent and avirulent *Burkholderia* species. *Genome Res*

688    2004;14(11):2295-307. DOI: 10.1101/gr.1608904

689    69.    Kim HS, Schell MA, Yu Y, Ulrich RL, Sarria SH, Nierman WC, *et al.* Bacterial

690    genome adaptation to niches: divergence of the potential virulence genes in three

691    *Burkholderia* species of different survival strategies. *BMC Genomics* 2005;6(1):174. DOI:

692    10.1186/1471-2164-6-174

693    70.    Majerczyk CD, Brittnacher MJ, Jacobs MA, Armour CD, Radey MC, Bunt R, *et al.*

694    Cross-species comparison of the *Burkholderia pseudomallei, Burkholderia thailandensis*, and

695    *Burkholderia mallei* quorum-sensing regulons. *J Bacteriol* 2014;196(22):3862-71. DOI:

696    10.1128/JB.01974-14

697    71.    Yu Y, Kim HS, Chua HH, Lin CH, Sim SH, Lin D, et al. Genomic patterns of

698    pathogen evolution revealed by comparison of *Burkholderia pseudomallei,* the causative

699    agent of melioidosis, to avirulent *Burkholderia thailandensis. BMC Microbiology*

700    2006;6(1):46. DOI: 10.1186/1471-2180-6-46

701    72.    Hizukuri Y, Yakushi T, Kawagishi I, Homma M. Role of the intramolecular disulfide

702    bond in FlgI, the flagellar P-ring component of *Escherichia coli*. *J Bacteriol*

703    2006;188(12):4190-7 DOI: 10.1128/JB.01896-05

704 73.    Allen RC, Popat R, Diggle SP, Brown SP. Targeting virulence: can we make

705 evolution-proof drugs? *Nat Rev Microbiol* 2014;12(4):300-8. DOI: 10.1038/nrmicro3232

706 74.    Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics*.

707 2014;30(14):2068-9 DOI: 10.1093/bioinformatics/btu153

708 75.    Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MT, *et al.* Roary: rapid

709 large-scale prokaryote pan genome analysis. *Bioinformatics* 2015;31(22):3691-3. DOI:

710 10.1093/bioinformatics/btv421

711 76.    Madeira F, Park YM, Lee J, Buso N, Gur T, Madhusoodanan N, *et al.* The EMBL-

712 EBI search and sequence analysis tools APIs in 2019. *Nucleic Acids Res*

713 2019;47(W1):W636-W41. DOI: 10.1093/nar/gkz268

714 77.    Almagro Armenteros JJ, Tsirigos KD, Sonderby CK, Petersen TN, Winther O,

715 Brunak S, *et al.* SignalP 5.0 improves signal peptide predictions using deep neural networks.

716 *Nat Biotechnol* 2019;37(4):420-3. DOI: 10.1038/s41587-019-0036-z

717 78.    Käll L, Krogh A, Sonnhammer EL. Advantages of combined transmembrane

718 topology and signal peptide prediction—the Phobius web server. *Nucleic Acids Res*

719 2007;35(Supplementary 2 Web server issue):W429-W32. DOI: 10.1093/nar/gkm256

720 79.    Gene Ontology Consortium. The Gene Ontology (GO) database and informatics

721 resource. *Nucleic Acids Res* 2004;32(Supplementary 1 Database issue):D258-D61. DOI:

722 10.1093/nar/gkh036

723 80.    Törönen P, Medlar A, Holm L. PANNZER2: a rapid functional annotation web

724 server. *Nucleic Acids Res* 2018;46(W1):W84-W8. DOI: 10.1093/nar/gky350

725 81.    Hayashi S, Abe M, Kimoto M, Furukawa S, Nakazawa T. The dsbA-dsbB disulfide

726 bond formation system of *Burkholderia cepacia* is involved in the production of protease and

727 alkaline phosphatase, motility, metal resistance, and multi-drug resistance. *Microbiol*

728 *Immunol* 2000;44(1):41-50. DOI: 10.1111/j.1348-0421.2000.tb01244.x

729     82.     Johnson M, Zaretskaya I, Raytselis Y, Merezhuk Y, McGinnis S, Madden TL. NCBI

730     BLAST: a better web interface. *Nucleic Acids Res* 2008;36(Supplementary 2 Web Server

731     issue):W5-9. DOI: 10.1093/nar/gkn201

732     83.     Zhou C, Chen Z, Zhang L, Yan D, Mao T, Tang K, *et al.* SEPPA 3.0—enhanced

733     spatial epitope prediction enabling glycoprotein antigens. *Nucleic Acids Res*

734     2019;47(W1):W388-W94. DOI: 10.1093/nar/gkz413

735     84.     Shey RA, Ghogomu SM, Esoh KK, Nebangwa ND, Shintouo CM, Nongley NF, *et al.*

736     In-silico design of a multi-epitope vaccine candidate against onchocerciasis and related

737     filarial diseases. *Sci Rep* 2019;9(1):1-18. DOI: 10.1038/s41598-019-40833-x

738

# Supporting Information

740     **S1 Fig.** Accession numbers for disease related genomes of *B. pseudomallei* used in this

741     analysis

742     **S2 Fig.** Core genome (4,496 gene products) of disease related B. pseudomallei (fasta format).

743     **S3 Fig** *B. pseudomallei* proteins from the core genome with a signal peptide (removed before

744     counting cysteines) and even number of cysteines (263 proteins, fasta format).

745     **S4 File Gene Ontology (GO) classification of the gene and gene-product descriptions.**

746     **S5 File Predicted virulence-associated substrates of DsbA**

747     **S6 File Predicted B-cell epitopes**

748

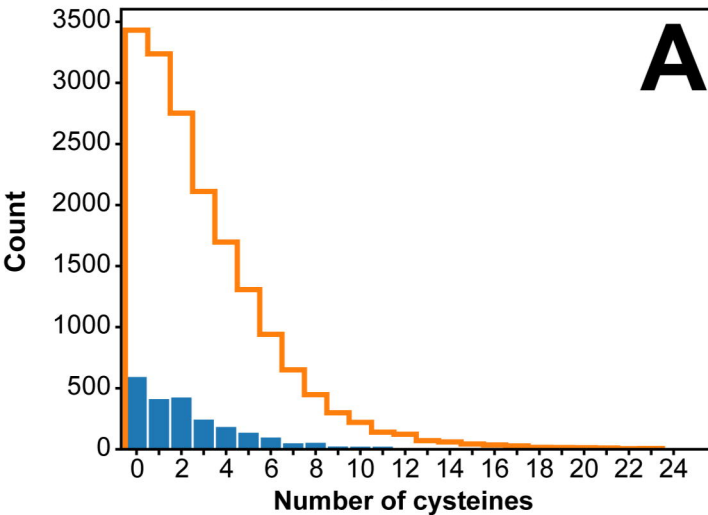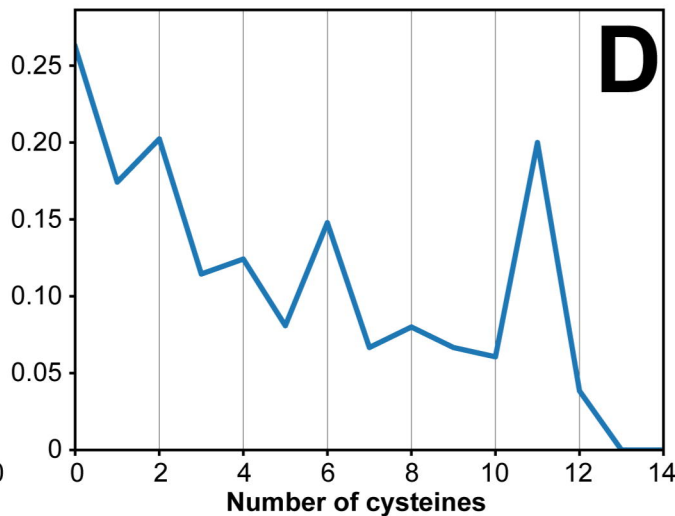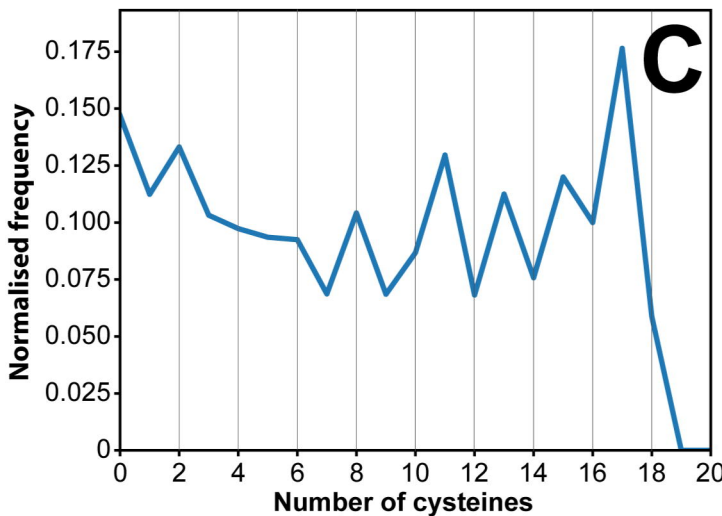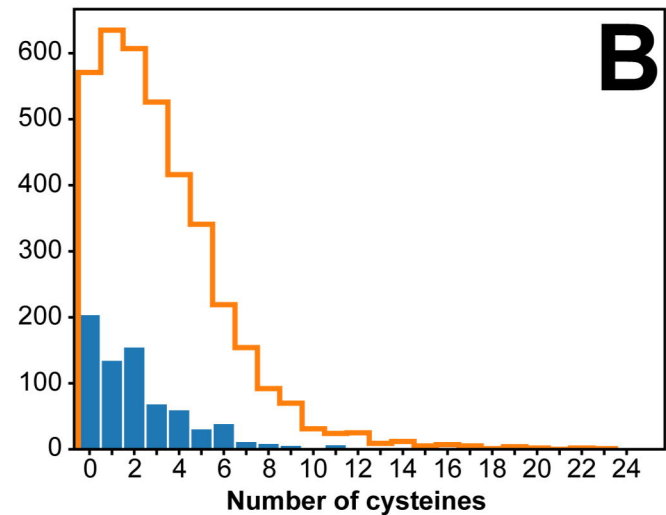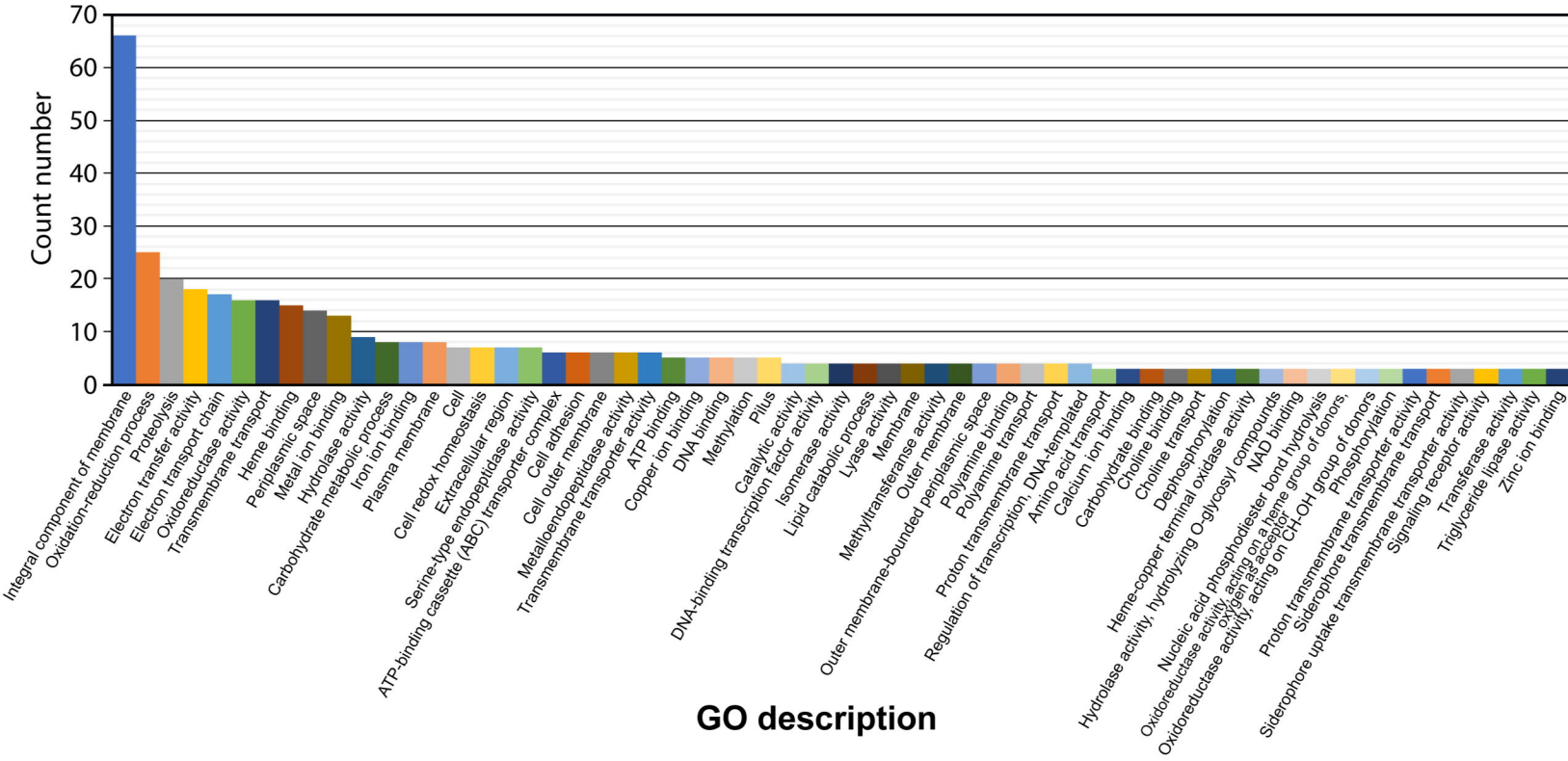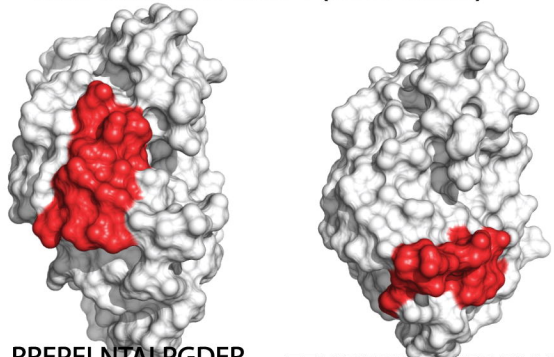| Total number of genomes found on NCBI | 1,577 *B. pseudomallei* genomes |
| Number of genomes associated with disease | 511 disease-associated genomes |
| Pan-genome of 511 disease-assocated genomes | 19,991 unique proteins |
| Core genome of 511 strains | 4,496 core proteins |
| Predicted proteins with signal sequence | 726 extra-cytoplasmic proteins |
| Putative DsbA substrates | 263 proteins with an even number of Cys |

**Pangenome**

**Core genome**

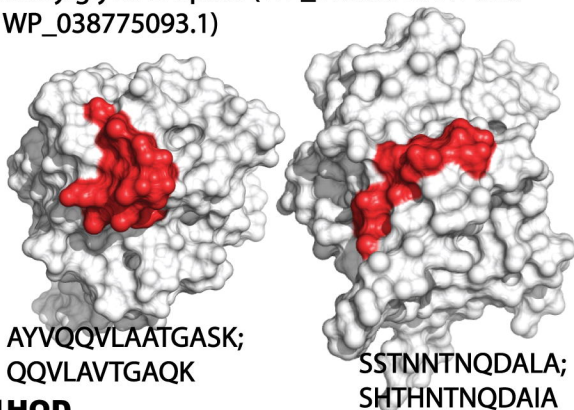Beta-lactamase Toho-1 (KGV04506.1)

RREPELNTALPGDER    TTMRNPNAQARDDVIA

3W4O

Type 1 fimbrial protein
(WP_063597677.1)

SSKAYTIAEGDNTF

5N2B

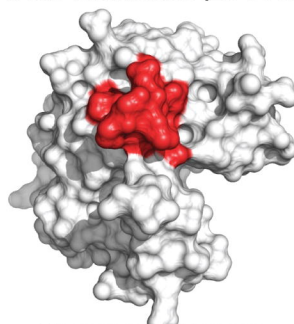Triacylglycerol lipase (WP_063597677.1 and
WP_038775093.1)

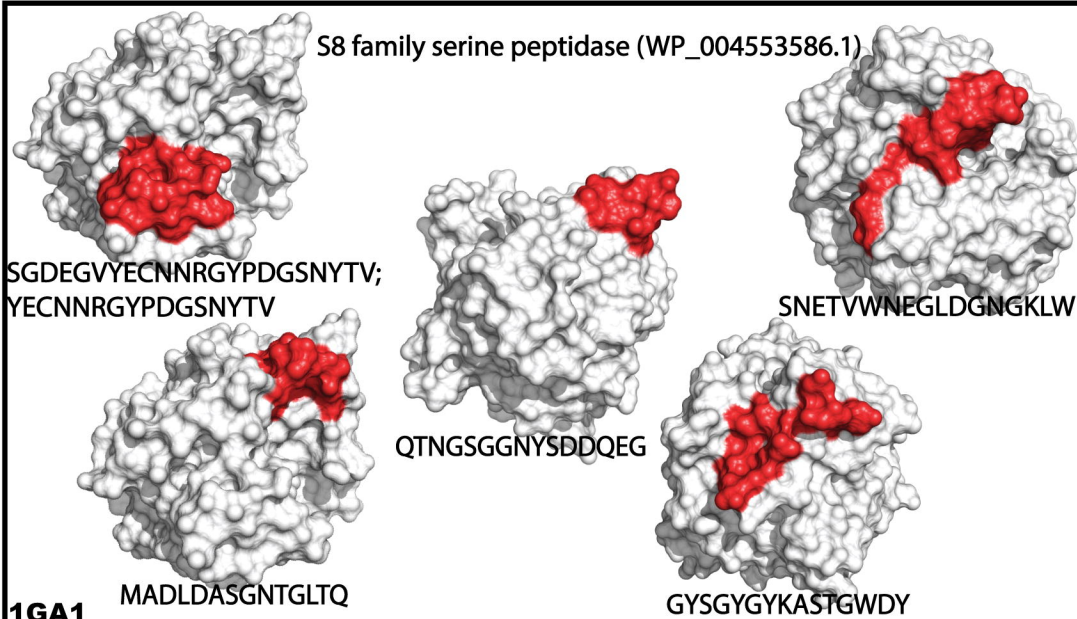AYVQQVLAATGASK;
QQVLAVTGAQK

SSTNNTNQDALA;
SHTHNTNQDAIA

1HQD

Class D beta-lactamase (EDO89205.1)

VSGDPGQNNGLDR

6NI0

S8 family serine peptidase (WP_004553586.1)

SGDEGVYECNNRGYPDGSNYTV;
YECNNRGYPDGSNYTV

SNETVWNEGLDGNGKLW

MADLDASGNTGLTQ

QTNGSGGNYSDDQEG

GYSGYGYKASTGWDY

1GA1