# Outer membrane 1 pore protein prediction in mycobacteria using genomic comparison

*Nancy Mah, Carolina Perez-Iratxeta, Miguel A. Andrade-Navarro*

# Outer membrane pore protein prediction in mycobacteria using genomic comparison

**Nancy Mah [1], Carolina Perez-Iratxeta [2], Miguel A. Andrade-Navarro [1]**

[1]Max Delbrück Center for Molecular Medicine, Robert-Rössle-Str. 10, 13125 Berlin, Germany

[2]Ottawa Health Research Institute, 501 Smyth Road, Ottawa, Ontario K1H 8L6, Canada

Correspondence: Nancy Mah, Max Delbrück Center for Molecular Medicine, Robert-Rössle-Str. 10, 13125 Berlin, Germany; Telephone: +49 30 9406-2831; Fax: +49 30 9406-4240; E-mail: nancy.mah@mdc-berlin.de

Running title: outer membrane pore protein prediction in mycobacteria

Contents category: Research Paper (Microbial Pathogenicity)

Abbreviations: OMP, outer membrane protein; Mt, *Mycobacterium tuberculosis*; Mb, *Mycobacterium bovis*; Ml, *Mycobacterium leprae*; Mu, *Mycobacterium ulcerans*; Ms, *Mycobacterium smegmatis*, Ma, *Mycobacterium avium*; Mm, *Mycobacterium marinum*, aa, amino acid; TM, transmembrane helix

## Summary

Proteins responsible for outer membrane transport across the unique membrane structure of *Mycobacterium* spp. are attractive drug targets in the treatment of human diseases caused by the mycobacterial pathogens, *M. tuberculosis*, *M. bovis*, *M. leprae* and *M. ulcerans*. In contrast to *E. coli*, relatively few outer membrane proteins (OMPs) have been identified in *Mycobacterium* spp., largely due to the difficulties in isolating mycobacterial membrane proteins and our incomplete understanding of secretion mechanisms and cell wall structure in these organisms. To further expand our knowledge of these elusive proteins in Mycobacterium, we have improved upon our previous method of OMP prediction in mycobacteria by taking advantage of genomic data from seven mycobacteria species. Our improved algorithm suggests 4333 sequences as putative OMPs in these seven species with varying degrees of confidence. The most virulent pathogenic mycobacterial species are slightly enriched in these selected sequences. We present examples of predicted OMPs involved in horizontal transfer and paralogy expansion. Analysis of local secondary structure content allowed identifying small domains predicted to perform as OMPs; some examples show their involvement in events of tandem duplication and domain rearrangements. We discuss the taxonomic distribution of these discovered families and architectures, often specific to mycobacteria or the wider taxonomic class of Actinobacteria. Our results suggest that OMP functionality in mycobacteria is richer than expected and provide a resource to guide future research of these understudied proteins.

## Introduction

Mycobacteria are responsible for some of the most terrible human diseases including leprosy and tuberculosis (Cosma *et al.*, 2003). However, not all mycobacteria are

49   pathogenic to humans despite their considerable genomic similarity. Part of their

50   variable properties in infective ability and specificity are likely related to their

51   variable cell wall (Brennan & Nikaido, 1995). Outer membrane proteins (OMPs) are

52   an important component of the mycobacterial cell wall, yet they are poorly studied in

53   mycobacteria (Niederweis *et al.*, 2010). OMPs are transmembrane proteins that form

54   a beta-barrel structure consisting of amphipathic beta strands and are secreted into the

55   periplasmic space and inserted into the outer membrane to act as channels (Faller *et*

56   *al.*, 2004). This type of protein is therefore an important target for antibacterial

57   therapy and an object of study to the elucidation of the mechanisms of pathogenicity

58   (Niederweis, 2008). However, currently there is evidence of just some mycobacterial

59   OMPs, in large degree for MspA (Stahl *et al.*, 2001), and less for another two Mt

60   proteins: Rv1973 (Song *et al.*, 2008) and Rv1698 (Siroy *et al.*, 2008; Song *et al.*,

61   2008). This is not only due to the difficulties of culturing mycobacteria, but also to the

62   difficulty of identifying these proteins.

63

64   Computational methods of OMP detection have been developed and applied to

65   *Mycobacterium tuberculosis* with relative success (Pajon *et al.*, 2006; Song *et al.*,

66   2008), but there is room for improvement, especially when it comes to prioritizing

67   targets for research. The increasing number of related mycobacterial genomes offers a

68   unique opportunity to support the predictions by addressing their coherence across

69   orthologs and to pinpoint OMP families specific to pathogenic mycobacteria.

70

71   In this work we accomplished parameter optimization of a previous method that

72   predicts OMPs based solely on their potential to be secreted and to form an

73   amphiphilic beta-barrel (Song *et al.*, 2008). We explored the predictive power of a set

74      of OMP-related properties by contrasting the robustness of the results on the complete

75      proteomes for seven mycobacteria: three obligate pathogens *M. tuberculosis*, *M. bovis*

76      and *M. leprae*, two facultative pathogens *M. marinum* and *M. ulcerans*, one

77      opportunistic pathogen *M. avium*, and the non-pathogenic *M. smegmatis*.

78

79      The relatedness between these species and their pathogenic properties are

80      heterogeneous. The closest genomes by far are those of *M. tuberculosis* and *M. bovis*,

81      but their host ranges are different (*M. bovis* can cause tuberculosis in several

82      mammals, whereas the natural hosts of *M. tuberculosis* are humans). Therefore, both

83      genomes were included in the analysis since we considered that a comparison of

84      OMPs between these two species can lead to insights that help to explain the

85      differences in ability to infect different hosts.

86

87      The predictions on these seven genomes directed us to a number of OMP domains

88      present in mycobacteria, most of them exclusive to actinobacteria and without

89      homologs in eukaryotes.


90      # Methods

91      **Calculation of parameters for outer membrane protein prediction in**
92      **mycobacteria**
93      Protein sequences were obtained for seven mycobacterial genomes (Table 1)

94      [ftp://ftp.ncbi.nlm.nih.gov/genomes/Bacteria/], including: *Mycobacterium avium* 104,

95      *Mycobacterium bovis* AF2122/97, *Mycobacterium leprae* TN, *Mycobacterium*

96      *marinum* M, *Mycobacterium smegmatis* str. MC2 155, *Mycobacterium tuberculosis*

97      H37Rv, and *Mycobacterium ulcerans* Agy99. All proteins were predicted to be OMPs

98      based on the two main properties: 1) the ability to be secreted to the outer membrane;

99      and 2) the ability to form beta-barrel structures.

100

101   Secreted proteins were predicted according to both classical secretion mechanisms

102   (SignalP-v3.0) and non-classical secretion such as twin arginine translocation (TatP-

103   v1.0) or leaderless secretion (SecretomeP-v1.0) (Bendtsen *et al.*, 2004a; Bendtsen *et*

104   *al.*, 2004b; Bendtsen *et al.*, 2005). Prediction of classically secreted proteins is well-

105   studied and most secreted bacterial proteins are exported in this manner (Malen *et al.*,

106   2007); therefore predictions determined by SignalP-v3.0 algorithm or the presence of

107   a single predicted transmembrane alpha helix in the first 70 aa of protein (TMHMM;

108   (Krogh *et al.*, 2001)) were considered to be superior to other prediction methods. TatP

109   prediction is specific to bacteria and was considered to be next most reliable.

110   SecretomeP prediction was the least specific of all methods, but was still considered

111   to be useful as we are aiming for high recall.

112   To demonstrate the efficacy of the secretion prediction methods, the algorithms were

113   run on positive and negative validation sets. A set of 53 experimentally verified,

114   classically secreted mycobacterial proteins was obtained from Leversen and co-

115   workers (Leversen *et al.*, 2009). A set of non-classically secreted mycobacterial

116   proteins, including proteins secreted by Tat or SecA2 systems, was assembled by

117   literature search. Negative controls for secretion were represented by 1725 reviewed

118   cytoplasmic proteins from Mycobacterium sp. (UniProtKB release 15.3).

119

120   In the second step of the OMP prediction, the propensity of the proteins to form beta-

121   barrel structures was determined using various beta-strand properties. Secondary

122   structure for all proteins was predicted using Jnet v1.0 (Cuff & Barton, 1999; Cuff &

123   Barton, 2000). Beta-strands of 5 or more residues (B5 strands) were evaluated for

124   amphiphilicity (FracB5) as previously described (Song *et al.*, 2008). As additional

125 measures of 'betaness', the overall proportion of residues in beta strands

126 (PercentBeta), the number of B5 strands (numB5) and the total number of residues in

127 B5 strands (numB5res) were recorded. Positive and negative control sets for beta-

128 barrels were used to demonstrate the utility of these parameters to predict beta-barrel

129 structure. Positive validation of beta-barrel structure was represented by 428 bacterial

130 and eukaryotic sequences taken from proteins containing Pfam motifs or bacterial

131 sequences with solved 3D structures annotated as forming a beta-barrel. Negative

132 validation of beta sheet prediction consisted of 90 actinobacterial proteins with solved

133 3D structure of low beta content. Protein structures were obtained from the Protein

134 Data Bank (PDB, www.pdb.org).

135

136 Additional parameters, including the number of cysteine residues (numcys) and the

137 isoelectric point (pI) (BioPerl pICalculator, www.bioperl.org) were also evaluated.

138 Initial OMP prediction (Method 1) was carried out with similar parameters and

139 thresholds as previously used by Song et al. (Song *et al.*, 2008), namely: FracB5 $\geq$

140 0.19, PercentBeta $\geq$ 0.10, Smean (from the signal peptide predictor SignalP) $\geq$ 0.50 or

141 numpredhel (number of predicted transmembrane helices) = 1 and firsthel (position in

142 amino acids of the most N-terminal predicted transmembrane helix) $\leq$ 70

143 (Supplementary Table S1). Here, the OMP prediction method was further refined by

144 using sequence homology information and optimization of the algorithm as described

145 below.

146 **Clustering of Mycobacterial sequences**
147 Clustering of the seven mycobacterium genomes listed in Table 1 was carried out

148 using a very strict sequence similarity criterion that enforces all proteins in a cluster to

149 be homologous along their full lengths, ensuring their domain content and

150    architectures are equivalent (Perez-Iratxeta *et al.*, 2007). About 1% of the sequences

151    were not used because they were too short (< 50 aa). The remaining sequences

152    (30,605) were distributed into 11,633 clusters. Pfam motif information (Finn *et al.*,

153    2008) of representative members of the clusters was retrieved to further characterize

154    them.


155    **Optimization of OMP prediction**
156    To optimize the OMP prediction, two training sets of clusters (with at least five

157    sequences each) were defined: 1) OMP-rich clusters which contained ≥ 80% OMP-

158    predicted sequences (1151 sequences in 168 clusters); and 2) OMP-poor clusters,

159    which contained ≤ 20% OMP-predicted sequences (981 sequences in 96 clusters).

160    First, optimization was carried out by applying varying thresholds for OMP prediction

161    on parameters that were not included in Method 1, namely pI, sequence length,

162    number of cysteines, numB5, numB5res, Tat-secretion score, and leaderless secretion

163    score. Next, the Method 1 thresholds on parameters including frac, PerBeta, and

164    Smean were optimized. The optimal cutoff values were defined to be those that most

165    reduced the fraction of predicted OMPs in OMP-poor clusters while retaining over

166    90% of the predicted OMPs in OMP-rich clusters. The new set of optimized criteria

167    was called Method 2 and used for the remainder of the analysis.

168

169    A scoring system (producing a score ranging from zero to 16) was used to indicate the

170    confidence of the OMP prediction. Propensity for secretion and beta-barrel formation

171    were given equal weight (maximum of 8 points each). Proteins predicted to be

172    secreted by general secretion or Tat mechanisms were awarded 8 points. In the

173    absence of these two predictions, proteins were awarded 3 points if the SecretomeP

174    score ≥ 0.574. Beta-sheet related parameters were assessed on: 1) the entire protein

175  length; and 2) within a 300 aa sliding window of protein sequence to detect local

176  regions of high beta content. A window size of 300 aa was chosen because this is the

177  size of the beta barrel domain in some known OMP structures (Faller *et al.*, 2004;

178  Song *et al.*, 1996); however, such domains can be formed by association of beta

179  strands from multiple protein monomers, and it can happen that an OMP protein or

180  region is much smaller than 300 aa. One point was awarded for each of the following

181  criteria satisfied as frac $\geq$ 0.28, PerBeta $\geq$ 0.11, numB5 $\geq$ 3, numB5res $\geq$ 17, for a

182  maximum of 8 points.

## 183  Results

### 184  Optimization of OMP prediction

185  To optimize our predictions in seven mycobacterial genomes, we clustered their

186  protein sequences and evaluated the fraction of proteins in each cluster predicted as

187  OMP (see Methods). Contrasting the parameters of OMP-rich and OMP-poor clusters

188  enabled us to further refine the OMP criteria, based on the assumption that all proteins

189  in a cluster should be predicted either as being OMPs or not; and therefore, that OMP-

190  poor clusters indirectly reflected false positives.

191

192  In the first instance, we tested nine parameters that could play a minor role in the

193  initial prediction of OMPs but were not previously used. The fraction of OMP

194  predicted sequences in OMP-rich and OMP-poor groups was monitored when

195  applying increasingly restrictive thresholds in each of these nine parameters

196  (Supplementary Fig. S1). Four parameters (pI, number of cysteines, first helix, dvalue

197  score from TatP) showed little improvement in reducing the number of potential false

198  negatives in the OMP-poor group of clusters. Potential improvements could be made

199  for the remaining five parameters: sequence length, number of residues in B5 sheets,

200 number of B5 sheets, number of predicted TM helices, and nnscore (prediction of

201 leaderless score from SecretomeP).

202

203 Sequence length was rejected as a constraining factor, as it was undesirable to

204 eliminate short (~100 aa) predicted OMPs potentially composing homologous

205 multimeric structures. Requiring the number of B5 sheets and the number of residues

206 in B5 sheets to be a minimum of 3 and 17, respectively, was successful in reducing

207 the number of OMPs in the OMP-poor clusters by >10%, while keeping 96% of the

208 OMPs in the OMP-rich clusters (Supplementary Fig. S1). A rather stringent threshold

209 was used for prediction of leaderless secretion. At nnscore >= 0.71, 56% of sequences

210 from the OMP-rich clusters were retained, while rejecting 74% of sequences from the

211 OMP-poor clusters as targets for leaderless secretion. This was not considered to be

212 overly stringent, since proteins with signal sequences, (which would be identified by

213 SignalP-3.0 or TatP- 1.0) were likely to have a high nnscore anyway (Bendtsen *et al.*,

214 2004a).

215

216 In the second instance, the five parameters initially used to predict OMPs were varied

217 and the fraction of predicted OMPs was recorded. This analysis suggested fine

218 adjustments in the cutoff values for amphiphilicity (frac), proportion of residues in

219 beta strands (perbeta), and general secretion score (Smean) (Supplementary Fig. S2).

220 The optimized criteria were applied to the dataset, expanding the number of predicted

221 OMPs, compared to the original method (Tables 1 and Supplementary Table S1). Up

222 to this point, the optimization had been carried out by varying parameters on an

223 individual basis. A scoring system was implemented to summarize the effect of all the

224 optimized parameters (with 8 points for beta-barrel formation and 8 points for

225  secretion, see Methods). Assuming that most of the sequences in OMP-rich clusters

226  should actually be OMPs, and that sequences in OMP-poor clusters should not be

227  OMPs, a threshold of OMP score = 12 to accept an OMP prediction was found to be

228  optimal (Fig. 1; Supplementary Fig. S3(a)). At this threshold, 94% of sequences from

229  OMP-rich clusters are classified as OMPs, while 89% of the sequences in the OMP-

230  poor clusters are rejected as OMPs.


231  **Validation of signal sequence and beta-barrel prediction in mycobacteria**

232  None of the secretion prediction programs were specifically designed to predict signal

233  sequences in *Mycobacterium* spp., although the SignalP neural net predictions were

234  based on Gram-positive bacteria, and the TatP server was trained on both Gram-

235  negative and Gram-positive sequences. Mycobacteria are classified as Gram-positive,

236  despite the fact that the mycobacterial outer membrane has distinct properties not

237  found in either functionally classified Gram-negative nor Gram-positive bacteria (Hett

238  & Rubin, 2008). Therefore, it was important to test these algorithms for their ability to

239  detect signal sequences in mycobacteria.

240

241  Using known cytoplasmic and known mycobacterial proteins secreted by the general

242  secretory pathway, it could be shown that the optimized cutoff (Smean = 0.54) was

243  sufficient to correctly predict secretion in 93% of the known GSP proteins and reject

244  98% of the cytoplasmic proteins (Supplementary Fig. S3(b)). Secretion by the non-

245  classical Tat system could be predicted at a cutoff of dvalue = 0.36 in 79% of the

246  known Tat-secreted mycobacterial sequences, while rejecting 98% of the cytoplasmic

247  proteins for secretion (Supplementary Fig. S3(c)).

248

249    Prediction of leaderless secretion in mycobacteria at the chosen cutoff nnscore ≥ 0.71

250    correctly identified 50% (6/12) known leaderless secreted proteins (Supplementary

251    Fig. S3(d)), which included secreted proteins by the recently described bacterial

252    export system ESX-1 and the SecA2 (Sec-independent) system. 81% of the known

253    cytoplasmic proteins were predicted as being secreted in this instance, making the

254    leaderless secretion prediction the least precise of all three secretion prediction

255    methods. As a result, less emphasis was placed on leaderless secretion scores in the

256    OMP prediction.

257    Prediction of beta-barrel structures was based on beta-sheet content and the

258    amphiphilicity of predicted beta strands (computed as in (Song *et al.*, 2008)). As a

259    measure of the protein's propensity to form beta-barrel structures, beta-barrel scores

260    were calculated globally (over whole protein) and locally (sliding window) for a

261    maximum score of 8 (see Methods). For a beta-barrel score ≥ 6, a total of 97% of

262    known bacterial OMP and 90% of annotated beta-barrel proteins were correctly

263    identified as containing beta-barrels, whereas non beta-barrel structures were

264    predicted in 74% of solved sequences lacking certain beta-barrel structure (Fig. 2).

265

266    After the optimization stage, Method 2 was able to correctly identify 90% (27/30) of

267    known bacterial OMPs with high scores (score ≥14; Supplementary Fig. S4)

268    corresponding to strong predictions. Among the three OMPs missed by Method 2,

269    there were two from *Rhodobacter* spp. Although one of them (PORI_RHOBL)

270    contained sufficient beta structure for a beta-barrel, it was predicted to be secreted by

271    leaderless secretion resulting in a weak OMP prediction (OMP score = 11). OmpG

272    from *E. coli* was as well not identified as OMP by this method, due to a lack of beta-

273    strand prediction from Jnet v1.0.

274

275 The selection criteria of both the optimized Method 2 and the previous method show

276 substantial differences between the results (Supplemental Table S1). When applying

277 them to seven mycobacterial genomes (see Methods) 3340 proteins are predicted to be

278 OMPs by both methods. A total of 993 proteins are newly predicted by Method 2

279 whereas 406 proteins predicted to be OMPs by Method 1, are now rejected as false

280 positives.

281 **Identification of OMPs**
282 Table 1 specifies the number of sequences and predicted OMPs. The seven genomes

283 analysed have a genomic size in the 4000-5000 range except for the small Ml genome

284 (1605 genes, an extreme case of genome downsizing (Cole *et al.*, 2001)) and the

285 larger Ms (6,716 genes). When considering the percentage of predicted OMPs it is

286 interesting to note that the three obligate pathogenic organisms (Mt, Mb, and Ml)

287 have the largest percentage (15.1 -15.8%) while the opportunistic pathogen Ma and

288 the non-pathogenic species Ms have the lowest values (12.5 -12.6%).

289

290 We showcase the results of our method with some examples of newly OMP-predicted

291 mycobacterial proteins. Because the taxonomic distribution of a protein can give an

292 indication of its functional relevance, we have categorized the examples by this

293 property. Our selection of examples was facilitated by the clustering used for the

294 optimization of the method, e.g. when searching for OMPs present in all seven

295 mycobacteria. The complete results are available in Supplementary Table S2.

296 **OMPs present in Mt but not in Mb**
297 Though Mt and Mb are very closely related (they both belong to the Mycobacterium

298 tuberculosis complex – Mt complex) and share many 100% identical proteins, they

299  have obvious differences regarding pathogenicity. It is therefore interesting to find

300  OMPs in Mt that have no equivalent in Mb. Two outstanding examples are Mt

301  proteins mce3c (Rv1968; 410 aa) and mce3e (Rv1970; 377 aa) encoded by two of six

302  genes from the putative mce3 operon Mce3A-F. Both Mce3C and Mce3E proteins

303  were found to react with antibodies from serum of TB patients (Ahmad *et al.*, 2004)

304  and have one predicted MCE domain each (at positions 38-114 and 36-112,

305  respectively).

306

307  The presence of the MCE domain in these sequences is relevant because many of the

308  genes with this domain have been shown to be expressed during natural infection of

309  Mt and it is thought that they are related to mycobacterial pathogenicity (Ahmad *et*

310  *al.*, 1999). Mt has a total of 24 genes with this domain arranged in four mce operons,

311  which contain two integral membrane proteins followed by six genes with the MCE

312  domain (Cole *et al.*, 1998). In our previous work (Song *et al.*, 2008) we predicted that

313  23/24 of Mt MCE genes were OMPs. The present method predicts all 24 MCE genes

314  in Mt H37Rv as OMPs with a score of 16 (the maximum possible). We speculate that

315  the MCE domain is actually a beta-barrel characteristic of OMPs, which is coherent

316  with their proposed role at the mycobacterial cell surface (Flesselles *et al.*, 1999).

317  **Present in the Mt complex but not in all mycobacteria**
318  We observed many OMP clusters with members in Mt/Mb but missing in all or some

319  of the five species outside the Mt complex. This can be either due to genes being

320  invented (or horizontally transferred) within the mycobacteria lineage, or to selective

321  gene loss (as in the massive pseudogenization that occurred in the Ml genome (Cole *et*

322  *al.*, 2001)). Here we show examples of each of those.

323 **An OMP unique to Mt/Mb**
324 Rv1351 is a Mt 109 aa protein that we predict to be an OMP. The Mb ortholog is

325 100% identical, and there are no homologous sequences in the other five

326 mycobacteria analysed in this work, or outside mycobacteria (no PSIBLAST hits with

327 E-value below 8.3). We also predict that the gene next to it, Rv1352 (encoding a 123

328 aa protein), is also a small OMP. According to the STRING database (Jensen *et al.*,

329 2009), these two predicted OMPs are in an operon conserved between Mt/Mb, which

330 includes upstream genes Rv1348 and Rv1349 (two uncharacterized ABC transporter

331 ATP-binding proteins) and fabG2/Rv1350 (predicted as 3-ketoacyl-(acyl-carrier-

332 protein) reductase). Therefore, these two predicted OMPs seem to form part of an

333 Mt/Mb specific operon and although they are rather small they could form a barrel by

334 multimerization, which would explain the need to co-express them in an operon. Such

335 operon could carry out a function inherent to the Mt complex. The three genes,

336 Rv1348, Rv1349 and Rv1350, are essential genes for growth of Mt as determined by

337 Sassetti et al. (Sassetti *et al.*, 2003).


338 **A mycobacterial OMP with horizontal transfer**
339 Mt Rv1914c (135 aa) is predicted as an OMP with orthologs in Mb/Mu/Mm but

340 without apparent equivalents in Ml/Ma/Ms. Curiously, the only match in the database

341 outside mycobacteria is a very clear hit (>50% identity) on a distant bacteria,

342 Proteobacteria *Geobacter uraniireducens* (sequence GI:148265072, 135 aa). This

343 suggests an event of horizontal transfer of this gene between mycobacteria and

344 geobacteria. One could speculate that the function of this OMP would not be

345 associated to pathogenicity given its presence both in pathogenic and non-pathogenic

346 mycobacteria (and in *G. uraniireducens*). Incidentally, Rv1914c was one of 224 genes

347 found to be deleted in one or more clinical isolates of a H37Rv strain from San

348 Francisco (Tsolaki *et al.*, 2004).

349 **C4: a novel putative OMP domain that occurs as a tandem repeat.**

350 Mt Rv2270 (175 aa) defines a family with orthologs in five of seven mycobacteria

351 tested (missing in Ms/Ml) and corynebacteria. This implies that the gene was invented

352 in Corynebacterineae and that there was a selective loss of this gene within some

353 members of the mycobacteria lineage indicating that it is not essential for them.

354

355 Sequence analysis indicated that the family contains a C-terminal 120 aa domain (that

356 we termed C4 for its conserved four cysteines, see Supplementary Fig. S5), which is

357 present in other two protein families, one where the domain is tandemly repeated

358 (with orthologs in all seven mycobacteria, e.g. Mt Rv3835), and another where it is

359 combined with an N-terminal Ser/Thr Kinase C domain (present exclusively in a

360 series of Actinomycetales species, e.g. *Stackebrandtia nassauensis* GI:229864975,

361 577 aa; see Supplementary Fig. S5).

362 The prediction of Mt Rv2270 as containing a lipoprotein anchor signal may invalidate

363 the OMP function, but the predicted C4 domain has high beta content and high

364 amphiphilicity; its involvement in variable domain architectures suggests that it can

365 be used as a biological module.

366 **Present in all seven mycobacterial genomes**

367 We found a total of 588 clusters with sequences from all seven mycobacteria, and 61

368 of these were predicted as OMP families. These families are likely to represent

369 proteins important for all mycobacteria but possibly not for pathogenicity since they

370 are present both in pathogenic and in non-pathogenic organisms. We present two

371 interesting cases below.

372 **ACT: an actinobacteria OMP domain greatly expanded in Corynebacterineae**

373 Rv0431 is an Mt predicted OMP with orthologs in all seven mycobacteria. Sequence

374 analysis indicated that the family contains a C-terminal domain of about 100 aa (that

375     we name ACT for the names given to the proteins where it is present: Alanine rich,

376     CpsA, Tuberculin related) present in five Mt sequences that define five families (see

377     Supplementary Fig. S6). In three of the five families the ACT domain is preceded by

378     a predicted domain of around 170 aa of unknown function (PFAM LytR_cpsA_psr).

379

380     The ACT domain is present in some genera outside but close to mycobacteria, chiefly

381     Nocardiodes and Corynebacterium, but not all species have the five sequences and

382     Ma has an extra copy of one of the five. These results suggest that the ACT domain

383     was invented before the divergence of mycobacteria, corynebacteria and nocardiodes.

384     Its high level of duplication and a number of gene losses and duplications in

385     mycobacteria suggest that it confers some kind of low-specific functional advantage.

386     **An OMP essential for Mycobacteria growth**
387     Rv0227c is another predicted OMP in a cluster with orthologs in all seven

388     mycobacterial genomes. The proteins in this cluster have no known function, and

389     closer analysis by PSI-BLAST revealed that there are distant homologs in nocardia

390     and corynebacteria. The protein itself is characterized by a by a signal peptide with a

391     predicted cleavage point before the first TM helix and a 300 aa beta-domain

392     surrounded by two TMs. Mutagenesis and comparative genomic analyses have

393     identified Rv0227c as being a 'core' mycobacterial gene, required for optimal growth

394     (Marmiesse *et al.*, 2004; Sassetti *et al.*, 2003).

395     **Example OMPs identified by new criteria**
396     Method 2 includes two predictive features that have not been used before: export

397     signals other than those reported by SignalP and a window analysis of secondary

398     structure. Those allowed the identification of many extra OMPs respect to our

399    previous work. Here we present two examples of OMPs detected by each of these

400    new criteria.

**An actinobacteria-specific protein with low global beta content**

401

402    One of our clusters represents a family with members in five of the mycobacteria

403    tested, four of which have OMP scores of 15 (Rv2345, MAV2041, Mb2374,

404    MMAR_3652; ~660 aa) and one with an OMP score of 13 (MSMEG_4484). Notably

405    absent are sequence homologs in Ml and Mu (confirmed using PSI-BLAST under

406    default parameters), but we found orthologs of this protein in a wide range of

407    Actinobacteria. These proteins contain a predicted Pfam domain of unknown function

408    (DUF477), followed by a predicted TM, and a very variable glycine-rich region at the

409    end. The percentage of beta structure of the whole sequence is well below the

410    threshold of 0.11 that we use for selection. However, the window analysis shows that

411    the DUF447 domain has a high percentage of beta content and high amphiphilicity,

412    potentially characterizing an OMP function (Fig. 3).

413

414    Similarity searches uncovered a much shorter second homolog in Ma (MAV_2102),

415    also present in *M. intracellulare*, which keeps the N-terminal domain, the predicted

416    TM following it, and a C-terminal domain, but lacks the middle region and the

417    Glycine-rich region (Fig. 3). We predict that both the long and the short families are

418    OMPs.

**Mycobacteria-specific OMPs secreted by the Tat system**

419

420    An example found using the predicted Tat-system secretion that would not have been

421    detected using SignalP was Rv2577 from Mt. This is an OMP predicted protein with

422    orthologs in Mm and Ma (all of them with maximum OMP score = 16), apparently

423    absent from Mu and Ml. The C-terminal end contains a predicted

424  metallophosphoesterase domain (similarity to COG1409 Predicted phosphohydrolases

425  according to database annotations) with clear homologs to other species outside

426  actinobacteria.

427

428  In Mb AF2122/97, the syntenic gene of Rv2577 (529 aa) is separated into two open

429  reading frames (Mb2607 and Mb2608) due to a base transversion (G-A), which

430  introduces a stop codon. Mb2607 (83 aa) contains the signal sequence and sufficient

431  beta strand structure for an OMP prediction of perfect score. Mb2608 (434 aa)

432  matches Rv2577 from position 96 on, so that just 12 amino acids of the Mt protein are

433  not represented in any of the two Mb proteins. The complete predicted

434  phosphoesterase domain is intact. The N-terminal region has high content of potential

435  amphiphilic beta-strand but this extends further to the region of homology to Mb2608.

436  In the absence of sequences with homology to Mb2607 but not to Mb2608, we cannot

437  support that Mb2607 can form an independent domain, although the gene split

438  suggests this possibility.

439

440  Both Mb2607/Mb2608 transcripts have been shown to be up-regulated in a virulent

441  strain of *M. bovis* during bacterial replication in macrophages (Blanco *et al.*, 2009).

442  The G-A transversion is absent in avirulent *M. bovis* strains used for human vaccine

443  development (*Mycobacterium bovis* BCG str. Tokyo 172, *Mycobacterium bovis* BCG

444  Pasteur 1173P2). The splitting of this gene may extend host-specific modular

445  functions of this protein in *M. bovis* AF2122/97, which is pathogenic to cattle

446  (Garnier *et al.*, 2003).

## Discussion

Outer membrane proteins (OMPs) act as gatekeepers to the external environment.
They are exposed as quorum sensors or acting in response to its environment, and are
likely to be essential for general survival of the cell. In pathogenic species, the
function of the OMPs may play important roles in host-cell interactions that enable
the persistence of mycobacterial infection. As such, OMPs are logical drug targets -
not only in tuberculosis and leprosy, but also in opportunistic infections in
immunocompromised patients, which in total kill millions of people world-wide every
year and are complicated by new problems like co-infection with HIV and resistance
to drugs (Meya & McAdam, 2007).

Existing computational methods to predict beta-barrel outer membrane proteins
primarily focus on known OMPs in Gram-negative bacteria (Berven *et al.*, 2004;
Bigelow *et al.*, 2004; Casadio *et al.*, 2003; Remmert *et al.*, 2009; Zhai & Saier, 2002).
Unfortunately, OMPs of the beta-barrel type are almost totally uncharacterized in
mycobacteria. Models for Gram-positive and Gram-negative OMPs can only partially
extend to mycobacteria due to its unique cell wall construction that eludes clear
functional classification as Gram-positive or Gram-negative. For example, the beta-
barrels of the OMPs of mycobacteria have to be longer than those typically known
from Gram-negative bacteria, in agreement with the greater thickness of the
mycobacterial wall (Alahari *et al.*, 2007; Hoffmann *et al.*, 2008; Zuber *et al.*, 2008);
this is the case for MspA (Faller *et al.*, 2004). As a result, the study of mycobacterial
OMPs has to rely on tools specific to it.

Computational methods can be used to predict OMPs based on two properties: OMPs
must form an amphipathic beta-barrel and be secreted. However, in view of the

473      current small number of mycobacterial OMPs with which to benchmark such

474      methods, and aware of the fact that mycobacterial OMPs are expected to be very

475      different to those known outside Actinobacteria, we resourced to benchmark our

476      method according to its ability to produce coherent predictions across proteins with

477      high similarity to Mt proteins.

478

479      Accordingly, the method we presented here uses clusters of homologous sequences

480      from seven mycobacterial genomes to optimize OMP prediction, based on the

481      assumption that cross-genomic sequences within the same cluster should share similar

482      properties, and therefore if the majority of sequences within a cluster were predicted

483      to be OMPs, then those sequences that escaped prediction were also likely to be OMP

484      sequences. In this manner, we were able to examine the effect of changing the

485      thresholds of different parameters on the number of predicted OMPs to set final

486      thresholds that reduced the number of spurious OMP predictions in clusters with low

487      OMP content while maintaining OMP predictions for clusters with initially high OMP

488      content. Moreover, we performed a sliding window analysis on all proteins to identify

489      local regions of beta-content within larger proteins with low overall propensity to

490      form beta-barrels. This method predicts practically all known mycobacterial OMPs

491      with close to maximum scores.

492

493      We computed a set of 4300 potential OMPs in seven genomes (+600 alone in Mt). It

494      is unlikely that all of them will be OMPs as current estimations of OMPs in Mt are in

495      the order of 100s (Niederweis *et al.*, 2010). We do not think that with the current

496      information on mycobacterial OMPs we can devise a more sensitive scoring system.

497      In any case we note that this dataset includes a higher proportion of sequences from

498     obligate pathogenic mycobacteria compared to opportunistic or non-pathogenic

499     mycobacteria (15% versus 13%) suggesting that the set is enriched in genes with a

500     function related to pathogenicity. Many of these proteins, as we have shown, define

501     families specific to Mycobacteria or Actinobacteria that remain yet to be functionally

502     characterized.

503

504     Our work proposes a number of putative OMP domains. Some of them are reused in

505     multiple domain architectures and duplicated in paralogs (e.g. the ACT domain) or

506     inside genes (e.g. the tandemly repeated C4 domain in Rv3835). Some of these

507     domains or even some of the entire OMP predicted proteins are probably too small to

508     form a beta-barrel by themselves (<150 aa). However, the many cases where such

509     proteins appear together in putative operons (e.g. the mce operons) suggest that they

510     may associate to form a barrel. OMP formation by oligomerization is already

511     suspected in the predicted OMP Rv1698. Rv1698 has been observed to dimerize and

512     the observation that channel complexes containing Rv1698 have variable conductance

513     states suggest that Rv1698 might form oligomers (Siroy *et al.*, 2008). The formation

514     of self-associations is also a possibility that has been reported. For example, both Mt

515     MspA and the alpha-hemolysin porin of *Staphylococcus aureus* (from different

516     phylum firmicutes) form a beta-barrel with each monomer contributing just a 50

517     amino acids loop to the beta-barrel associating as homo-octamer (Faller *et al.*, 2004)

518     or homo-heptamer (Song *et al.*, 1996), respectively.

## 519  Conclusions

520 In summary, our results suggest that potential OMPs are a large contributor to the

521 protein baggage of mycobacteria, possibly of Actinobacteria. Should a large fraction

522 of our predictions be demonstrated experimentally to be OMPs, this would point to

523 this function as an important factor for shaping the evolution, variability, and

524 adaptability of these organisms. Using genomic information we have been able to

525 tune an OMP prediction algorithm and produced a set of OMP predictions for more

526 than 4300 mycobacterial proteins. Their profiles of taxonomic conservation can be

527 used to hypothesize the functional importance and pathogenicity relevance.

528

529 We note that while this manuscript was under review, one of our predicted OMPs,

530 Rv0899, has been the focus of an experimental effort to characterize it as an OMP

531 (Teriete et al., 2010). Although the result was negative, this indicates that our method

532 produces targets that align well with those that the researchers in the field choose

533 using their intuition and knowledge. As new experimental evidence accumulates, we

534 will be able to refine our algorithm. In addition, the expected sequencing of novel

535 mycobacterial genomes will allow us to further complete the picture of the

536 evolutionary history of OMPs and to pinpoint their association to pathogenicity,

537 hopefully leading to new strategies to combat a number of terrible diseases.

538

543

## **References**
545 **Ahmad, S., Akbar, P. K., Wiker, H. G., Harboe, M. & Mustafa, A. S. (1999).**
546 Cloning, expression and immunological reactivity of two mammalian cell entry
547 proteins encoded by the mce1 operon of Mycobacterium tuberculosis. *Scand J*
548 *Immunol* **50**, 510-518.

549

550 **Ahmad, S., El-Shazly, S., Mustafa, A. S. & Al-Attiyah, R. (2004).** Mammalian
551 cell-entry proteins encoded by the mce3 operon of Mycobacterium tuberculosis are
552 expressed during natural infection in humans. *Scand J Immunol* **60**, 382-391.
553
554 **Alahari, A., Saint, N., Campagna, S., Molle, V., Molle, G. & Kremer, L. (2007).**
555 The N-terminal domain of OmpATb is required for membrane translocation and pore-
556 forming activity in mycobacteria. *J Bacteriol* **189**, 6351-6358.
557
558 **Bendtsen, J. D., Jensen, L. J., Blom, N., Von Heijne, G. & Brunak, S. (2004a).**
559 Feature-based prediction of non-classical and leaderless protein secretion. *Protein*
560 *Eng Des Sel* **17**, 349-356.
561
562 **Bendtsen, J. D., Nielsen, H., von Heijne, G. & Brunak, S. (2004b).** Improved
563 prediction of signal peptides: SignalP 3.0. *J Mol Biol* **340**, 783-795.
564
565 **Bendtsen, J. D., Nielsen, H., Widdick, D., Palmer, T. & Brunak, S. (2005).**
566 Prediction of twin-arginine signal peptides. *BMC Bioinformatics* **6**, 167.
567
568 **Berven, F. S., Flikka, K., Jensen, H. B. & Eidhammer, I. (2004).** BOMP: a
569 program to predict integral beta-barrel outer membrane proteins encoded within
570 genomes of Gram-negative bacteria. *Nucleic Acids Res* **32**, W394-399.
571
572 **Bigelow, H. R., Petrey, D. S., Liu, J., Przybylski, D. & Rost, B. (2004).** Predicting
573 transmembrane beta-barrels in proteomes. *Nucleic Acids Res* **32**, 2566-2577.
574
575 **Blanco, F. C., Nunez-Garcia, J., Garcia-Pelayo, C. & other authors (2009).**
576 Differential transcriptome profiles of attenuated and hypervirulent strains of
577 Mycobacterium bovis. *Microbes Infect* **11**, 956-963.
578
579 **Brennan, P. J. & Nikaido, H. (1995).** The envelope of mycobacteria. *Annu Rev*
580 *Biochem* **64**, 29-63.
581
582 **Casadio, R., Fariselli, P., Finocchiaro, G. & Martelli, P. L. (2003).** Fishing new
583 proteins in the twilight zone of genomes: the test case of outer membrane proteins in
584 Escherichia coli K12, Escherichia coli O157:H7, and other Gram-negative bacteria.
585 *Protein Sci* **12**, 1158-1168.
586
587 **Cole, S. T., Brosch, R., Parkhill, J. & other authors (1998).** Deciphering the
588 biology of Mycobacterium tuberculosis from the complete genome sequence. *Nature*
589 **393**, 537-544.
590
591 **Cole, S. T., Eiglmeier, K., Parkhill, J. & other authors (2001).** Massive gene decay
592 in the leprosy bacillus. *Nature* **409**, 1007-1011.
593
594 **Cosma, C. L., Sherman, D. R. & Ramakrishnan, L. (2003).** The secret lives of the
595 pathogenic mycobacteria. *Annu Rev Microbiol* **57**, 641-676.
596
597 **Cuff, J. A. & Barton, G. J. (1999).** Evaluation and improvement of multiple
598 sequence methods for protein secondary structure prediction. *Proteins* **34**, 508-519.
599

600 **Cuff, J. A. & Barton, G. J. (2000).** Application of multiple sequence alignment
601 profiles to improve protein secondary structure prediction. *Proteins* **40**, 502-511.
602
603 **Faller, M., Niederweis, M. & Schulz, G. E. (2004).** The structure of a mycobacterial
604 outer-membrane channel. *Science* **303**, 1189-1192.
605
606 **Finn, R. D., Tate, J., Mistry, J. & other authors (2008).** The Pfam protein families
607 database. *Nucleic Acids Res* **36**, D281-288.
608
609 **Flesselles, B., Anand, N. N., Remani, J., Loosmore, S. M. & Klein, M. H. (1999).**
610 Disruption of the mycobacterial cell entry gene of Mycobacterium bovis BCG results
611 in a mutant that exhibits a reduced invasiveness for epithelial cells. *FEMS Microbiol*
612 *Lett* **177**, 237-242.
613
614 **Garnier, T., Eiglmeier, K., Camus, J. C. & other authors (2003).** The complete
615 genome sequence of Mycobacterium bovis. *Proc Natl Acad Sci U S A* **100**, 7877-
616 7882.
617
618 **Hett, E. C. & Rubin, E. J. (2008).** Bacterial growth and cell division: a
619 mycobacterial perspective. *Microbiol Mol Biol Rev* **72**, 126-156, table of contents.
620
621 **Hoffmann, C., Leis, A., Niederweis, M., Plitzko, J. M. & Engelhardt, H. (2008).**
622 Disclosure of the mycobacterial outer membrane: cryo-electron tomography and
623 vitreous sections reveal the lipid bilayer structure. Proc Natl Acad Sci U S A 105,
624 3963-3967.
625
626 **Jensen, L. J., Kuhn, M., Stark, M. & other authors (2009).** STRING 8--a global
627 view on proteins and their functional interactions in 630 organisms. *Nucleic Acids Res*
628 **37**, D412-416.
629
630 **Krogh, A., Larsson, B., von Heijne, G. & Sonnhammer, E. L. (2001).** Predicting
631 transmembrane protein topology with a hidden Markov model: application to
632 complete genomes. *J Mol Biol* **305**, 567-580.
633
634 **Leversen, N. A., de Souza, G. A., Malen, H., Prasad, S., Jonassen, I. & Wiker, H.**
635 **G. (2009).** Evaluation of signal peptide prediction algorithms for identification of
636 mycobacterial signal peptides using sequence data from proteomic methods.
637 *Microbiology* **155**, 2375-2383.
638
639 **Malen, H., Berven, F. S., Fladmark, K. E. & Wiker, H. G. (2007).** Comprehensive
640 analysis of exported proteins from Mycobacterium tuberculosis H37Rv. *Proteomics* **7**,
641 1702-1718.
642
643 **Marmiesse, M., Brodin, P., Buchrieser, C., Gutierrez, C., Simoes, N., Vincent, V.,**
644 **Glaser, P., Cole, S. T. & Brosch, R. (2004).** Macro-array and bioinformatic analyses
645 reveal mycobacterial 'core' genes, variation in the ESAT-6 gene family and new
646 phylogenetic markers for the Mycobacterium tuberculosis complex. *Microbiology*
647 **150**, 483-496.
648

649    **Meya, D. B. & McAdam, K. P. (2007).** The TB pandemic: an old problem seeking
650    new solutions. *J Intern Med* **261**, 309-329.

651

652    **Niederweis, M. (2008).** Nutrient acquisition by mycobacteria. *Microbiology* **154**,
653    679-692.

654

655    **Niederweis, M., Danilchanka, O., Huff, J., Hoffmann, C. & Engelhardt, H.**
656    **(2010).** Mycobacterial outer membranes: in search of proteins. *Trends Microbiol*.

657

658    **Pajon, R., Yero, D., Lage, A., Llanes, A. & Borroto, C. J. (2006).** Computational
659    identification of beta-barrel outer-membrane proteins in Mycobacterium tuberculosis
660    predicted proteomes as putative vaccine candidates. *Tuberculosis (Edinb)* **86**, 290-
661    302.

662

663    **Perez-Iratxeta, C., Palidwor, G. & Andrade-Navarro, M. A. (2007).** Towards
664    completion of the Earth's proteome. *EMBO Rep* **8**, 1135-1141.

665

666    **Remmert, M., Linke, D., Lupas, A. N. & Soding, J. (2009).** HHomp--prediction
667    and classification of outer membrane proteins. *Nucleic Acids Res* **37**, W446-451.

668

669    **Sassetti, C. M., Boyd, D. H. & Rubin, E. J. (2003).** Genes required for
670    mycobacterial growth defined by high density mutagenesis. *Mol Microbiol* **48**, 77-84.

671

672    **Senaratne, R. H., Mobasheri, H., Papavinasasundaram, K. G., Jenner, P., Lea, E.**
673    **J. & Draper, P. (1998).** Expression of a gene for a porin-like protein of the OmpA
674    family from Mycobacterium tuberculosis H37Rv. *J Bacteriol* **180**, 3541-3547.

675

676    **Siroy, A., Mailaender, C., Harder, D., Koerber, S., Wolschendorf, F.,**
677    **Danilchanka, O., Wang, Y., Heinz, C. & Niederweis, M. (2008).** Rv1698 of
678    Mycobacterium tuberculosis represents a new class of channel-forming outer
679    membrane proteins. *J Biol Chem* **283**, 17827-17837.

680

681    **Song, H., Sandie, R., Wang, Y., Andrade-Navarro, M. A. & Niederweis, M.**
682    **(2008).** Identification of outer membrane proteins of Mycobacterium tuberculosis.
683    *Tuberculosis (Edinb)* **88**, 526-544.

684

685    **Song, L., Hobaugh, M. R., Shustak, C., Cheley, S., Bayley, H. & Gouaux, J. E.**
686    **(1996).** Structure of staphylococcal alpha-hemolysin, a heptameric transmembrane
687    pore. *Science* **274**, 1859-1866.

688

689    **Stahl, C., Kubetzko, S., Kaps, I., Seeber, S., Engelhardt, H. & Niederweis, M.**
690    **(2001).** MspA provides the main hydrophilic pathway through the cell wall of
691    Mycobacterium smegmatis. *Mol Microbiol* **40**, 451-464.

692

693    **Teriete, P., Yao, Y., Kolodzik, A., Yu, J., Song, H., Niederweis, M. & Marassi, F.**
694    **M. (2010).** Mycobacterium tuberculosis Rv0899 adopts a mixed alpha/beta-structure
695    and does not form a transmembrane beta-barrel. *Biochemistry* **49**, 2768-2777.

696

697 **Tsolaki, A. G., Hirsh, A. E., DeRiemer, K. & other authors (2004).** Functional and
698 evolutionary genomics of Mycobacterium tuberculosis: insights from genomic
699 deletions in 100 strains. *Proc Natl Acad Sci U S A* **101**, 4865-4870.
700
701 **Zhai, Y. & Saier, M. H., Jr. (2002).** The beta-barrel finder (BBF) program, allowing
702 identification of outer membrane beta-barrel proteins encoded within prokaryotic
703 genomes. *Protein Sci* **11**, 2196-2207.
704
705 **Zuber, B., Chami, M., Houssin, C., Dubochet, J., Griffiths, G. & Daffe, M.**
706 **(2008).** Direct visualization of the outer membrane of mycobacteria and
707 corynebacteria in their native state. J Bacteriol 190, 5672-5680.
708

# Tables

709

710 **Table 1  - Predictions for seven genomes.**

| Genome (NCBI Taxon ID) | Habitat | Annotated proteins | Suggested OMP proteins with score >=12 (% of total) | Suggested OMPs unique to a genome* |
|---|---|---|---|---|
| Mycobacterium tuberculosis H37Rv (83332) | obligate pathogen | 3991 | 629 (15.8 %) | 35 |
| Mycobacterium bovis AF2122/97 (233413) | obligate pathogen | 3920 | 617(15.7%) | 39 |
| Mycobacterium leprae TN (272631) | obligate pathogen | 1605 | 242 (15.1%) | 77 |
| Mycobacterium marinum M (216594) | environmental, facultative pathogen | 5462 | 799 (14.6%) | 184 |
| Mycobacterium ulcerans Agy99 (362242) | environmental, facultative pathogen | 4160 | 561 (13.5%) | 111 |
| Mycobacterium smegmatis str. MC2 155 (246196) | environmental, not pathogenic | 6716 | 844 (12.6 %) | 459 |
| Mycobacterium avium 104 (243243) | environmental, facultative opportunistic pathogen | 5120 | 641 (12.5%) | 251 |
| Total | | 30974 | 4333 (14%) | |

711 * numbers of predicted OMPs in within clusters in other genomic patterns: present in
712 all genomes: 189; present only in obligate pathogens (Mt, Mb, Ml): 48; present only
713 in facultative pathogens (Mm, Mu, Ma): 66; all other combinations: 2874
714

715

# Figure legends

**Fig. 1 - OMP scores for sequences in OMP-rich and OMP-poor clusters.**

The minimum score for OMP prediction was set to >= 12 (dotted vertical blue line).

At this threshold, 94% of sequences from OMP-rich clusters (black line) are classified

as OMPs, while 85% of the sequences in the OMP-poor clusters (red line) are rejected

as OMPs.

**Fig. 2 - Beta-barrel prediction scores for control sequences.**

Positive controls for beta-barrels include 428 bacterial or eukaryotic proteins from

Pfam or PDB with annotated beta-barrel structures and solved structure information.

Negative controls include 90 actinobacterial sequence fragments with low beta

content, as determined from solved structures in PDB. At beta-barrel score $\geq 6$, 97%

and 90% of known bacterial OMPs and annotated beta-barrels, respectively, are

predicted to be beta-barrels, and 74% of low beta sheet content sequences are

predicted to be without beta-barrel structure.

**Fig. 3 - Frac and PerBeta in a sliding window for Rv2345.**

Rv2345 defines a family conserved in Actinobacteria and present in the mycobacteria

tested with the exception of Ml and Mu. Top: average on a 300 aa window of

percentage of beta sheet (PercentBeta) and amphiphilicity of beta strands (FracB5) for

Rv2345. The horizontal lines indicate the thresholds used for these two parameters.

The plot suggests that the N-terminal of Rv2345 contains a highly amphiphilic beta

structure. Bottom: the N-terminal end of Rv2345 and orthologs contains a predicted

Pfam domain of unknown function (DUF477). Ma protein MAV_102 represents a

different architecture but is potentially a shorter OMP as it keeps the N-terminal

domain. Other predicted sequence features for these proteins are: transmembrane

740    helices (blue boxes), a 300 aa domain (blue hexagon), a C-terminal domain (yellow

741    oval), and a G-rich amino acid biased region (orange bar).

742

**Supplementary material legends**

744 **Supplementary Fig. S1– Fraction of OMPs remaining with increased restriction**
745 **of non-Method 1 parameters.**
746 This figure shows the fractions of OMPs remaining (y-axis) from the OMP-rich (S1,

747 solid line) and OMP-poor (S2, dotted line) groups of clusters, as the parameter

748 thresholds become increasingly restrictive (x-axis).


749 **Supplementary Fig. S2- Effect of changing original parameters on OMP**
750 **prediction in OMP-rich (S1) and OMP-poor (S2) clusters.**
751 Cutoff criteria for the parameters frac, percent beta, and general secretion score

752 (Smean) were varied and the fraction of predicted OMPs relative to the Method 1

753 prediction, was recorded. Optimal cutoffs (shown in red) eliminated 5-25% of the

754 OMPs in S2 clusters while maintaining at least 94% of the OMPs in S1 clusters.

755


756 **Supplementary Fig. S3– Validation of OMP prediction and signal sequence**
757 **prediction in mycobacteria.**
758 **(a)** Recall-precision curve for predicted OMPs. This figure shows the recall and

759 precision curve, based on the assumptions that OMP-rich clusters (S1) contained true-

760 positives and OMP-poor clusters (S2) contained true negatives. An OMP score of 12

761 was chosen as the threshold for an OMP prediction. **(b)** General secretion scores for

762 mycobacterial proteins. This figure shows the Smean scores, as determined from

763 SignalP-v3.0, for 1723 cytoplasmic and 58 proteins known to be secreted by the

764 general secretion pathway. Cytoplasmic proteins were taken from annotated, reviewed

765 proteins in UniProtKB. Experimentally verified GSP secreted proteins were taken

766 from the literature. Proteins with Smean$\geq$0.54 (vertical dotted blue line) were

767 considered to be secreted. At this cutoff, 93% of known GSP-secreted proteins are

768 correctly predicted to be secreted, while 98% of the cytoplasmic proteins are not

769 predicted to be secreted. **(c)** Twin arginine translocation scores for mycobacterial

770 proteins. In this figure, experimentally verified Mycobacterium proteins secreted by

771    Tat system (19) were found by literature search. At the selected cutoff (TatP

772    dvalue=0.36, vertical dotted blue line), the TatP algorithm correctly predicts 79% of

773    mycobacterial positive validation proteins to be Tat-secreted. 98% of the cytoplasmic

774    proteins were not predicted to be Tat-secreted. **(d)** Leaderless secretion prediction for

775    mycobacterial sequences. In this figure, 50% of known leaderless secreted proteins

776    (12) are correctly predicted to be secreted at nnscore ≥0.71 (SecretomeP-v1.0). This

777    includes 5/6 ESX-1 secreted proteins and 1/5 SecA2 secreted proteins. A single

778    protein (GLNA1_MYCTU), whose secretion mechanism is unknown, was not

779    predicted to be secreted. 81% of known cytoplasmic proteins were correctly predicted

780    to not be secreted.

781

782    **Supplementary Fig. S4– OMP scores for known bacterial OMPs.**
783    In this figure, each horizontal bar summarizes the points awarded to each bacterial

784    OMP (indicated by UniProtKB accession). Points were awarded for prediction of

785    secretion by signal sequence prediction (general or Tat secretion, 8 points; black bar)

786    or leaderless secretion (3 points; red bar). For the beta-barrel structure, one point was

787    awarded for 4 parameters (frac, perbeta, numB5, resB5) over the whole sequence

788    (green bar) or for a sliding window of 300 aa (blue bar), for a maximum of 8 beta-

789    barrel points.

790

791    **Supplementary Fig. S5- C4: a tandemly repeated OMP domain**
792    We found a domain that occurs in Mt proteins Rv2770 and Rv3835, tandemly

793    repeated in the latter. In some Actinomycetales the domain occurs with a Ser/Thr

794    protein kinase catalytic domain at the N-terminal and a predicted TM helix in

795    between.

796

**Supplementary Fig. S6- ACT. A domain duplicated and lost many times.**

798 We identified a novel domain (ACT) as a candidate OMP domain occurring C-

799 terminal in Mt proteins Rv0431, Rv2700, Rv0822c, Rv3267 and Rv3484. In three of

800 them it is combined with an N-terminal extracellular domain of unknown function

801 (LytR) found in a number of putative membrane-bound proteins. Left: phylogenetic

802 tree from an alignment of instances of the domain in the seven mycobacterial species

803 analyzed and in Corynebacterineae species: *Corynebacterium amycolatum* SK46

804 (Ca), *Rhodococcus opacus* B4 (Ro) and *Nocardia farcinica* (Nf). Right: sequence

805 features of the five Mt sequences. Trans-membrane alpha helix (TM, blue) and signal

806 peptide (SP, red) were predicted using TMHMM and SignalP-v3.0, respectively. The

807 TM in Rv0822c was under the default cut-off of TMHMM and was not predicted but

808 the 18 aa region reported displays a maximum of probability of being a TM (with

809 scores above 0.6).

810

811 **Supplementary Table S1– Criteria for OMP prediction.**
812 This table shows the criteria used to predict OMPs in Method 1 and Method 2 (this

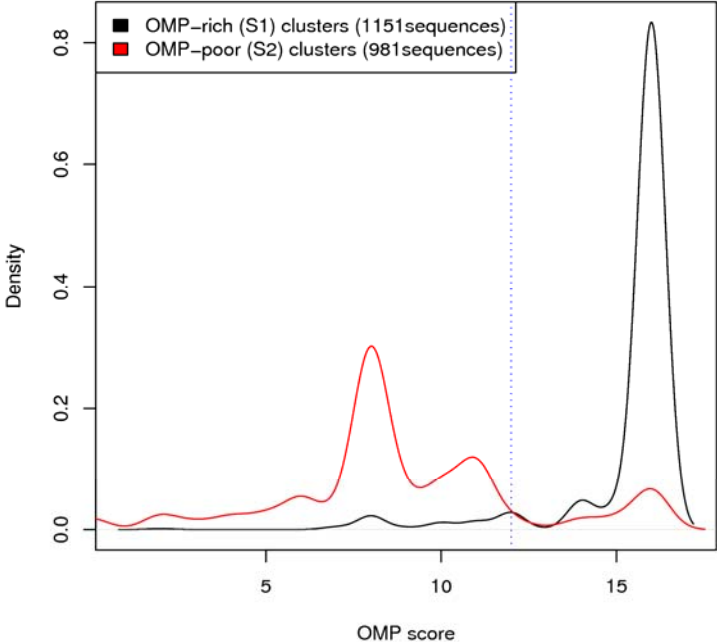813 study). The parameter values for 30 known OMPs are included for comparison.

814

815 **Supplementary Table S2– OMP score for mycobacterial sequences.**
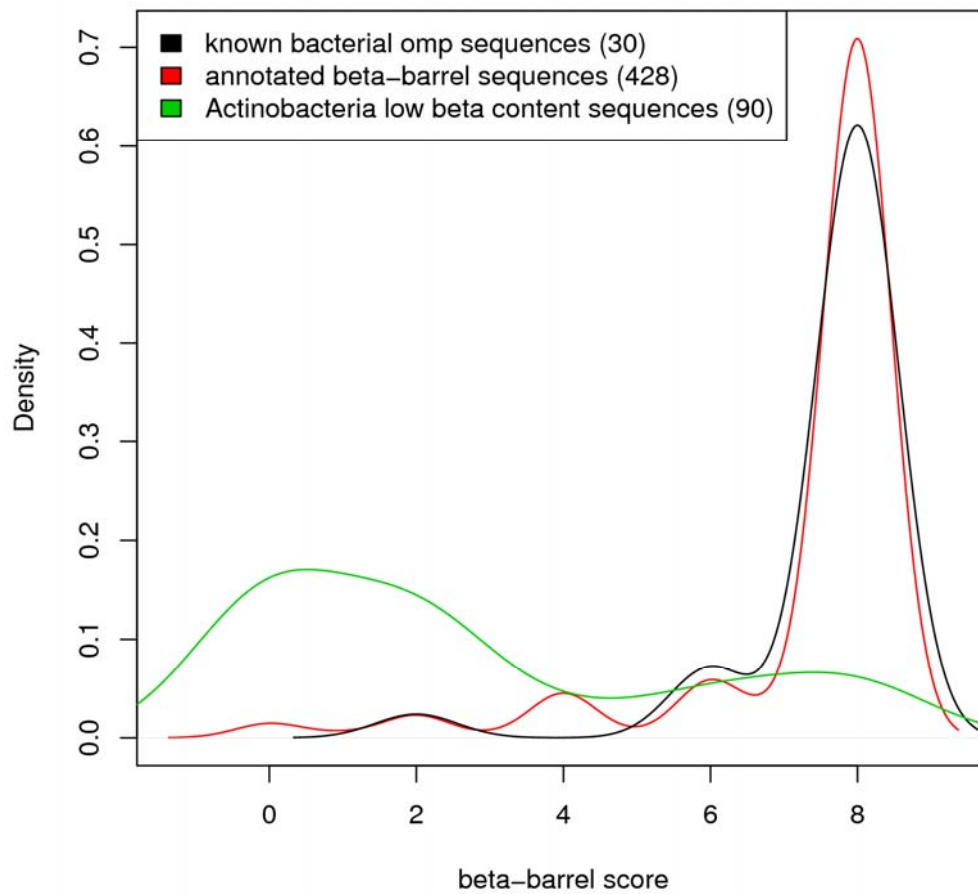816 This file provides the results of the different scores and tests for protein sequences of

817 seven mycobacteria used in this manuscript (30,605 sequences of length 50 amino

818 acids or more). Each row represents a sequence. The columns indicate (1) species, (2)

819 gene identifier, (3) OMP score, (4) Frac, (5) PerBeta, (6) Smean, (7) Dval (from tat),

820 (8) number of predicted beta strands of five residues or longer and (9) number of

821 residues in those, (10) number of predicted transmembrane helices, (11) position of

822 the first transmembrane helix, (12) length of the sequence, (13) computed pI, (14)

823 number of cysteines. Columns 15-20 regard the properties found on a 300 amino acid

824     window whose position was selected as indicated in Methods: (15) window left start

825     position, (16) number of beta strands of length five residues or more, (17) residues on

826     those, (18) percentage of beta structure predicted and (19) Frac inside the 300 amino

827     acid window, and (20) window beta score.

828

**OMP scores for S1 and S2 (Group.size>5)**

**Beta-barrel scores for control sequences**

Legend:
- known bacterial omp sequences (30)
- annotated beta-barrel sequences (428)
- Actinobacteria low beta content sequences (90)

Y-axis: Density

X-axis: beta-barrel score

Measures of perbeta and frac over a 300aa Window (Rv2345; CAB06160.1)