

1 Updated functional annotation of the *Mycobacterium bovis* 2 AF2122/97 reference genome.

3
4 Damien Farrell¹, Joseph Crispell¹, Stephen V. Gordon^{1,2,3,4†}.

5
6 ¹UCD School of Veterinary Medicine, University College Dublin. Ireland

7 ²UCD Conway Institute of Biomolecular and Biomedical Research, University College
8 Dublin, Ireland

9 ³UCD School of Medicine, University College Dublin, Dublin 4, Ireland.

10 ⁴UCD School of Biomolecular and Biomedical Science, University College Dublin, Dublin
11 4, Ireland.

12

13

14 **†Corresponding Author**

15

16 Damien Farrell, E-mail: damien.farrell@ucd.ie

17

18

19 **Keywords:**

20 *Mycobacterium bovis*; annotation; genome; TB

21

22 **Abstract**

23 *Mycobacterium bovis* AF2122/97 is the reference strain for the bovine tuberculosis
24 bacillus. We here report an update to the *M. bovis* AF2122/97 genome annotation to
25 reflect 616 new protein identifications which replace many of the old hypothetical coding
26 sequences and proteins of unknown function in the genome. These changes integrate
27 information from functional assignments of orthologous coding sequences in the
28 *Mycobacterium tuberculosis* H37Rv genome. We have also added 69 additional new gene
29 names.

30 **Background**

31 *Mycobacterium bovis* (*M. bovis*) is a causative agent of bovine tuberculosis (bTB) and the
32 most widely studied animal-adapted member of the *Mycobacterium tuberculosis* complex
33 (MTBC). The genome of the *M. bovis* AF2122/97 strain was first sequenced in 2003
34 [1] and is considered to be the reference sequence for this species. *Mycobacterium bovis*
35 has considerable importance as the basis for comparative studies into animal- and human-
36 adapted species of the MTBC. This genome was revised in 2017 [2] with an updated

37 sequence that included the previously missing RD900 region and added 42 new coding
38 sequences. These revisions brought the annotation in line with updates to the *M.*
39 *tuberculosis* H37Rv genome, to which *M. bovis* shares high identity.

40 **Hypothetical and unknown proteins**

41 Proteins encoded by genes with no clear functional activity have been traditionally
42 annotated with designations such as ‘unknown protein’, ‘hypothetical protein’ or ‘conserved
43 hypothetical’. These labels are usually placed in the `/product` field of the annotation file (in
44 this context we refer to the genbank file format). This `/product` qualifier is not automatically
45 updated when new protein functions are discovered, and hence genome annotation files
46 become out of date over time if not regularly updated. Cross references to databases with
47 functional information are automatically added in a `/db_xref` qualifier during updates on
48 the DDBJ/EMBL/GenBank system; however these are not human readable. Therefore, it is
49 a valuable and necessary exercise to update the protein product information directly in
50 genome annotations of reference species.

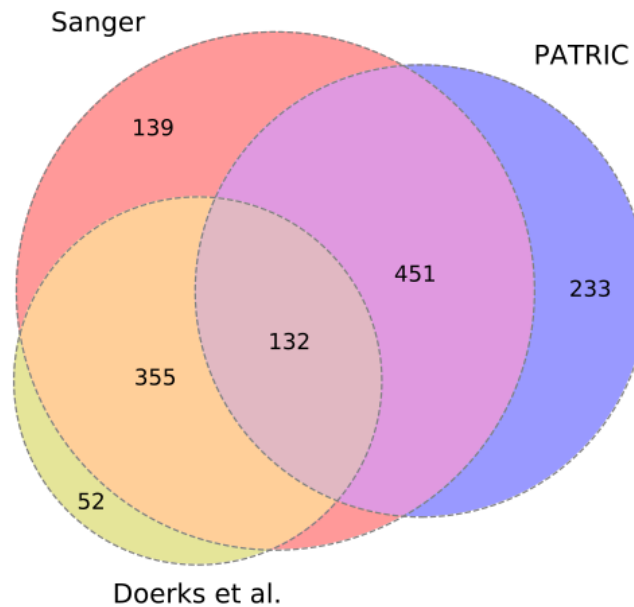
51 **New sources of annotation**

52 Since the original annotation of *M. bovis* AF2122/97, the function of many hypothetical and
53 newly identified proteins has been recognised. These data were collated by Doerks et al.
54 in 2012 [3] who added approximately 620 new functional assignments to the remaining
55 unknowns in the *M. tuberculosis* H37Rv reference genome. This was done using orthology
56 and genomic context evidence. If a hypothetical protein was a member of a known
57 orthologous group in the eggNOG database, this annotation was transferred to the
58 corresponding *M. tuberculosis* protein. The STRING tool was also used, combining gene
59 fusion events and significant co-occurrence to predict links with known proteins. The
60 annotations vary in specificity; those predicted through association are more ‘functional
61 hints’ as to the nature of the protein function rather than ascribing a specific function.

62 The PATRIC database [4] have used their own pipeline based on RASTtk to re-annotate
63 H37Rv. The PATRIC annotation has some additional functional assignments and small
64 genes that were not in the *M. tuberculosis* H37Rv reference. Some of the genes found by
65 Doerks et al. are also assigned in PATRIC and a few have since been assigned to the
66 reference. Figure 1 shows the overlap in these two sets with the existing unknown proteins
67 in the H37Rv reference.

68 UniProt [5] stores up-to-date protein annotations for the *M. tuberculosis* H37Rv strain,
69 several of which are recent and were not present in the either PATRIC or Doerks et al.
70 data. These three sources were used to make an integrated table that was used to update
71 the protein product field of the *M. bovis* AF2122/97 genome annotation.

72



73 Figure 1: Overlap between the hypothetical/unknown proteins in the reference (Sanger)
74 and PATRIC H37Rv annotations and those found by Doerks et al.

75

76 **Methods**

77 Data from three sources was used: 1) Doerks et al. (2012); 2) the PATRIC H37Rv
78 annotation (<https://www.patricbrc.org/view/Genome/83332.12>); and 3) the UniProt *M.*
79 *tuberculosis* H37Rv proteome (<https://www.uniprot.org/proteomes/UP000001584>) were
80 downloaded as csv files and combined together by matching corresponding entries on the
81 H37Rv locus tags (/locus_tag field). For all the unknown proteins in the H37Rv genome we
82 selected updated annotations from these combined sources: if an annotation was present
83 in the Doerk dataset this was used preferentially, then PATRIC and finally UniProt. The
84 order of preference was not significant. The majority of the data was provided from the
85 Doerks dataset though 62 proteins from this dataset were excluded due to lack of a
86 specific functional description.

87 The resulting table with new protein products was then matched to the *M. bovis*
88 orthologous genes using a mapping between *M. tuberculosis* H37Rv and *M. bovis* locus
89 tags. In a final step we performed a BLAST [6] of the remaining unknowns to the Protein
90 Data Bank and found five additional proteins that have structures and function ascribed
91 which were also added. Analysis was done in Python, utilising the Biopython [7] and
92 pandas [8] libraries.

93

94 Results

95 The current *M. bovis* genome annotation contains 3989 protein coding genes [2]. Of these
96 1097 were marked as hypothetical, conserved or unknown proteins. These data are
97 summarized in Table 1, with two *M. tuberculosis* H37Rv annotations for comparison. We
98 have now added a total of 616 new protein product annotations to the genome which
99 includes five products from a PDB search. 69 new gene names were added from the
100 UniProt data. There are now 488 hypothetical/unknowns remaining in the updated *M.*
101 *bovis* AF2122/97 annotation. This revised annotation has been submitted to
102 DDBJ/ENA/GenBank and is available under the accession no. [LT708304](https://www.ncbi.nlm.nih.gov/nuccore/LT708304).

103

	Coding sequences	CDS with gene names	Hypothetical	Pseudogenes
Mbovis AF2122/97	3989	1964 (2026)	1097 (488)	11
MTB-H37Rv (Reference)	4018	1953	1097	27
MTB-H37Rv (Broad)	4143	16	843	92

104

105 Table 1: Current annotation statistics for the *M. bovis* and H37Rv genomes. The revised
106 numbers for the updated annotations for AF2122/97 are shown in brackets.

107

108 Conclusion

109 Reference genomes stored on the INSDC (International Nucleotide Sequence Database
110 Collaboration) provide the primary sources of annotation for virtually all bacterial species.
111 Though there are now multiple alternative information sources for bacterial genomes, most
112 are specialist databases and not universally known. The proliferation of data sources risks
113 fragmentation of genome annotation and linked functional information; it is vitally important
114 that reference sequence annotation in GenBank and Ensembl, the first port of call for the

115 majority of researchers, are as up to date as possible. The *M. bovis* AF2122/97 strain is
116 well established as a reference for *M. bovis* and the MTBC, and will remain a research
117 cornerstone for the foreseeable future; maintaining an updated genome annotation to the
118 research community drove our current work. We note that the latest reference annotation
119 of *M. tuberculosis* H37Rv also lacks many of the updates we have added here to
120 hypothetical proteins, underlining the need for constant curation of reference sequence
121 annotation.

122

123 **Data Availability**

124 All data sources, output files and the Jupyter notebook used to produce this analysis are
125 stored in a github repository at the following url: [https://github.com/dmnfarrell/gordon-](https://github.com/dmnfarrell/gordon-group/tree/master/mbovis_annotation)
126 [group/tree/master/mbovis_annotation](https://github.com/dmnfarrell/gordon-group/tree/master/mbovis_annotation).

127

128 **Acknowledgements**

129

130 This work was supported by the Irish Department of Agriculture Food and the Marine grant
131 15/S/651 (NEXUSMAP) and Science Foundation Ireland (SFI) grant 16/BBSRC/3390, part
132 of the SFI–Biotechnology and Biological Sciences Research Council joint funding
133 partnership.

134

135 **References**

136

- 137 [1] T. Garnier *et al.*, “The complete genome sequence of *Mycobacterium bovis*,” *Proc.*
138 *Natl. Acad. Sci. U. S. A.*, vol. 100, no. 13, pp. 7877–82, Jun. 2003.
- 139 [2] K. M. Malone, D. Farrell, T. P. Stuber, and O. T. Schubert, “Updated Reference
140 Genome Sequence and Annotation of *Mycobacterium bovis* AF2122/97,” *Am. Soc.*
141 *Microbiol.*, vol. 5, no. 14, pp. 17–18, 2017.
- 142 [3] T. Doerks, V. van Noort, P. Minguéz, and P. Bork, “Annotation of the *M. tuberculosis*
143 hypothetical orfeome: adding functional information to more than half of the
144 uncharacterized proteins,” *PLoS One*, vol. 7, no. 4, p. e34302, Jan. 2012.
- 145 [4] A. R. Wattam *et al.*, “PATRIC, the bacterial bioinformatics database and analysis
146 resource,” *Nucleic Acids Res.*, 2014.
- 147 [5] A. Bateman *et al.*, “UniProt: The universal protein knowledgebase,” *Nucleic Acids*
148 *Res.*, 2017.

- 149 [6] S. F. Altschul, W. Gish, W. Miller, E. W. Myers, and D. J. Lipman, "Basic local
150 alignment search tool," *J. Mol. Biol.*, no. 215, pp. 403–410, 1990.
- 151 [7] P. J. A. Cock *et al.*, "Biopython: Freely available Python tools for computational
152 molecular biology and bioinformatics," *Bioinformatics*, 2009.
- 153 [8] W. Mckinney, "Pandas, Python Data Analysis Library," 2015. [Online]. Available:
154 <http://pandas.pydata.org/>.
155
156
157
158