



Provided by the author(s) and University College Dublin Library in accordance with publisher policies., Please cite the published version when available.

Title	Genome-level analyses of Mycobacterium bovis lineages reveal the role of SNPs and antisense transcription in differential gene expression
Authors(s)	Golby, Paul; Nunez, Javier; Witney, Adam; Gordon, Stephen V.; et al.
Publication date	2013-10
Publication information	BMC Genomics, 14 (1): 710
Publisher	BioMed Central
Item record/more information	http://hdl.handle.net/10197/5310
Publisher's version (DOI)	10.1186/1471-2164-14-710

Downloaded 2019-03-24T07:24:37Z

The UCD community has made this article openly available. Please share how this access benefits you. Your story matters! (@ucd_oa)



Some rights reserved. For more information, please see the item record link above.



Genome-level analyses of *Mycobacterium bovis* lineages reveal the role of SNPs and antisense transcription in differential gene expression

Paul Golby¹, Javier Nunez¹, Adam Witney², Jason Hinds², Michael A. Quail³, Stephen Bentley³, Simon Harris³, Noel Smith¹, R. Glyn Hewinson¹, and Stephen V. Gordon⁴.

¹Animal Health and Veterinary Laboratories Agency, Woodham Lane, New Haw, Addlestone, Surrey KT15 3NB.

²Bacterial Microarray Group, Centre for Infection & Immunity, Division of Clinical Sciences, St George's, University of London, Cranmer Terrace, London SW17 0RE, UK.

³The Wellcome Trust Sanger Institute, Hinxton, Cambridgeshire, UK.

⁴UCD School of Veterinary Medicine and UCD Conway Institute, University College Dublin, Dublin 4, Ireland.

1 **Abstract**

2 Bovine tuberculosis (bTB) is a disease with major implications for animal welfare and
3 productivity, as well as having the potential for zoonotic transmission. In Great Britain (GB)
4 alone, controlling bTB costs in the region of £100 million annually, with the current control
5 scheme seemingly unable to stop the inexorable spread of infection. One aspect that may
6 be driving the epidemic is evolution of the causative pathogen, *Mycobacterium bovis*. To
7 understand the underlying genetic changes that may be responsible for this evolution, we
8 performed comprehensive genome-level analyses of 4 *M. bovis* strains that encompass the
9 main molecular types of the pathogen circulating in GB. We show that while these strains
10 show extensive similarities in their genetic make-up and gene expression profiles, they
11 exhibit distinct expression of a subset of genes. We provide genomic, transcriptomic and
12 functional data to show that synonymous point mutations (sSNPs) on the coding strand can
13 lead to the expression of antisense transcripts on the opposing strand, a finding with
14 implications for how we define a 'silent' nucleotide change. Furthermore, we show that
15 transcriptomic data based solely on amplicon arrays can generate spurious results in terms
16 of gene expression profiles due to hybridisation of antisense transcripts. Overall our data
17 suggest that subtle genetic differences, such as sSNPs, may have important consequences
18 for gene expression and subsequent phenotype.

19

20

21 Introduction

22 *Mycobacterium bovis* is the causative agent of bovine tuberculosis (bTB), an endemic
23 disease of cattle in Great Britain (GB) with the potential for zoonotic transmission to
24 humans. In GB the primary control of bTB is through ‘test and slaughter’ surveillance,
25 whereby cattle that are positive to the tuberculin skin test [1] are removed from the herd
26 and slaughtered. In spite of this approach, which has been in place since the 1950s, the
27 number of TB-positive cattle slaughtered is increasing year on year - approximately 30,000
28 cattle were tested and slaughtered between 2012-2013, compared to 300 between 1995-
29 1996 (<http://www.defra.gov.uk/animal-diseases/a-z/bovine-tb/>). The UK (GB and Northern
30 Ireland) governments currently spend approximately £100 million per year collectively on
31 control measures and compensation to farmers for slaughtered cattle. The failure of the
32 test-and-slaughter policy to control the spread of infection in large parts of GB suggests that
33 we need a much greater understanding of the TB disease dynamic, including the role of
34 pathogen diversity as a potential driver of this process.

35 *M. bovis* isolates that are cultured from skin test-reactor animals are currently genetically
36 typed using a combination of spoligotyping [2] and VNTR [3]. Spoligotyping exploits a
37 polymorphic region of the genome called the DR locus which consists of multiple, identical
38 36bp repeats interspersed with unique sequences known as spacers. Isolates of *M. bovis*
39 differ in the presence or absence of spacers and adjacent DRs, allowing a ‘barcode’ to be
40 generated for each molecular type. Spoligotypes 9 and 17 are the dominant molecular types
41 in the UK, with more than one third of all isolates corresponding to Type 9 and a quarter to
42 Type 17. VNTR measures the variation at repeat sequences in 6 regions of the genome.
43 There are 6 major VNTR types for Type 9, while all others show only one dominant profile,
44 suggesting that *M. bovis* Type 9 strains are more genetically variable compared with other
45 spoligotypes. Integration of molecular typing with geographical information systems allows
46 temporal and spatial distribution of molecular types to be mapped across GB. Type 9
47 isolates are widely distributed across GB, while type 17 is an emerging clone which has
48 expanded out of foci around Gloucester, Hereford and Worcester. Similarly, Types 25 and 35
49 have expanded out of Staffordshire/Shropshire and Hereford/Worcester, respectively.
50 Between them, types 25, 35, 9 and 17 encompass the diversity of the major clonal lineages
51 of *M. bovis* circulating in the UK.

52 An analysis of molecular typing data from ~11,500 *M. bovis* isolates revealed that the
53 population structure of *M. bovis* in GB could not be explained by random mutation and drift
54 and instead, it appeared that certain strains were increasing at a faster rate relative to
55 others [4]. One suggestion for the 'clonal expansion' of GB *M. bovis* genotypes was that
56 certain genotypes had a selective advantage over others leading to an increase in their
57 frequency in the population [4]. Supportive of this hypothesis, several lines of evidence have
58 suggested that *M. bovis* isolates show phenotypic differences to each other. Fourier-
59 Transform Infrared Spectroscopy (FT-IR) has been used to generate metabolic profiles of the
60 10 major spoligotype groups of *M. bovis* isolates circulating in GB. Clustering analysis of the
61 resulting spectra showed that the spectra could be differentiated according to spoligotype,
62 indicating that strains of different spoligotypes possess phenotypically distinct traits [5]. In
63 addition, it has also been shown that type 17 isolates have lower incorporation rates of
64 propionate into membrane lipid components compared to other field strains, suggesting a
65 degree of metabolic remodelling in the type 17 lineage [6]. Hence it appears that genetic
66 differences across *M. bovis* lineages may impact on phenotypic traits. This latter finding may
67 have important implications for vaccine and diagnostic test development, in terms of which
68 experimental challenge strains to test vaccines against or on influencing diagnostic test
69 performance.

70 In an attempt to better define genetic differences across the major *M. bovis* lineages
71 circulating in GB that may give rise to phenotypic differences of practical importance, we
72 have used a combination of genome sequencing, transcriptome analyses, and recombinant
73 DNA technology. The genomes of three *M. bovis* field isolates were sequenced using
74 Illumina sequencing technology and strain specific differences in gene expression were
75 measured during in vitro growth and in ex vivo bovine alveolar macrophages (M ϕ) using a
76 whole genome amplicon microarray. Recent discoveries of small non coding RNA within
77 mycobacteria [7] [8] prompted us to assess differences in sRNA expression across the
78 isolates using a whole genome tiled oligonucleotide microarray. SNP/small base pair
79 insertion and deletions (INDELs) and gene expression data were overlaid onto the genomic
80 sequence of the fully sequenced strain of *M. bovis* 2122/97 to link observed strain specific
81 genomic differences with differences in RNA expression.

82

84 Results

85 *Comparative genomics of M. bovis field isolates using whole genome sequencing and* 86 *microarrays*

87 The strains for this study were chosen to reflect the genomic diversity of the *M.*
88 *bovis* population circulating in GB, and are listed in Table 1. *M. bovis* strains were typed
89 using a combination of spoligotyping and VNTR. For each spoligotype group, an isolate
90 which possessed the most common VNTR profile was selected, so that each chosen strain
91 was the most representative of each spoligotype group (Table 1). Of the four studied strains,
92 2451/01 and 1307/01/01 diverged earliest during descent from the most recent common
93 ancestor of *M. bovis* in GB and are more distant to strains 1121/01 and 1307/01 (Smith, N.
94 personal communication). All 4 strains were isolated from diseased cows belonging to herds
95 which were taken from farms in geographically diverse areas of the country.

96 The genomes of the three *M. bovis* strains 1121/01, 2451/01 and 1307/01 were
97 paired-end sequenced using Illumina sequencing technology. Processed sequence reads
98 were mapped to the genome of the fully sequenced and annotated strain 2122/97 [9] to
99 identify SNPs. The detection of INDELS using short read Illumina sequence data is
100 challenging and requires the application of complex algorithms. We have therefore focussed
101 our attentions on the identification of SNPs only. In total, 1031 polymorphic sites were
102 identified associated with single nucleotide changes. SNPs were identified across all four
103 sequenced *M. bovis* strains, and their positions, together with their SNP class, are listed in
104 Table S1. The numbers of SNPs between *M. bovis* 2122/97 and the other three strains were
105 found to be consistent with their predicted evolutionary distances from each other. Strain
106 1121/01 (type 17) is most closely related to the original genome sequenced strain 2122/97
107 (type 9) with 118 SNPs, whereas the more distantly related strains 2451/01 (type 25) and
108 1307/01 (type 35) have 485 and 618 SNP differences respectively.

109 For each strain, the percentages of non-synonymous (nSNP), synonymous (sSNP) and
110 intergenic SNPs (iSNP) are shown in Fig. 1B. The ratios of each class of SNP for strains
111 2451/01 and 1307/01 are almost identical to each other at 4:2:1 (nSNP: sSNP: iSNP). The

112 SNP profile for strain 1121/01 was, however, different to that of 2451/01 and 1307/01, with
113 the ratio of SNPs in each class to be approximately 3.2:2.5:1.

114 Large sequence polymorphisms (LSPs) across the *M. bovis* field strains were
115 identified using a combination of *in silico* comparisons of the genome data as well as
116 microarray technology, the latter being achieved by isolation of genomic DNA from all four
117 strains, labelling with fluorescent dyes and hybridisation to a whole genome *M.*
118 *tuberculosis/M. bovis* amplicon microarray (see Methods). Table 2 lists several LSPs that
119 were detected across the strains. The large 6.8kb deletion (RDbovis(d)_0173) which appears
120 to be specific to UK strains belonging to Type 17, has been described in a previous study [10]
121 and encompasses genes Mb1963-Mb1971. Several of these gene products are predicted to
122 encode proteins involved in lipid metabolism, but the lipid composition of several type 17
123 isolates was found to be no different to other *M. bovis* strains, although their ability to
124 incorporate propionate into mycolic acids was found to be lower [6]. A smaller 1.6 kb
125 deletion specific 1307/01 was detected that comprises the 3' end of Mb2056c, Mb2055c,
126 and the 5' end of *pkfB* (Mb2054c). Due to a single base deletion, Mb2056c and Mb2055c are
127 pseudogenes in 2122/97, but the two genes exist as one intact functional gene in 2145 and
128 H37Rv (Rv2030c). The *pkfB* gene encodes a phosphofructokinase homologue and is strongly
129 immunogenic in human TB patients, while Rv2030c encodes an erythromycin esterase. Both
130 *pkfB* and Mb2056c are members of the DosR regulon, which are highly upregulated under
131 anaerobic conditions and have been implicated in bacterial persistence *in vivo* [11]. Other
132 smaller deletions detected include a deletion of a probable lipid transfer protein encoding
133 gene Mb1699c, which is specific to 1307/01, and an aldo/keto reductase encoding gene,
134 Mb2320 that is specific to 1121/01.

135

136 *Linking SNPs to genes that show differential expression amongst M bovis strains grown*
137 *under vitro conditions and in ex vivo macrophages*

138 The four *M. bovis* field strains were grown to mid-logarithmic phase in pyruvate-
139 containing Middlebrook 7H9 liquid media, and then used to infect bovine alveolar M ϕ using
140 a multiplicity of infection (MOI) of 10:1 (bacilli: M ϕ). Mycobacterial RNA was recovered from
141 infected host cells 4 and 24 hrs post infection using a differential lysis procedure and

142 amplified using a modified procedure similar to that described by van Gelder et al ([12]; see
143 Methods). As a control, RNAs were also extracted from strains that had been incubated
144 statically in RPMI cell culture media for a period of 4 hrs. To eliminate potential skewing
145 effects on the transcriptome resulting from the amplification process, comparisons were
146 made only between amplified RNA generated from samples collected at the same time
147 point and biological replicate. For the in vitro growth condition, total RNAs were extracted
148 from the four strains grown in a pyruvate-containing Middlebrook 7H9 liquid media and
149 rolled during incubation.

150 The RNAs extracted from each of the growth conditions were converted to cyanine
151 labelled cDNA using reverse transcriptase and hybridised to whole genome amplicon
152 microarrays. Using only those genes that are common to all four strains, we found a total
153 set of 70 genes that showed a 2.5-fold or more difference in expression in one or more
154 strains when pairwise comparisons were made between the transcriptomes of 2122/97 and
155 1121/01, 2451/01 or 1307/01 (Table S2). A subset of these 70 genes is shown in Table 3
156 where key examples of alterations in metabolic processes are shown. The numbers of genes
157 that were found to show differential expression across the four strains reflected the
158 evolutionary distances between 2122/97 and the other 3 strains. Thus, 1121/01, which is
159 closely related to 2122/97, shows only 5 differentially expressed genes, while the most
160 distantly related strain 1307/01 shows 56 genes differently expressed compared to
161 2122/97. Of these 56 genes, 5 were specific to the in vitro condition, while 19 were specific
162 to the M ϕ . Ten genes were common to both conditions, which serves to validate the
163 technical reproducibility of the RNA amplification process.

164 Using the genome sequencing information determined for each of the four strains,
165 we attempted to correlate the observed strain-specific differences in gene expression with
166 the presence or absence of mutations within the coding regions or promoters of those
167 genes that show differential gene expression, or in genes that are known to regulate the
168 activity of those genes. Mb1749c and Mb1750c are two genes that are specifically
169 upregulated in 1307/01 and encode a toxin and antitoxin (TA) pair, respectively, belonging
170 to type II TA systems of the VapBC family [13]. Members of VapB type toxins contain PIN
171 domains that cleave RNA and thus function to control translation of mRNA transcripts [14].
172 The homologous genes from strain 1307/01 show up to 19- and 10-fold higher levels of

173 expression, respectively, than those of the other three strains. An analysis of the coding
174 sequences of Mb1749c across all four strains revealed that the 1307/01 homologue has a
175 unique nSNP at position 1932704 (wrt 2122/97 genomic sequence), a C-T transition that
176 results in the nonconservative substitution of Gly19 to Asp. Research has shown that TA
177 gene pairs negatively regulate their own expression through binding of the TA protein
178 complex to the promoter region of the TA gene pair, thus preventing access to RNA
179 polymerase [15]. The G19D mutation in Mb1749c could therefore impair the ability of the
180 complex to bind to the promoter resulting in the deregulation of the TA gene pair.

181 Mb2007c, which shows a 4-fold higher expression in 1307/01 only, encodes a
182 transcriptional regulator of the LysR class. There are two SNPs present in the coding
183 sequence of Mb2007c in 1307/01 which are absent in the homologues of the other three
184 strains: the first is a nSNP which results in the conservative substitution of Arg137 to Gln,
185 while the second is a more debilitating nonsense SNP, which ultimately leads to a protein
186 whose length is only 60% that of the wild-type. Many regulators belonging to the LysR
187 family regulate their own expression through a negative autoregulatory mechanism similar
188 to that described above for VapBC TA systems [16]. A loss in protein integrity could,
189 therefore, result in the regulator being unable to bind the regulatory region, leading to the
190 observed upregulation in the expression of this gene in 1307/01. As the product of this gene
191 is predicted to be a transcriptional regulator, it was speculated that the regulation of gene(s)
192 controlled by regulator could be affected in 1307/01 due to the severely truncated form of
193 this protein. As LysR regulators are often found to regulate genes that are divergently
194 transcribed from the lysR gene, it was surprisingly to find that expression of the Mb2008
195 homologue in 1307/01, which is predicted to encode a lysine transporter, does not show
196 any difference in expression in 1307/01 to 2122/97. To define the regulon of this regulator,
197 we first compared the transcriptomes of 2122/97 transformed with a multicopy plasmid
198 expressing the truncated copy of *mb2007c* against a vector only control. No differences in
199 expression were found (data not shown), which could indicate that the regulator does not
200 control any other genes apart from itself, or that experimental conditions did not favour the
201 active form of the regulator. LysR regulators regulate expression of their regulon through
202 binding of a co-inducer to the C-terminal domain, and the failure to observe any changes in
203 gene expression could therefore be due to the absence of the co-inducer during the

204 experiment. A further experiment to compare the profiles of 2122/97 expressing either the
205 truncated or wild type forms of the protein also showed no differences in expression (data
206 not shown).

207 Nitrite reductase catalyses the reduction of nitrite to ammonia and is strongly
208 expressed during growth in the presence of nitrate or nitrite, but repressed in the presence
209 of ammonia [17]. The gene encoding the large subunit of the nitrite reductase, *nirB*
210 (Mb0258), shows approximately 9-fold higher expression in 1307/01 compared to the other
211 3 strains in our standard ammonia containing 7H9 growth media, suggesting that the strain
212 has lost regulatory control of this gene. Expression of *nirB* in *M. tuberculosis* has been
213 shown to be controlled by the response regulator GlnR [18], but an analysis of the sequence
214 of the *glnR* orthologue from all four strains revealed no differences in either the coding or
215 promoter sequences. A comparison of the *nirB* sequence from all 4 strains did, however,
216 reveal the presence of a single base (C to T) transition leading to a sSNP that is specific to
217 1307/01. It was not readily apparent why a sSNP in the coding sequence of a gene should
218 lead to an increase in expression of that gene, but there are several reports that show sSNPs
219 leading to changes in stability of mRNA transcripts [19] [20]. Rv0987 and Rv0988 of *M.*
220 *tuberculosis* H37Rv encode part of an ABC transporter and a putative secreted hydrolase,
221 respectively. In 2122/97, a single base transition (G-A) introduces a stop codon that splits
222 Rv0987 into the two pseudogenes, Mb1013 and Mb1014. Previous microarray based gene
223 expression studies by our group have shown that Rv0987 and Rv0988 in *M. tuberculosis*
224 show higher levels of expression than the orthologous Mb1013/Mb1014 and Mb1015,
225 respectively, in *M. bovis* 2122/97 [21], and in the present study the Mb1013/Mb1014 and
226 Mb1015 homologues in 2451/01 and 1307/01 also showed higher expression (up to 10-fold)
227 than the homologues in 2122/97 and 1121/01. Comparing the sequences of
228 Mb1013/Mb1014 and Mb1015 across all 4 strains indicated that strains that show high
229 expression have the 'G' allele.

230 Mb3477c encodes an ATP binding membrane protein, part of the Esx4 secretion
231 system [22], and gene shows up to 10-fold higher expression in 2451/01 and 1307/01
232 compared to 2122/97. The gene also contains an A to C transition at position 3812465, a
233 nSNP at position resulting in the non-conservative substitution of a serine to a glycine
234 residue.

235 Of the 19 genes that show specific differential expression in the M ϕ , the most
236 notable are Mb1914c and *echA21*, which show upregulation in 2451/01 only (up to 6- and
237 23-fold, respectively). Both genes encode proteins that could be involved in lipid
238 metabolism, and both genes contain single sSNPs that are present in 2451/01, but absent in
239 the other three strains.

240 Real time RT-PCR was used to verify a selection of genes that showed differential
241 gene expression as predicted by the microarray analysis. Figure 2 compares the fold changes
242 in the expression levels of 4 genes as measured by microarray and real time RT-PCR. The
243 *nirB* and Mb1749c genes were selected because they showed strong upregulation in
244 1307/01 in both *in vitro* and *ex vivo* M ϕ while Mb1914c and *echA21* were chosen because
245 the array data predicted them to be specifically upregulated in 2451/01 and only in *ex vivo*
246 M ϕ . For each of the 4 genes, the strain dependent pattern of expression as measured by
247 real time RT-PCR was consistent with that measured by microarray, although the fold
248 changes measured by real time RT-PCR were higher than those measured by microarray.

249

250 *Functional analysis of SNP role in differential gene expression*

251 The above data showed that many of the strain specific differences in gene
252 expression were linked to the presence of synonymous or non-synonymous SNPs located
253 within the coding regions of the genes that show variable expression. Non-synonymous
254 SNPs lead to changes in amino acid sequence which can lead to changes in protein function.
255 The C to T transition at position 1932704 (wrt 2122/97) in the coding sequence of the
256 1307/01 Mb1749c homologue leads to the non-conservative substitution of Gly19 to Asp,
257 and this nSNP appears to be linked to the upregulation of both Mb1749c and Mb1750c in
258 that strain. In order to confirm this, a 0.9 kb DNA fragment containing the *Mb1749c-*
259 *Mb1750c-MB1751c* region of 1307/01 (containing the 'T' allele) and the equivalent region
260 from 2122/97 (with the 'C' allele) were PCR amplified and the fragments were cloned
261 separately into the mycobacterial shuttle vector pKINT (see Methods) to create the
262 constructs pPG107 and pPG106, respectively. The two constructs were introduced into
263 *Mycobacterium smegmatis* mc²155, separately, and then the expression of Mb1749c and
264 Mb1750c in *M. smegmatis* pPG101 was compared to that of *M. smegmatis* pPG102 using

265 real time RT-PCR. Table 4 shows that the expression levels of Mb1750c and Mb1749c in the
266 strain expressing the mutated forms of Mb1749c/Mb1750c are 13- and 9-fold higher,
267 respectively, compared to the strain expressing the wild type forms, confirming that this
268 SNP is responsible for the observed up-regulation of the two genes in 1307/01.

269 Synonymous substitutions do not lead to changes in protein sequence and have
270 generally been considered to be 'silent' or benign. Recent studies, however, have suggested
271 that sSNPs can have functional effects, such as decreased mRNA stability and translation
272 [19] [20]. In our own studies, we have found several genes whose expression levels
273 correlate with the presence of sSNPs in the coding regions of those genes. For example, a C-
274 T transversion at position 303227 (wrt 2122/97) within the coding sequence of *nirB* of
275 1307/01 is a sSNP that appears to be linked with the upregulation in expression of *nirB*
276 within that strain. To confirm that this is the case, we PCR amplified 3.5 kb DNA fragments
277 containing the *hsp-nirB-nirD-cobU* region of strain 1307/01 (containing the 'C' allele) and the
278 equivalent region from strain 2122/97 (with the 'T' allele) and cloned them separately into
279 the integrating vector pKINT to create the constructs pPG108 and pPG109, respectively.
280 These constructs were introduced into 2122/97 and the expression levels of *nirB* and *nirD*
281 were found to be 30- and 2-fold, respectively, higher in the strain expressing the mutated
282 *nirBD* locus compared to the strain expressing the wild-type form (Table 5). This confirms
283 that this mutation is responsible for the upregulation of the two genes in this strain.

284

285 *Use of a high density tiled oligonucleotide microarray to detect differentially expressed small* 286 *RNA transcripts in M. bovis isolates*

287 The *M. tuberculosis/M. bovis* amplicon arrays used in the present study were
288 specifically designed to measure expression levels of genes annotated in the genomic
289 sequence of *M. bovis* 2122/97 [9]. They were not, however, designed to monitor the
290 expression of non-coding RNA such as small RNA within intergenic regions or antisense
291 sRNA. Hence, a high density tiled oligonucleotide microarray consisting of approximately
292 180,000 partially-overlapping (10-base overlap) short 60 mer oligonucleotides was designed
293 that offered an unbiased approach to the detection of strand specific transcripts encoded
294 over the entire *M. bovis* 2122/97 chromosome. Total RNA enriched for small sized (<100 nt)

295 RNA species was extracted from the four *M. bovis* strains that had been grown in liquid
296 media and hybridised to the oligonucleotide microarray. To avoid potential secondary
297 strand synthesis during cDNA synthesis, which could be interpreted as sRNA, the RNA was
298 directly labelled with cyanine based dyes. After pairwise comparisons were performed
299 between 2122/97 and 1121/01, 2451/01 or 1307/01, 220 oligonucleotide probes were
300 identified that detected differentially expressed transcripts (2.5 fold cut off) in one or more
301 of the three strains (Table S4). Only transcripts detected by multiple (2 or more) overlapping
302 probes were regarded as genuine transcripts as those detected by single probes could be
303 due to cross-hybridisation effects or represent spurious transcripts. Using these criteria, 26
304 transcripts, designated T1-T26, were found to show differential expression in one or more of
305 the strains (Table 6), and those transcripts can be divided into those that are encoded within
306 intergenic regions and those encoded within the genomic co-ordinates encompassing
307 annotated coding sequences. Comparison of the differentially expressed gene lists identified
308 using amplicon vs. oligonucleotide arrays (Table 5), it is clear that many of the transcripts
309 detected using the amplicon arrays are not necessarily encoded on the sense gene strand,
310 as had been previous interpreted. For example, the amplicon array data had appeared to
311 suggest that Mb1914c and *echA21* were upregulated in 2451/01, but the oligo array data
312 indicates that transcripts 11 and 25, which are encoded within the co-ordinates encoded by
313 those two genes, are actually encoded on the antisense strands. This apparent discrepancy
314 can be rationalised once we consider that double stranded amplicon microarray probes
315 cannot discriminate between transcripts encoded on the sense or antisense strands.
316 Transcripts 11 and 25 can therefore be considered as potential antisense sRNAs (asRNA),
317 which could be involved in translational or post-transcriptional control of the sense
318 transcript. Other potential cis-encoded sRNAs detected using the arrays include T6, T14, and
319 T15/T16 which are encoded on the antisense strands to Mb1618c, Mb2117 and Mb2607,
320 respectively, and for each of these transcripts, their expression appears to be linked to the
321 presence of a single SNP within the co-ordinates of the genes. The approximate boundaries
322 of these transcripts can be derived from the genomic co-ordinates of the oligonucleotide
323 probes that detect the expression of the transcript. Thus, the transcripts appear to be
324 between 100-300 nt in size and the positions of the linked SNPs appear to be positioned
325 either just upstream or within the predicted 5' end of the transcripts (Figure 3). Three of the
326 transcripts (T11, T14 and T25) are antisense to the central part of the sense encoded gene,

327 while T6 is encoded antisense to the 5' end of Mb1618c. As well as antisense transcripts, we
328 also saw the differential expression of sense transcripts. The amplicon microarray data
329 (confirmed by real time RT-PCR) indicated that *nirB* is strongly upregulated specifically in
330 1307/01 in both in vitro and ex-vivo M ϕ . An analysis of the oligonucleotide array data,
331 however, indicates that there are two short transcripts, T1 and T2 (sense and antisense,
332 respectively) that are encoded within the genomic co-ordinates of the *nirB* gene. T2 is the
333 longer in size (255 vs. 155nt) and more highly expressed (5 vs. 3-fold) than T1, and both
334 transcripts appear to be linked to the presence of a SNP that is located within the middle of
335 T1 and approximately 50nt upstream of T2.

336 Some of the transcripts are *bone fide* gene sense strand mRNA transcripts, such as
337 T9 and T10 which are encoded by Mb1749c and Mb1750c, respectively. Although it would
338 appear that the two genes are transcribed separately, it is probable that the two transcripts
339 are co-transcribed as the stop codon of Mb1750c overlaps the start codon of Mb1749c.
340 Eight of the transcripts listed in Table 6 are encoded within intergenic regions, 7 of which
341 are encoded within the polymorphic direct repeat (DR) locus. The DR locus of strains
342 belonging to the *M. tuberculosis* complex has been suggested to constitute a CRISPR locus
343 which have been shown in many species of bacteria to be involved in protection against
344 exogenous foreign DNA such as plasmids and phage [23]. All the DR encoded transcripts are
345 short (approx. 100 nt), straddle contiguous repeat and spacer sequences and show
346 approximately 5-fold higher levels of expression in 2451/01/01 and 1307/01 compared with
347 2122/97 and 1121/01.

348

349 *Characterisation of differentially expressed cis asRNA*

350 The genomic co-ordinates of the oligonucleotide probes that detected the antisense
351 species described above can only serve as approximate estimations as to their start and end
352 points. Thus, we used 5' RLM-RACE (RNA Ligase Mediated Rapid Amplification of cDNA
353 Ends) in an attempt to accurately define the transcriptional start sites (TSS) for the short
354 sense transcript T2, and the antisense transcripts T6, T11 and T25 described in the above
355 section (see Methods). These transcripts were chosen as their expression levels are high and
356 their transcript lengths were considered to be sufficiently long to enable the RLM-RACE

357 methodology to work. Table 6 details the sizes of the PCR products obtained after RLM-
358 RACE was performed using oligonucleotide primers designed to sequences predicted for
359 transcripts T6, T11 and T25. No PCR product was obtained for transcript T2. For each of the
360 three transcripts, the TSS was determined to be a G residue, which is the most commonly
361 used residue type for mycobacterial TSS's [24]. For each of the T6, T11 and T25 transcripts,
362 expression of the asRNAs was linked to the presence of a SNP (C to T) proximal to the 5' end
363 of the asRNA. Strains exhibiting the 'C' allele showed no expression of the asRNA, whilst the
364 strain that showed expression had the 'T' allele. An analysis of the nucleotide sequence in
365 the vicinity of the SNPs reveals that for each of the three transcripts the SNP constitutes the
366 6th residue of a motif that has strong homology to the consensus sequence for the -10
367 element of Group A mycobacterial promoters [24]. The finding that a 'T' residue is
368 associated with expression is consistent with the consensus sequence which indicates that
369 86% of all -10 elements have a 'T' residue at the 6th residue position. Several residues that
370 flank the -10 motif also show a degree of conservation. Sequence motifs which show strong
371 homology to group A -35 elements are present 18-19 bp upstream of the putative -10
372 elements, and the distances between the -35, -10 and TSS elements are consistent with
373 those elements of the consensus sequence. No protein encoding open reading frames were
374 detected within the T6, T11 and T25 transcripts.

375 In a parallel study, high density oligonucleotide microarrays were also used to
376 interrogate the transcriptomes of *M. tuberculosis* H37Rv, *M. bovis* BCG Pasteur,
377 *Mycobacterium caprae* and *M. bovis* AN5 that had been grown in Middlebrook 7H9 media.
378 As a result of these experiments, two asRNA species were found to be expressed within the
379 antisense strands of the *ino1* and *narH* genes of *M. tuberculosis* H37Rv, but not in any of the
380 other 3 strains tested (data not shown). A comparison of nucleotide sequences of the
381 orthologous genes across the species suggested that expressions of the as_sRNAs correlated
382 with the presence of a sSNP (C to T transition at positions 50555 and 1292100 wrt H37Rv
383 genomic sequence for *as_ino1* and *as_narH*, respectively) upstream of the asRNAs.
384 Approximate information regarding the transcriptional start site was deduced from the
385 binding co-ordinates of the probes that detected the transcripts. As with the *M. bovis*
386 antisense sRNAs described above, the T residue associated with the expression of the *M.*
387 *tuberculosis* asRNAs is part of a putative -10 element. A -35 element with an appropriate

388 spacing to the -10 element was identified for as_ino1, but not for as_narH, suggesting that
389 the as_narH promoter may belong to group B mycobacterial promoters that have a
390 conserved -10, but no -35 motif [24].

391

392 Discussion

393 The aim of this work was to define possible phenotypic variation across *M. bovis* field
394 isolates through a combination of genome sequencing, comparative genomics and
395 transcriptome analyses from both in vitro and ex vivo conditions. Using these approaches
396 we uncovered a range of novel findings, the most striking of which was the realisation that
397 genes that had been predicted to be differentially expressed based on amplicon-microarray
398 data were in fact not upregulated, and that instead it was an antisense transcript that was
399 showing differential expression. Analysing both transcriptome and genome sequence data
400 allowed us to identify SNPs responsible for the transcription of antisense RNAs, with
401 generation of a consensus -10 promoter sequence the likely mechanism. Our results suggest
402 that data generated from amplicon arrays in the past may need to be revisited, as it is
403 possible that some coding-sequences identified as being differentially expressed were
404 instead antisense transcripts.

405 With the growth of technologies such as high density tiled oligonucleotide
406 microarrays and next generation sequencing there has been a rapid increase in the number
407 of reports describing the existence of non-coding RNAs (ncRNAs) in bacteria. Non-coding
408 RNAs broadly consist of two types, cis- and trans-encoded RNA. Trans RNA includes
409 intergenic encoded RNA, while cis-encoded RNA includes 5' and 3' untranslated regions of
410 mRNA and antisense RNA. To study the expressions of both cis- and trans encoded ncRNA
411 we used a high density oligonucleotide tiled microarray since our amplicon microarray was
412 unable to detect intergenic transcripts or differentiate between sense and antisense
413 transcripts. Previous studies using *M. tuberculosis* have identified substantial amounts of
414 ncRNA encoded in both intergenic and intragenic regions [7] [8]. We detected substantial
415 amounts of ncRNAs in *M. bovis*, including many instances of cis-antisense RNA species. Due
416 to their perfect complementarity, cis asRNA form a duplex with the sense strand encoded

417 transcript resulting in either degradation [25] or translation inhibition [26] of the sense
418 mRNA. Antisense RNAs vary in length, ranging from 10s to 1000's of nucleotides and can be
419 classified according to their encoded position with respect to the opposite sense encoded
420 gene. Thus, they can be classified as 5' or 3' overlapping, while others are classified as
421 internally located. The 1121/01 specific as_Mb1618c is an example of a 5' overlapping
422 asRNA, which is encoded antisense to gene Mb1618c which is predicted to express a
423 secretory lipase. The location of the asRNA transcript suggests it may function to prevent
424 translation of Mb1618c mRNA by steric hindrance of the ribosome binding site.

425 In the work presented here strain 2451/01 expressed two asRNAs, as_Mb1914c and
426 as_echA21, that are not expressed by any of the other 3 strains. They are encoded within
427 the central part of the opposite genes and are therefore likely to modulate the stability of
428 the transcripts. Mb1914c encodes a short chain dehydrogenase while *echA21* encodes an
429 enoyl-CoA hydratase. Short chain dehydrogenases catalyse a wide range of functions so the
430 precise function and identity of the substrate is difficult to deduce from sequence alone.
431 Enoyl-CoA hydratases hydrate double carbon-carbon bonds of macromolecules and are vital
432 in the metabolism of fatty acids. Both gene products would therefore appear to be involved
433 in the metabolism of a macromolecule and their similar expression profiles in this strain
434 could indicate involvement in the metabolism of the same molecule, or molecules that are
435 of the same pathway.

436 In many instances, upregulation of asRNA negatively correlates with the
437 transcription of the antisense gene [25], but in many cases expression of the antisense
438 transcript has no effect on the transcription of the opposite gene. In our studies,
439 expressions of as_Mb1618c, as_Mb1914c and as_echA21 did not appear to have any effect
440 on the expressions of the opposite sense encoded genes (data not shown). We have shown
441 that the expressions of the asRNAs are associated with the presence of SNPs, which are
442 either synonymous or non-synonymous with respect to the sense transcript, but upstream
443 of the asRNA transcriptional start site. This highlights the fact that mutations can potentially
444 affect expression of transcripts on both strands, and that the classification of a SNP is strand
445 dependent. For each of the three asRNAs, the associated SNP was found to be located
446 within a putative -10 promoter motif of group A mycobacterial promoters. The sixth residue
447 of the -10 hexamer motif consensus sequence is a strongly conserved 'T' residue, which is

448 present in 81% of all group A mycobacterial promoter elements. Its importance is
449 underlined by the finding that the strains that exhibit a 'C' residue at this position show no
450 detectable expression of the asRNA, while strains having a 'T' residue at this position exhibit
451 expression.

452 Single nucleotide polymorphisms were found to be the most frequent form of
453 genetic variation that exists between the isolates, with a total of 1013 SNPs detected across
454 the three strains 1121/01, 2451/01 and 1307/01 compared to the reference strain *M. bovis*
455 2122/97. Non-synonymous SNPs, which include both non-sense and missense SNPs, are a
456 class of SNPs most likely to impact on protein function and contribute to phenotypic
457 variation. Non-sense SNPs, which results in the expression of a truncated polypeptide due to
458 the introduction of a premature stop codon, were identified in five genes across the strains.
459 Of these, we focussed our attentions on the non-sense SNP present in a gene encoding the
460 LysR regulator, Mb2007c as mutations affecting regulators are likely to impact on the
461 expression of one or more genes that are part of the regulon of the regulator and are
462 therefore more likely to result in phenotypic variation. Experiments to compare the
463 transcriptomes of a strain that exhibited the mutation with a strain overexpressing a
464 functional regulator did not, however, reveal any differences. The reason for this unclear,
465 but could reflect a requirement of the regulator for a co-inducer that was absent under the
466 conditions of the experiment. The consequences of missense SNPs are more difficult to
467 predict, as substitutions of one amino acid for another in a protein sequence do not
468 necessarily lead to a change in protein function. However, for genes that are controlled by
469 an autoregulatory mechanism, a mutation that affects the ability of the product of the gene
470 to regulate itself will result in a change in expression of the gene. In our studies, we have
471 shown that the presence of a missense mutation in a VapB type toxin encoding gene
472 Mb1749c in strain 1307/01/01 results in the upregulation in expression of the toxin-
473 antitoxin encoding pair of genes Mb1750c-Mb1749c due to the inability of the encoded
474 proteins to self regulate themselves. Toxin-antitoxin systems have a variety of proposed
475 cellular functions including general regulation of mRNA stability levels in the cell [27].
476 Further experiments are required to fully understand the consequences of this mutation.

477 Genomic deletions have played an important role in the evolution of strains
478 belonging to the mycobacterial complex [28], and in the derivation of the tuberculosis

479 vaccine strain *M. bovis* BCG [29]. In addition to the previously described 6.8 kb gene
480 deletion that is specific to strains having a spoligotype 17 pattern [6] , we have identified a
481 1.6 kb multi-gene deletion that is specific to strain 1307/01 and encompasses genes that are
482 part of the DosR regulon. However, one of the deleted genes exists as a pseudo gene in
483 DosR in strain 2122/97, so its importance to the biology of *M. bovis* is unlikely to be
484 significant. Several other genes with internal deletions were detected but none of the
485 encoded proteins have any significant similarity to any protein with a defined function.

486 In conclusion we have performed a comprehensive analysis of 4 *M. bovis* strains of
487 the most common molecular types circulating in GB. We show that while these strains show
488 extensive similarities in their genetic make-up and gene expression profiles, they show
489 distinct differences in the expression of a subset of genes. We provide functional data to
490 show that SNPs can lead to the expression of antisense RNA, a finding with implications for
491 how we define a 'silent' nucleotide change. Furthermore, we show that the interpretation
492 of transcriptome data based solely on amplicon arrays could lead to artefacts due to
493 expression of antisense transcripts, a caveat that needs to be kept in mind for previous
494 studies of global expression analysis in bacteria.

495

496

497 **Materials and Methods**

498

499 *Bacterial strains, media and growth conditions*

500

501 For the M ϕ infection experiments, bovine alveolar M ϕ were cultivated in tissue culture
502 media R10, which consisted of RPMI (Invitrogen) media plus 2 mM glutamine, 10 % calf fetal
503 serum and 1 % amphotericin. Where used, antibiotics gentamycin and ampicillin were
504 added at concentrations of 50 and 100 μ g / ml, respectively. *M. bovis* field strains were pre-
505 grown in Middlebrook 7H9 broth supplemented with 10 % albumin-dextrose catalase (ADC,
506 Difco), 0.05 % Tween and 10 mM pyruvate. Cultures were harvested in mid-logarithmic
507 phase (OD₆₀₀ of 0.3-0.8), washed and then resuspended in RPMI containing 0.05 % Tween
508 80.

509

510 *Isolation of bovine alveolar macrophages and infection with mycobacteria*

511 The lungs of a 6-8 week old male Holstein-Friesian calf were removed and a whole lung
512 lavage procedure was performed to washout the alveolar M ϕ . Briefly, 4-5 x 500 ml aliquots
513 of Hanks' Balanced Sterile Salts solution (HBSS) were used to infuse the lungs via the
514 trachea, and the washings were pooled in a sterile beaker. The M ϕ cells contained in the
515 washes were pelleted by centrifugation at 500 x g for 10 mins at 4 °C, washed and then
516 resuspended in R10 growth media supplemented with antibiotics (R10+) to a concentration
517 of 1-2 x 10⁷ / ml. Approximately 0.5-1.5 x 10⁹ M ϕ were isolated per calf lung.

518 Vented 225 cm² tissue culture flasks containing R10+ media were seeded with 3-4 x
519 10⁷ alveolar M ϕ and placed in a humidified 37°C incubator containing 5% CO₂. Typically, 2-4
520 flasks were used per strain and time point. After 24 hrs, the growth media was decanted to
521 remove non-adherent cells and then replaced with fresh R10+ media. After a further 24 hrs,
522 the growth media was discarded and the M ϕ monolayer was washed with RPMI to remove
523 traces of the antibiotic containing growth media. The monolayer was then covered with R10
524 media without antibiotics (R10-) and then infected with mid-logarithmic phase grown

525 mycobacteria using an MOI of approximately 10:1 (bacilli: M ϕ). The AlvM ϕ were incubated
526 with mycobacteria for 4 hrs, after which the cell monolayers were washed with RPMI and
527 then either processed for RNA extraction (4 hr time point) or incubated in fresh R10+ media
528 for a further 20 hrs before being processed for RNA extraction (24 hr time point).

529

530 *Extraction and amplification of mycobacterial RNA from infected macrophages*

531 M ϕ cell monolayers were lysed using a guanidinium thiocyanate (GTC) containing solution.
532 The lysed M ϕ 's were vortexed and passed twice through a 21G blunt ended needle to sheer
533 host genomic DNA and thereby reduce the viscosity of the solution. Mycobacterial cells
534 were then pelleted by centrifugation at 4600 rpm for 20 mins at room temperature and
535 washed with GTC solution to remove host genomic DNA. Cells were then resuspended in
536 Trizol and RNA was extracted using the protocol outlined in Bacon et al [30]. The amount of
537 purified DNase-treated RNA recovered was of the order 100-500 ng per time point. RNA was
538 amplified using the 'MessageAmp II-Bacteria RNA Amplification Kit' (Ambion) according to
539 the manufacturers' instructions. Using an input of 100 ng of unamplified RNA, 20-100 μ g of
540 amplified RNA was recovered.

541

542 *Amplicon microarray analysis*

543 For the in vitro growth experiments, three independent experiments (biological replicates)
544 were carried out, and for each strain in each experiment two microarrays (technical
545 replicates) were performed. Thus, for each strain 6 microarrays were performed. Three
546 independent AlvM ϕ infection experiments were carried out and for each experiment two
547 microarrays were performed for each of the control RPMI samples, and the 4- and 24 hrs
548 post-infection samples. Cy5 and Cy3 fluorescently-labelled probes were synthesised from
549 RNA and genomic DNA, respectively, and hybridised to whole genome *M. bovis* / *M.*
550 *tuberculosis* microarrays. The array design is available in B μ G@Sbase (accession number A-
551 BUGS-31;
552 <http://bugs.sgul.ac.uk/A-BUGS-31>) and also ArrayExpress (accession number

553 A-BUGS-31). Details of probe synthesis, hybridization conditions and manufacture of the
554 microarray can be found in Golby et al 2008. Microarrays were scanned using an Affymetrix
555 428 Microarray scanner and scanned images were quantified using BlueFuse for Microarrays
556 v3.2 software (BlueGnome). See Golby et al for further details.

557 Normalisation was performed by dividing the log ratio of the Cy5 to Cy3 signal for
558 every spot by the median of the log ratios for all spots, except control spots. A median
559 absolute (MAD) scale transformation was applied to the normalised data from the pas an
560 additional normalisation step. For every microarray, duplicate spots were averaged, and
561 then the average expression of every gene across all technical replicate microarrays was
562 calculated. Averages of the three biological replicates were used to compare gene
563 expression between strains. For each gene, a moderate t-test was applied and those genes
564 with a P- value less than 0.05 were selected. From this gene list, those genes whose average
565 expression differed by more than 2.5-fold between strains were selected. Fully annotated
566 microarray data have been deposited in BμG@Sbase (accession number: E-BUGS-150;
567 http://bugs.sgul.ac.uk/E_BUGS-150) and also ArrayExpress (accession number: E-BUGS-
568 150).

569

570 *Oligonucleotide microarray analysis*

571 Experiments were performed in a similar way to that described for the in vitro amplicon
572 array experiments, except that the RNA was purified using the mirVana miRNA Isolation kit
573 (Ambion), which is designed to capture small (>20 nt) RNA species. RNA and genomic DNA
574 were directly labelled with Cy5- and Cy3, respectively, using the ULS microRNA labelling kit
575 (Kreatech), according to the manufacturer's instructions. Purified Cy5- and Cy3-labelled
576 probes were co-purified and applied to an Agilent 40K custom made tiled (10 nt overlap) 60-
577 mer oligonucleotide microarray designed to the genomic sequence of *M. bovis* 2122/97/97.
578 The array design is available in BμG@Sbase (accession number A-BUGS-52;
579 <http://bugs.sgul.ac.uk/A-BUGS-52>) and also ArrayExpress (accession number A-BUGS-52).
580 Microarrays were hybridised at 65°C for 18 hrs and then washed in a solution containing ...
581 at room temperature for 1 minute. The slides were then washed in a second solution
582 containing ... at 37°C for 1 minute, dried and then scanned at 2 mm using an Agilent DNA

583 microarray scanner.

584 Tiling array data was analysed using the Limma package of R/Bioconductor [31] . The
585 signal median was quantile normalised between arrays followed by a LOESS normalisation
586 within arrays. Differential expression analysis was performed by pairwise comparison using
587 linear models and empirical Bayes methods, and P values adjusted using the Benjamini and
588 Hochberg's method to control for multiple testing. Fully annotated microarray data have
589 been deposited in BμG@Sbase (accession number: E-BUGS-150; [http://](http://bugs.sgul.ac.uk/E_BUGS-150)
590 bugs.sgul.ac.uk/E_BUGS-150) and also ArrayExpress (accession number: E-BUGS-150).

591

592 *Whole genome sequencing of M. bovis field isolates*

593 Whole genome paired end sequencing was performed using an Illumina Genome Analyser II
594 at the Wellcome Trust Sanger Institute (Hinxton, Cambridge). Raw sequence data was
595 uploaded to the European Nucleotide Archive (ENA) and can be downloaded at
596 <http://www.ebi.ac.uk/ena/data/view/ERX006616-ERX06617,ERX012284-ERX012286>. Raw
597 sequence data was processed to remove adapter sequences and low quality reads. Filtered
598 reads were aligned to the *M. bovis* reference strain 2122/97 [9] with the MAQ program
599 (<http://maq.sourceforge.net/maq-man.shtml>) using default settings to identify SNPs. For all
600 sequenced strains, more than 99.9% of the genomes were covered by reads.

601

602 *Real time RT-PCR*

603 Quantitative real-time syber green based PCR (qRT-PCR) experiments were performed using
604 a RotorGene 3000 (Corbett research) as described by Golby et al. [21]. Fold changes were
605 calculated using relative standard curve method and pcr controls included no template and
606 no reverse transcriptase. Primer pair sequences are given in Supplementary Table S3,
607 available with the online version of this paper.

608

609 *Construction of the Rv1749c-Rv1750c and nirBD overexpressing plasmids*

610 The Rv1749c-Rv1750c overexpressing plasmids pPG106 and pPG107 were constructed by
611 PCR amplification of a 907bp fragment encompassing *Rv1749c-Rv1750c-Rv1751c* using
612 primers *tox_f* and *tox_r*. For pPG106, the fragment was amplified using 2122/97/97
613 genomic DNA as a template, while pPG107 was amplified using 1307/01 gDNA. Both
614 fragments were digested with *SpeI* and cloned into the *SpeI* cut mycobacterial attP-
615 integrating shuttle vector pKINT (a gift from Douglas Young, Imperial College, London).
616 Plasmids pPG108 and pPG109, which contain a 3.5 kb *hsp'-nirB-nirD-cobU* fragment was
617 constructed in several steps. Firstly, two 1.6 kb *hsp'-nirB'* PCR fragments were PCR amplified
618 separately using primers *nirB1_f* and *nirB1_r* and 2122/97 and 1307/01 as genomic
619 templates. Similarly, two 1.8kb '*nirB-nirD-cobU* fragments were amplified using primers
620 *nirB2_f* and *nirB2_r* and genomic DNAs 2122/97 and 1307/01. The two PCR products were
621 digested with *SpeI* and *BamHI* and then co-ligated into pKINT. Details concerning the
622 nucleotide sequences of the pcr primer pairs are given in the Supplementary Table S3.

623

624 *5'-RLM-RACE PCR*

625 Transcriptional start site mapping of was determined using the First Choice RNA ligase-
626 mediated rapid amplification of cDNA ends (RLM-RACE) kit (Ambion) as per manufacturers
627 instructions. Briefly, 10 ug of total RNA was treated with calf intestinal phosphatase (CIP)
628 and tobacco acid pyrophosphatase (TAP) before ligation of an RNA Adapter oligonucleotide
629 to the 5' ends of the mRNA transcripts. A random-primed reverse transcription reaction was
630 carried out to generate cDNA and then nested PCR reactions were performed on the cDNA
631 using combinations of adapter and gene specific primers. Details concerning the nucleotide
632 sequences of the 5' outer and inner adapter sequences as well as 3' outer and inner gene
633 specific primers are given in the Supplementary Table S3. PCR products generated using the
634 5' inner adapter and 3' inner gene specific primers were sequenced by Sanger sequencing
635 using the 3' inner gene specific primer.

636

637

638 Figure Legends

639 **Fig. 1.** SNPs present in the genome sequences of *M. bovis* field strains 1121/01, 2451/01 and
640 1307/01 compared with 2122/97. (a) A Venn diagram to illustrate specific and shared SNPs
641 across the strains. Numbers in parenthesis after the strain ID indicate the total number of
642 SNPs compared to 2122/97. (b) Distribution of each SNP category per strain. Vertical bars
643 show the percentage of non-synonymous (red), synonymous (blue) and intergenic SNPs
644 (black) shown by each strain.

645 **Fig. 2.** Venn diagram to show genes differentially expressed in vivo and ex vivo M ϕ .

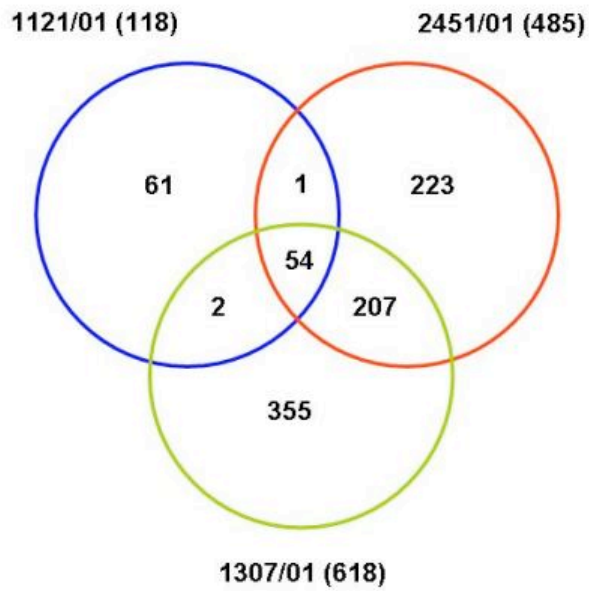
646 **Fig. 3.** Confirmation of amplicon microarray results with real time RT-PCR. The fold changes
647 in gene expression for (a) Mb1750c, (b) *nirB*, (c) *echA21* and (d) *Mb1914c* measured by
648 microarray (open bars) in each of the four strains were compared to that measured by real
649 time RT-PCR (closed bars).

650 **Fig. 4.** Expressions and schematic representation of genomic locations of selected cis-
651 encoded antisense sRNAs identified using a tiled oligonucleotide microarray. Three asRNAs
652 (open arrows) are (a) T6, (b) T14 and (c) T25. For each asRNA, a histogram plots the fold
653 changes for each of the oligonucleotide probes that detected the asRNA, and for each probe
654 the binding position relative to the 2122/97 genome is indicated. Closed and open arrows
655 indicate lengths and direction of transcription of genes and asRNAs, respectively.

656 **Fig. 5.** Promoters of anti sense RNAs. A. Promoters of the asRNAs as_mb1618c, as_1914c
657 and as_echA21. -10 and -35 elements are indicated in bold and italics. Transcriptional start
658 sites are indicated by large font G characters, while SNP residue that leads to the expression
659 of the asRNA is indicated by a large font red T residue. The consensus sequence for group A
660 mycobacterial promoters is indicated. Numerical subscripts indicate the percentage of the
661 total number of promoters for which a transcriptional start site has been experimentally
662 determined that show the indicated residue. B. Promoters of the differentially expressed
663 asRNAs as_ino1 and as_narH in *M. tuberculosis*. -10 and -35 elements are indicated in bold
664 italics. The red residue indicates SNP responsible for differential expression.

Fig.1. SNPs present in the genome sequences of *M. bovis* field strains 1121/01, 2451/01 and 1307/01 compared with 2122/97.

(a)



(b)

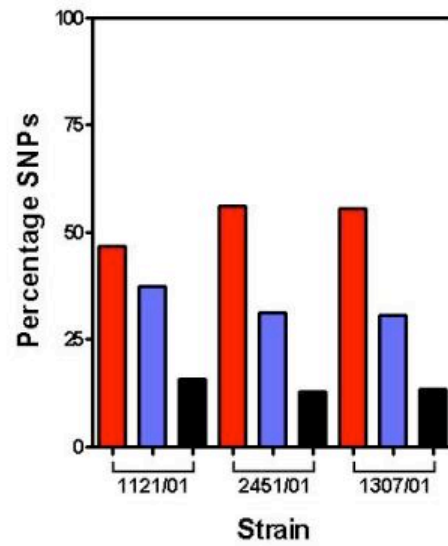


Fig. 2. Venn diagram to show genes differentially expressed in vivo and ex vivo M ϕ .

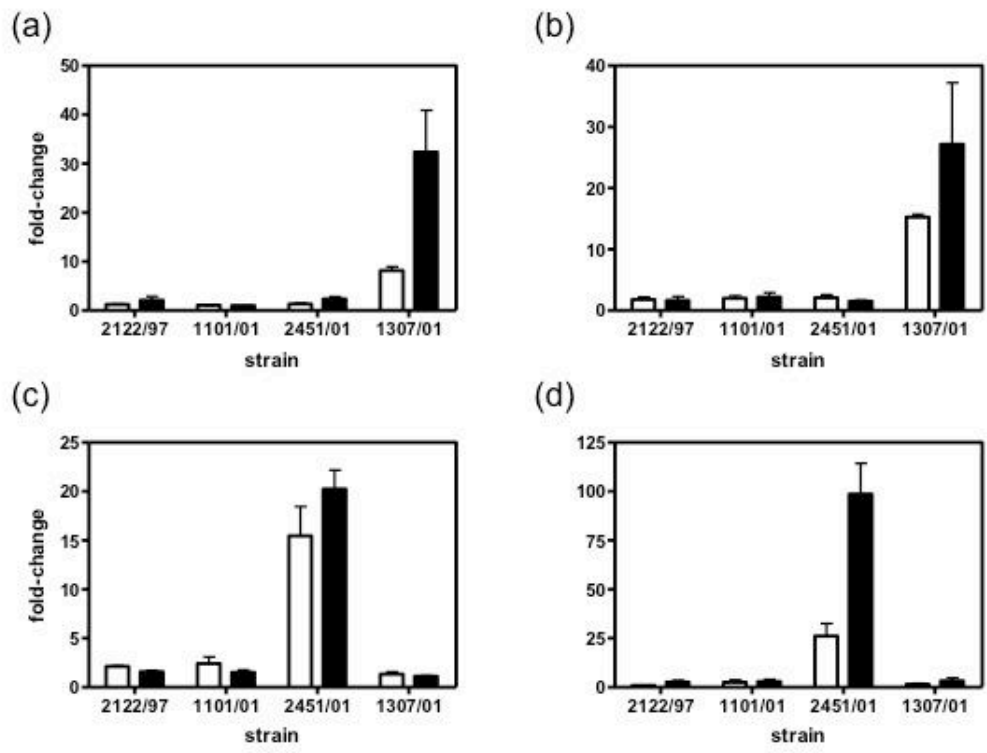


Fig. 3. Confirmation of amplicon microarray results with real time RT-PCR.

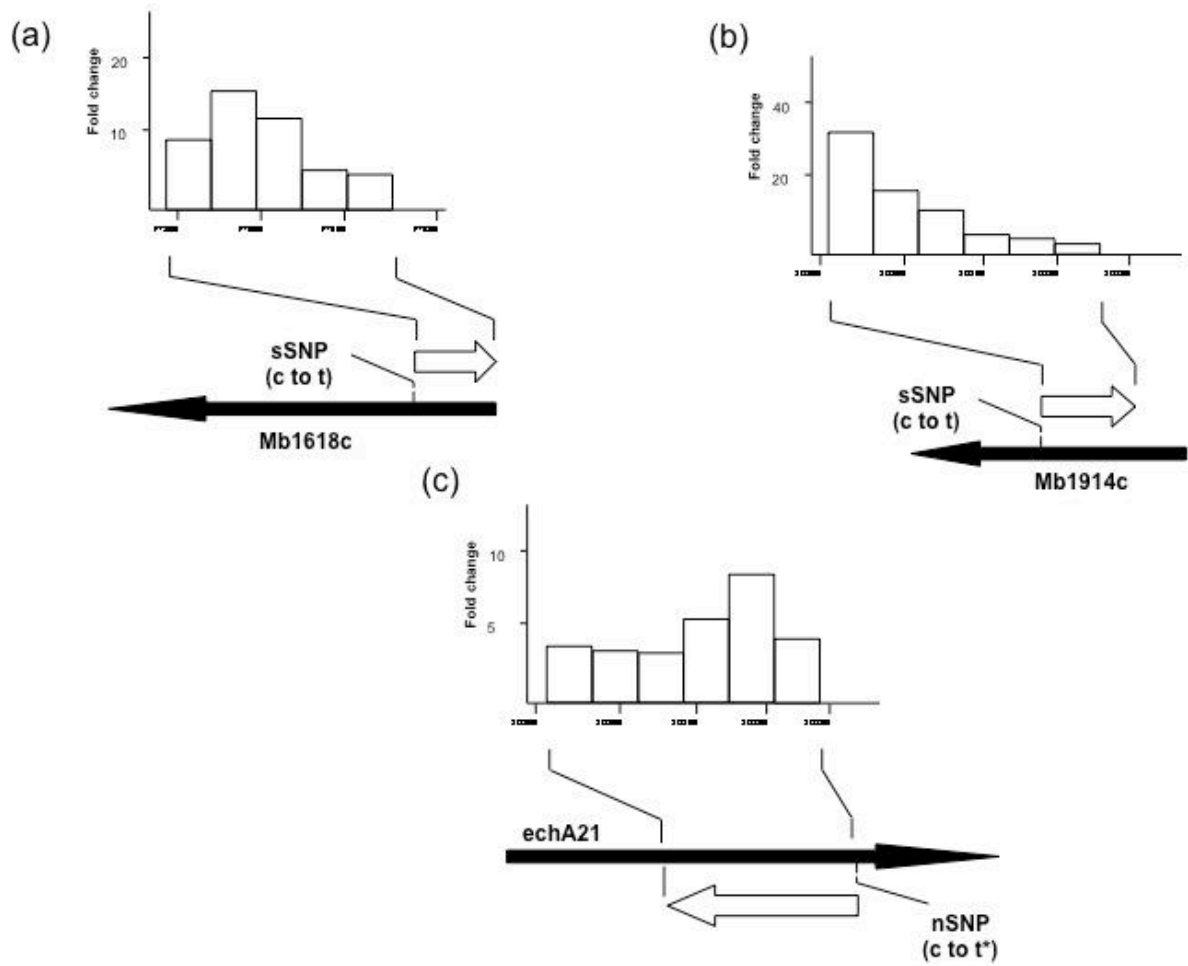


Fig. 4. Expressions and schematic representation of genomic locations of selected cis-encoded antisense sRNAs identified using a tiled oligonucleotide microarray.

