

Deciphering the biology of *Mycobacterium tuberculosis* from the complete genome sequence

S. T. Cole*, R. Brosch*, J. Parkhill, T. Garnier*, C. Churcher, D. Harris, S. V. Gordon*, K. Eiglmeier*, S. Gas*, C. E. Barry III†, F. Tekaia‡, K. Badcock, D. Basham, D. Brown, T. Chillingworth, R. Connor, R. Davies, K. Devlin, T. Feltwell, S. Gentles, N. Hamlin, S. Holroyd, T. Hornsby, K. Jagels, A. Krogh§, J. McLean, S. Moule, L. Murphy, K. Oliver, J. Osborne, M. A. Quail, M.-A. Rajandream, J. Rogers, S. Rutter, K. Seeger, J. Skelton, R. Squares, S. Squares, J. E. Sulston, K. Taylor, S. Whitehead & B. G. Barrell

Sanger Centre, Wellcome Trust Genome Campus, Hinxton CB10 1SA, UK

* Unité de Génétique Moléculaire Bactérienne, and ‡ Unité de Génétique Moléculaire des Levures, Institut Pasteur, 28 rue du Docteur Roux, 75724 Paris Cedex 15, France

† Tuberculosis Research Unit, Laboratory of Intracellular Parasites, Rocky Mountain Laboratories, National Institute of Allergy and Infectious Diseases, National Institutes of Health, Hamilton, Montana 59840, USA

§ Center for Biological Sequence Analysis, Technical University of Denmark, Lyngby, Denmark

Countless millions of people have died from tuberculosis, a chronic infectious disease caused by the tubercle bacillus. The complete genome sequence of the best-characterized strain of *Mycobacterium tuberculosis*, H37Rv, has been determined and analysed in order to improve our understanding of the biology of this slow-growing pathogen and to help the conception of new prophylactic and therapeutic interventions. The genome comprises 4,411,529 base pairs, contains around 4,000 genes, and has a very high guanine + cytosine content that is reflected in the biased amino-acid content of the proteins. *M. tuberculosis* differs radically from other bacteria in that a very large portion of its coding capacity is devoted to the production of enzymes involved in lipogenesis and lipolysis, and to two new families of glycine-rich proteins with a repetitive structure that may represent a source of antigenic variation.

Despite the availability of effective short-course chemotherapy (DOTS) and the Bacille Calmette-Guérin (BCG) vaccine, the tubercle bacillus continues to claim more lives than any other single infectious agent¹. Recent years have seen increased incidence of tuberculosis in both developing and industrialized countries, the widespread emergence of drug-resistant strains and a deadly synergy with the human immunodeficiency virus (HIV). In 1993, the gravity of the situation led the World Health Organisation (WHO) to declare tuberculosis a global emergency in an attempt to heighten public and political awareness. Radical measures are needed now to prevent the grim predictions of the WHO becoming reality. The combination of genomics and bioinformatics has the potential to generate the information and knowledge that will enable the conception and development of new therapies and interventions needed to treat this airborne disease and to elucidate the unusual biology of its aetiological agent, *Mycobacterium tuberculosis*.

The characteristic features of the tubercle bacillus include its slow growth, dormancy, complex cell envelope, intracellular pathogenesis and genetic homogeneity². The generation time of *M. tuberculosis*, in synthetic medium or infected animals, is typically ~24 hours. This contributes to the chronic nature of the disease, imposes lengthy treatment regimens and represents a formidable obstacle for researchers. The state of dormancy in which the bacillus remains quiescent within infected tissue may reflect metabolic shutdown resulting from the action of a cell-mediated immune response that can contain but not eradicate the infection. As immunity wanes, through ageing or immune suppression, the dormant bacteria reactivate, causing an outbreak of disease often many decades after the initial infection³. The molecular basis of dormancy and reactivation remains obscure but is expected to be genetically programmed and to involve intracellular signalling pathways.

The cell envelope of *M. tuberculosis*, a Gram-positive bacterium with a G + C-rich genome, contains an additional layer beyond the peptidoglycan that is exceptionally rich in unusual lipids, glycolipids and polysaccharides^{4,5}.

Novel biosynthetic pathways generate cell-wall components such as mycolic acids, mycocerosic acid, phenolthiocerol, lipoarabinomannan and arabinogalactan, and several of these may contribute to mycobacterial longevity, trigger inflammatory host reactions and act in pathogenesis. Little is known about the mechanisms involved in life within the macrophage, or the extent and nature of the virulence factors produced by the bacillus and their contribution to disease.

It is thought that the progenitor of the *M. tuberculosis* complex, comprising *M. tuberculosis*, *M. bovis*, *M. bovis* BCG, *M. africanum* and *M. microti*, arose from a soil bacterium and that the human bacillus may have been derived from the bovine form following the domestication of cattle. The complex lacks interstrain genetic diversity, and nucleotide changes are very rare⁶. This is important in terms of immunity and vaccine development as most of the proteins will be identical in all strains and therefore antigenic drift will be restricted. On the basis of the systematic sequence analysis of 26 loci in a large number of independent isolates⁶, it was concluded that the genome of *M. tuberculosis* is either unusually inert or that the organism is relatively young in evolutionary terms.

Since its isolation in 1905, the H37Rv strain of *M. tuberculosis* has found extensive, worldwide application in biomedical research because it has retained full virulence in animal models of tuberculosis, unlike some clinical isolates; it is also susceptible to drugs and amenable to genetic manipulation. An integrated map of the 4.4 megabase (Mb) circular chromosome of this slow-growing pathogen had been established previously and ordered libraries of cosmids and bacterial artificial chromosomes (BACs) were available^{7,8}.

Organization and sequence of the genome

Sequence analysis. To obtain the contiguous genome sequence, a combined approach was used that involved the systematic sequence analysis of selected large-insert clones (cosmids and BACs) as well as

random small-insert clones from a whole-genome shotgun library. This culminated in a composite sequence of 4,411,529 base pairs (bp) (Figs 1, 2), with a G + C content of 65.6%. This represents the second-largest bacterial genome sequence currently available (after that of *Escherichia coli*)⁹. The initiation codon for the *dnaA* gene, a hallmark for the origin of replication, *oriC*, was chosen as the start point for numbering. The genome is rich in repetitive DNA, particularly insertion sequences, and in new multigene families and duplicated housekeeping genes. The G + C content is relatively constant throughout the genome (Fig. 1) indicating that horizontally transferred pathogenicity islands of atypical base composition are probably absent. Several regions showing higher than average G + C content (Fig. 1) were detected; these correspond to sequences belonging to a large gene family that includes the polymorphic G + C-rich sequences (PGRSs).

Genes for stable RNA. Fifty genes coding for functional RNA molecules were found. These molecules were the three species produced by the unique ribosomal RNA operon, the 10Sa RNA involved in degradation of proteins encoded by abnormal messenger RNA, the RNA component of RNase P, and 45 transfer RNAs. No 4.5S RNA could be detected. The *rrn* operon is situated unusually as it occurs about 1,500 kilobases (kb) from the putative *oriC*; most eubacteria have one or more *rrn* operons near to *oriC* to exploit the gene-dosage effect obtained during replication¹⁰. This arrangement may be related to the slow growth of *M. tuberculosis*. The genes encoding tRNAs that recognize 43 of the 61 possible sense codons were distributed throughout the genome and, with one

exception, none of these uses A in the first position of the anticodon, indicating that extensive wobble occurs during translation. This is consistent with the high G + C content of the genome and the consequent bias in codon usage. Three genes encoding tRNAs for methionine were found; one of these genes (*metV*) is situated in a region that may correspond to the terminus of replication (Figs 1, 2). As *metV* is linked to defective genes for integrase and excisionase, perhaps it was once part of a phage or similar mobile genetic element.

Insertion sequences and prophages. Sixteen copies of the promiscuous insertion sequence IS6110 and six copies of the more stable element IS1081 reside within the genome of H37Rv⁸. One copy of IS1081 is truncated. Scrutiny of the genomic sequence led to the identification of a further 32 different insertion sequence elements, most of which have not been described previously, and of the 13E12 family of repetitive sequences which exhibit some of the characteristics of mobile genetic elements (Fig. 1). The newly discovered insertion sequences belong mainly to the IS3 and IS256 families, although six of them define a new group. There is extensive similarity between IS1561 and IS1552 with insertion sequence elements found in *Nocardia* and *Rhodococcus* spp., suggesting that they may be widely disseminated among the actinomycetes.

Most of the insertion sequences in *M. tuberculosis* H37Rv appear to have inserted in intergenic or non-coding regions, often near tRNA genes (Fig. 1). Many are clustered, suggesting the existence of insertional hot-spots that prevent genes from being inactivated, as has been described for *Rhizobium*¹¹. The chromosomal distribution of the insertion sequences is informative as there appears to have been a selection against insertions in the quadrant encompassing *oriC* and an overrepresentation in the direct repeat region that contains the prototype IS6110. This bias was also observed experimentally in a transposon mutagenesis study¹².

At least two prophages have been detected in the genome sequence and their presence may explain why *M. tuberculosis* shows persistent low-level lysis in culture. Prophages phiRv1 and phiRv2 are both ~10 kb in length and are similarly organized, and some of their gene products show marked similarity to those encoded by certain bacteriophages from *Streptomyces* and saprophytic mycobacteria. The site of insertion of phiRv1 is intriguing as it corresponds to part of a repetitive sequence of the 13E12 family that itself appears to have integrated into the biotin operon. Some strains of *M. tuberculosis* have been described as requiring biotin as a growth supplement, indicating either that phiRv1 has a polar effect on expression of the distal *bio* genes or that aberrant excision, leading to mutation, may occur. During the serial attenuation of *M. bovis* that led to the vaccine strain *M. bovis* BCG, the phiRv1 prophage was lost¹³. In a systematic study of the genomic diversity of prophages and insertion sequences (S.V.G. *et al.*, manuscript in preparation), only IS1532 exhibited significant variability, indicating that most of the prophages and insertion sequences are currently stable. However, from these combined observations, one can conclude that horizontal transfer of genetic material into the free-living ancestor of the *M. tuberculosis* complex probably occurred in nature before the tubercle bacillus adopted its specialized intracellular niche.

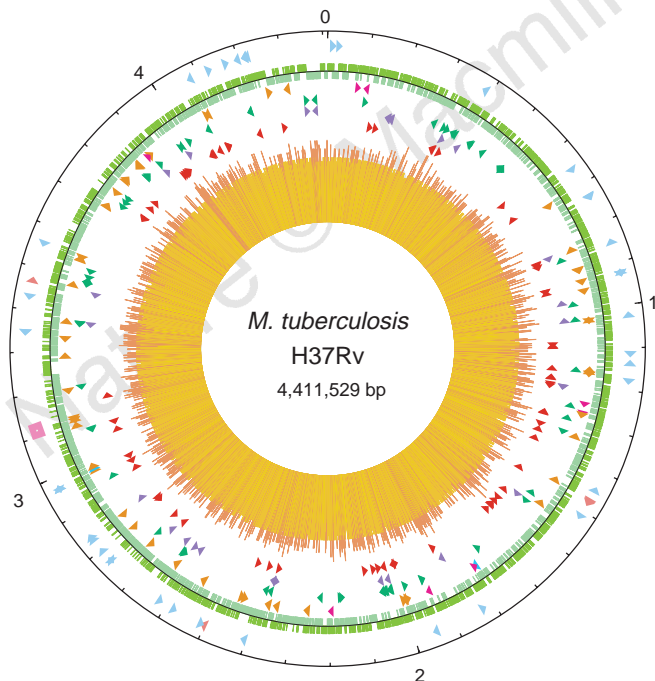


Figure 1 Circular map of the chromosome of *M. tuberculosis* H37Rv. The outer circle shows the scale in Mb, with 0 representing the origin of replication. The first ring from the exterior denotes the positions of stable RNA genes (tRNAs are blue, others are pink) and the direct repeat region (pink cube); the second ring inward shows the coding sequence by strand (clockwise, dark green; anticlockwise, light green); the third ring depicts repetitive DNA (insertion sequences, orange; 13E12 REP family, dark pink; prophage, blue); the fourth ring shows the positions of the PPE family members (green); the fifth ring shows the PE family members (purple, excluding PGRS); and the sixth ring shows the positions of the PGRS sequences (dark red). The histogram (centre) represents G + C content, with <65% G + C in yellow, and >65% G + C in red. The figure was generated with software from DNASTAR.

Figure 2 Linear map of the chromosome of *M. tuberculosis* H37Rv showing the position and orientation of known genes and coding sequences (CDS). We used the following functional categories (adapted from ref. 20): lipid metabolism (black); intermediary metabolism and respiration (yellow); information pathways (pink); regulatory proteins (sky blue); conserved hypothetical proteins (orange); proteins of unknown function (light green); insertion sequences and phage-related functions (blue); stable RNAs (purple); cell wall and cell processes (dark green); PE and PPE protein families (magenta); virulence, detoxification and adaptation (white). For additional information about gene functions, refer to <http://www.sanger.ac.uk>.

Genes encoding proteins. 3,924 open reading frames were identified in the genome (see Methods), accounting for ~91% of the potential coding capacity (Figs 1, 2). A few of these genes appear to have in-frame stop codons or frameshift mutations (irrespective of the source of the DNA sequenced) and may either use frameshifting during translation or correspond to pseudogenes. Consistent with the high G + C content of the genome, GTG initiation codons (35%) are used more frequently than in *Bacillus subtilis* (9%) and *E. coli* (14%), although ATG (61%) is the most common translational start. There are a few examples of atypical initiation codons, the most notable being the ATC used by *infC*, which begins with ATT in both *B. subtilis* and *E. coli*^{9,14}. There is a slight bias in the orientation of the genes (Fig. 1) with respect to the direction of replication as ~59% are transcribed with the same polarity as replication, compared with 75% in *B. subtilis*. In other bacteria, genes transcribed in the same direction as the replication forks are believed to be expressed more efficiently^{9,14}. Again, the more even distribution in gene polarity seen in *M. tuberculosis* may reflect the slow growth and infrequent replication cycles. Three genes (*dnaB*, *recA* and *Rv1461*) have been invaded by sequences encoding inteins (protein introns) and in all three cases their counterparts in *M. leprae* also contain inteins, but at different sites¹⁵ (S.T.C. *et al.*, unpublished observations).

Protein function, composition and duplication. By using various database comparisons, we attributed precise functions to ~40% of the predicted proteins and found some information or similarity for another 44%. The remaining 16% resembled no known proteins and may account for specific mycobacterial functions. Examination of the amino-acid composition of the *M. tuberculosis* proteome by correspondence analysis¹⁶, and comparison with that of other microorganisms whose genome sequences are available, revealed a statistically significant preference for the amino acids Ala, Gly, Pro, Arg and Trp, which are all encoded by G + C-rich codons, and a comparative reduction in the use of amino acids encoded by A + T-rich codons such as Asn, Ile, Lys, Phe and Tyr (Fig. 3). This approach also identified two groups of proteins rich in Asn or Gly that belong to new families, PE and PPE (see below). The fraction of the proteome that has arisen through gene duplication is similar to that seen in *E. coli* or *B. subtilis* (~51%; refs 9, 14), except that the level of sequence conservation is considerably higher, indicating that there may be extensive redundancy or differential production of the corresponding polypeptides. The apparent lack of divergence following gene duplication is consistent with the hypothesis that *M. tuberculosis* is of recent descent⁶.

General metabolism, regulation and drug resistance

Metabolic pathways. From the genome sequence, it is clear that the tubercle bacillus has the potential to synthesize all the essential amino acids, vitamins and enzyme co-factors, although some of the pathways involved may differ from those found in other bacteria. *M. tuberculosis* can metabolize a variety of carbohydrates, hydrocarbons, alcohols, ketones and carboxylic acids^{2,17}. It is apparent from genome inspection that, in addition to many functions involved in lipid metabolism, the enzymes necessary for glycolysis, the pentose phosphate pathway, and the tricarboxylic acid and glyoxylate cycles are all present. A large number (~200) of oxidoreductases, oxygenases and dehydrogenases is predicted, as well as many oxygenases containing cytochrome P450, that are similar to fungal proteins involved in sterol degradation. Under aerobic growth conditions, ATP will be generated by oxidative phosphorylation from electron transport chains involving a ubiquinone cytochrome *b* reductase complex and cytochrome *c* oxidase. Components of several anaerobic phosphorylative electron transport chains are also present, including genes for nitrate reductase (*narGHJI*), fumarate reductase (*frdABCD*) and possibly nitrite reductase (*nirBD*), as well as a new reductase (*narX*) that results from a rearrangement of a homologue of the *narGHJI* operon. Two genes encoding haemoglobin-like

proteins, which may protect against oxidative stress or be involved in oxygen capture, were found. The ability of the bacillus to adapt its metabolism to environmental change is significant as it not only has to compete with the lung for oxygen but must also adapt to the microaerophilic/anaerobic environment at the heart of the burgeoning granuloma.

Regulation and signal transduction. Given the complexity of the environmental and metabolic choices facing *M. tuberculosis*, an extensive regulatory repertoire was expected. Thirteen putative sigma factors govern gene expression at the level of transcription initiation, and more than 100 regulatory proteins are predicted (Table 1). Unlike *B. subtilis* and *E. coli*, in which there are >30 copies of different two-component regulatory systems¹⁴, *M. tuberculosis* has only 11 complete pairs of sensor histidine kinases and response regulators, and a few isolated kinase and regulatory genes. This relative paucity in environmental signal transduction pathways is probably offset by the presence of a family of eukaryotic-like serine/threonine protein kinases (STPKs), which function as part of a phosphorelay system¹⁸. The STPKs probably have two domains: the well-conserved kinase domain at the amino terminus is predicted to be connected by a transmembrane segment to the carboxy-terminal region that may respond to specific stimuli. Several of the predicted envelope lipoproteins, such as that encoded by *lppR* (Rv2403), show extensive similarity to this putative receptor domain of STPKs, suggesting possible interplay. The STPKs probably function in signal transduction pathways and may govern important cellular decisions such as dormancy and cell division, and although their partners are unknown, candidate genes for phosphoprotein phosphatases have been identified.

Drug resistance. *M. tuberculosis* is naturally resistant to many antibiotics, making treatment difficult¹⁹. This resistance is due mainly to the highly hydrophobic cell envelope acting as a permeability barrier⁴, but many potential resistance determinants are also encoded in the genome. These include hydrolytic or drug-modifying enzymes such as β -lactamases and aminoglycoside acetyl transferases, and many potential drug-efflux systems, such as 14 members of the major facilitator family and numerous ABC transporters. Knowledge of these putative resistance mechanisms will promote better use of existing drugs and facilitate the conception of new therapies.

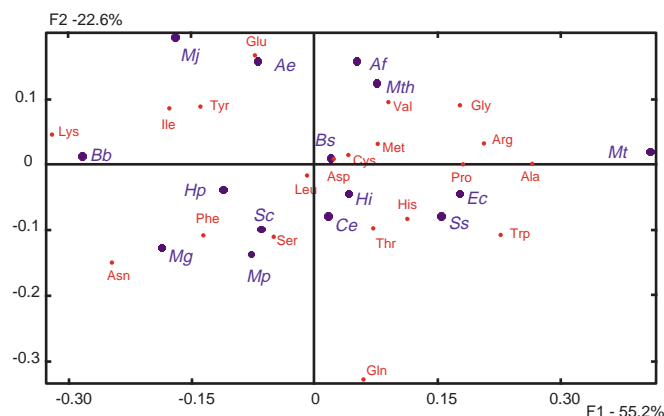


Figure 3 Correspondence analysis of the proteomes from extensively sequenced organisms as a function of amino-acid composition. Note the extreme position of *M. tuberculosis* and the shift in amino-acid preference reflecting increasing G + C content from left to right. Abbreviations used: Ae, *Aquifex aeolicus*; Af, *Archaeoglobus fulgidis*; Bb, *Borrelia burgdorferi*; Bs, *B. subtilis*; Ce, *Caenorhabditis elegans*; Ec, *E. coli*; Hi, *Haemophilus influenzae*; Hp, *Helicobacter pylori*; Mg, *Mycoplasma genitalium*; Mj, *Methanococcus jannaschi*; Mp, *Mycoplasma pneumoniae*; Mt, *M. tuberculosis*; Mth, *Methanobacterium thermoautotrophicum*; Sc, *Saccharomyces cerevisiae*; Ss, *Synechocystis* sp. strain PCC6803. F1 and F2, first and second factorial axes¹⁶.

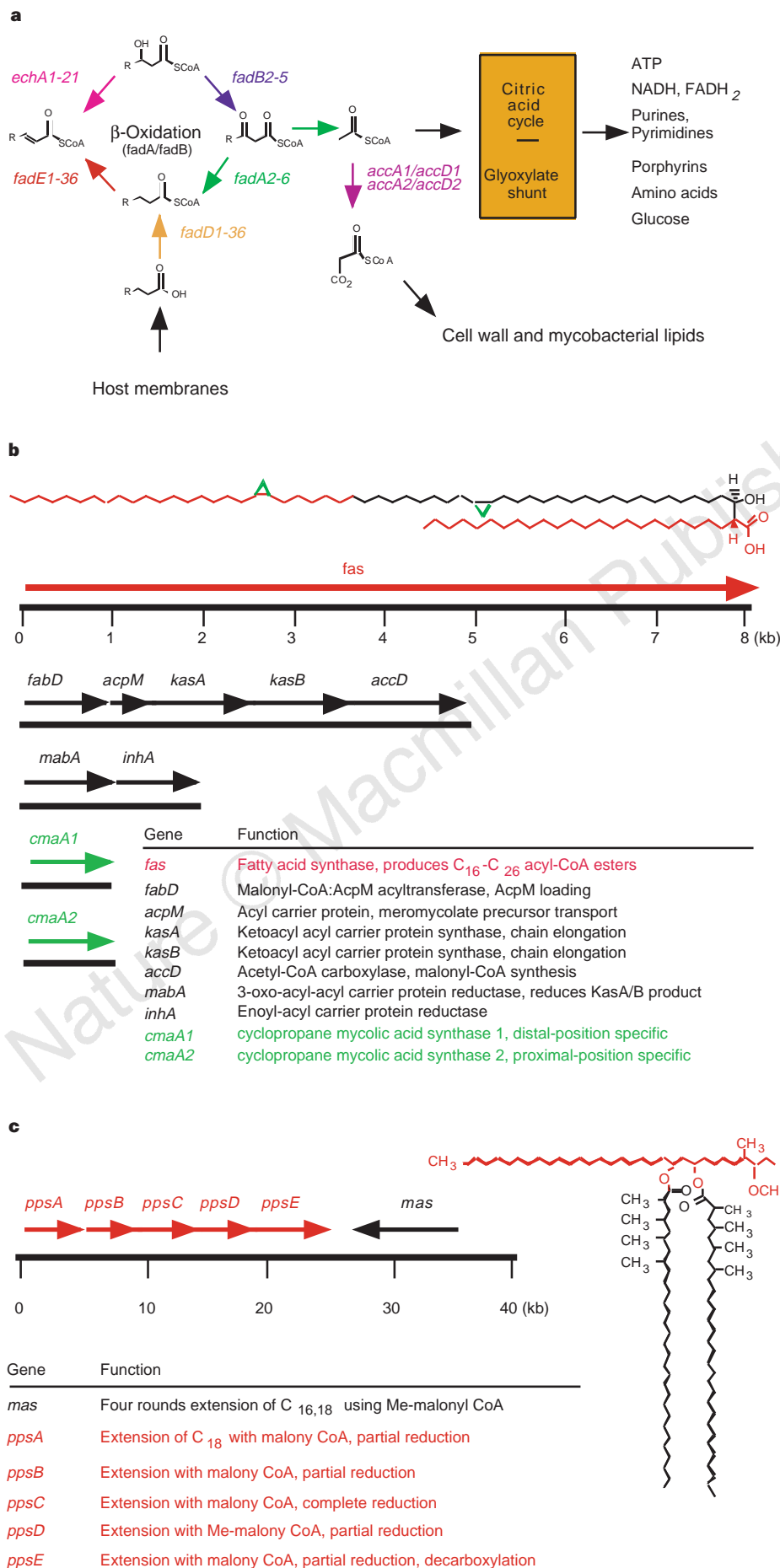


Figure 4 Lipid metabolism. **a**, Degradation of host-cell lipids is vital in the intracellular life of *M. tuberculosis*. Host-cell membranes provide precursors for many metabolic processes, as well as potential precursors of mycobacterial cell-wall constituents, through the actions of a broad family of β -oxidative enzymes encoded by multiple copies in the genome. These enzymes produce acetyl CoA, which can be converted into many different metabolites and fuel for the bacteria through the actions of the enzymes of the citric acid cycle and the glyoxylate shunt of this cycle. **b**, The genes that synthesize mycolic acids, the dominant lipid component of the mycobacterial cell wall, include the type I fatty acid synthase (*fas*) and a unique type II system which relies on extension of a precursor bound to an acyl carrier protein to form full-length (~80-carbon) mycolic acids. The *cma* genes are responsible for cyclopropanation. **c**, The genes that produce phthiocerol dimycocerosate form a large operon and represent type I (*mas*) and type II (the *pps* operon) polyketide synthase systems. Functions are colour coordinated.

Lipid metabolism

Very few organisms produce such a diverse array of lipophilic molecules as *M. tuberculosis*. These molecules range from simple fatty acids such as palmitate and tuberculostearate, through isoprenoids, to very-long-chain, highly complex molecules such as mycolic acids and the phenolphthiocerol alcohols that esterify with mycocerosic acid to form the scaffold for attachment of the mycosides. Mycobacteria contain examples of every known lipid and polyketide biosynthetic system, including enzymes usually found in mammals and plants as well as the common bacterial systems. The biosynthetic capacity is overshadowed by the even more remarkable radiation of degradative, fatty acid oxidation systems and, in total, there are ~250 distinct enzymes involved in fatty acid metabolism in *M. tuberculosis* compared with only 50 in *E. coli*²⁰.

Fatty acid degradation. *In vivo*-grown mycobacteria have been suggested to be largely lipolytic, rather than lipogenic, because of the variety and quantity of lipids available within mammalian cells and the tubercle² (Fig. 4a). The abundance of genes encoding components of fatty acid oxidation systems found by our genomic approach supports this proposition, as there are 36 acyl-CoA synthases and a family of 36 related enzymes that could catalyse the first step in fatty acid degradation. There are 21 homologous enzymes belonging to the enoyl-CoA hydratase/isomerase superfamily of enzymes, which rehydrate the nascent product of the acyl-CoA dehydrogenase. The four enzymes that convert the 3-hydroxy fatty acid into a 3-keto fatty acid appear less numerous, mainly because they are difficult to distinguish from other members of the short-chain alcohol dehydrogenase family on the basis of primary sequence. The five enzymes that complete the cycle by thiolysis of the β -ketoester, the acetyl-CoA C-acetyltransferases, do indeed appear to be a more limited family. In addition to this extensive set of dissociated degradative enzymes, the genome also encodes the canonical FadA/FadB β -oxidation complex (Rv0859 and Rv0860). Accessory activities are present for the metabolism of odd-chain and multiply unsaturated fatty acids.

Fatty acid biosynthesis. At least two discrete types of enzyme system, fatty acid synthase (FAS) I and FAS II, are involved in fatty acid biosynthesis in mycobacteria (Fig. 4b). FAS I (Rv2524, *fas*) is a single polypeptide with multiple catalytic activities that generates several shorter CoA esters from acetyl-CoA primers⁵ and probably creates precursors for elongation by all of the other fatty acid and polyketide systems. FAS II consists of dissociable enzyme components which act on a substrate bound to an acyl-carrier protein (ACP). FAS II is incapable of *de novo* fatty acid synthesis but instead elongates palmitoyl-ACP to fatty acids ranging from 24 to 56 carbons in length^{17,21}. Several different components of FAS II may be targets for the important tuberculosis drug isoniazid, including the enoyl-ACP reductase *InhA*²², the ketoacyl-ACP synthase *KasA* and the ACP *AcpM*²¹. Analysis of the genome shows that there are only three potential ketoacyl synthases: *KasA* and *KasB* are highly related, and their genes cluster with *acpM*, whereas *KasC* is a more distant homologue of a ketoacyl synthase III system. The number of ketoacyl synthase and ACP genes indicates that there is a single FAS II system. Its genetic organization, with two clustered ketoacyl synthases, resembles that of type II aromatic polyketide biosynthetic gene clusters, such as those for actinorhodin, tetracycline and tetracenomycin in *Streptomyces* species²³. *InhA* seems to be the sole enoyl-ACP reductase and its gene is co-transcribed with a *fabG* homologue, which encodes 3-oxoacyl-ACP reductase. Both of these proteins are probably important in the biosynthesis of mycolic acids.

Fatty acids are synthesized from malonyl-CoA and precursors are generated by the enzymatic carboxylation of acetyl (or propionyl)-CoA by a biotin-dependent carboxylase (Fig. 4b). From study of the genome we predict that there are three complete carboxylase systems, each consisting of an α - and a β -subunit, as well as three β -subunits without an α -counterpart. As a group, all of the carboxylases seem to be more related to the mammalian homo-

logues than to the corresponding bacterial enzymes. Two of these carboxylase systems (*accA1*, *accD1* and *accA2*, *accD2*) are probably involved in degradation of odd-numbered fatty acids, as they are adjacent to genes for other known degradative enzymes. They may convert propionyl-CoA to succinyl-CoA, which can then be incorporated into the tricarboxylic acid cycle. The synthetic carboxylases (*accA3*, *accD3*, *accD4*, *accD5* and *accD6*) are more difficult to understand. The three extra β -subunits might direct carboxylation to the appropriate precursor or may simply increase the total amount of carboxylated precursor available if this step were rate-limiting.

Synthesis of the paraffinic backbone of fatty and mycolic acids in the cell is followed by extensive postsynthetic modifications and unsaturations, particularly in the case of the mycolic acids^{24,25}. Unsaturation is catalysed either by a FabA-like β -hydroxyacyl-ACP dehydrase, acting with a specific ketoacyl synthase, or by an aerobic terminal mixed function desaturase that uses both molecular oxygen and NADPH. Inspection of the genome revealed no obvious candidates for the FabA-like activity. However, three potential aerobic desaturases (encoded by *desA1*, *desA2* and *desA3*) were evident that show little similarity to related vertebrate or yeast enzymes (which act on CoA esters) but instead resemble plant desaturases (which use ACP esters). Consequently, the genomic data indicate that unsaturation of the meromycolate chain may occur while the acyl group is bound to AcpM.

Much of the subsequent structural diversity in mycolic acids is generated by a family of S-adenosyl-L-methionine-dependent enzymes, which use the unsaturated meromycolic acid as a substrate to generate *cis* and *trans* cyclopropanes and other mycolates. Six members of this family have been identified and characterized²⁵ and two clustered, convergently transcribed new genes are evident in the genome (*umaA1* and *umaA2*). From the functions of the known family members and the structures of mycolic acids in *M. tuberculosis*, it is tempting to speculate that these new enzymes may introduce the *trans* cyclopropanes into the meromycolate precursor. In addition to these two methyltransferases, there are two other unrelated lipid methyltransferases (*Ufa1* and *Ufa2*) that share homology with cyclopropane fatty acid synthase of *E. coli*²⁵. Although cyclopropanation seems to be a relatively common modification of mycolic acids, cyclopropanation of plasma-membrane constituents has not been described in mycobacteria. Tuberculostearic acid is produced by methylation of oleic acid, and may be synthesized by one of these two enzymes.

Condensation of the fully functionalized and preformed meromycolate chain with a 26-carbon α -branch generates full-length mycolic acids that must be transported to their final location for attachment to the cell-wall arabinogalactan. The transfer and subsequent transesterification is mediated by three well-known immunogenic proteins of the antigen 85 complex²⁶. The genome encodes a fourth member of this complex, antigen 85C' (*fbpC2*, Rv0129), which is highly related to antigen 85C. Further studies are needed to show whether the protein possesses mycolyltransferase activity and to clarify the reason behind the apparent redundancy.

Polyketide synthesis. Mycobacteria synthesize polyketides by several different mechanisms. A modular type I system, similar to that involved in erythromycin biosynthesis²³, is encoded by a very large operon, *ppsABCDE*, and functions in the production of phenolphthiocerol⁵. The absence of a second type I polyketide synthase suggests that the related lipids phthiocerol A and B, phthiodiolone A and phthiotriol may all be synthesized by the same system, either from alternative primers or by differential postsynthetic modification. It is physiologically significant that the *pps* gene cluster occurs immediately upstream of *mas*, which encodes the multifunctional enzyme mycocerosic acid synthase (MAS), as their products phthiocerol and mycocerosic acid esterify to form the very abundant cell-wall-associated molecule phthiocerol dimycocerosate (Fig. 4c).

Members of another large group of polyketide synthase enzymes are similar to MAS, which also generates the multiply methyl-branched fatty acid components of mycosides and phthiocerol dimycocerosate, abundant cell-wall-associated molecules⁵. Although some of these polyketide synthases may extend type I FAS CoA primers to produce other long-chain methyl-branched fatty acids such as mycolipenic, mycolipodienic and mycolipanic acids or the phthioceranic and hydroxyphthioceranic acids, or may even show functional overlap⁵, there are many more of these enzymes than there are known metabolites. Thus there may be new lipid and polyketide metabolites that are expressed only under certain conditions, such as during infection and disease.

A fourth class of polyketide synthases is related to the plant enzyme superfamily that includes chalcone and stilbene synthase²³. These polyketide synthases are phylogenetically divergent from all other polyketide and fatty acid synthases and generate unreduced polyketides that are typically associated with anthocyanin pigments and flavonoids. The function of these systems, which are often linked to apparent type I modules, is unknown. An example is the gene cluster spanning *pk10*, *pk7*, *pk8* and *pk9*, which includes two of the chalcone-synthase-like enzymes and two modules of an apparent type I system. The unknown metabolites produced by these enzymes are interesting because of the potent biological activities of some polyketides such as the immunosuppressor rapamycin.

Siderophores. Peptides that are not ribosomally synthesized are made by a process that is mechanistically analogous to polyketide synthesis^{23,27}. These peptides include the structurally related iron-scavenging siderophores, the mycobactins and the exochelins^{2,28}, which are derived from salicylate by the addition of serine (or threonine), two lysines and various fatty acids and possible polyketide segments. The *mbt* operon, encoding one apparent salicylate-activating protein, three amino-acid ligases, and a single module of a type I polyketide synthase, may be responsible for the biosynthesis of the mycobacterial siderophores. The presence of only one non-ribosomal peptide-synthesis system indicates that this pathway may generate both siderophores and that subsequent modification of a single ϵ -amino group of one lysine residue may account for the different physical properties and function of the siderophores²⁸.

Immunological aspects and pathogenicity

Given the scale of the global tuberculosis burden, vaccination is not only a priority but remains the only realistic public health intervention that is likely to affect both the incidence and the prevalence of the disease²⁹. Several areas of vaccine development are promising, including DNA vaccination, use of secreted or surface-exposed proteins as immunogens, recombinant forms of BCG and rational attenuation of *M. tuberculosis*²⁹. All of these avenues of research will benefit from the genome sequence as its availability will stimulate more focused approaches. Genes encoding ~90 lipoproteins were identified, some of which are enzymes or components of transport systems, and a similar number of genes encoding preproteins (with type I signal peptides) that are probably exported by the Sec-dependent pathway. *M. tuberculosis* seems to have two copies of *secA*. The potent T-cell antigen Esat-6 (ref. 30), which is probably secreted in a Sec-independent manner, is encoded by a member of a multigene family. Examination of the genetic context reveals several similarly organized operons that include genes encoding large ATP-hydrolysing membrane proteins that might act as transporters. One of the surprises of the genome project was the discovery of two extensive families of novel glycine-rich proteins, which may be of immunological significance as they are predicted to be abundant and potentially polymorphic antigens.

The PE and PPE multigene families. About 10% of the coding capacity of the genome is devoted to two large unrelated families of acidic, glycine-rich proteins, the PE and PPE families, whose genes are clustered (Figs 1, 2) and are often based on multiple copies of the polymorphic repetitive sequences referred to as PGRSs, and major polymorphic tandem repeats (MPTRs), respectively^{31,32}. The names PE and PPE derive from the motifs Pro-Glu (PE) and Pro-Pro-Glu (PPE) found near the N terminus in most cases³³. The 99 members of the PE protein family all have a highly conserved N-terminal domain of ~110 amino-acid residues that is predicted to have a globular structure, followed by a C-terminal segment that varies in size, sequence and repeat copy number (Fig. 5). Phylogenetic analysis separated the PE family into several subfamilies. The largest of these is the highly repetitive PGRS class, which contains 61 members; members of the other subfamilies, share very limited sequence similarity in their C-terminal domains (Fig. 5). The predicted molecular weights of the PE proteins vary considerably as a few members contain only the N-terminal domain, whereas most have C-terminal extensions ranging in size from 100 to 1,400 residues. The PGRS proteins have a high glycine content (up to 50%), which is the result of multiple tandem repetitions of Gly-Gly-Ala or Gly-Gly-Asn motifs, or variations thereof.

The 68 members of the PPE protein family (Fig. 5) also have a conserved N-terminal domain that comprises ~180 amino-acid residues, followed by C-terminal segments that vary markedly in sequence and length. These proteins fall into at least three groups, one of which constitutes the MPTR class characterized by the presence of multiple, tandem copies of the motif Asn-X-Gly-X-Gly-Asn-X-Gly. The second subgroup contains a characteristic,

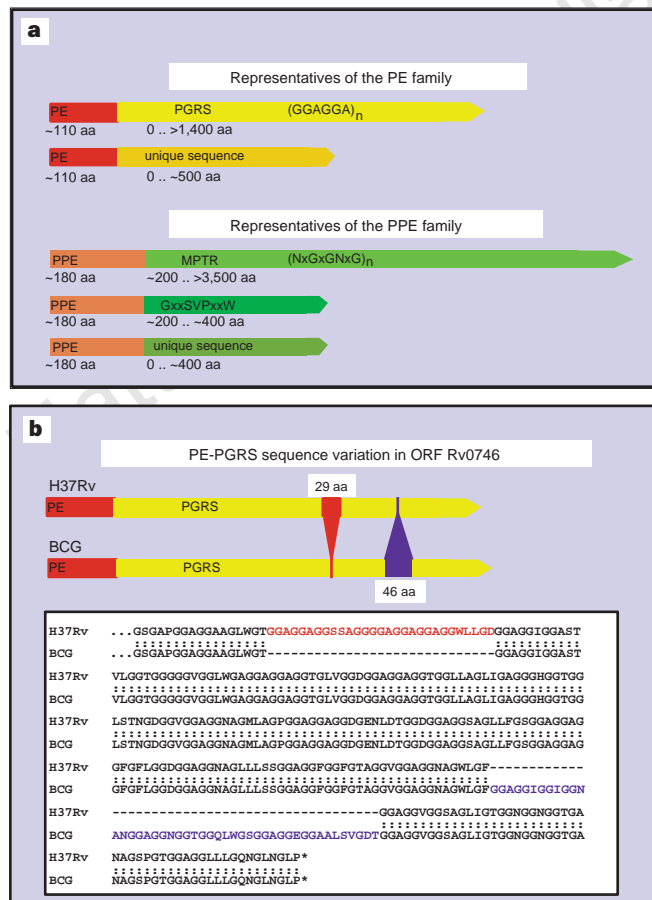


Figure 5 The PE and PPE protein families. **a**, Classification of the PE and PPE protein families. **b**, Sequence variation between *M. tuberculosis* H37Rv and *M. bovis* BCG-Pasteur in the PE-PGRS encoded by open reading frame (ORF) Rv0746.

well-conserved motif around position 350, whereas the third contains proteins that are unrelated except for the presence of the common 180-residue PPE domain.

The subcellular location of the PE and PPE proteins is unknown and in only one case, that of a lipase (Rv3097), has a function been demonstrated. On examination of the protein database from the extensively sequenced *M. leprae*¹⁵, no PGRS- or MPTR-related polypeptides were detected but a few proteins belonging to the non-MPTR subgroup of the PPE family were found. These proteins include one of the major antigens recognized by leprosy patients, the serine-rich antigen³⁴. Although it is too early to attribute biological functions to the PE and PPE families, it is tempting to speculate that they could be of immunological importance. Two interesting possibilities spring to mind. First, they could represent the principal source of antigenic variation in what is otherwise a genetically and antigenically homogeneous bacterium. Second, these glycine-rich proteins might interfere with immune responses by inhibiting antigen processing.

Several observations and results support the possibility of antigenic variation associated with both the PE and the PPE family proteins. The PGRS member Rv1759 is a fibronectin-binding protein of relative molecular mass 55,000 (ref. 35) that elicits a variable antibody response, indicating either that individuals mount different immune responses or that this PGRS protein may vary between strains of *M. tuberculosis*. The latter possibility is supported by restriction fragment length polymorphisms for various PGRS and MPTR sequences in clinical isolates³³. Direct support for genetic variation within both the PE and the PPE families was obtained by comparative DNA sequence analysis (Fig. 5). The gene for the PE-PGRS protein Rv0746 of BCG differs from that in H37Rv by the deletion of 29 codons and the insertion of 46 codons. Similar variation was seen in the gene for the PPE protein Rv0442 (data not shown). As these differences were all associated with repetitive sequences they could have resulted from intergenic or intragenic recombinational events or, more probably, from strand slippage during replication³². These mechanisms are known to generate antigenic variability in other bacterial pathogens³⁶.

There are several parallels between the PGRS proteins and the Epstein-Barr virus nuclear antigens (EBNAs). Members of both polypeptide families are glycine-rich, contain extensive Gly-Ala repeats, and exhibit variation in the length of the repeat region between different isolates. The Gly-Ala repeat region of EBNA1 functions as a *cis*-acting inhibitor of the ubiquitin/proteasome antigen-processing pathway that generates peptides presented in the context of major histocompatibility complex (MHC) class I molecules^{37,38}. MHC class I knockout mice are very susceptible to *M. tuberculosis*, underlining the importance of a cytotoxic T-cell response in protection against disease³⁹. Given the many potential effects of the PPE and PE proteins, it is important that further studies are performed to understand their activity. If extensive antigenic variability or reduced antigen presentation were indeed found, this would be significant for vaccine design and for understanding protective immunity in tuberculosis, and might even explain the varied responses seen in different BCG vaccination programmes⁴⁰.

Pathogenicity. Despite intensive research efforts, there is little information about the molecular basis of mycobacterial virulence⁴¹. However, this situation should now change as the genome sequence will accelerate the study of pathogenesis as never before, because other bacterial factors that may contribute to virulence are becoming apparent. Before the completion of the genome sequence, only three virulence factors had been described⁴¹: catalase-peroxidase, which protects against reactive oxygen species produced by the phagocyte; *mce*, which encodes macrophage-colonizing factor⁴²; and a sigma factor gene, *sigA* (aka *rpoV*), mutations in which can lead to attenuation⁴¹. In addition to these single-gene virulence

factors, the mycobacterial cell wall⁴ is also important in pathology, but the complex nature of its biosynthesis makes it difficult to identify critical genes whose inactivation would lead to attenuation.

On inspection of the genome sequence, it was apparent that four copies of *mce* were present and that these were all situated in operons, comprising eight genes, organized in exactly the same manner. In each case, the genes preceding *mce* code for integral membrane proteins, whereas *mce* and the following five genes are all predicted to encode proteins with signal sequences or hydrophobic stretches at the N terminus. These sets of proteins, about which little is known, may well be secreted or surface-exposed; this is consistent with the proposed role of Mce in invasion of host cells⁴². Furthermore, a homologue of *smpB*, which has been implicated in intracellular survival of *Salmonella typhimurium*, has also been identified⁴³. Among the other secreted proteins identified from the genome sequence that could act as virulence factors are a series of phospholipases C, lipases and esterases, which might attack cellular or vacuolar membranes, as well as several proteases. One of these phospholipases acts as a contact-dependent haemolysin (N. Stoker, personal communication). The presence of storage proteins in the bacillus, such as the haemoglobin-like oxygen captors described above, points to its ability to stockpile essential growth factors, allowing it to persist in the nutrient-limited environment of the phagosome. In this regard, the ferritin-like proteins, encoded by *bfrA* and *bfrB*, may be important in intracellular survival as the capacity to acquire enough iron in the vacuole is very limited. □

Methods

Sequence analysis. Initially, ~3.2 Mb of sequence was generated from cosmids⁸ and the remainder was obtained from selected BAC clones⁷ and 45,000 whole-genome shotgun clones. Sheared fragments (1.4–2.0 kb) from cosmids and BACs were cloned into M13 vectors, whereas genomic DNA was cloned in pUC18 to obtain both forward and reverse reads. The PGRS genes were grossly underrepresented in pUC18 but better covered in the BAC and cosmid M13 libraries. We used small-insert libraries⁴⁴ to sequence regions prone to compression or deletion and, in some cases, obtained sequences from products of the polymerase chain reaction or directly from BACs⁷. All shotgun sequencing was performed with standard dye terminators to minimize compression problems, whereas finishing reactions used dRhodamine or BigDye terminators (<http://www.sanger.ac.uk>). Problem areas were verified by using dye primers. Thirty differences were found between the genomic shotgun sequences and the cosmids; twenty of which were due to sequencing errors and ten to mutations in cosmids (1 error per 320 kb). Less than 0.1% of the sequence was from areas of single-clone coverage, and <0.2% was from one strand with only one sequencing chemistry.

Informatics. Sequence assembly involved PHRAP, GAP4 (ref. 45) and a customized perl script that merges sequences from different libraries and generates segments that can be processed by several finishers simultaneously. Sequence analysis and annotation was managed by DIANA (B.G.B. *et al.*, unpublished). Genes encoding proteins were identified by TB-parse⁴⁶ using a hidden Markov model trained on known *M. tuberculosis* coding and non-coding regions and translation-initiation signals, with corroboration by positional base preference. Interrogation of the EMBL, TrEMBL, SwissProt, PROSITE⁴⁷ and in-house databases involved BLASTN, BLASTX⁴⁸, DOTTER (<http://www.sanger.ac.uk>) and FASTA⁴⁹. tRNA genes were located and identified using tRNAscan and tRNAscan-SE⁵⁰. The complete sequence, a list of annotated cosmids and linking regions can be found on our website (<http://www.sanger.ac.uk>) and in MycDB (<http://www.pasteur.fr/mycdb/>).

Received 15 April; accepted 8 May 1998.

1. Snider, D. E. Jr, Ravignione, M. & Kochi, A. in *Tuberculosis: Pathogenesis, Protection, and Control* (ed. Bloom, B. R.) 2–11 (Am. Soc. Microbiol., Washington DC, 1994).
2. Wheeler, P. R. & Ratledge, C. in *Tuberculosis: Pathogenesis, Protection, and Control* (ed. Bloom, B. R.) 353–385 (Am. Soc. Microbiol., Washington DC, 1994).
3. Chan, J. & Kaufmann, S. H. E. in *Tuberculosis: Pathogenesis, Protection, and Control* (ed. Bloom, B. R.) 271–284 (Am. Soc. Microbiol., Washington DC, 1994).
4. Brennan, P. J. & Draper, P. in *Tuberculosis: Pathogenesis, Protection, and Control* (ed. Bloom, B. R.) 271–284 (Am. Soc. Microbiol., Washington DC, 1994).

5. Kolattukudy, P. E., Fernandes, N. D., Azad, A. K., Fitzmaurice, A. M. & Sirakova, T. D. Biochemistry and molecular genetics of cell-wall lipid biosynthesis in mycobacteria. *Mol. Microbiol.* **24**, 263–270 (1997).
6. Sreevatsan, S. *et al.* Restricted structural gene polymorphism in the *Mycobacterium tuberculosis* complex indicates evolutionarily recent global dissemination. *Proc. Natl Acad. Sci. USA* **93**, 9869–9874 (1997).
7. Brosch, R. *et al.* Use of a *Mycobacterium tuberculosis* H37Rv bacterial artificial chromosome library for genome mapping, sequencing and comparative genomics. *Infect. Immun.* **66**, 2221–2229 (1998).
8. Philipp, W. J. *et al.* An integrated map of the genome of the tubercle bacillus, *Mycobacterium tuberculosis* H37Rv, and comparison with *Mycobacterium leprae*. *Proc. Natl Acad. Sci. USA* **93**, 3132–3137 (1996).
9. Blattner, F. R. *et al.* The complete genome sequence of *Escherichia coli* K-12. *Science* **277**, 1453–1462 (1997).
10. Cole, S. T. & Saint-Girons, I. Bacterial genomics. *FEMS Microbiol. Rev.* **14**, 139–160 (1994).
11. Freiberg, C. *et al.* Molecular basis of symbiosis between *Rhizobium* and legumes. *Nature* **387**, 394–401 (1997).
12. Bardarov, S. *et al.* Conditionally replicating mycobacteriophages: a system for transposon delivery to *Mycobacterium tuberculosis*. *Proc. Natl Acad. Sci. USA* **94**, 10961–10966 (1997).
13. Mahairas, G. G., Sabo, P. J., Hickey, M. J., Singh, D. C. & Stover, C. K. Molecular analysis of genetic differences between *Mycobacterium bovis* BCG and virulent *M. bovis*. *J. Bacteriol.* **178**, 1274–1282 (1996).
14. Kunst, F. *et al.* The complete genome sequence of the gram-positive bacterium *Bacillus subtilis*. *Nature* **390**, 249–256 (1997).
15. Smith, D. R. *et al.* Multiplex sequencing of 1.5 Mb of the *Mycobacterium leprae* genome. *Genome Res.* **7**, 802–819 (1997).
16. Greenacre, M. *Theory and Application of Correspondence Analysis* (Academic, London, 1984).
17. Ratledge, C. R. in *The Biology of the Mycobacteria* (eds Ratledge, C. & Stanford, J.) 53–94 (Academic, San Diego, 1982).
18. Av-Gay, Y. & Davies, J. Components of eukaryotic-like protein signaling pathways in *Mycobacterium tuberculosis*. *Microb. Comp. Genomics* **2**, 63–73 (1997).
19. Cole, S. T. & Telenti, A. Drug resistance in *Mycobacterium tuberculosis*. *Eur. Resp. Rev.* **8**, 701S–713S (1995).
20. Riley, M. & Labeledan, B. in *Escherichia coli and Salmonella* (ed. Neidhardt, F. C.) 2118–2202 (ASM, Washington, 1996).
21. Mdulki, K. *et al.* Inhibition of a *Mycobacterium tuberculosis* β -ketoacyl ACP synthase by isoniazid. *Science* **280**, 1607–1610 (1998).
22. Banerjee, M. A., Stachelhaus, T. & Mootz, H. D. Modular peptide synthetases involved in nonribosomal peptide synthesis. *Chem. Rev.* **97**, 2651–2673 (1997).
23. Hopwood, D. A. Genetic contributions to understanding polyketide synthases. *Chem. Rev.* **97**, 2465–2497 (1997).
24. Minnikin, D. E. in *The Biology of the Mycobacteria* (eds Ratledge, C. & Stanford, J.) 95–184 (Academic, London, 1982).
25. Barry, C. E. III *et al.* Mycolic acids: structure, biosynthesis, and physiological functions. *Prog. Lipid Res.* (in the press).
26. Belisle, J. T. *et al.* Role of the major antigen of *Mycobacterium tuberculosis* in cell wall biogenesis. *Science* **276**, 1420–1422 (1997).
27. Marahiel, M. A., Stachelhaus, T. & Mootz, H. D. Modular peptide synthetases involved in nonribosomal peptide synthesis. *Chem. Rev.* **97**, 2651–2673 (1997).
28. Gobin, J. *et al.* Iron acquisition by *Mycobacterium tuberculosis*: isolation and characterization of a family of iron-binding exochelins. *Proc. Natl Acad. Sci. USA* **92**, 5189–5193 (1995).
29. Young, D. B. & Fruth, U. in *New Generation Vaccines* (eds Levine, M., Woodrow, G., Kaper, J. & Cobon, G. S.) 631–645 (Marcel Dekker, New York, 1997).
30. Sorensen, A. L., Nagai, S., Houen, G., Andersen, P. & Anderson, A. B. Purification and characterization of a low-molecular-mass T-cell antigen secreted by *Mycobacterium tuberculosis*. *Infect. Immun.* **63**, 1710–1717 (1995).
31. Hermans, P. W. M., van Soolingen, D. & van Embden, J. D. A. Characterization of a major polymorphic tandem repeat in *Mycobacterium tuberculosis* and its potential use in the epidemiology of *Mycobacterium kansasii* and *Mycobacterium goodii*. *J. Bacteriol.* **174**, 4157–4165 (1992).
32. Poulet, S. & Cole, S. T. Characterisation of the polymorphic GC-rich repetitive sequence (PGRS) present in *Mycobacterium tuberculosis*. *Arch. Microbiol.* **163**, 87–95 (1995).
33. Cole, S. T. & Barrell, B. G. in *Genetics and Tuberculosis* (eds Chadwick, D. J. & Cardew, G., *Novartis Foundation Symp.* 217) 160–172 (Wiley, Chichester, 1998).
34. Vega-Lopez, F. *et al.* Sequence and immunological characterization of a serine-rich antigen from *Mycobacterium leprae*. *Infect. Immun.* **61**, 2145–2153 (1993).
35. Abou-Zeid, C. *et al.* Genetic and immunological analysis of *Mycobacterium tuberculosis* fibronectin-binding proteins. *Infect. Immun.* **59**, 2712–2718 (1991).
36. Robertson, B. D. & Meyer, T. F. Genetic variation in pathogenic bacteria. *Trends Genet.* **8**, 422–427 (1992).
37. Levitskaya, J. *et al.* Inhibition of antigen processing by the internal repeat region of the Epstein-Barr virus nuclear antigen-1. *Nature* **375**, 685–688 (1995).
38. Levitskaya, J., Sharipo, A., Leonchiks, A., Ciechanover, A. & Masucci, M. G. Inhibition of ubiquitin/proteasome-dependent protein degradation by the Gly-Ala repeat domain of the Epstein-Barr virus nuclear antigen 1. *Proc. Natl Acad. Sci. USA* **94**, 12616–12621 (1997).
39. Flynn, J. L., Goldstein, M. A., Treibold, K. J., Koller, B. & Bloom, B. R. Major histocompatibility complex class-I restricted T cells are required for resistance to *Mycobacterium tuberculosis* infection. *Proc. Natl Acad. Sci. USA* **89**, 12013–12017 (1992).
40. Bloom, B. R. & Fine, P. E. M. in *Tuberculosis: Pathogenesis, Protection, and Control* (ed. Bloom, B. R.) 531–557 (Am. Soc. Microbiol., Washington DC, 1994).
41. Collins, D. M. In search of tuberculosis virulence genes. *Trends Microbiol.* **4**, 426–430 (1996).
42. Arruda, S., Bomfim, G., Knights, R., Huima-Byron, T. & Riley, L. W. Cloning of an *M. tuberculosis* DNA fragment associated with entry and survival inside cells. *Science* **261**, 1454–1457 (1993).
43. Baumler, A. J., Kusters, J. G., Stojkovic, I. & Heffron, F. *Salmonella typhimurium* loci involved in survival within macrophages. *Infect. Immun.* **62**, 1623–1630 (1994).
44. McMurray, A. A., Sulston, J. E. & Quail, M. A. Short-insert libraries as a method of problem solving in genome sequencing. *Genome Res.* **8**, 562–566 (1998).
45. Bonfield, J. K., Smith, K. F. & Staden, R. A new DNA sequence assembly program. *Nucleic Acids Res.* **24**, 4992–4999 (1995).
46. Krogh, A., Mian, I. S. & Haussler, D. A hidden Markov model that finds genes in *E. coli* DNA. *Nucleic Acids Res.* **22**, 4768–4778 (1994).
47. Bairoch, A., Bucher, P. & Hofmann, K. The PROSITE database, its status in 1997. *Nucleic Acids Res.* **25**, 217–221 (1997).
48. Altschul, S., Gish, W., Miller, W., Myers, E. & Lipman, D. A basic local alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990).
49. Pearson, W. & Lipman, D. Improved tools for biological sequence comparisons. *Proc. Natl Acad. Sci. USA* **85**, 2444–2448 (1988).
50. Lowe, T. M. & Eddy, S. R. tRNAscan-SE: a program for improved detection of transfer RNA genes in genomic DNA. *Nucleic Acids Res.* **25**, 955–964 (1997).

Acknowledgements. We thank Y. Av-Gay, F.-C. Bange, A. Danchin, B. Dujon, W. R. Jacobs Jr, L. Jones, M. McNeil, I. Moszer, P. Rice and J. Stephenson for advice, reagents and support. This work was supported by the Wellcome Trust. Additional funding was provided by the Association Française Raoul Follereau, the World Health Organisation and the Institut Pasteur. S.V.G. received a Wellcome Trust travelling research fellowship.

Correspondence and requests for materials should be addressed to B.G.B. (barrell@sanger.ac.uk) or S.T.C. (stcole@pasteur.fr). The complete sequence has been deposited in EMBL/GenBank/DBJ as MTBH37RV, accession number AL123456.

YOURS TO HAVE AND TO HOLD BUT NOT TO COPY

The publication you are reading is protected by copyright law. Photocopying copyright material without permission is no different from stealing a magazine from a newsagent, only it doesn't seem like theft.

If you take photocopies from books, magazines and periodicals at work your employer should be licensed with CLA.

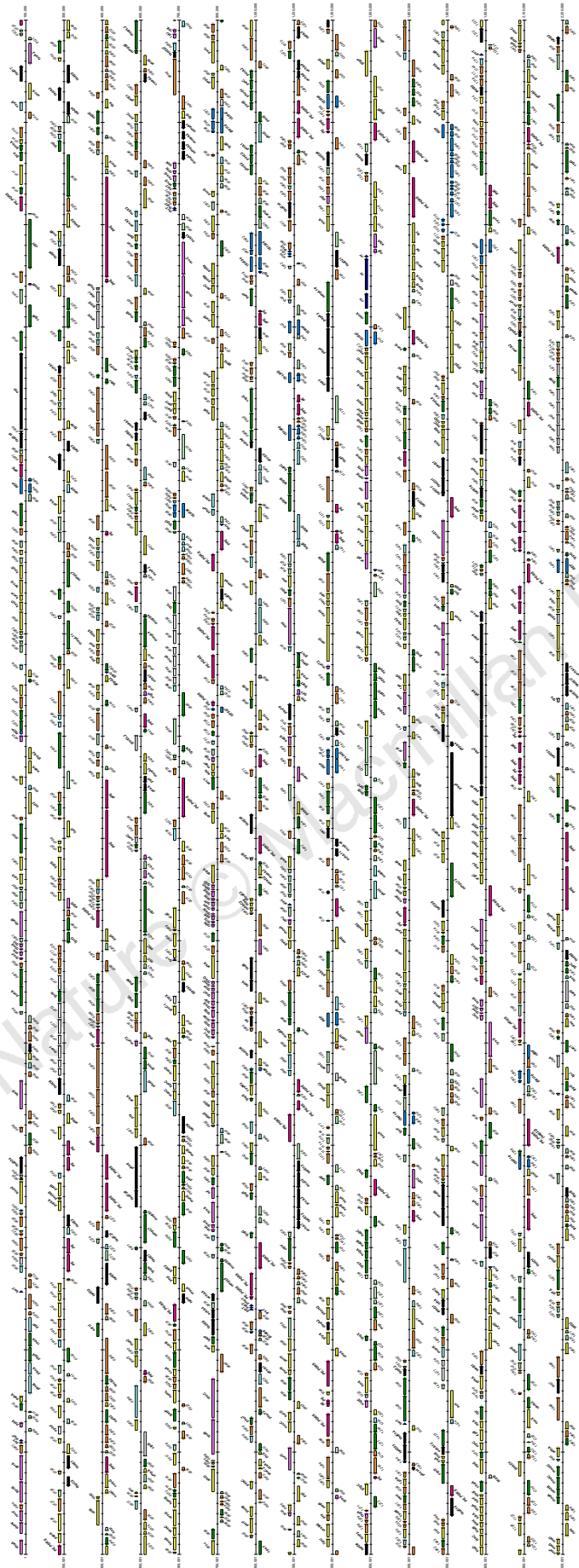
Make sure you are protected by a photocopying licence.



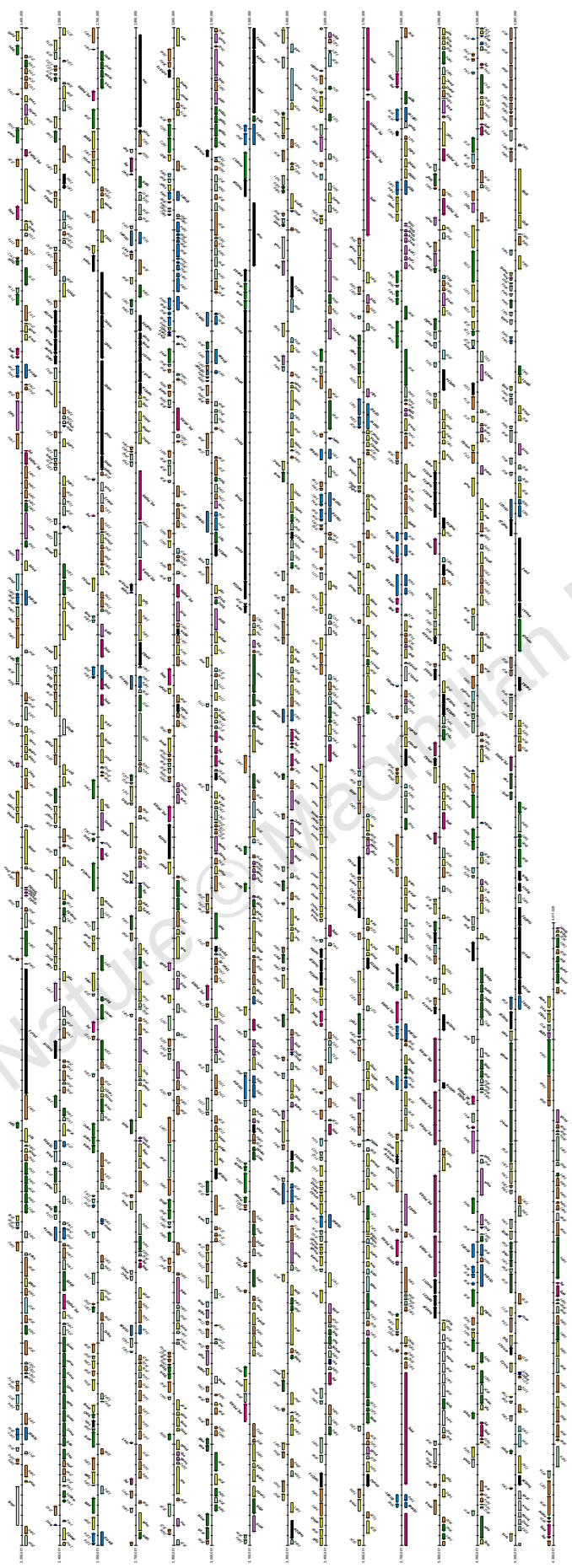
The Copyright Licensing Agency Limited
90 Tottenham Court Road, London W1P 0LP
Telephone: 0171 436 5931 Fax: 0171 436 3986

Rv2436	<i>rbsK</i>	ribokinase	Rv3250c	<i>rubB</i>	rubredoxin B	Rv1878	<i>glnA3</i>	probable glutamine synthase
Rv1408	<i>rpe</i>	ribulose-phosphate 3-epimerase				Rv2860c	<i>glnA4</i>	probable glutamine synthase
Rv2465c	<i>rpi</i>	phosphopentose isomerase	7. Miscellaneous oxidoreductases and oxygenases 171			Rv2918c	<i>glnD</i>	uridylyltransferase
Rv1448c	<i>tal</i>	transaldolase				Rv2221c	<i>glnE</i>	glutamate-ammonia-ligase
Rv1449c	<i>tkk</i>	transketolase						adenylyltransferase
Rv1121	<i>zwf</i>	glucose-6-phosphate 1-dehydrogenase	8. ATP-proton motive force			Rv3859c	<i>gltB</i>	ferredoxin-dependent glutamate synthase
Rv1447c	<i>zwf2</i>	glucose-6-phosphate 1-dehydrogenase	Rv1308	<i>atpA</i>	ATP synthase chain	Rv3858c	<i>gltD</i>	small subunit of NADH-dependent glutamate synthase
			Rv1304	<i>atpB</i>	ATP synthase chain			
			Rv1311	<i>atpC</i>	ATP synthase ϵ chain	Rv3704c	<i>gshA</i>	possible -glutamylcysteine synthase
			Rv1310	<i>atpD</i>	ATP synthase chain			
			Rv1305	<i>atpE</i>	ATP synthase c chain	Rv2427c	<i>proA</i>	-glutamyl phosphate reductase
			Rv1306	<i>atpF</i>	ATP synthase b chain	Rv2439c	<i>proB</i>	glutamate 5-kinase
			Rv1309	<i>atpG</i>	ATP synthase chain	Rv0500	<i>proC</i>	pyrroline-5-carboxylate reductase
			Rv1307	<i>atpH</i>	ATP synthase chain			
6. Respiration								
<i>a. aerobic</i>								
Rv0527	<i>ccsA</i>	cytochrome <i>c</i> -type biogenesis protein	C. Central intermediary metabolism					
			1. General					
Rv0529	<i>ccsB</i>	cytochrome <i>c</i> -type biogenesis protein	Rv2589	<i>gabT</i>	4-aminobutyrate aminotransferase	2. Aspartate family		
Rv1451	<i>ctaB</i>	cytochrome <i>c</i> oxidase assembly factor	Rv3432c	<i>gabB</i>	glutamate decarboxylase	Rv3708c	<i>asd</i>	aspartate semialdehyde dehydrogenase
Rv2200c	<i>ctaC</i>	cytochrome <i>c</i> oxidase chain II	Rv1832	<i>gcvB</i>	glycine decarboxylase	Rv3709c	<i>ask</i>	aspartokinase
Rv3043c	<i>ctaD</i>	cytochrome <i>c</i> oxidase polypeptide I	Rv1826	<i>gcvH</i>	glycine cleavage system H protein	Rv2201	<i>asnB</i>	asparagine synthase B
			Rv2211c	<i>gcvT</i>	T protein of glycine cleavage system	Rv3565	<i>aspB</i>	aspartate aminotransferase
Rv2193	<i>ctaE</i>	cytochrome <i>c</i> oxidase polypeptide III	Rv1213	<i>glgC</i>	glucose-1-phosphate adenylyltransferase	Rv0337c	<i>aspC</i>	aspartate aminotransferase
Rv1542c	<i>glnN</i>	hemoglobin-like, oxygen carrier	Rv3842c	<i>glpQ1</i>	glycerophosphoryl diester phosphodiesterase	Rv2753c	<i>dapA</i>	dihydrodipicolinate synthase
Rv2470	<i>glnO</i>	hemoglobin-like, oxygen carrier	Rv0317c	<i>glpQ2</i>	glycerophosphoryl diester phosphodiesterase	Rv2773c	<i>dapB</i>	dihydrodipicolinate reductase
Rv2249c	<i>glpD1</i>	glycerol-3-phosphate dehydrogenase	Rv3566c	<i>nhoA</i>	N-hydroxyarylamino <i>o</i> -acetyltransferase	Rv1202	<i>dapE</i>	succinyl-diaminopimelate desuccinylase
Rv3302c	<i>glpD2</i>	glycerol-3-phosphate dehydrogenase	Rv0155	<i>pntAA</i>	pyridine transhydrogenase subunit 1	Rv2141c	<i>dapE2</i>	ArgE/DapE/Acy1/Cpg2/yscS family
Rv0694	<i>lldD1</i>	L-lactate dehydrogenase (cytochrome)	Rv0156	<i>pntAB</i>	pyridine transhydrogenase subunit 2	Rv2726c	<i>dapF</i>	diaminopimelate epimerase
			Rv0157	<i>pntB</i>	pyridine transhydrogenase subunit	Rv1293	<i>lysA</i>	diaminopimelate decarboxylase
Rv1872c	<i>lldD2</i>	L-lactate dehydrogenase	Rv1127c	<i>ppdK</i>	similar to pyruvate, phosphate dikinase	Rv3341	<i>metA</i>	homoserine <i>o</i> -acetyltransferase
Rv1854c	<i>ndh</i>	probable NADH dehydrogenase				Rv1079	<i>metB</i>	cystathionine -synthase
Rv3145	<i>nuoA</i>	NADH dehydrogenase chain A				Rv3340	<i>metC</i>	cystathionine -lyase
Rv3146	<i>nuoB</i>	NADH dehydrogenase chain B				Rv1133c	<i>metE</i>	5-methyltetrahydropteroyltrimethylglutamate-homocysteine methyltransferase
Rv3147	<i>nuoC</i>	NADH dehydrogenase chain C						
Rv3148	<i>nuoD</i>	NADH dehydrogenase chain D				Rv2124c	<i>metH</i>	5-methyltetrahydrofolate-homocysteine methyltransferase
Rv3149	<i>nuoE</i>	NADH dehydrogenase chain E				Rv1392	<i>metK</i>	S-adenosylmethionine synthase
Rv3150	<i>nuoF</i>	NADH dehydrogenase chain F				Rv0391	<i>metZ</i>	<i>o</i> -succinylhomoserine sulfhydrylase
Rv3151	<i>nuoG</i>	NADH dehydrogenase chain G	2. Gluconeogenesis					
Rv3152	<i>nuoH</i>	NADH dehydrogenase chain H	Rv0211	<i>pckA</i>	phosphoenolpyruvate carboxylase	Rv1294	<i>thrA</i>	homoserine dehydrogenase
Rv3153	<i>nuoI</i>	NADH dehydrogenase chain I				Rv1296	<i>thrB</i>	homoserine kinase
Rv3154	<i>nuoJ</i>	NADH dehydrogenase chain J	Rv0069c	<i>sdaA</i>	L-serine dehydratase 1	Rv1295	<i>thrC</i>	homoserine synthase
Rv3155	<i>nuoK</i>	NADH dehydrogenase chain K						
Rv3156	<i>nuoL</i>	NADH dehydrogenase chain L						
Rv3157	<i>nuoM</i>	NADH dehydrogenase chain M	3. Sugar nucleotides					
Rv3158	<i>nuoN</i>	NADH dehydrogenase chain N	Rv1512	<i>epiA</i>	nucleotide sugar epimerase	3. Serine family		
Rv2195	<i>qcrA</i>	Rieske iron-sulphur component of <i>ubiQ</i> - <i>cytB</i> reductase	Rv3784	<i>epiB</i>	probable UDP-galactose 4-epimerase	Rv0815c	<i>cysA2</i>	thiosulfate sulfurtransferase
			Rv1511	<i>gmdA</i>	GDP-mannose 4,6 dehydratase	Rv3117	<i>cysA3</i>	thiosulfate sulfurtransferase
Rv2196	<i>qcrB</i>	cytochrome component of <i>ubiQ</i> - <i>cytB</i> reductase	Rv0334	<i>rmlA</i>	glucose-1-phosphate thymidyltransferase	Rv2335	<i>cysE</i>	serine acetyltransferase
			Rv3264c	<i>rmlA2</i>	glucose-1-phosphate thymidyltransferase	Rv0511	<i>cysG</i>	uroporphyrin-III <i>c</i> -methyltransferase
Rv2194	<i>qcrC</i>	cytochrome <i>b/c</i> component of <i>ubiQ</i> - <i>cytB</i> reductase	Rv3464	<i>rmlB</i>	dTDP-glucose 4,6-dehydratase	Rv2847c	<i>cysG2</i>	multifunctional enzyme, siroheme synthase
			Rv3634c	<i>rmlB2</i>	dTDP-glucose 4,6-dehydratase	Rv2334	<i>cysK</i>	cysteine synthase A
			Rv3468c	<i>rmlB3</i>	dTDP-glucose 4,6-dehydratase	Rv1336	<i>cysM</i>	cysteine synthase B
			Rv3465	<i>rmlC</i>	dTDP-4-dehydrothiamine	Rv1077	<i>cysM2</i>	cystathionine -synthase
						Rv0848	<i>cysM3</i>	putative cysteine synthase
			Rv3266c	<i>rmlD</i>	3,5-epimerase	Rv1093	<i>glyA</i>	serine hydroxymethyltransferase
			Rv0322	<i>udgA</i>	dTDP-4-dehydrothiamine reductase	Rv0070c	<i>glyA2</i>	serine hydroxymethyltransferase
						Rv2996c	<i>serA</i>	D-3-phosphoglycerate dehydrogenase
			Rv3265c	<i>wbbL</i>	UDP-glucose dehydrogenase/GDP-mannose 6-dehydrogenase	Rv0505c	<i>serB</i>	probable phosphoserine phosphatase
			Rv1525	<i>wbb2</i>	dTDP-rhamnosyl transferase	Rv3042c	<i>serB2</i>	C-term similar to phosphoserine phosphatase
			Rv3400	-	dTDP-rhamnosyl transferase probable -phosphoglucomutase	Rv0884c	<i>serC</i>	phosphoserine aminotransferase
			4. Amino sugars					
			Rv3436c	<i>glmS</i>	glucosamine-fructose-6-phosphate aminotransferase	4. Aromatic amino acid family		
						Rv3227	<i>aroA</i>	3-phosphoshikimate
			5. Sulphur metabolism			Rv2538c	<i>aroB</i>	1-carboxyvinyl transferase
			Rv0711	<i>atsA</i>	arylsulfatase	Rv2537c	<i>aroD</i>	3-dehydroquininate dehydrogenase
			Rv3299c	<i>atsB</i>	probable arylsulfatase	Rv2552c	<i>aroE</i>	3-dehydroquininate dehydrogenase
			Rv0663	<i>atsD</i>	probable arylsulfatase	Rv2540c	<i>aroF</i>	shikimate 5-dehydrogenase
			Rv3077	<i>atsF</i>	probable arylsulfatase	Rv2178c	<i>aroG</i>	chorismate synthase
			Rv0296c	<i>atsG</i>	probable arylsulfatase	Rv2539c	<i>aroK</i>	DAHPh synthase
			Rv3796	<i>atsH</i>	probable arylsulfatase	Rv3838c	<i>pheA</i>	shikimate kinase I
			Rv1285	<i>cysD</i>	ATP:sulphurylase subunit 2	Rv1613	<i>trpA</i>	prephenate dehydratase
			Rv1286	<i>cysN</i>	ATP:sulphurylase subunit 1	Rv1612	<i>trpB</i>	tryptophan synthase chain
			Rv2131c	<i>cysQ</i>	homologue of <i>M.leprae</i> <i>cysQ</i>	Rv1611	<i>trpC</i>	tryptophan synthase chain
			Rv3248c	<i>sahH</i>	adenosylhomocysteinase			
			Rv3283	<i>sseA</i>	thiosulfate sulfurtransferase	Rv2192c	<i>trpD</i>	anthranilate phosphoribosyltransferase
			Rv2291	<i>sseB</i>	thiosulfate sulfurtransferase	Rv1609	<i>trpE</i>	anthranilate synthase
			Rv3118	<i>sseC</i>	thiosulfate sulfurtransferase			
			Rv0814c	<i>sseC2</i>	thiosulfate sulfurtransferase	Rv2386c	<i>trpE2</i>	anthranilate synthase component I
			Rv3762c	-	probable alkyl sulfatase	Rv3754	<i>tyrA</i>	prephenate dehydrogenase
			D. Amino acid biosynthesis					
			1. Glutamate family					
			Rv1654	<i>argB</i>	acetylglutamate kinase	5. Histidine		
			Rv1652	<i>argC</i>	N-acetyl- -glutamyl-phosphate reductase	Rv1603	<i>hisA</i>	phosphoribosylformimino-5-aminoimidazole carboxamide
			Rv1655	<i>argD</i>	acetylornithine aminotransferase	Rv1601	<i>hisB</i>	ribonucleotide isomerase
			Rv1656	<i>argF</i>	ornithine carbamoyltransferase	Rv1600	<i>hisC</i>	imidazole glycerol-phosphate dehydratase
			Rv1658	<i>argG</i>	arginosuccinate synthase			
			Rv1659	<i>argH</i>	arginosuccinate lyase	Rv3772	<i>hisC2</i>	histidinol-phosphate aminotransferase
			Rv1653	<i>argJ</i>	glutamate N-acetyltransferase			
			Rv2220	<i>glnA1</i>	glutamine synthase class I	Rv1599	<i>hisD</i>	histidinol dehydrogenase
			Rv2222c	<i>glnA2</i>	glutamine synthase class II			

Rv1605	<i>hisF</i>	imidazole glycerol-phosphate synthase	Rv3048c	<i>urdG</i>	subunit ribonucleoside-diphosphate small subunit	Rv3119	<i>moaE</i>	subunit 1 molybdopterin-converting factor subunit 2
Rv2121c	<i>hisG</i>	ATP phosphoribosyltransferase	Rv3053c	<i>urdH</i>	glutaredoxin electron transport component of NrdEF system	Rv0866	<i>moaE2</i>	molybdopterin-converting factor subunit 2
Rv1602	<i>hisH</i>	amidotransferase	Rv3052c	<i>urdI</i>	NrdI/YgaO/YmaA family thymidylate kinase	Rv3322c	<i>moaE3</i>	molybdopterin-converting factor subunit 2
Rv2122c	<i>hisI</i>	phosphoribosyl-AMP cyclohydro-lase	Rv3247c	<i>tmk</i>	thymidylate kinase	Rv0994	<i>moaA</i>	molybdopterin biosynthesis
Rv1606	<i>hisI2</i>	probable phosphoribosyl-AMP 1,6 cyclohydrolyase	Rv2764c	<i>thyA</i>	thymidylate synthase	Rv3116	<i>moaB</i>	molybdopterin biosynthesis
Rv0114	-	similar to HisB	Rv0570	<i>urdZ</i>	ribonucleotide reductase, class II	Rv2338c	<i>moaW</i>	molybdopterin biosynthesis
6. Pyruvate family			Rv3752c	-	probable cytidine/deoxycytidylate deaminase	Rv1681	<i>moaX</i>	weak similarity to <i>E. coli</i> MoaA
Rv3423c	<i>alr</i>	alanine racemase	4. Salvage of nucleosides and nucleotides			Rv1355c	<i>moaY</i>	probably involved in molybdopterin biosynthesis
7. Branched amino acid family			Rv3313c	<i>add</i>	probable adenosine deaminase	Rv3206c	<i>moaZ</i>	molybdopterin biosynthesis
Rv1559	<i>ilvA</i>	threonine deaminase	Rv2584c	<i>apt</i>	adenine phosphoribosyltransferases	Rv0865	<i>mog</i>	molybdopterin biosynthesis
Rv3003c	<i>ilvB</i>	acetolactate synthase I large sub-unit	Rv3315c	<i>cdt</i>	probable cytidine deaminase	5. Pantothenate		
Rv3470c	<i>ilvB2</i>	acetolactate synthase large sub-unit	Rv3314c	<i>deoA</i>	thymidine phosphorylase	Rv1092c	<i>coaA</i>	pantothenate kinase
Rv3001c	<i>ilvC</i>	ketol-acid reductoisomerase	Rv0478	<i>deoC</i>	deoxyribose-phosphate aldolase	Rv2225	<i>panB</i>	3-methyl-2-oxobutanoate hydroxymethyltransferase
Rv0189c	<i>ilvD</i>	dihydroxy-acid dehydratase	Rv3307	<i>deoD</i>	probable purine nucleoside phosphorylase	Rv3602c	<i>panC</i>	pantoate- α -alanine ligase
Rv2210c	<i>ilvE</i>	branched-chain-amino-acid transaminase	Rv3624c	<i>hpt</i>	probable hypoxanthine-guanine phosphoribosyltransferase	Rv3601c	<i>panD</i>	aspartate 1-decarboxylase
Rv1820	<i>ilvG</i>	acetolactate synthase II	Rv3393	<i>iunH</i>	probable inosine-uridine preferring nucleoside hydrolase	6. Pyridoxine		
Rv3002c	<i>ilvN</i>	acetolactate synthase I small sub-unit	Rv0535	<i>pnp</i>	phosphorylase from Pnp/MtaP family 2	Rv2607	<i>pdxH</i>	pyridoxamine 5'-phosphate oxidase
Rv3509c	<i>ilvX</i>	probable acetohydroxyacid synthase I large subunit	Rv3309c	<i>upp</i>	uracil phosphoribosyltransferase	7. Pyridine nucleotide		
Rv3710	<i>leuA</i>	-isopropyl malate synthase	5. Miscellaneous nucleoside/nucleotide reactions			Rv1594	<i>nadA</i>	quinolinate synthase
Rv2995c	<i>leuB</i>	3-isopropylmalate dehydrogenase	Rv0733	<i>adk</i>	probable adenylate kinase	Rv1595	<i>nadB</i>	L-aspartate oxidase
Rv2988c	<i>leuC</i>	3-isopropylmalate dehydratase large subunit	Rv2364c	<i>bex</i>	GTP-binding protein of Era/ThdF family	Rv1596	<i>nadC</i>	nicotinate-nucleotide pyrophosphatase
Rv2987c	<i>leuD</i>	3-isopropylmalate dehydratase small subunit	Rv1712	<i>cmk</i>	cytidylate kinase	Rv0423c	<i>thiC</i>	thiamine synthesis, pyrimidine moiety
<i>E. Polyamine synthesis</i>			Rv2344c	<i>dgt</i>	probable deoxyguanosine triphosphate hydrolase	8. Thiamine		
Rv2601	<i>speE</i>	spermidine synthase	Rv2404c	<i>lepA</i>	GTP-binding protein LepA	Rv0422c	<i>thiD</i>	phosphomethylpyrimidine kinase
<i>F. Purines, pyrimidines, nucleosides and nucleotides</i>			Rv2727c	<i>miaA</i>	tRNA (2)-isopentenylpyrophosphate transferase	Rv0414c	<i>thiE</i>	thiamine synthesis, thiazole moiety
1. Purine ribonucleotide biosynthesis			Rv2445c	<i>ndkA</i>	nucleoside diphosphate kinase	Rv0417	<i>thiG</i>	thiamine synthesis, thiazole moiety
Rv1389	<i>gmk</i>	putative guanylate kinase	Rv2440c	<i>obg</i>	Obg GTP-binding protein	Rv2977c	<i>thiL</i>	probable thiamine-monophosphate kinase
Rv3396c	<i>guaA</i>	GMP synthase	Rv2583c	<i>relA</i>	(p)ppGpp synthase I	9. Riboflavin		
Rv1843c	<i>guaB1</i>	inosine-5'-monophosphate dehydrogenase	<i>G. Biosynthesis of cofactors, prosthetic groups and carriers</i>			Rv1940	<i>ribA</i>	GTP cyclohydrolase II
Rv3411c	<i>guaB2</i>	inosine-5'-monophosphate dehydrogenase	1. Biotin			Rv1415	<i>ribA2</i>	probable GTP cyclohydrolase II
Rv3410c	<i>guaB3</i>	inosine-5'-monophosphate dehydrogenase	Rv1568	<i>bioA</i>	adenosylmethionine-8-amino-7-oxononanoate aminotransferase	Rv1412	<i>ribC</i>	riboflavin synthase chain
Rv1017c	<i>prsA</i>	ribose-phosphate pyrophosphokinase	Rv1589	<i>bioB</i>	biotin synthase	Rv2671	<i>ribD</i>	probable riboflavin deaminase
Rv0357c	<i>purA</i>	adenylosuccinate synthase	Rv1570	<i>bioD</i>	dethiobiotin synthase	Rv2786c	<i>ribF</i>	riboflavin kinase
Rv0777	<i>purB</i>	adenylosuccinate lyase	Rv1569	<i>bioF</i>	8-amino-7-oxononanoate synthase	Rv1409	<i>ribG</i>	riboflavin biosynthesis
Rv0780	<i>purC</i>	phosphoribosylaminoimidazole-succinocarboxamide synthase	Rv0032	<i>bioF2</i>	C-terminal similar to <i>B. subtilis</i> BioF	Rv1416	<i>ribH</i>	riboflavin synthase chain
Rv0772	<i>purD</i>	phosphoribosylamine-glycine ligase	Rv3279c	<i>birA</i>	biotin apo-protein ligase	Rv3300c	-	probable deaminase, riboflavin synthesis
Rv3275c	<i>purE</i>	phosphoribosylaminoimidazole carboxylase	Rv1442	<i>bisC</i>	biotin sulfoxide reductase	10. Thioredoxin, glutaredoxin and mycothiol		
Rv0808	<i>purF</i>	amidophosphoribosyltransferase	Rv0089	-	possible <i>bioC</i> biotin synthesis gene	Rv0773c	<i>ggtA</i>	putative -glutamyl transpeptidase
Rv0957	<i>purH</i>	phosphoribosylaminoimidazole-carboxamide formyltransferase	2. Folic acid			Rv2394	<i>ggtB</i>	-glutamyltranspeptidase precursor
Rv3276c	<i>purK</i>	phosphoribosylaminoimidazole carboxylase ATPase subunit	Rv2763c	<i>dfrA</i>	dihydrofolate reductase	Rv2855	<i>gorA</i>	glutathione reductase homologue
Rv0803	<i>purL</i>	phosphoribosylformylglycinamide synthase II	Rv2447c	<i>folC</i>	folypolyglutamate synthase	Rv0816c	<i>thiX</i>	equivalent to <i>M. leprae</i> ThiX
Rv0809	<i>purM</i>	5'-phosphoribosyl-5-aminoimidazole synthase	Rv3356c	<i>folD</i>	methylene tetrahydrofolate dehydrogenase	Rv1470	<i>trxA</i>	thioredoxin
Rv0956	<i>purN</i>	phosphoribosylglycinamide formyltransferase I	Rv3609c	<i>folE</i>	GTP cyclohydrolase I	Rv1471	<i>trxB</i>	thioredoxin reductase
Rv0788	<i>purQ</i>	phosphoribosylformylglycinamide synthase I	Rv3606c	<i>folK</i>	7,8-dihydro-6-hydroxymethylpterin pyrophosphokinase	Rv3913	<i>trxB2</i>	thioredoxin reductase
Rv0389	<i>purT</i>	phosphoribosylglycinamide formyltransferase II	Rv3608c	<i>folP</i>	dihydropterate synthase	Rv3914	<i>trxC</i>	thioredoxin
Rv2964	<i>purU</i>	formyltetrahydrofolate deformylase	Rv1207	<i>folP2</i>	dihydropterate synthase	11. Menaquinone, PQQ, ubiquinone and other terpenoids		
2. Pyrimidine ribonucleotide biosynthesis			Rv3607c	<i>folX</i>	may be involved in folate biosynthesis	Rv2682c	<i>dxs</i>	1-deoxy-D-xylulose 5-phosphate synthase
Rv1383	<i>carA</i>	carbamoyl-phosphate synthase subunit	Rv0013	<i>pabA</i>	<i>p</i> -aminobenzoate synthase	Rv0562	<i>grcC1</i>	heptaprenyl diphosphate synthase II
Rv1384	<i>carB</i>	carbamoyl-phosphate synthase subunit	Rv1005c	<i>pabB</i>	glutamine amidotransferase	Rv0989c	<i>grcC2</i>	heptaprenyl diphosphate synthase II
Rv1380	<i>pyrB</i>	aspartate carbamoyltransferase	Rv0812	<i>pabC</i>	<i>p</i> -aminobenzoate synthase	Rv3398c	<i>idsA</i>	geranylgeranyl pyrophosphate synthase
Rv1381	<i>pyrC</i>	dihydroorotate dehydrogenase	3. Lipote			Rv2173	<i>idsA2</i>	geranylgeranyl pyrophosphate synthase
Rv2139	<i>pyrD</i>	dihydroorotate dehydrogenase	Rv2218	<i>lipA</i>	lipote biosynthesis protein A	Rv3383c	<i>idsB</i>	transfergeranyl, similar geranyl pyrophosphate synthase
Rv1385	<i>pyrF</i>	orotidine 5'-phosphate decarboxylase	Rv2217	<i>lipB</i>	lipote biosynthesis protein B	Rv0534c	<i>menA</i>	4-dihydroxy-2-naphthoate octaprenyltransferase
Rv1699	<i>pyrG</i>	CTP synthase	4. Molybdopterin			Rv0548c	<i>menB</i>	naphthoate synthase
Rv2883c	<i>pyrH</i>	uridylylate kinase	Rv3109	<i>moaA</i>	molybdenum cofactor biosynthesis, protein A	Rv0553	<i>menC</i>	<i>o</i> -succinylbenzoate-CoA synthase
Rv0382c	<i>umpA</i>	probable uridine 5'-monophosphate synthase	Rv0869c	<i>moaA2</i>	molybdenum cofactor biosynthesis, protein A	Rv0555	<i>menD</i>	2-succinyl-6-hydroxy-2,4-cyclohexadiene-1-carboxylate synthase
3. 2'-deoxyribonucleotide metabolism			Rv0438c	<i>moaA3</i>	molybdenum cofactor biosynthesis, protein A	Rv0542c	<i>menE</i>	<i>o</i> -succinylbenzoic acid-CoA ligase
Rv0321	<i>dcd</i>	deoxycytidine triphosphate deaminase	Rv3110	<i>moaB</i>	molybdenum cofactor biosynthesis, protein B	Rv3853	<i>menG</i>	S-adenosylmethionine: 2-demethylmenaquinone phytoene synthase
Rv2697c	<i>dut</i>	deoxyuridine triphosphatase	Rv0984	<i>moaB2</i>	molybdenum cofactor biosynthesis, protein B	Rv3397c	<i>phyA</i>	coenzyme PQQ synthesis
Rv0233	<i>nrdB</i>	ribonucleoside-diphosphate reductase B2 (eukaryotic-like)	Rv3111	<i>moaC</i>	molybdenum cofactor biosynthesis, protein C	Rv0693	<i>pqqE</i>	protein E
Rv3051c	<i>nrdE</i>	ribonucleoside diphosphate reductase chain	Rv0864	<i>moaC2</i>	molybdenum cofactor biosynthesis, protein C	Rv0558	<i>ubiE</i>	ubiquinone/menaquinone biosynthesis methyltransferase
Rv1981c	<i>nrdF</i>	ribonucleotide reductase small subunit	Rv3324c	<i>moaC3</i>	molybdenum cofactor biosynthesis, protein C	12. Heme and porphyrin		
			Rv3112	<i>moaD</i>	molybdopterin converting factor subunit 1	Rv0509	<i>hemA</i>	glutamyl-tRNA reductase
			Rv0868c	<i>moaD2</i>	molybdopterin converting factor	Rv0512	<i>hemB</i>	-aminolevulinic acid dehydratase
						Rv0510	<i>hemC</i>	porphobilinogen deaminase
						Rv2678c	<i>hemE</i>	uroporphyrinogen decarboxylase



Nature © Macmillan Publishers Ltd 1998



Nature © Macmillan Publishers Ltd 1998

Rv0932c	<i>pstS</i>	phosphate transport system PstS component of phosphate uptake	Rv1821	<i>secA2</i>	unit SecA, preprotein translocase sub-unit	Rv3500c	-	part of <i>mce4</i> operon
Rv2400c	<i>subI</i>	sulphate binding precursor	Rv2587c	<i>secD</i>	protein-export membrane protein	Rv3501c	-	part of <i>mce4</i> operon
Rv0143c	-	probable chloride channel	Rv0638	<i>secE</i>	SecE preprotein translocase	Rv3896c	-	putative p60 homologue
Rv1707	-	probable sulphate permease	Rv2586c	<i>secF</i>	protein-export membrane protein	Rv3922c	-	possible hemolysin
Rv1739c	-	possible sulphate transporter	Rv1440	<i>secG</i>	protein-export membrane protein	<i>B. IS elements, Repeated sequences, and Phage</i>		
Rv3679	-	possible anion transporter	Rv0732	<i>secY</i>	SecY subunit of preprotein translocase	1. IS elements		
Rv3680	-	probable anion transporter				IS6110	16 copies	
5. Fatty acid transport			Rv2462c	<i>tig</i>	chaperone protein, similar to trigger factor	IS1081	6 copies	
Rv2790c	<i>lip1</i>	non-specific lipid transport protein	Rv2813	-	probable general secretion pathway protein	Others	37 copies	
Rv3540c	<i>lip2</i>	non-specific lipid transport protein				2. REP13E12 family 7 copies		
6. Efflux proteins			<i>E. Adaptations and atypical conditions</i>			3. Phage-related functions		
Rv2936	<i>draA</i>	similar daunorubicin resistance ABC-transporter	Rv1901	<i>cinA</i>	competence damage protein	Rv2894c	<i>xerC</i>	integrase/recombinase
Rv2937	<i>draB</i>	similar daunorubicin resistance transmembrane protein	Rv3648c	<i>cspA</i>	cold shock protein, transcriptional regulator	Rv1701	<i>xerD</i>	integrase/recombinase
Rv2938	<i>draC</i>	similar daunorubicin resistance transmembrane protein	Rv0871	<i>cspB</i>	probable cold shock protein	Rv1054	-	integrase-a
Rv2846c	<i>efpA</i>	putative efflux protein	Rv3063	<i>cstA</i>	starvation-induced stress response protein	Rv1055	-	integrase-b
Rv3065	<i>emrE</i>	resistance to ethidium bromide	Rv0871	<i>cspB</i>	probable cold shock protein	Rv1573	-	phiRV1 phage related protein
Rv0783c	-	multidrug resistance protein	Rv3490	<i>otsA</i>	probable , -trehalose-phosphate synthase	Rv1574	-	phiRV1 phage related protein
Rv0849	-	possible quinolone efflux pump	Rv2006	<i>otsB</i>	trehalose-6-phosphate phosphatase	Rv1575	-	phiRV1 phage related protein
Rv1145	-	probable drug transporter	Rv3372	<i>otsB2</i>	trehalose-6-phosphate phosphatase	Rv1576c	-	phiRV1 phage related protein
Rv1146	-	probable drug transporter	Rv3758c	<i>proV</i>	osmoprotection ABC transporter	Rv1577c	-	phiRV1 phage related protein
Rv1250	-	probable drug efflux protein	Rv3757c	<i>proW</i>	transport system permease	Rv1578c	-	phiRV1 phage related protein
Rv1258c	-	probable multidrug resistance pump	Rv3759c	<i>proX</i>	similar to osmoprotection proteins	Rv1579c	-	phiRV1 phage related protein
Rv1410c	-	probable drug efflux protein	Rv3756c	<i>proZ</i>	transport system permease	Rv1580c	-	phiRV1 phage related protein
Rv1634	-	probable drug efflux protein	Rv1026	-	probable pppGpp-5-phosphohydro-lase	Rv1581c	-	phiRV1 phage related protein
Rv1819c	-	probable multidrug resistance pump				Rv1582c	-	phiRV1 phage related protein
Rv2136c	-	putative bacitracin resistance protein	<i>F. Detoxification</i>			Rv1583c	-	phiRV1 phage related protein
Rv2209	-	probable drug efflux protein	Rv2428	<i>ahpC</i>	alkyl hydroperoxide reductase	Rv1584c	-	phiRV1 phage related protein
Rv2333c	-	probable tetracycline C resistance protein	Rv2429	<i>ahpD</i>	member of AhpC/TSA family	Rv1585c	-	phiRV1 phage related protein
Rv2994	-	probable fluoroquinolone efflux protein	Rv2238c	<i>ahpE</i>	member of AhpC/TSA family	Rv1586c	-	phiRV1 phage related protein
Rv1877	-	probable drug efflux protein	Rv2521	<i>bcp</i>	bacterioferritin comigratory protein	Rv2309c	-	integrase
Rv2459	-	probable drug efflux protein	Rv1608c	<i>bcpB</i>	probable bacterioferritin comigratory protein	Rv2310	-	excisionase
<i>B. Chaperones/Heat shock</i>			Rv3473c	<i>bpoA</i>	probable non-heme bromoperoxidase	Rv2646	-	phiRV2 integrase
Rv0384c	<i>clpB</i>	heat shock protein	Rv1123c	<i>bpoB</i>	probable non-heme bromoperoxidase	Rv2647	-	phiRV2 phage related protein
Rv0352	<i>dnaJ</i>	acts with GrpE to stimulate DnaK ATPase	Rv0554	<i>bpoC</i>	probable non-heme bromoperoxidase	Rv2650c	-	phiRV2 phage related protein
Rv2373c	<i>dnaJ2</i>	DnaJ homologue	Rv3617	<i>epaA</i>	probable epoxide hydrolase	Rv2651c	-	phiRV2 phage related protein
Rv0350	<i>dnaK</i>	70 kD heat shock protein, chromosome replication	Rv1938	<i>epaB</i>	probable epoxide hydrolase	Rv2652c	-	phiRV2 phage related protein
Rv3417c	<i>groEL1</i>	60 kD chaperonin 1	Rv1124	<i>epaC</i>	probable epoxide hydrolase	Rv2653c	-	phiRV2 phage related protein
Rv0440	<i>groEL2</i>	60 kD chaperonin 2	Rv2214c	<i>epaD</i>	probable epoxide hydrolase	Rv2654c	-	phiRV2 phage related protein
Rv3418c	<i>groES</i>	10 kD chaperone	Rv3670	<i>epaE</i>	probable epoxide hydrolase	Rv2655c	-	phiRV2 phage related protein
Rv0351	<i>grpE</i>	stimulates DnaK ATPase activity	Rv0134	<i>epaF</i>	probable epoxide hydrolase	Rv2656c	-	phiRV2 phage related protein
Rv2374c	<i>hrcA</i>	heat-inducible transcription repressor	Rv3171c	<i>hpx</i>	probable non-heme haloperoxidase	Rv2657c	-	similar to gp36 of mycobacteriophage L5
Rv0251c	<i>hsp</i>	possible heat shock protein	Rv1908c	<i>katG</i>	catalase-peroxidase	Rv2658c	-	phiRV2 phage related protein
Rv0353	<i>hspR</i>	heat shock regulator	Rv3846	<i>sodA</i>	superoxide dismutase	Rv2659c	-	phiRV2 integrase
Rv2031c	<i>hspX</i>	14kD antigen, heat shock protein Hsp20 family	Rv0432	<i>sodC</i>	superoxide dismutase precursor - (Cu-Zn)	Rv2830c	-	similar to phage P1 <i>phd</i> gene excisionase
Rv2299c	<i>htpG</i>	heat shock protein Hsp90 family	Rv1932	<i>tpx</i>	thiol peroxidase	Rv3750c	-	putative integrase
Rv0563	<i>htpX</i>	probable (transmembrane) heat shock protein	Rv0634c	-	putative glyoxylase II	Rv3751	-	
Rv2701c	<i>shhB</i>	putative extragenic suppressor protein	Rv2581c	-	putative glyoxylase II	<i>C. PE and PPE families</i>		
Rv3269	-	probable heat shock protein	Rv3177	-	probable non-heme haloperoxidase	1. PE family		
<i>C. Cell division</i>			IV. Other			PE subfamily 38 members		
Rv3641c	<i>fic</i>	possible cell division protein	<i>A. Virulence</i>			PE_PGRS subfamily 61 members		
Rv3102c	<i>ftsE</i>	membrane protein	Rv0169	<i>mce1</i>	cell invasion protein	2. PPE family 68 members		
Rv3610c	<i>ftsH</i>	inner membrane protein, chaperone	Rv0589	<i>mce2</i>	cell invasion protein	<i>D. Antibiotic production and resistance</i>		
Rv2748c	<i>ftsK</i>	chromosome partitioning	Rv1966	<i>mce3</i>	cell invasion protein	Rv2068c	<i>blaC</i>	class A -lactamase
Rv2151c	<i>ftsQ</i>	ingrowth of wall at septum	Rv3499c	<i>mce4</i>	cell invasion protein	Rv3290c	<i>lat</i>	lysine-ε aminotransferase
Rv2154c	<i>ftsW</i>	membrane protein (shape determination)	Rv3100c	<i>smfB</i>	probable small protein b	Rv2043c	<i>pncA</i>	pyrazinamide resistance/sensitivity
Rv3101c	<i>ftsX</i>	membrane protein	Rv1694	<i>tlyA</i>	cytotoxin/hemolysin homologue	Rv0133	-	possible puromycin N-acetyltransferase
Rv2921c	<i>ftsY</i>	cell division protein FtsY	Rv0024	-	putative p60 homologue	Rv0262c	-	aminoglycoside 2'-N-acetyltransferase
Rv2150c	<i>ftsZ</i>	circumferential ring, GTPase	Rv0167	-	part of <i>mce1</i> operon	Rv0802c	-	acetyltransferase
Rv3919c	<i>gid</i>	glucose inhibited division protein B	Rv0168	-	part of <i>mce1</i> operon	Rv1082	-	similar to <i>S. lincolnsensis lmbE</i>
Rv3625c	<i>mesJ</i>	probable cell cycle protein	Rv0170	-	part of <i>mce1</i> operon	Rv1170	-	similar to <i>S. lincolnsensis lmbE</i>
Rv3917c	<i>parA</i>	chromosome partitioning; DNA - binding	Rv0171	-	part of <i>mce1</i> operon	Rv1347c	-	possible aminoglycoside 6'-N-acetyltransferase
Rv3918c	<i>parB</i>	possibly involved in chromosome partitioning	Rv0172	-	part of <i>mce1</i> operon	Rv2036	-	similar to lincomycin production genes
Rv2922c	<i>smc</i>	member of Smc1/Cut3/Cut14 family	Rv0174	-	part of <i>mce1</i> operon	Rv2303c	-	similar to <i>S. griseus</i> macrotetrolide resistance protein
Rv0012	-	possible cell division protein	Rv0587	-	part of <i>mce2</i> operon	Rv3225c	-	probable aminoglycoside 3'-phosphotransferase
Rv0435c	-	ATPase of AAA-family	Rv0588	-	part of <i>mce2</i> operon	Rv3700c	-	probable acetyltransferase
Rv2115c	-	ATPase of AAA-family	Rv0590	-	part of <i>mce2</i> operon	Rv3817	-	probable aminoglycoside 3'-phosphotransferase
Rv3213c	-	possible role in chromosome segregation	Rv0591	-	part of <i>mce2</i> operon	<i>E. Bacteriocin-like proteins</i> 3		
Rv1708	-	possible role in chromosome partitioning	Rv0592	-	part of <i>mce2</i> operon	<i>F. Cytochrome P450 enzymes</i> 22		
<i>D. Protein and peptide secretion</i>			Rv0594	-	part of <i>mce2</i> operon	<i>G. Coenzyme F420-dependent enzymes</i> 3		
Rv2916c	<i>ffh</i>	signal recognition particle protein	Rv1085c	-	possible hemolysin	<i>H. Miscellaneous transferases</i> 61		
Rv2903c	<i>lepB</i>	signal peptidase I	Rv1477	-	putative exported p60 protein homologue	<i>I. Miscellaneous phosphatases, lyases, and hydrolases</i> 18		
Rv1614	<i>lgt</i>	prolipoprotein diacylglycerol transferase	Rv1478	-	putative exported p60 protein homologue	<i>J. Cyclases</i> 6		
Rv1539	<i>lspA</i>	lipoprotein signal peptidase	Rv1566c	-	putative exported p60 protein homologue	<i>K. Chelatases</i> 2		
Rv0379	<i>sec</i>	probable transport protein	Rv1964	-	part of <i>mce3</i> operon	<i>V. Conserved hypotheticals</i> 912		
Rv3240c	<i>secA</i>	SecA, preprotein translocase sub-	Rv1965	-	part of <i>mce3</i> operon	<i>VI. Unknowns</i> 606		
			Rv1966	-	part of <i>mce3</i> operon	TOTAL 3924		
			Rv1967	-	part of <i>mce3</i> operon			
			Rv1968	-	part of <i>mce3</i> operon			
			Rv1969	-	part of <i>mce3</i> operon			
			Rv1971	-	part of <i>mce3</i> operon			
			Rv2190c	-	putative p60 homologue			
			Rv3494c	-	part of <i>mce4</i> operon			
			Rv3496c	-	part of <i>mce4</i> operon			
			Rv3497c	-	part of <i>mce4</i> operon			
			Rv3498c	-	part of <i>mce4</i> operon			