



Research Journal of **Microbiology**

ISSN 1816-4935



Academic
Journals Inc.

www.academicjournals.com

Genome Wide Single Nucleotide Polymorphism Analysis of *Mycobacterium* Species and Subspecies

S.K. Srivastava, M. Agrawal and M. Grover
Amity Institute of Biotechnology, Amity University, Sec-125, Expressway,
Noida 210304, Uttar Pradesh, India

Abstract: In this study we report the reannotation of the genome of seven *Mycobacterial* species and subspecies. We have used bioinformatics tools for annotation and reevaluated each of the Protein-Coding Sequences (CDS) previously annotated and presented the combined results of recent database searches. We have also used comparative genomics tools to focus on comparative analysis as an effective strategy. Pair wise comparison between the various *Mycobacterium* strains was performed so as to predict the relationships between them. Among the wide variety of mycobacterium strains present, we selected seven and showed how their genome is interrelated by studying synteny with the genomes of various strains studied. The genome wide SNP analysis in the seven genomes of *Mycobacterium* sp. was also done in this study and the base by base changes in the genome of these seven subspecies were identified. The gene based SNPs were further classified into the marker SNPs (SNPs which are unique amongst all the seven studied species). Out of a total of 2073 SNPs, 966 were identified as marker SNPs. This study may be used for further analysis of host pathogen interactions at the pathway and product level. The present investigation will also be useful for study of evolutionary relationship.

Key words: Gene annotation, *Mycobacterium tuberculosis*, SNPs, synteny

INTRODUCTION

One third of the world's population-nearly two billion people are infected with the bacterium that causes tuberculosis. Two million people die from it each year. Effective drugs to treat and cure the disease have been available for more than 50 years, yet every 15 sec, someone in the world dies from TB. In addition to this, a person is reported to be newly infected with TB every second of every day. Left untreated, a person with active TB will infect an average of 10 to 15 other people every year. (<http://www.worldlungfoundation.org/tuberculosis.php>). The etiologic agent of tuberculosis in humans is bacterium *Mycobacterium tuberculosis*.

The sequenced genomes of pathogenic bacteria provide useful information for understanding genomic traits. The complete genome sequence of *M. tuberculosis* strain H37Rv was obtained by Cole *et al.* (1998). The genome comprises 4,411,529 base pairs, containing around 4,000 genes. The sequence of the genome and its comparison to sequences of other microorganisms reported in several databases has allowed the assignment of precise functions to 40% of the predicted proteins and the identification of 44% of orthologues (genes with very similar functions in a different species), leaving 16% as unique unknown proteins.

Some unique characteristics of the biology of the tubercle bacillus, such as its characteristic slow growth, the nature of its complex cell wall, certain genes related to its virulence and persistence and the apparent stability of its genome have been elucidated with the help of genome sequence of

M. tuberculosis H37Rv. The availability of the complete genome sequences has also led to the development of downstream sciences that take advantage of this sequence information. It is critical, therefore, that genome annotations are frequently updated if the information they contain is to remain accurate, relevant and useful. Several other *Mycobacterial* genomes have thus been reannotated recently (Dandekar *et al.*, 2000; Serres *et al.*, 2001; Gaasterland and Opera, 2001; Bocs *et al.*, 2002).

In this study we report the reannotation of seven *Mycobacterial* species and subspecies genomes. We have used bioinformatic tools for annotation and reevaluated each of the Protein-Coding Sequences (CDS) previously annotated and presented the combined results of recent database searches.

In the current investigation we have also used the bioinformatic tools for comparative genomics. Pair wise comparison between the various mycobacterium strains was performed so as to predict the relationships between them. Among the wide variety of *Mycobacterium* strains present, we chose seven and showed how their genome is interrelated by studying synteny with the genomes of various strains taken. The genome wide SNP analysis in the seven genomes of *Mycobacterium* sp. was also done in this study and the base by base changes in the genome of these seven subspecies were identified. The gene based SNPs were further classified into the marker SNPs (SNPs which are unique amongst all the seven studied species). Out of a total of 2073 SNPs, 966 were identified as marker SNPs. This study may be used for further analysis of host pathogen interactions at the pathway and product level. The present investigation will also be useful for study of evolutionary relationships.

MATERIALS AND METHODS

The present investigation was carried out in June 2008.

Sequences

The following sets of *Mycobacterium* subspecies were used in this study:

- *Mycobacterium tuberculosis* H37Rv, Accession No. AL123456 (Cole *et al.*, 1998, 2001)
- *Mycobacterium leprae* TN, Accession No. AL450380 (Geluk *et al.*, 2005)
- *Mycobacterium bovis* BCG Pasteur, Accession No. AM408590 (Garnier *et al.*, 2003)
- *Mycobacterium bovis* subsp. *bovis* AF2122/97, Accession No. BX248333 (Garnier *et al.*, 2003)
- *Mycobacterium* CDC 1551, Accession No. AE000516 (Fleischmann *et al.*, 2002)
- *Mycobacterium avium* subsp. *Paratuberculosis*, Accession No. AE016958 (Li *et al.*, 2005)
- *Mycobacterium smegmatis* str. MC2 155, Accession No. CP000480 (TIGR)

The whole genome sequence of each of the above sub-species of *Mycobacterium* made available by National Center for Biotechnology Information (NCBI) was obtained from <http://www.ncbi.nlm.nih.gov/> in a FASTA format, a text-based format for representing either nucleic acid sequences or peptide sequences, in which base pairs or amino acids were represented using single-letter codes. The format also allowed for sequence names and comments to precede the sequences.

Sequence Analysis and Re-Annotation

Each of the downloaded sequence was split up into parts by using JR Split File (<http://www.spadixbd.com/freetools/jsplit.htm>). The File Split function allowed us to split a large file into smaller files and create a standard bat file that could be used to reconstruct a copy of the original file. Then each of the split sequence was re-annotated. Re-annotation was supported by FGENESB software (<http://www.softberry.com/berry.phtml?topic=fgenesb&group=programs&subgroup=gfindb>), a tool for Bacterial Operon and Gene Prediction. Each of the *mycobacterium* sub-species coding genes were predicted and annotated. Further analysis of putative

proteins was performed, based on the results of similarity searches using BLASTP with the default parameters and Fasta sequence comparisons using non-redundant (nr) data from the GenBank database. Once we determined the total number of hits to nr, the best top BLAST hit for each gene was obtained.

Identification of SNPs

Computational strategy was used to identify SNPs present in seven *Mycobacterium* genomes. The sequence data sets for each sub-species were aligned using SNPs Finder software (<http://snpsfinder.lanl.gov/>) to detect polymorphic sites (taking *Mycobacterium tuberculosis* H37Rv as a reference sequence) which were differentiated as SNPs and cSNPs. SNPs were differentiated from cSNPs by marking the coding and non coding sequence in the reference sequence (*Mycobacterium tuberculosis* 37 Rv).

RESULTS AND DISCUSSION

The comparative genomic studies focuses primarily on biotechnological areas and potential biomedical benefits include insights into the specialized, shared systems used by disease-causing organisms (pathogens) to disable or destroy human cells. Comparing these genomic data with those of other microbes may help understand a diverse range of pathogens that have remarkably similar methods for infiltrating organisms with protein-coding genes capable of sneaking past human defense systems. Hence, comparative genomics and related technologies are helping to unravel the molecular basis of the pathogenesis, host range, evolution and phenotypic differences of the slow-growing mycobacteria.

Selection of Species

In this post-genomic era, an understanding of strain-to-strain variation is expected to fast-forward the rational design and development of effective diagnostic reagents, vaccines and drugs in the case of the difficult-to-crack TB pathogen. Therefore, not surprisingly the present era of *Mycobacterium* genomics is rapidly evolving into a period of comparative genomics.

Several genomes including two strains of *M. tuberculosis* (i.e., H37Rv and CDC1551), *M. bovis*, *M. leprae*, *M. smegmatis*, *M. ulcerans*, *Mycobacterium* sp. MCS, *Mycobacterium* sp. KMS, *M. vanbaalenii* PYR-1, *M. avium* and *M. avium* subsp. *paratuberculosis* have been sequenced and annotated completely (www.ncbi.nlm.nih.gov/genomes/lproks.cgi). Fifteen additional *Mycobacterial* genome sequencing projects have either been completed and awaiting annotation or are underway.

On the basis of this, we chose seven sub-species of *Mycobacterium* for our whole genome comparative study. A combined approach of genomics and bio-informatics was adopted for this study, which was undertaken with the widely used reference strain *M. tuberculosis* H37Rv. This comparative genome analysis revealed the similarities and differences across various genomes.

In silico Gene Prediction and Re-Annotation

The initial genome annotation in *M. tuberculosis* H37Rv strain soon became out of date. For this reason, a re-annotation of that genome sequence was published in 2002. This re-annotation incorporated eighty-two new Protein-Coding Sequences (CDS) and 22 of these have a predicted function. According to the above study, it is possible to assign a function to 2058 proteins (52% of the 3995 proteins predicted) and only 376 putative proteins share no homology with known proteins and thus could be unique to *M. tuberculosis*.

Here, in this paper complete re-annotation of the genome sequence of *Mycobacterium tuberculosis* sub-species is presented. The tool (FGENESB) was used on seven different organisms' viz., *Mycobacterium tuberculosis* H37Rv, *Mycobacterium leprae* TN, *Mycobacterium bovis* BCG Pasteur, *Mycobacterium bovis* subsp. *bovis* AF2122/97, *Mycobacterium* CDC 1551, *Mycobacterium avium* subsp. *paratuberculosis* and *Mycobacterium smegmatis* MC2 155.

The input submitted to the tool was a set of protein sequences in FASTA format, which were made available from NCBI. The tool was run at default parameters and against the closest organism *Mycobacterium tuberculosis* H37Rv, taken as a reference. The tool identifies genes based on similarities with the genes present in the database. The results obtained deciphered the number of genes predicted (CDS), along with their starting and ending positions on the genome, position on the chain and the sequence of each gene annotated in each of the species taken. The scores were also obtained due to matching of predicted genes with the database.

Revising the Number of Genes in the Genome

When the results were compared with the first annotation of the genome, it was found that there is a rise in the number of annotated putative genes in all the organisms taken (Table 1). All of the new genes are believed to encode polypeptides. The numbering of the new CDS has not interfered with the labeling of the existing genes.

The gene sequences obtained after gene annotation were re-analysed based on the results of similarity searches using BLASTP with the default parameters. The tool identifies essential genes based on similarities with the genes present in the database. The tool matches all the genes in the input data set with the database using BLAST and displayed hits at the specified E-value and a bit score, using top hit.

Genome Wide SNP Analysis

The genome wide SNP analysis in the seven genomes of *Mycobacterium* sp. was done in this study and the base by base changes in the genome of these seven subspecies were identified which is shown in Table 2. The gene based SNPs were further classified into the marker SNPs (SNPs which are unique amongst all the seven studied species). The genome wide SNP analysis was extended to find out the variations at the gene level. The annotated genes of *Mycobacterium tuberculosis* Rv were taken as a reference sequence and the variations falling in the genes of Rv strain were coded as gene SNPs (Table 2).

On analysis of SNP discovery, we observed that a total of 2073 SNPs amongst seven *mycobacterium* subspecies are spread throughout the genome out of which 966 represent the unique/ marker SNPs, which means that marker SNPs (46.6% of the total SNPs) have a greater importance in the genomic studies.

The new genomic annotation of *M. tuberculosis* H37Rv as a reference sequence and other species are used in this study to know the various changes at the nucleotide level. The above study was used to know the base by base changes in the genome of seven different species in their genic and intergenic region. In a similar study a combination of previously and newly identified SNPs based on whole-genome comparisons of *M. tuberculosis* strains H37Rv, CDC1551 and 210 and *M. bovis* AF2122/97

Table 1: Total numbers of genes found in each of the *Mycobacterium* species and subspecies and probable SNP markers in numbers with reference to *Mycobacterium tuberculosis* H37Rv, control sequence

Species	Accession No.	Size of genome (Mbp)	Total No. of genes	Total SNPs	Marker SNPs
<i>Mycobacterium tuberculosis</i> H37Rv (Ref. Sequence)	AL123456	4.40000	4356	0	0
<i>Mycobacterium leprae</i> TN	AL450380	3.26820	5253	160	108
<i>Mycobacterium bovis</i> BCG Pasteur	AM408590	4.40000	4293	1238	202
<i>Mycobacterium bovis</i> subsp. <i>bovis</i> AF2122/97	BX248333	4.34549	4280	1235	237
<i>Mycobacterium</i> CDC 1551	AE000516	4.40384	4352	616	302
<i>Mycobacterium avium</i> subsp. <i>Paratuberculosis</i>	AE016958	4.80000	1578	160	117
<i>Mycobacterium smegmatis</i> str. MC2 155	CP000480	7.00000	7005		

Table 2: The annotated genes and SNPs in *Mycobacterium* species and sub species

Gene No	Gene Coordinate	CNS-96					CNS-96					CNS-96												
		hdsr	Mb-ahn	Mb-hss	Mb-hssM22	Mb-DK100	hdsr	Mb-ahn	Mb-hss	Mb-hssM22	Mb-DK100	hdsr	Mb-ahn	Mb-hss	Mb-hssM22	Mb-DK100								
1	467 A	G	A	A	A	A	141	145865 G	G	A	G	G	313	338020 A	C	A	595321 T	T	T	G	T			
1	1057 G	A	G	G	G	G	141	147455 C	C	A	C	C	313	338100 T	C	T	560	597553 G	G	C	G	G		
2	2347 A	G	A	A	A	A	142	147919 A	A	A	G	A	319	346034 A	A	A	A	560	599165 A	A	A	C	A	
2	2532 T	C	T	T	T	T	142	147985 C	C	C	T	C	328	357278 C	T	C	C	C	567	606216 C	G	C	C	C
	3751 T	G	T	T	T	T	145	150327 A	A	A	C	A	331	361089 G	G	G	C	G	570	609926 G	G	A	G	G
	4480 C	T	C	C	C	C	145	150334 C	C	C	G	C	332	362867 C	C	C	T	C	585	621236 G	G	C	G	G
5	5752 G	A	G	G	G	G	145	150338 G	G	G	C	G	334	362513 G	G	A	G	G	588	623931 G	G	G	T	G
5	6406 C	T	C	C	C	C	148	155293 G	G	G	A	G	334	363757 C	C	T	C	C	593	628895 G	G	G	A	G
5	6446 G	T	G	G	G	G	153	159179 G	G	A	G	G	337	364975 C	C	T	C	C		631689 A	A	A	G	A
6	8285 C	T	C	C	C	C	154	160335 T	T	T	C	T	337	364978 C	C	T	C	C	608	643002 G	G	C	G	G
6	8741 C	T	C	C	C	C	154	160344 C	C	T	C	C	349	382904 G	G	G	A	G	608	643319 T	C	T	T	T
6	9143 T	C	T	T	T	T	155	161351 C	A	C	C	C	365	384808 C	C	C	T	C	611	646331 G	G	G	A	G
6	9217 A	C	A	A	A	A	156	162309 G	G	A	G	G	367	397275 G	G	G	C	G	613	648667 C	C	T	C	C
9	10727 A	G	A	A	A	A	158	163573 A	A	A	C	A	371	399660 T	G	T	T	T	621	657269 A	A	A	G	A
	11370 C	C	T	T	T	T	160	165609 A	A	G	A	A	374	405400 A	A	A	G	A	625	659808 G	A	G	G	G
17	13197 G	C	G	G	G	G	161	166506 G	G	A	G	G	374	405750 T	T	T	G	T	629	663451 G	G	G	G	G
17	13460 A	G	A	A	A	A		166852 C	C	C	T	C	376	407507 T	C	T	T	T	629	664045 G	T	G	G	G
19	14401 A	G	A	A	A	A	169	172901 C	C	A	C	C	377	408515 G	A	G	G	G	639	672637 C	C	T	C	C
20	15117 C	G	C	C	C	C	170	174684 C	C	C	A	C	379	409965 C	C	C	G	C	639	674741 A	G	A	A	A
22	17857 T	T	C	T	T	T	174	178946 C	C	C	T	C	390	423289 A	G	A	A	A	639	674744 C	G	C	C	C
27	24679 A	A	G	A	A	A	184	187738 G	G	G	A	G	393	424424 C	C	C	A	C	639	674830 C	C	C	G	C
27	25306 C	C	G	C	C	C	187	191887 G	G	A	G	G	393	424425 C	C	C	A	C	646	679173 T	T	G	T	T
31	27714 C	T	C	C	C	C	187	191988 G	G	A	G	G	394	426909 A	C	C	C	C	647	680171 G	G	G	C	G
34	30519 C	C	C	T	T	T	187	193014 T	C	T	T	T	394	428221 C	C	C	G	C	648	682667 C	C	C	T	C
34	30688 T	T	T	G	T	T	188	193381 C	C	C	G	C	394	428228 G	G	G	T	G	648	682712 G	T	G	G	G
34	30943 C	C	C	T	T	T	196	202369 G	G	G	A	G	394	428231 A	A	A	G	A	648	682964 T	C	T	T	T
	31081 C	C	C	T	T	T	201	208305 A	G	A	A	A	394	429060 T	T	C	T	T	648	683826 C	C	T	C	C
	33804 G	G	A	G	G	G	211	217779 T	G	T	T	T	398	437465 G	G	C	G	G	651	686655 C	C	C	T	C
41	35342 C	T	C	C	C	C	212	219144 C	T	C	C	C	402	440745 T	T	C	T	T	657	692588 C	C	T	C	C
45	39093 G	A	G	G	G	G	214	220050 C	C	C	T	C	408	445415 C	C	T	C	C		697855 T	C	T	T	T
45	39273 G	A	G	G	G	G	217	224441 T	C	T	T	T	408	446405 C	C	C	G	C	688	700084 C	C	G	C	C
46	40478 C	C	C	A	A	A	220	227022 A	A	G	A	A	412	448905 C	C	T	C	C	688	700141 G	G	A	G	G
	41244 C	C	A	C	C	C	223	231994 C	C	G	C	C	416	453394 T	T	T	G	T	689	700943 C	T	C	C	C
47	41453 C	C	T	C	C	C	224	232574 G	G	G	T	G	416	453992 C	C	C	T	C	684	711074 C	C	T	C	C
	50005 G	G	A	G	G	G	224	233979 C	T	C	C	C	416	453993 C	C	C	T	C	685	712082 C	C	G	C	C
57	50786 A	G	A	A	A	A	230	241776 A	C	A	A	A	421	458949 A	A	A	G	A	697	721924 A	A	G	A	A
58	51649 G	G	G	C	C	C	237	251466 G	G	G	T	G	426	466326 T	T	C	T	T	697	723473 C	T	C	C	C
59	51954 T	T	T	A	A	A	238	253032 C	T	C	C	C	426	466385 A	A	C	A	A	698	725734 G	G	A	G	G
61	53785 C	C	C	G	G	G	239	254012 G	G	G	A	G	427	466923 C	T	C	C	C	698	726816 G	G	G	C	G
61	55552 G	T	G	G	G	G	242	258128 T	G	T	T	T	445	483292 A	A	A	T	A	699	728729 C	T	C	C	C
61	55557 C	C	C	T	T	T	245	261235 G	G	G	A	G	447	484874 G	T	G	G	G	713	737857 C	T	C	C	C
63	57393 A	A	A	T	T	T	247	263195 C	C	C	T	C	448	487878 T	T	G	T	T	718	746123 G	G	G	A	G
	59871 G	G	A	G	G	G	247	264021 T	C	T	T	T	448	489073 C	C	C	G	C	738	763575 G	C	G	G	G
69	62657 G	G	G	A	A	A	249	265554 A				C	448	489859 A	A	G	A	A	738	765150 G	G	G	A	G
	65786 A	C	A	A	A	A	249	265988 C				G	452	494270 C	T	C	C	C	754	781588 C	T	C	C	C
74	65867 G	G	A	G	G	G	250	266859 G	G	A	G	G	458	501615 C	C	C	G	C	755	782246 G	G	G	A	G
74	66285 C	C	C	G	G	G	250	266860 G	G	A	G	G	458	501620 C	C	T	C	C	760	786259 G	G	G	T	G
77	69072 C	T	C	C	C	C	254	271665 C	C	T	C	C	463	505302 G	A	G	G	G		800219 T	T	T	C	T
77	69871 C	C	C	T	T	T		277885 C	C	C	G	C	506910 C	C	T	C	C	C	778	803175 C	C	T	C	C
77	69923 C	C	T	C	C	C	264	284623 G	G	G	A	G	469	514245 C	C	C	T	C	785	806630 T	C	T	T	T
77	70092 A	A	G	A	A	A	267	288952 C	C	T	C	C	473	518393 C	C	T	C	C	785	806779 G	G	T	G	G
77	71396 C	C	C	G	G	G	270	290454 C	C	T	C	C	477	519805 G	G	A	G	G		809915 A	A	C	A	A
	75264 G	C	C	T	T	T	274	294127 C	T	C	C	C	480	524327 C	C	C	G	C	787	810887 A	A	G	A	A
88	83982 A	A	A	G	G	G	274	294128 C	T	C	C	C	487	531550 C	G	C	C	C	805	822822 T	C	T	T	T
89	84972 A	G	A	A	A	A	277	297194 C	C	T	C	C	490	534394 C	T	C	C	C	809	826181 G	G	C	G	G
92	87079 T	C	T	T	T	T	283	304679 G	G	G	T	G	495	540552 C	C	T	C	C		835012 A	A	T	A	A
96	89124 C	T	C	C	C	C	287	308167 C	C	C	A	C	498	543181 T	T	C	T	T	822	835375 C	T	C	C	C
99	90929 C	C	T	C	C	C	287	308337 G	G	G	A	G	499	544461 C	C	C	T	C	823	836487 C	C	C	A	C
103	95591 G	G	G	T	T	T	291	313086 C	T	C	C	C		546914 G	G	G	A	G	823	837033 A				A
104	97388 G	G	G	T	T	T	293	315502 C	C	C	A	C	504	549144 T	T	T	C	C	823	837434 G	G	G	A	G
119	114093 C	A	C	C	C	C	295	316586 T	C	T	T	T	512	555330 C	C	T	C	C	823	837580 G	G	A	G	G
127	126042 G	C	G	G	G	G	296	318088 C	T	C	C	C	514	558720 C	T	C	C	C	824	838990 C				C
127	126043 G																							

Table 2: Continued

Gene No.	Gene Coordinate	Anchor	Min. width	Min. length	Min. width	Min. length	Gene No.	Gene Coordinate	Anchor	Min. width	Min. length	Min. width	Min. length	Gene No.	Gene Coordinate	Anchor	Min. width	Min. length	Min. width	Min. length						
1475615	T	C	T	T	T	T	1475697	C	C	C	C	C	C	1589	1997137	C	T	C	C	1753	1781177	C	T	C	C	
1475638	C	T	C	C	C	C	1476608	C	G	C	C	C	C	1596	1605149	C	C	T	C	1758	1783650	G	A	G	G	
1475642	T	C	T	T	T	T	1476610	C	G	C	C	C	C	1606	16157	A	G	A	A	1761	1785534	C	G	C	C	
1475645	C	T	C	C	C	C	1476631	G	G	G	G	A	G	1606280	G	G	G	A	G	1775	1798570	C	A	C	C	
1475649	T	C	T	T	T	T	1476633	G	G	G	G	A	G	1606281	G	G	G	C	G	1777	1800521	A	A	G	A	
1475656	A	A	A	A	A	G	1476684	C	C	C	C	T	C	1598	1608201	C	C	C	T	1789	1810730	C	C	C	T	
1475659	G	G	G	G	G	A	1476686	C	C	C	C	T	C	1605	1616069	T	C	T	T	1792	1813260	T	T	C	T	
1475666	G	G	G	G	G	A	1476689	A	T	A	A	A	A	1610	1618416	G	G	A	G	G	1815	1815222	T	C	T	T
1475700	C	C	C	C	T	C	1476691	A	T	A	A	A	A	1610	1618624	G	G	T	G	G	1796	1816848	G	G	T	G
1475707	C	C	C	C	T	C	1476693	G	A	G	G	G	G	1610	1618999	C	C	C	T	C	1802	1825327	G	G	G	A
1475754	C	C	C	C	T	C	1476695	G	A	G	G	G	G	1622	162763	C	T	C	C	C	1804	1827304	C	T	C	C
1475761	C	C	C	C	T	C	1476697	A	T	T	T	T	T	1623	162098	A	G	A	A	A	1808	1830284	G	A	G	G
1475882	T	C	T	T	T	T	1476716	A	T	A	A	A	A	1623	1623101	C	G	C	C	C	1813	1838287	A	T	A	A
1475889	T	C	T	T	T	T	1476718	A	T	A	A	A	A	1622	1631596	C	G	C	C	C	1822	1839204	G	C	G	G
1475941	C	C	C	C	T	C	1476733	T	C	T	T	T	T	1622	1631826	A	A	A	G	A	1817	1843555	G	A	G	G
1475948	C	C	C	C	T	C	1476735	T	C	T	T	T	T	1622	1632140	T	T	G	T	T	1821	1849729	G	T	G	G
1475989	A	G	A	A	A	A	1476741	G	G	G	G	A	G	1622	1632737	C	C	C	T	C	1822	1852372	G	T	G	G
1475996	A	G	A	A	A	A	1476743	G	G	G	G	A	G	1622	1632825	G	G	G	A	G	1832	1863504	A	A	A	G
1475998	A	A	A	A	A	G	1476760	T	C	T	T	T	T	1632	1643934	C	C	C	A	C	1835	1868578	C	C	C	T
1476025	G	T	G	G	G	G	1476762	T	C	T	T	T	T	1633	1645334	C	C	C	T	C	1840	1872906	G	G	T	G
1476026	G	T	G	G	G	G	1476778	C	C	C	C	T	C	1634	1646777	A	T	A	A	A	1841	1875207	C	C	T	C
1476031	C	G	C	C	C	C	1476780	C	C	C	C	T	C	1635	1647117	T	T	C	T	T	1842	1878234	G	G	A	G
1476032	C	G	C	C	C	C	1476787	A	A	A	A	G	A	1635	1647777	C	C	C	T	C	1843	1884697	G	G	A	G
1476041	C	C	C	C	T	C	1476789	A	A	A	A	G	A	1635	1649211	G	C	G	G	G	1850	1896367	T	T	C	T
1476042	C	C	C	C	T	C	1478883	C	C	T	C	C	C	1635	1649217	C	G	C	C	C	1851	1896864	A	G	A	A
1476046	G	T	G	G	G	G	1487084	C	C	C	T	C	C	1635	1649226	G	C	G	G	G	1859	1902867	C	C	C	T
1476047	G	T	G	G	G	G	1479	1488444	G	G	A	G	G	1635	1649232	T	C	T	T	T	1861	1905117	C	G	C	C
1476048	G	T	G	G	G	G	1479	1488445	G	G	A	G	G	1635	1649235	G	A	G	G	G	1908	1908312	T	C	T	T
1476049	G	T	G	G	G	G	1479	1488448	G	G	T	G	G	1635	1649239	C	T	C	C	C	1879	1920911	A	A	A	G
1476079	C	C	C	C	T	C	1479	1488449	G	G	T	G	G	1635	1649241	G	T	G	G	G	1882	1923938	C	C	C	T
1476080	C	C	C	C	T	C	1479	1488457	G	G	T	G	G	1635	1649278	G	A	G	G	G	1882	1924702	C	C	G	C
1476086	G	G	G	G	A	G	1479	1488458	G	G	T	G	G	1635	1649301	C	G	C	C	C	1883	1925710	C	C	G	C
1476087	G	G	G	G	A	G	1479	1488956	C	C	T	C	C	1635	1649331	C	T	C	C	C	1883	1926185	G	A	G	G
1476090	T	C	T	T	T	T	1479	1489562	C	C	T	C	C	1635	1649337	C	G	C	C	C	1887	1929580	C	T	C	C
1476091	T	C	T	T	T	T	1479	1489571	C	C	A	G	A	1635	1649394	C	G	C	C	C	1888	1930350	C	C	T	C
1476096	G	C	C	C	C	C	1489971	C	T	C	C	C	C	1635	1649466	G	C	G	G	G	1888	1931083	C	T	C	C
1476097	C	G	C	C	C	C	1480	1491645	G	G	T	G	G	1635	1649514	G	A	G	G	G	1889	1931564	G	G	A	G
1476104	C	C	C	C	G	C	1480	1491673	G	G	C	G	G	1637	1650880	C	T	C	C	C	1889	1932650	T	G	T	T
1476105	C	C	C	C	G	C	1482	1495326	A	A	G	A	A	1650	1660652	G	G	T	G	G	1900	1943222	G	G	A	G
1476107	G	T	G	G	G	G	1484	1499488	C	C	T	C	C	1650	1664286	G	A	G	G	G	1909	1950464	G	A	G	G
1476112	G	G	G	G	G	G	1486	1501225	T	T	G	T	T	1654	1669358	C	C	C	A	C	1912	1952032	G	C	G	G
1476113	G	G	G	G	G	G	1492	1505535	G	T	G	G	G	1656	1672136	A	G	A	A	A	1912	1952601	C	C	C	T
1476117	T	C	T	T	T	T	1494	1507043	G	G	A	G	G	1671	1686305	C	T	C	C	C	1913	1953335	A	T	A	A
1476118	T	C	T	T	T	T	1513	1525390	G	C	G	G	G	1671	1686314	G	G	G	A	G	1917	1957629	G	G	T	G
1476119	G	C	G	G	G	G	1526703	G	T	G	G	G	G	1675	1689571	C	C	C	G	C	1920	1960391	G	G	A	G
1476120	G	C	G	G	G	G	1518	1530731	G	G	A	G	G	1677	1690984	C	G	C	C	C	1921	1960782	T	T	C	T
1476202	A	A	A	A	A	T	1523	1534929	G	T	G	G	G	1678	1692685	C	C	A	G	C	1932	1968172	G	G	T	G
1476203	A	A	A	A	A	T	1526	1538659	T	T	T	G	G	1682	1695336	A	A	A	G	A	1933	1969170	T	T	A	T
1476206	C	C	C	C	T	C	1539508	G	A	G	G	G	G	1685	1698911	G	A	A	A	A	1936	1971849	A	A	G	A
1476207	C	C	C	C	T	C	1534	1544460	A	A	A	A	A	1689	1704554	C	T	A	A	A	1939	1977247	C	A	C	C
1476226	A	G	A	A	A	A	1535	1545049	A	A	A	A	A	1691	1706119	T	C	A	A	A	1942	1980935	G	C	G	G
1476227	A	G	A	A	A	A	1535	1545125	A	A	A	A	A	1697	1711336	G	A	G	G	G	1945	1985034	T	C	T	T
1476237	A	A	A	A	A	G	1536	1546544	A	A	C	A	A	1698	1712193	C	C	C	T	C	1952	1990099	C	T	C	C
1476238	A	A	A	A	A	G	1541	1548796	G	G	G	G	G	1700	1714773	C	C	T	C	C	1955	1995944	C	T	C	C
1476253	T	T	T	T	T	C	1544	1551885	G	G	T	G	G	1700	1717003	T	T	C	T	T	1961	1999264	A	G	A	A
1476254	T	T	T	T	T	C	1552630	C	C	T	C	C	C	1705	1723697	C	A	C	C	C	1964	2001232	A	C	A	A
1476256	A	A	A	A	A	C	1546	1553633	T	T	T	C	T	1705	1724722	T	T	T	G	T	1965	2003091	C	T	C	C
1476259	A	A	A	A	A	C	1550	1557796	C	C	T	C	C	1728	1728622	C	C	C	G	G	1965	2003853	G	G	A	G
1476263	T	C	T	T	T	T	1561411	G	A	G	G	G	G	1728	1728664	C	C	T	C	C	1966	2003901	C	C	A	C
1476264	T	C	T	T	T	T	1563632	C	C	T	C	T	C	1707	1730367	A	G	A	A	A	1969	2007194	G	A	G	G
1476299	C	C	C	C	T	C	1556	1564898	C	G	C	C	C	1708	1731798	G	A	G	G	G	1969	2007502	G	G	A	G
1476300	C	C	C	C	T	C	1559	1567002	A	A	G	A	A	1708	1731990	G	G	G	A	A	1970	2008870	C	C	G	C
1476301	C	C	C	C	T	C	1559	1567109	C	T	C	C	C	1708	1732361	C	T	C	C	C	1975	2012436	G	C	G	G
1476302	C	C	C	C	T	C	1559	1567408	G	A	G	G	G	1725	1749146	T	C	T	T	T	1975	2012767	C	T	C	C
1476425	G	G	G	G	A	G	1563	1572322	G	G	A	G	G	1726	1753131	A	A	G	A	A	1975	2014185	C	C	T	C
1476427	G	G	G	G	A	G	1566	1575674	G	C	G	G	G	1730	1757541	A	G	A	A	A	1978	2016802	G	G	A	G
1476540	C	C	C	C	T	C	1569	1578626	C	C	A	G	A	1749	1778259	C	T	C	C	C	1978	2017443	C	T	C	C
1476542	C	C	C	C	T	C	1571	1580680	C	C	T	C	C	1780	1780627	G	G	C	C	C	1980	2020144	G			

Table 2: Continued

Gene No.	Exon	Coordinate	Number	Substitution	Gene No.	Exon	Coordinate	Number	Substitution
			Substitution	Gene No.			Substitution	Substitution	Gene No.
1984	2026148	C	C	T	2131	2165479	G	A	G
1985	2026776	G	G	A	2131	2165500	G	A	G
1986	2027609	T	T	G	2131	2165503	T	A	T
1988	2027617	G	A	G	2132	2168319	G	G	G
1988	2030026	C	C	T	2132	2168921	C	A	C
1991	2030855	T	T	T	2150	2185906	G	G	A
1991	2030862	T	T	C	2152	2187274	C	C	A
1991	2030942	G	G	A		2195922	T	C	T
1993	2033025	A	A	G	2161	2197230	G	A	G
2000	2039901	C	C	T	2161	2197237	A	C	A
2000	2040746	A	G	A	2161	2197239	A	C	A
2003	2043086	C	T	C	2161	2197260	G	A	G
2004	2044585	C	C	T	2161	2197263	G	A	G
2005	2045118	C	C	A	2161	2197271	T	C	T
2005	2045122	G	G	A	2161	2197273	G	C	G
2005	2045128	C	C	G	2161	2197279	T	C	T
2005	2045174	C	C	T	2161	2197288	A	G	A
2005	2045183	C	C	T	2161	2197291	C	A	C
2005	2045196	C	C	A	2161	2197294	C	A	C
2005	2045201	C	C	T	2172	2202833	G	G	A
2005	2045210	C	C	T		2204260	A	A	A
2005	2045228	C	C	T	2185	2211826	A	G	A
2005	2045288	C	C	T	2189	2216443	C	A	C
2005	2045294	C	C	A	2194	2220512	T	G	T
2008	2049495	C	T	C	2195	2221335	A	G	A
2013	2054307	G	G	A	2196	2222308	T	C	T
2013	2054805	G	G	A	2197	2223682	G	G	C
2016	2056450	T	T	C	2198	2225175	A	C	A
2017	2056805	G	A	G	2200	2225456	A	A	T
2017	2057141	G	G	A	2202	2226901	G	A	G
2020	2060423	A	A	C	2205	2229801	C	G	C
2020	2060438	A	A	C		2231486	A	G	A
2020	2060516	C	C	G	2212	2233751	G	G	C
2020	2060529	A	A	G	2221	2239349	G	G	A
	2062680	A	A	G	2222	2240572	A	C	A
2022	2064376	A	A	G	2227	2247905	G	A	G
2024	2065047	G	G	A	2227	2248179	A	A	A
2025	2068884	G	T	G	2232	2253453	T	T	T
2034	2074219	C	C	T	2232	2254221	G	A	G
2039	2080644	T	T	C	2232	2255735	T	T	T
2047	2088981	C	C	G	2242	2263760	C	C	T
2048	2090265	A	A	C	2251	2269220	A	A	A
2049	2090921	C	C	T		2271820	A	A	A
2050	2093582	G	G	A	2268	2287482	G	G	G
2052	2095504	A	G	A	2275	2294846	C	T	C
2062	2102106	C	G	C	2275	2294848	C	A	C
2062	2102193	G	G	A	2275	2294876	G	A	G
2064	2104425	C	C	T	2275	2294894	G	C	G
2064	2104479	C	T	C	2275	2294896	G	A	G
2074	2114738	T	C	T	2275	2294903	G	A	G
2074	2115210	A	A	C	2275	2294911	A	G	A
2076	2119208	A	G	A	2275	2295026	T	C	T
2077	2120686	A	G	A	2275	2295029	T	G	T
2081	2123477	G	G	A	2275	2295031	C	G	C
2082	2124333	C	C	G	2275	2295038	T	A	T
2086	2128844	G	G	A	2275	2295048	A	C	A
2089	2132152	T	T	C	2275	2296424	A	A	A
2094	2135450	A	G	A	2275	2297703	G	A	G
2096	2136619	A	G	A	2275	2298095	G	T	G
2109	2144704	T	C	T	2275	2298193	C	C	C
	2145917	C	C	G	2275	2300207	T	T	T
2113	2147126	C	C	T		2307701	G	G	A
2117	2151178	C	C	T	2279	2311302	G	G	A
2123	2156486	G	A	G	2279	2311322	C	C	A
2126	2158582	G	G	A	2279	2312264	G	G	A
2127	2159209	C	T	C	2288	2316774	C	C	T
	2162803	T	T	C	2289	2318692	G	A	G
2131	2165286	A	C	T	2295	2325009	C	C	T
2131	2165428	T	A	T		2326813	G	A	G

were used to investigate *M. tuberculosis* evolution and phylogeny. These authors examined 212 SNPs (159 synonymous SNPs [sSNPs], 35 nonsynonymous SNPs [nsSNPs] and 18 intergenic SNPs [igSNPs]) discovered through pairwise comparisons of the *M. tuberculosis* H37Rv, CDC1551 and strain 210 and *M. bovis* AF2122/97 genomes (Filliol *et al.*, 2006). In another study, a polymerase chain reaction and sequencing strategy was used to identify single nucleotide polymorphisms (SNPs) in 25 genes in the sheep (S) and cattle (C) strains of *Mycobacterium avium* subsp. *paratuberculosis* (*M. a. paratuberculosis*) and between *M. a. paratuberculosis* and *M. a. avium*. From 12,117 bp of sequence representing 26 loci across 25 genes, 11 SNPs were identified between the S and C strains in eight genes: *hsp65*, *sodA*, *dnaA*, *dnaN*, *recF*, *gyrB*, *inhA* and *pks8*. An *in silico* comparison of these *M. a. paratuberculosis* sequences and the *M. a. avium* 104 genome revealed 86 SNPs. (Marsh and Whittington 2007). Hershberg *et al.* (2008) discovered a total of 488 SNPs in 108 strains of *Mycobacterium tuberculosis* complex. The *M. canettii* strain differed from each of the other *Mycobacterium tuberculosis* complex strains at 129–145 sites (0.2% of the examined sites), while the maximum number of SNPs between any two other *Mycobacterium tuberculosis* complex strains was 46 (0.07% of the examined sites).

In this study we have generated useful information, relevant to analysis of the variations of protein formation at genic level, which directly affect the pathogenicity of organisms to the host. This study may be used for further analysis of host pathogen interactions at the pathway and product level. The whole genome SNP information can also be used for analysis of evolutionary relationships.

REFERENCES

- Bocs, S., A. Danchin and C. Medigue, 2002. Re-annotation of genome microbial coding sequences: Finding new genes and inaccurately annotated genes. BMC. Bioinform., 3: 1-5.
- Cole, S.T., 1998. Comparative mycobacterial genomics. Curr. Opin. Microbiol., 1: 567-571.
- Cole, S.T., R. Brosch, J. Parkhill, T. Garnier and C. Churcher *et al.*, 1998. Deciphering the biology of *Mycobacterium tuberculosis* from the complete genome sequence. Nature, 393: 537-544.
- Cole, S.T., K. Eiglmeier, J. Parkhill, K.D. James and N.R. Thomson *et al.*, 2001. Massive gene decay in the leprosy bacillus. Nature, 409: 1007-1011.
- Dandekar, T., M. Huynen, J.T. Regula, B. Ueberle and C.U. Zimmermann *et al.*, 2000. Re-annotating the *Mycoplasma pneumoniae* genome sequence: Adding value, function and reading frames. Nucl. Acids Res., 28: 3278-3288.
- Filliol, I., A.S. Motiwala, M. Cavatore, W. Qi and M.H. Hazbón *et al.*, 2006. Global phylogeny of *Mycobacterium tuberculosis* based on single nucleotide polymorphism (SNP) analysis: Insights into tuberculosis evolution, phylogenetic accuracy of other DNA fingerprinting systems and recommendations for a minimal standard SNP set. J. Bacteriol., 188: 759-772.
- Fleischmann, R.D., D. Alland, J.A. Eisen, L. Carpenter and O. White *et al.*, 2002. Whole-genome comparison of *Mycobacterium tuberculosis* clinical and laboratory strains. J. Bacteriol., 184: 5479-5490.
- Gaasterland, T. and M. Opera, 2001. Whole genome analysis: Annotations and updates. Curr. Opin. Struct. Biol., 11: 377-381.
- Garnier, T., K. Eiglmeier, J.C. Camus, M. Medina and H. Mansoor *et al.*, 2003. The complete genome sequence of *Mycobacterium bovis*. Proc. Nat. Acad. Sci. USA., 100: 7877-7882.
- Geluk, A., M.R. Klein, K.L. Franken, K.E. van Meijgaarden and B. Wieles *et al.*, 2005. Postgenomic approach to identify novel *Mycobacterium leprae* antigens with potential to improve immunodiagnosis of infection. Infect. Immunol., 73: 5636-5644.
- Hershberg, R., M. Lipatov, P.M. Small, H. Sheffer and S. Niemann, 2008. High functional diversity in *Mycobacterium tuberculosis* driven by genetic drift and human demography. PLoS Biol., 6: e311-e311.

- Li, L., J.P. Bannantine, Q. Zhang, A. Amonsin and B.J. May *et al.*, 2005. The complete genome sequence of *Mycobacterium avium* subspecies *paratuberculosis*. *Proc. Nat. Acad. Sci. USA.*, 102: 12344-12349.
- Marsh, I.B. and R.J. Whittington, 2007. Genomic diversity in *Mycobacterium avium*: Single nucleotide polymorphisms between the S and C strains of *M. avium* subsp. *paratuberculosis* and with *M. a. avium*. *Mol. Cell. Probes*, 21: 66-75.
- Serres, M.H., S. Gopal, L.A. Nahum, P. Liang, T. Gaasterland and M. Riley, 2001. A functional update of the *Escherichia coli* K-12 genome. *Genome Biol.*, 2: 0035.1-0035.7.