

Application of DNA Markers to Estimate Genetic Diversity of *Mycobacterium tuberculosis* Strains

KAROL KORZEKWA*¹, KORNELIA POŁOK² and ROMAN ZIELIŃSKI²

¹ Department of Environmental Microbiology, Faculty of Environmental Sciences and Fisheries,
University of Warmia and Mazury in Olsztyn, Prawocheńskiego street 1, 10-720 Olsztyn

² Department of Genetics, Faculty of Biology, University of Warmia and Mazury in Olsztyn,
Łódzki square 3, 10-967 Olsztyn

Received 10 October 2005, received in revised form 21 December 2005, accepted 22 December 2005

This article is devoted to the memory of the late Prof. W.J.H. Kunicki-Goldfinger
on the tenth anniversary of his passing away

Abstract

The obligatory human pathogen, *Mycobacterium tuberculosis*, is the most important etiological factor of tuberculosis. Unfortunately, there is little information about genetic diversity of this pathogen. The main aim of this research was the estimation of genetic diversity of *M. tuberculosis* on the basis of various categories of DNA markers. The genome of 32 strains were scanned by DNA markers such RAPD, IS6110 and catalase-peroxidase *katG* gene. All 162 identified loci were polymorphic. The genetic diversity coefficient (H_T) of *M. tuberculosis* was 0.32 for RAPD and 0.27 for IS6110. There were 14 alleles in *katG* gene. All strains were characterised by the individual molecular pattern. Genetic similarity varied from 0.13 to 0.94 (RAPD markers) and from 0 to 1 for (IS6110). *M. tuberculosis* strains did not represent a clonal structure, single source of transmission and epidemiological relationships as well. The applied DNA markers proved to be highly efficient for analysis of genetic structure of *M. tuberculosis*.

Key words: *M. tuberculosis*, genetic diversity (H), RAPD, IS6110, *katG*

Introduction

Tubercle bacillus, *Mycobacterium tuberculosis* is gram positive, acid fast and slow growing bacteria, regarded as main etiological factor of tuberculosis over the world (WHO, 2005). Research in biology and genetics of *M. tuberculosis*, resulted in very good description of ecology, genome construction and its sequencing in 1998 (Cole, 1999). However, in the last decade research focused on polymorphism of tubercle bacillus appeared. Papers concerned mutations in drug resistance genes in *M. tuberculosis* (Jou *et al.*, 2005), selected sequences related to virulence (Gao *et al.*, 2005), insertion sequences (Kulaga *et al.*, 2004) and tandem repeats, as well (Fabre *et al.*, 2004). The data of polymorphism in *M. tuberculosis* strains concern selected genes or sequence of medical and epidemiological meaning (Kulaga *et al.*, 2004). It is not possible to generalize these data on the whole genome, because they are not representative (Tazi *et al.*, 2004). Existing paradox is caused simultaneously by a very good recognition of the *M. tuberculosis* genome and the lack of data related to the level of this bacterium's genetic variability. The use of all high-tech methods of molecular biology for structure and function research of chosen sequences in tubercle bacillus (Kaduma *et al.*, 2003) should be correlated with its biological proprieties, especially with the type of reproduction and pathogenicity (Korzekwa, 2004). Therefore, we should get to know the level of *M. tuberculosis* genetic diversity to qualify genetic similarity among its strains and to estimate the efficiency of different classes of DNA markers in such investigations.

* Corresponding author: Karol Korzekwa, UWM, Department of Environmental Microbiology, Prawocheńskiego 1, 10-720 Olsztyn, mail: idefix@moskit.uwm.edu.pl.

Experimental

Materials and Methods

Thirty two strains (including 7 drug resistant) were isolated from sputum of patients hospitalised in Warmia and Mazury Centre of Tuberculosis and Pulmonary Diseases in Olsztyn. All isolates were cultured on Löwenstein-Jensen medium at 37°C and identified as *Mycobacterium tuberculosis* by Tuberculosis and Lung Diseases Institute in Warsaw. Finally, from each patient three samples were collected and joined as one bulk. Colonies from L-J medium were harvested then boiled for 10 minutes at 100°C. The DNA was extracted by CTAB (cetyl-trimethyl-ammonium bromide) method described previously by Chen and Ronald (1999) with modifications (Polok, unpublished data). Average amount of DNA obtained in each sample was measured spectrophotometrically and was 122 µg. Random Amplified Polymorphic DNA (RAPD) technique was performed according to Williams *et al.* (1990) with modifications (Polok, unpublished data). Finally short 10 nucleotides (nt) and long 18 nt scanning primers were included (Table I). For each of 8 primers, the reproducibility of the patterns was tested three times for each stock. All bands obtained on RAPD gels were numbered and its presence was estimated. IS6110-PCR reaction was performed according to Ross and Dwyer (1993) with modifications (Polok, unpublished data). Two primers were included for right and left amplification of IS6110 region, divergently. For each of the 2 primers used (Table I) the reproducibility of the patterns was tested and all present bands obtained on gels were scored as locus. PCR for *katG* was performed according to Heym *et al.* (1995) with modifications (Polok, unpublished data). Total 24 primers (Table I) were included for amplification of *katG* that was divided into 12 fragments. The whole complete analysis was performed for 7 susceptible and 7 resistant strains to estimate the number of alleles and their frequencies in *katG* locus. Mutation sites were found with first fragment of *katG* analysis based on all 32 strains with first and second primers only. For each primer that was used the reproducibility of the patterns was tested and all present bands obtained on gels were scored as “amplification” or “mutation” site.

Taking into consideration that *M. tuberculosis* is haploid the obtained molecular phenotypes corresponded to genotypes. In every locus allele “1” or “0” were observed. Allele frequencies were calculated by POPGENE 1.32 software (Yeh *et al.*, 2000). All commonly used population genetics parameters of genetic diversity were calculated as: expected heterozygosity

$$H = 1 - \sum_{i=1}^m p_i^2$$

where p is frequency of i allele in locus in population with a mean value, through loci in population (H_S) and through loci in species (H_T) according to Nei and Kumar (2000). Additionally expected heterozygosity through loci for each starter was calculated to estimate primers efficiency. Coefficient of genetics similarity (I) was calculated on the basis of shareable bands between strains ($I = 2 \times X_{1,2} / X_1 + X_2$; where: X_1 and X_2 are number of bands from two different strains, $X_{1,2}$ are shared bands). Population was grouped by UPGMA algorithm (Nei and Kumar, 2000). Dendrograms were prepared by POPGENE 1.32.

Table I
Sequences of the primers used in this study

Primer	Sequence
RAPD primers	
ISJ3 _{10nt}	5' TGCAGGTCAG 3'
OPD-01	5' ACCGCGAGGG 3'
OPD-02	5' GGACCCAACC 3'
OPD-03	5' GTCGCCGTCA 3'
OPD-05	5' TGAGCGGACA 3'
OPD-08	5' GTGTGCCCA 3'
ISJ-2	5' ACTTACCTGAGGCGCCAC 3'
ISJ-4	5' GTCGGCGGACAGGTAATG 3'
IS6110 primers	
IS-L	5' ACCCATCCTTTCCAAGAAC 3'
IS-R	5' GGCTGAGGTCTCAGATCAG 3'
<i>katG</i> primers	
katG1-L	5' GACTACGCCCAACAGCTCC 3'
katG1-R	5' GCGATAACCCCGCAAGACC 3'
katG2-L	5' GCGGGGTATCGCCGATG 3'
katG2-R	5' GCCCTCGACGGGGTATTTTC 3'
katG3-L	5' AACGGCTGTCCCGTCGTG 3'
katG3-R	5' GTCGTGGATGCGGTAGGTG 3'

Primer	Sequence
katG4-L	5' TCGACTTGACGCCCTGACG 3'
katG4-R	5' CAGGTCCGCCCATGAGAG 3'
katG5-L	5' CGACAACGCCAGCTTGGAC 3'
katG5-R	5' GGTTACCGTAGATCAGCCCC 3'
katG6-L	5' GCAGATGGGGCTGATCTACG 3'
katG6-R	5' ACCTCGATGCCGCTGGTG 3'
katG7-L	5' GCTGGAGCAGATGGGCTTG 3'
katG7-R	5' ATCCACCCGCAGCGAGAG 3'
katG8-L	5' GTCACTGACCTCTCGCTG 3'
katG8-R	5' CGCCCATGCGGTCGAAAC 3'
katG9-L	5' GCGAAGCAGATTGCCAGCC 3'
katG9-R	5' ACAGCCACCGAGCACGAC 3'
katG10-L	5' CAAAGTGTCTTCGCCGACC 3'
katG10-R	5' CACCTACCAGCACCGTCATC 3'
katG11-L	5' TGCTCGACAAGGAGAACCCTG 3'
katG11-R	5' TCCGAGTTGGACCCGAAGAC 3'
katG12-L	5' TACCAGGGCAAGGATGGCAG 3'
katG12-R	5' GCAAACACCAGCACCCCG 3'

Results

RAPD primers that scanned genomes of all strains showed 32 different molecular phenotypes. Short and long primers identified an average 32 and 30 phenotypes, respectively. Considered separately, short primers obtained 19 up to 30 different phenotypes while long – 17 and 29. All RAPD loci considered were polymorphic ($P = 100\%$). Short RAPD primers disclosed 111 loci and long primers 32 loci. Both types of primers showed all polymorphic loci. On a locus, an average number of allele (A) for both types of primers was 2.0. In most of 143 studied loci “0” allele appeared with higher frequency than “1”. Mean total genetic diversity (\bar{H}_T) of *M. tuberculosis* strains obtained by RAPD marker was 0.34. Values of \bar{H}_T for individual loci ranged from 0.06 to 0.50. The parameter of average genetic diversity on locus in susceptible population (H_S) based on RAPD marker was 0.28 for short primers and 0.34 for long ones. Within RAPD primers more efficient were short primers, like the most efficient OPD-02 and 03, which revealed the highest number of loci. Genetic diversity for a given primer (H_p) was the highest in the case of OPD-01 (0.37). The best discrimination revealed OPD-02 which alone distinguished 30 from 32 strains. Summary RAPD analysis of 32 strains revealed similarity between them on level 0.88. Values of this parameter ranged from 0.13 to 0.88. The strains were divided into three main clusters and pseudo-clusters for RAPD markers (Fig. 1). Short and long RAPD primers showed a wide range of I values between 0.13–0.87 and 0.07–0.98, respectively. Both types of primers showed the same type of clustering with differences in particular branches.

Conterminal regions of insertion sequence revealed 31 different phenotypes. IS6110 marker showed 18 polymorphic loci ($P = 100\%$). In all loci “0” alleles appeared with frequency of 0.80 and “1” with 0.20. Genetic diversity value (\bar{H}_T) for IS6110 primers was 0.30 and H_T ranged from 0.06 up to 0.50. Mean H value for susceptible population (H_S) was 0.27 and H value ranged from 0.0 to 0.49. Efficiency of IS6110 primers revealed by H_p value (0.29) was lower even in comparison to long RAPD primers. Genetic similarity among 32 strains in relation to the presence of IS6110 was very high (0.94). This parameter ranged between 0 and 1. Each strain received zero value at least once. IS6110 mobile element divided 32 strains into three clusters (Fig. 1). The second cluster consisted of two A and B subclusters. One strain was classified into pseudo-cluster.

Estimation of genetic diversity in *katG* locus in selected 14 strains of *M. tuberculosis* revealed 14 different molecular phenotypes and particular primers identified: 2 (katG8) up to 11 (katG12) strains, with average 6. The first two primers that amplified a putative promoter fragment of *katG* revealed differences in molecular phenotypes between susceptible and resistant strains. The catalase-peroxidase gene amplified by katG1-katG12 primers at 14 selected strains revealed polymorphism (100%) and the presence of 14 different alleles. Average frequency of given allele for this gene was 0.071. About 10 sites without amplification were found and named as mutation sites in the area of the first 560 bp region of *katG*. Mean frequency of such sites was about 0.517. Average genetic diversity in *katG* locus for species (H_T) based on katG1-katG12 primers in 14 strains was 0.93 with mean genetic diversity in locus on population $H = 0.86$. The most effective primers amplified subterminal and terminal region of *katG* (e.g. katG10 and 12). Comparative analysis of 14 unique alleles for 14 strains showed an average genetic similarity value on the level 0.95 with a fluctuation in range of 0.38–0.99. The catalase-peroxidase gene that was analysed by 12 primers divided 14 strains into 2 clusters (Fig. 1).

Discussion

Research of genetic diversity and H parameter value estimation of tubercle bacillus complex were conducted by enzymatic methods in the middle of '90 (Feizabadi *et al.*, 1997). Within 135 analysed strains only 8 were *M. tuberculosis* and the value of H_S for all complex was 0.10. Up to 2004 there was lack of any information about genetic diversity of *M. tuberculosis* species estimated on the basis of DNA markers (Korzekwa, 2004; Tazi *et al.*, 2004). Although still the lack of many parameters such as (P , A , H_S) is present, especially for tubercle bacillus. Four times lower value of H_S for *M. tuberculosis* complex was reported by Feizabadi *et al.* (1997) in enzymatic analysis than by DNA markers (Korzekwa, 2004; Tazi *et al.*, 2004).

Parameters of genetic diversity for tubercle bacillus that we revealed in this paper were high: P 100%; A 2.0, \bar{H}_T 0.32, and G (genotypic diversity) 100%. The level of genetic variation increased three times from 0.10 for enzymatic markers up to 0.32 for DNA markers. Similar, but not exactly the same data, was obtained by Lewandowska (2001) during analysis of grasses. However, Tazi *et al.* (2004) noted mean genetic diversity of *M. tuberculosis* on the level 0.4 based on RAPD marker analysis. He compared it with a very high value

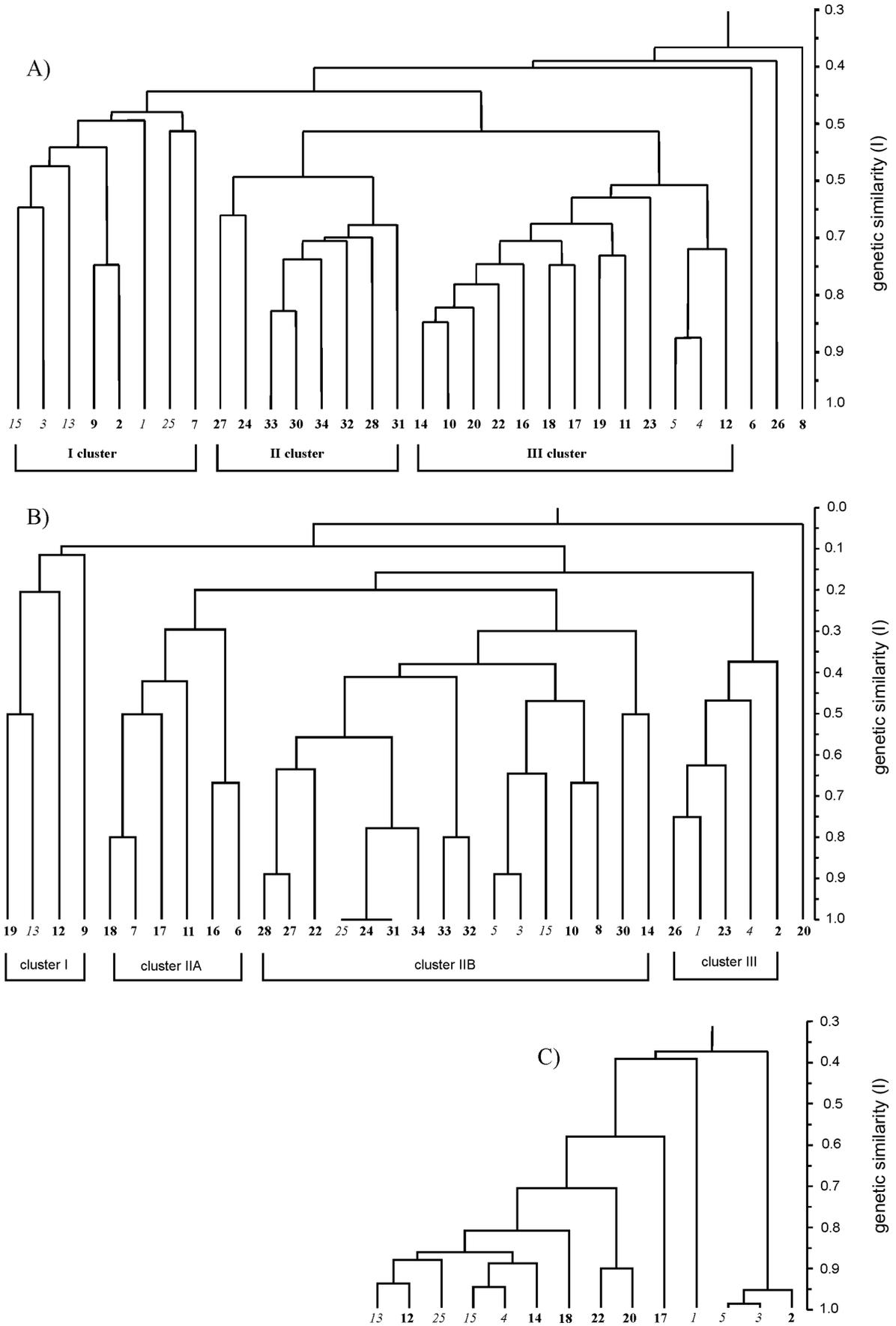


Fig. 1. Thirty two *M. tuberculosis* strains grouped by UPGMA method based on genetic similarity (I) obtained by DNA markers. Bold – susceptible strains, italic – resistant strains, a) – RAPD markers, B) – IS6110 sequence, C) – *katG* gen

of this parameter in *E. coli* (0.85) and other clonal species. In comparison to high H_S value obtained by enzymatic data for *E. coli* (0.27–0.52) (Woodward *et al.*, 1993) together with data obtained by Korzekwa (2004) and Tazi *et al.* (2004) concerned in *M. tuberculosis*, it was suggested that probably rule of three up to four times higher level of genetic diversity was revealed by DNA markers than enzymatic ones. Genetic diversity analysis of tubercle bacillus based on IS6110 marker confirmed its significant level (P 100%, A 2 and \bar{H}_T 0.30). This type of mobile element states about 3.4% of whole *M. tuberculosis* genome but it is dispersed throughout this genome simultaneously (Hatfull and Jacobs, 2000). Another problem is the number of IS6110 elements in a given strain (from zero up to several); (Hatfull and Jacobs, 2000). It means that in some cases information about IS elements in given genomes is essential. High polymorphism of analysed strains based on IS6110 marker revealed its efficiency for population genetics research. Hatfull and Jacobs (2000) joined transposition of the above mentioned genetic element with fluctuations of gene expression in bacteria. This hypothesis is not confirmed for *M. tuberculosis*. Within 18 IS6110 loci, one of them (IS6110-16) may preliminary pretend as resistant or susceptible strain type recognition and selection marker. In this locus amplification allele is the most frequent at resistant strains and is rare in susceptible ones. However, another research carried with more strains is needed to confirm this hypothesis. Moreover, high polymorphism of *katG* (Hatfull and Jacobs, 2000) was confirmed in this paper. Fourteen alleles and ten polymorphic sites in the first fragment of *katG* revealed in this paper were confirmed by the results of Saint-Joanis *et al.* (1999). Suerbaum *et al.* (2001) analysed 280 bp long fragments of selected genes in 33 *Campylobacter jejuni* strains and obtained 9 to 15 alleles (average 11). They proved relations between category of sequence and its polymorphism level. In our paper polymorphism of *katG* depends on an analysed fragment, as well.

Studies concerned genetic similarity were performed extensively for *Streptococcus* (Majewski *et al.*, 2000), *Pasteurella* (Blackall *et al.*, 1998), *Vibrio* (Farfán *et al.*, 2000), *Bacteroides* (Gutacker *et al.*, 2000). Based on enzymatic data Feizabadi *et al.* (1997) estimated Nei's genetic similarity (I) between strains of *Mycobacterium avium* complex (0.50). Belonging to one cluster *M. avium* subsp. *paratuberculosis* and grouped in 3 different clusters *M. scrofulaceum* strains revealed genetic similarity on the level of nonsibling species (I 0.3 and I 0.19, respectively). A wide range of genetic similarity (I , 0.0–1.0) of tubercle bacillus from Warmia and Mazury pointed at its allochthonic origin. Strains belonging to the same cluster were dispersed all over the voivodship.

RAPD markers that we used for *M. tuberculosis* studies generate moderate number of bands (average 6). Similar results were revealed by Gordon (1997) during genetic structure analysis of *E. coli*. The author obtained about 47 loci by only two RAPD primers what agrees with our 42 loci revealed by OPD02 and POD03. Every analysed locus during RAPD genome scanning analysis was polymorphic. The reproach of not enough reproducibility of RAPD technique between and even in the same laboratory can be eliminated by counting only strong and trusted bands. Moreover, RAPD-PCR needs high purity DNA with uniform concentration, trusted chemicals and equipment (Meunier and Grimont, 1993). The patterns that we obtained were stable and all reproducible in time and agreed with other experiences (Zervakis *et al.*, 2001). Further observation revealed an interesting methodical fact. Long RAPD primers generated similar polymorphism as 10 nt primers (two times shorter). IS6110, OPD02 and OPD01 were the best for strains identification. Generated polymorphism revealed greater differences between strains.

There are three hypotheses about *M. tuberculosis* evolution and genetic diversity level: homogeneity conception, moderate diversity and genome heterogeneity (Hatfull and Jacobs, 2000). The results presented in the recent paper confirm the theory about *M. tuberculosis* heterogeneity. Parameter of $G = 100\%$ means that absolutely the same strains were not found and shows the lack of a clonal structure of population and suggest invasive structure.

Acknowledgements. The research was financed in part by grant 6 PO4C 096 14.

Literature

- Blackall P.J., N. Fegan, G.I.T. Chew and D.J. Hampson. 1998. Population structure and diversity of avian isolates of *Pasteurella multocida* from Australia. *Microbiology* **144**: 279–289.
- Chen D.-H. and P.C. Ronald. 1999. A rapid DNA minipreparation method suitable for AFLP and other PCR applications. *Plant Mol. Biol. Rep.* **17**: 53–57.
- Cole T.S. 1999. Learning from the genome sequence of *Mycobacterium tuberculosis* H37Rv. *FEBS Lett.* **452**: 7–10.
- Fabre M., J.-L. Koeck, P. Le Fléche, F. Simon, V. Hervé, G. Vergnaud and C. Pourcel. 2004. High genetic diversity revealed by variable-number tandem repeat genotyping and analysis of *hsp65* gene polymorphism in a large collection of "*Mycobacterium canettii*" strains indicates that the *M. tuberculosis* complex is a recently emerged clone of "*M. canettii*". *J. Clin. Microb.* **42**: 3248–3255.

- Farfán M., D. Minana, M.C. Fuste and J.G. Loren. 2000. Genetic relationships between clinical and environmental *Vibrio cholerae* isolates based on multilocus enzyme electrophoresis. *Microbiology* **146**: 2613–2626.
- Feizabadi M.M., I.D. Robertson, D.V. Cousins, D.J. Dawson and D.J. Hampson. 1997. Use of multilocus enzyme electrophoresis to examine genetic relationships amongst isolates of *Mycobacterium intracellulare* and related species. *Microbiology* **143**: 1461–1469.
- Gao Q., K.E. Kripke, A.J. Saldanha, W. Yan, S. Holmes and P.M. Small. 2005. Gene expression diversity among *Mycobacterium tuberculosis* clinical isolates. *Microbiology* **151**: 5–14.
- Gordon D.M. 1997. The genetic structure of *Escherichia coli* population in feral house mice. *Microbiology*, **143**: 2039–2045.
- Gutacker M., C. Valsangiacomo and J-C. Piffaretti. 2000. Identification of two genetic groups in *Bacteroides fragilis* by multilocus enzyme electrophoresis: distribution of antibiotic resistance (*cfiA*, *cepA*) and enterotoxin (*bft*) encoding genes. *Microbiology* **146**: 1241–1254.
- Hatfull G.F. and W.R. Jacobs. 2000. Molecular genetics of mycobacteria. ASM Press, Washington.
- Heym B., P.M. Alzari, N. Honore and S.T. Cole. 1995. Missense mutations in the catalase-peroxidase gene, *katG*, are associated with isoniazid resistance in *Mycobacterium tuberculosis*. *Mol. Microbiol.* **15**: 235–245.
- Jou R., H.-Y. Chen, C.-Y. Chiang, M.-C. Yu and I.-J. Su. 2005. Genetic diversity of multidrug-resistant *Mycobacterium tuberculosis* isolates and identification of 11 novel *rpoB* alleles in Taiwan. *J. Clin. Microb.* **43**: 1390–1394.
- Kaduma E., T.D. McHugh and S.H. Gillespie. 2003. Molecular methods for *Mycobacterium tuberculosis* strain typing: a user's guide. *J. Appl. Microbiol.* **94**: 781–791.
- Korzekwa K. 2004. Ph.D. thesis. Warmia and Mazury University, Olsztyn.
- Kulaga S., M. Behr, D. Nguyen, J. Brinkman, J. Westley, D. Menzies, P. Brassard, T. Tannenbaum, L. Thibert, J.-F. Boivin, L. Joseph and K. Schwartzman. 2004. Diversity of *Mycobacterium tuberculosis* isolates in an immigrant population: evidence against a founder effect. *Am. J. Epidemiol.* **159**: 507–513.
- Lewandowska K. 2001. Ph.D. thesis. Warmia and Mazury University, Olsztyn.
- Majewski J., P. Zawadzki, P. Pickerill, F.M. Cohan and C.G. Dowson. 2000. Barriers to genetic exchange between bacterial species: *Srteptococcus pneumoniae* transformation. *J. Bacteriol.* **182**: 1016–1023.
- Meunier J.R. and P.A. Grimont. 1993. Factors affecting reproducibility of random amplified polymorphic DNA fingerprinting. *Res. Microbiol.* **144**: 373–379.
- Nei M. and S. Kumar. 2000. Molecular Evolution and Phylogenetics. Oxford Univ. Press, New York.
- Ross B.C. and B. Dwyer. 1993. Rapid, simple method for typing isolates of *Mycobacterium tuberculosis* by using the polymerase chain reaction. *J. Clin. Microbiol.* **31**: 329–334.
- Saint-Joanis B., H. Souchon, M. Wilming, K. Johnsson, P.M. Alzari and S.T. Cole. 1999. Use of site-directed mutagenesis to probe the structure, function and isoniazid activation of catalase/peroxidase, KatG, from *Mycobacterium tuberculosis*. *Biochem. J.* **338**: 753–760.
- Suerbaum S., M. Lohrengel, A. Sonnevend, F. Ruberg and M. Kist. 2001. Allelic diversity and recombination in *Campylobacter jejuni*. *J. Bacteriol.* **183**: 2553–2559.
- Tazi L., J. El Baghdadi, S. Lesjean, C. Loch, P. Supply, M. Tibayrenc and A.-L. Banuls. 2004. Genetic diversity and population structure of *Mycobacterium tuberculosis* in Casablanca, a Moroccan city with high incidence of tuberculosis. *J. Clin. Microbiol.* **42**: 461–466.
- Williams J.G., A.R. Kubelik, K.J. Livak, J.A. Rafalski and S.V. Tingey. 1990. DNA polymorphisms amplified by arbitrary primers are useful as genetic markers. *Nucl. Acids Res.* **18**: 6531–6535.
- Woodward J.M., I.D. Connaughton, V.A. Fahy, A.J. Lymbery and D.J. Hampson. 1993. Clonal analysis of *Escherichia coli* of serogroups O9, O20, and O101 isolated from Australian pigs with neonatal diarrhea. *J. Clin. Microbiol.* **31**: 1185–1188.
- World Health Organization. 2005. Global tuberculosis control: surveillance, planning, financing. WHO report 2005. Geneva.
- Yeh F., R. Yang and T. Boyle. 2000. Popgene version 1.32: A Microsoft Window-based freeware for population genetic analysis. Disp. in: <http://www.ualberta.ca/~fyeh/info.htm>. Accessed: 20 Nov. 2002.
- Zervakis G.I., G. Venturella and K. Papadopoulou. 2001. Genetic polymorphism and taxonomic infrastructure of the *Pleurotus eryngii* species-complex as determined by RAPD analysis, isozyme profiles and ecomorphological characters. *Microbiology* **147**: 3183–3194.