# Development of Genus- and Species-Specific Probe Design System for Pathogen Detection Based on 23S rDNA

PARK, JUN HYUNG[1], HEE KYUNG PARK[4], BYEONG CHUL KANG[5], EUN SIL SONG[4], HYUN JUNG JANG[4], AND CHEOL MIN KIM[1,2,3]*

[1]Busan Genome Center, [2]Medical Research Institute, and [3]Department of Biochemistry, College of Medicine, Pusan National University, Busan, South Korea
[4]Institute for Genomic Medicine, GeneIn. Co., Ltd., Busan, South Korea
[5]Division of Applied Bioengineering, Dongseo University, Busan, South Korea

**Abstract** Amplification by universal consensus sequences in pathogenic bacterial DNA would allow rapid identification of pathogenic bacteria, and amplification of genus-specific and species-specific sequences of pathogenic bacterial DNA might be used for genotyping at the genus and species levels. For design of probes for molecular diagnostics, several tools are available as stand-alone programs or as Web application. However, since most programs can design only a few probe sets at one time, they are not suitable for large-scale and automatic probes design. Therefore, for high-throughput design of specific probes in diagnostic array development, an automated design tool is necessary. Thus, we developed a Web-based automatic system for design of genus-specific and species-specific probes for pathogen detection. The system is available at http://www.miprobe.com.

**Key words:** Pathogenic bacteria, pathogen detection, PCR, oligonucleotide array, probe design, 23S rDNA

The reemergence of bacterial infections and antibiotics resistance makes it clinically important to quickly identify the pathogenic bacteria. Therefore, a considerable effort has been devoted to develop rapid, convenient, and accurate methods for the detection of these unfavorable organisms. Development of new diagnostics methods, using powerful molecular tools such as PCR and oligonucleotide array, has potential to advance traditional diagnostics to modern molecular diagnostics [15]. The modern molecular diagnostics in clinical microbiology include microbial genotyping and drug-resistance detection, based on sequence diversity and specificity. These methods also permit simultaneous monitoring

and analysis of a large number of target genes, depending on sequence diversity. Accurate probes design is actually one of the most important factors in successful PCR or oligonucleotide array development.

Several target genes, such as 16S rDNA, rpoB, hsp65, internal transcribed spacer (ITS), and gyrB, have been used as targets for the genotyping of microorganisms [4, 19]. The important features of these target genes are that these genes are present in all bacteria and contain both conserved and polymorphic regions [20, 13].

Specifically, rDNA possesses highly conserved regions that are suitable as sites for PCR primers that recognize a large group of organisms as well as variable regions that provide signatures for specific identification.

For many years, 16S rDNA has been an important target gene for determining the phylogenetic relationship between bacteria. The features of this molecular target that made it a useful phylogenetic tool would also make it a useful molecular target for bacterial detection and identification in the clinical laboratory. Sequence analysis of the 16S rDNA is a powerful method for identifying new pathogens in patients with suspected bacterial infection [8, 14]. The ITS regions between rDNAs are highly variable within many species, and this variation has been used for typing clinical isolates [7]. Park et al. [16] introduced the usefulness of the ITS region as a genetic marker for identification and differentiation of mycobacteria at the species level.

Recently, sequence data for 23S rDNA have become available for bacterial species. This region shows more variation than 16S rDNA between important pathogenic species, and therefore, probe design using 23S rDNA might be highly useful for clinical diagnosis. Anthony [18] introduced a rapid detection and identification system that uses universal PCR primers to amplify a variable region of bacterial 23S rDNA, followed by reverse hybridization of

*Corresponding author
Phone: 82-51-240-7725; Fax: 82-51-248-1118;
E-mail: kimcm@pusan.ac.kr

the products to a panel of oligonucleotides [18]: Amplication of a variable region of bacterial 23S rDNA by universal sequences would allow rapid identification of pathogenic bacteria. Furthermore, additional analysis of genus-specific and species-specific sequences of the variable region might be used for genotyping of pathogenic bacteria at the genus and species levels [12].

Although there are a number of software tools to design primers or probes, most of them require laborious manual work and are not systemically integrated. As basic bioinformatics tools to design probes, one needs various tools including databases such as NCBI GenBank [2], multiple alignment tools such as ClustalW [22], phylogenetic analysis tools such as PhyloDraw [10], and electronic sequence comparison tools such as BLAST [1]. As specific bioinformatics tools to design primers or probes, several programs are presently available as stand-alone or as Web application, such as PrimerMaster [17], PrimeArray [5], OligoArray [9], PRIDE [6], and Web Primer (http://genome-www2.stanford.edu/cgi-bin/SGD/web-primer).

However, most of these programs can design only a few probe sets at one time and are therefore less suitable for large-scale probes design [21, 3]. Furthermore, none is capable of automatically performing high-throughput design of probes for diagnostic oligonucleotide array development.

Therefore, to overcome time-consuming, laborious, and error-prone work in the high-throughput design of specific probes for diagnostic array development, we developed a Web-based integrated system for large-scale genus-specific and species-specific probes design for pathogen identification.

## Implementation

The system is implemented by the following four major components: 1) collection of target gene sequences; 2) construction of consensus sequences; 3) collection of probe candidates; and 4) *in-silico* evaluation of candidates by BLAST.

## METHODS

### Collection of Target Gene Sequences

In large-scaled design of probes, one of the most time-consuming tasks is collection of the target gene sequences. We can obtain target gene sequences from NCBI GenBank and by our own direct DNA sequencing. Sometimes, we use advanced search in order to correct more exact sequences.

Nevertheless, NCBI Entrez shows interest target gene sequences and non-interest target gene sequences at the same time. For example, we searched "*Helicobacter pylori* 23S rRNA" in the NCBI Entrez. Searched sequences included a total of 121 sequences related to the "*Helicobacter pylori* 23S rRNA*." Among them, only 52 sequences were directly related to the query. In this case, a researcher needs to manually check the reports one by one. The query might be different, like "23S rRNA," "23S rDNA," and "23S ribosomal RNA" and so on, and therefore, it is extremely laborious work. Sometimes, the number of searched sequences may be more than one thousand. In such a case, it is impossible to collect the target gene sequences manually through NCBI GenBank. Therefore, we developed NGSS (NCBI GenBank Search System) in order to get target gene sequences easily and correctly. NGSS consists of a customized GenBank nucleotide sequence list, filtering system, and Entrez search robot that is made by bioperl modules.

Using the automated and integrated Web-based tool, we have designed genus-specific and species-specific probes from 23S rDNA sequences of 41 genera including 133 pathogenic bacteria.

### Construction of Consensus Sequences

In molecular biology and bioinformatics, a consensus sequence is a way of representing the results of a multiple sequence alignment, where related sequences are compared with each other. The consensus sequence shows which
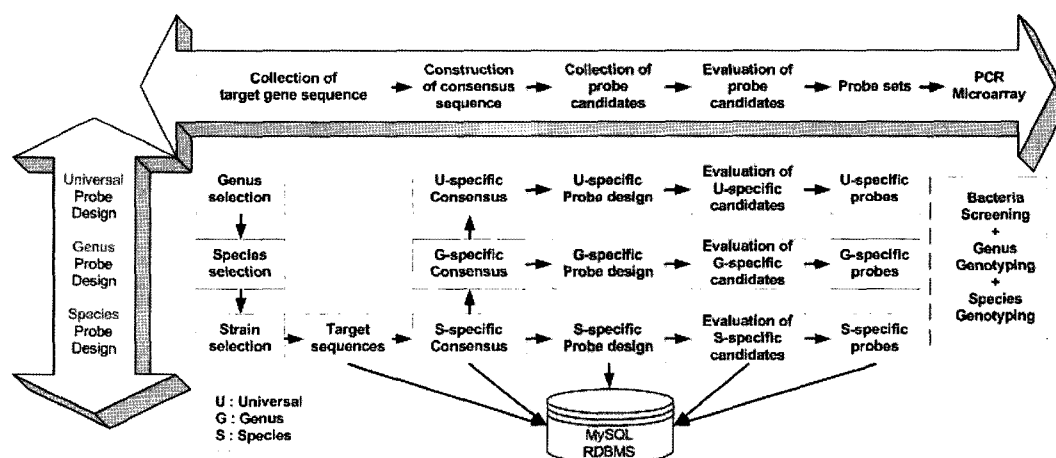


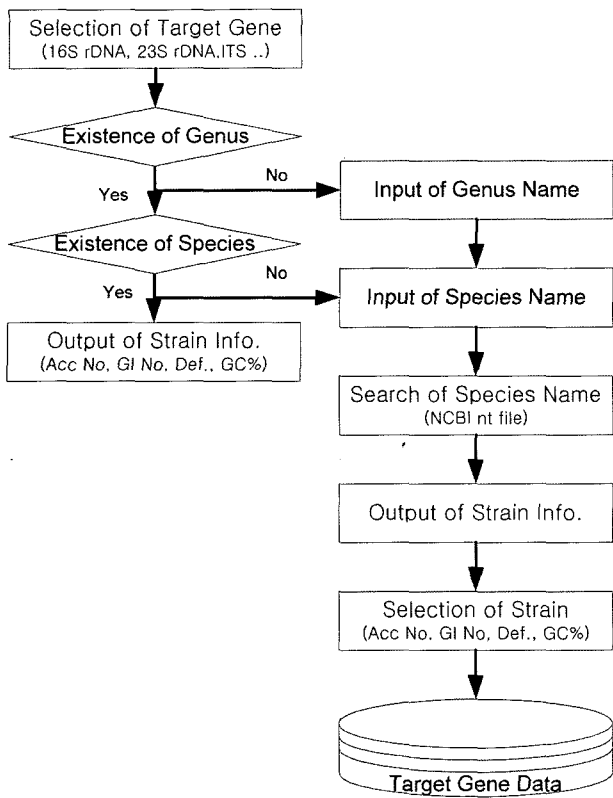**Fig. 1.** Concept and structure of developing the probe design system.

**Fig. 2.** A simplified flow chart, describing the collection process of a target gene sequence.



**Fig. 4.** Screenshot of the genus consensus sequence construction interface.

sequences (or residues) are conserved, and which sequences (or residues) are variable.

Using multiple sequence alignment, we were able to get various result files. Each result file has an extension tag like ".txt," ".aln," ".dnd," and .con." The file that has the ".txt" extension tag contains the result of multiple sequence alignment. The extension ".aln" means the alignment output file, ".dnd" means a "guide tree" for building the cladogram and phylogram, and ".con" means consensus sequence file.
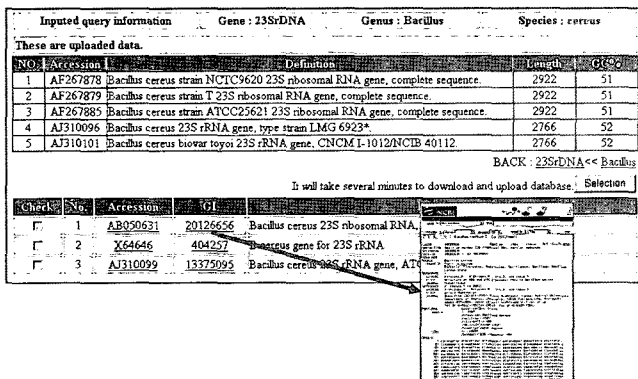
We can check the correlationship of target gene sequences using "Phylodendron" (http://iubio.bio.indiana.edu/treeapp/treeprint-form.html), which is linked to a consensus sequence viewer. Phylodendron is an application for drawing phylogenetic trees, which are used in evolutionary biology: It reads Newick format, and then displays graphical views of the phylogenetic tree.

A species-consensus sequence is the result of multiple sequence alignment of strain sequences, and a genus-consensus sequence is the result of multiple sequence alignment of species-consensus sequences. It is linked to each other in the vertical relation.

## Collection of Probe Candidates

The algorithm for collection of probe candidates is divided into two distinct parts. In the first step, it scans along the consensus sequence in a sliding window scheme in one-nucleotide shifts. It isolates conserved regions and nonconserved regions that include "N," uncertain sequence, in the consensus sequence. In the second step, the conserved regions of the consensus sequence are used for computation of all of the relevant parameters, including melting temperature (Tm), GC contents, length, thermodynamic stability, and self-complementarity.

Melting temperature is often used as an input for a primer or probe design program, because the researcher requires a



**Fig. 3.** Screenshot of the target gene selection interface.
The screenshot demonstrates how to get the target gene sequences.



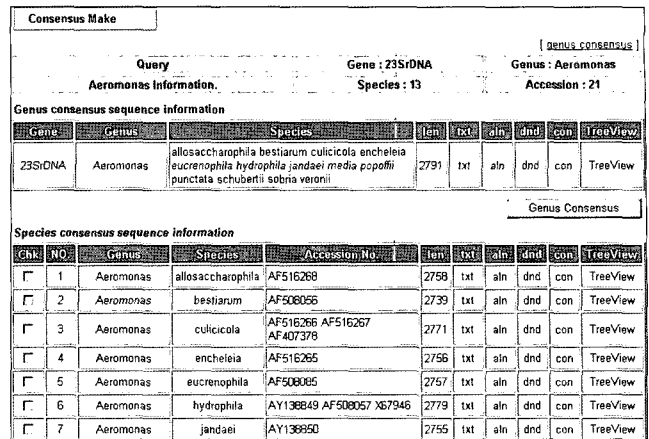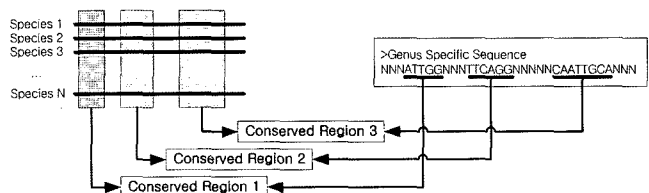**Fig. 5.** Scheme of consensus sequence construction and selection of conserved regions.

**Fig. 6.** Scheme of candidates selection by window sliding algorithm.



**Fig. 7.** Screenshot of the probe design interface.

primer or probe that will work under specified reaction conditions. One method for calculating melting temperature is the nearest-neighbor method. Melting temperature is calculated as a function of the sums of the entropy and enthalpy of the consecutive pairs of amino acids [11].

The stability of the primer-template DNA duplex is important for primer design, because it will affect the efficiency of priming. The GC content describes the stability of the primer-template duplex, because different energies are required to break apart GC pairs that have three hydrogen bonds, and AT pairs that have only two.

Developed system accepts minimum and maximum oligonucleotide length and GC contents instead of a fixed
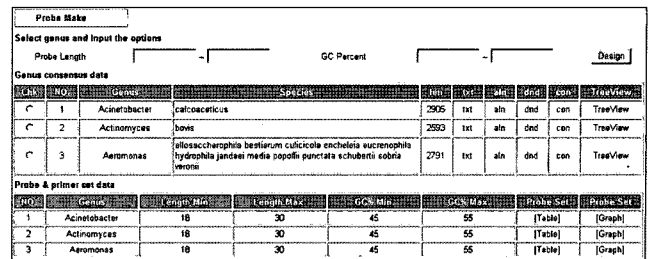
length and GC contents. Within the specified range, it can adjust the length of oligonucleotides to achieve greater specificity and uniformity among all oligonucleotides. It is accomplished by a built-in parameters setting, including Tm range, sequence complexity check, and difference on Tm, *etc.*

### *In-Silico* Evaluation of Candidates by BLAST

Once the probe candidates are collected, the next step is evaluation of candidate specificity between these probe candidates and known sequences in the NCBI. We have used two kinds of BLAST system: One is a stand-alone BLAST, which allows us to create custom-searchable
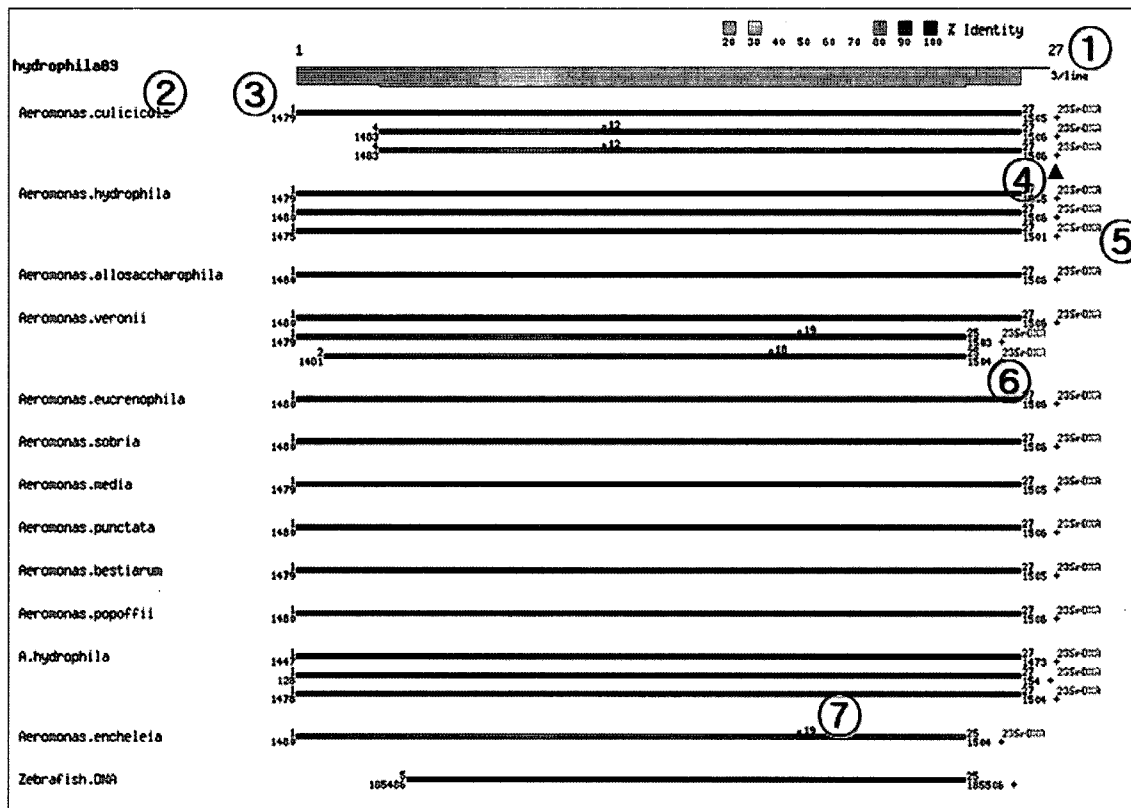


**Fig. 8.** Screenshot of the graph taxonomy BLAST result.
(1) Probe length. (2) Blast hit grouping of taxonomic browser type. (3) Hit position of query and subject sequence. (4) Strand of subject sequence. (5) Expression of 23S rDNA hit. (6) Expression of mismatch part in the 3' part. (7) Expression of mismatch part in the middle position.

**Table 1.** Primers for genus-specific identification.

| Genus | Primer name | Sequence | Size (bp) |
|---|---|---|---|
| *Enterococcus* | Entc-310F | 5'-TTGGGGTTGAGGACTCC(G/A)A-3' | 599 |
| | Entc-909R | 5'-GTGCTCTACCTCCATCATTCT-3' | |
| *Mycobacteria* | MB-2089F | 5'-TCCGTGCGAAGTCGCAAGACG-3' | 962 |
| | MB-3015R | 5'-AGGTTTCCCGCTTAGATGCT-3' | |
| *Streptococcus* | Str-791F | 5'-CAGGGCACGTTGAAAAGTGCTT-3' | 804 |
| | Str-1595R | 5'-GTGTGACATCACTAACTTCGC-3' | |

database of the probe candidates, and another is remote BLAST, which is run against NCBI GenBank through Internet. It is made by Bioperl modules and CGI scripts. Bioperl (http://www.bioperl.org) is an open-source community of bioinformatics professionals that develop and maintain code libraries and applications written in the Perl programming language. Results of the stand-alone BLAST search are parsed and the candidates are filtered according to its design type; e.g., species-specific design and genus-specific design. Candidates passed through the stand-alone BLAST results are transferred to remote BLAST and searched against NCBI GenBank. Results of the remote BLAST search are parsed and updated into the database. We have designed a graphical summary of a BLAST report using Thomas Boutell's GD graphical library (http://www.boutell.com/gd), which has been ported to Perl by Lincoln Stein (http://stein.cshl.org/WWW/software/GD). In order to graphically summarize a BLAST report, we have parsed the BLAST result in the NCBI tabular format. The fields that are separated in the tabular format are "qurey id," "subject id," "percent identity," "alignment length," "mismatches," "gap openings," "query start," "query end," "subject start," "subject end," "e-value," and "bit score." The graphic summary sorts the BLAST hits according to the species of the target sequence, so that all of the hits to the same organism will appear together. Within each species, the BLAST hits are sorted by score (as for the normal BLAST output). The species themselves are sorted by the strength of their strongest BLAST hit scores. Furthermore it shows the BLAST hits by sense strand, presence of mismatch regions, mismatch position, and presence of target gene in the case of complete sequences hit.
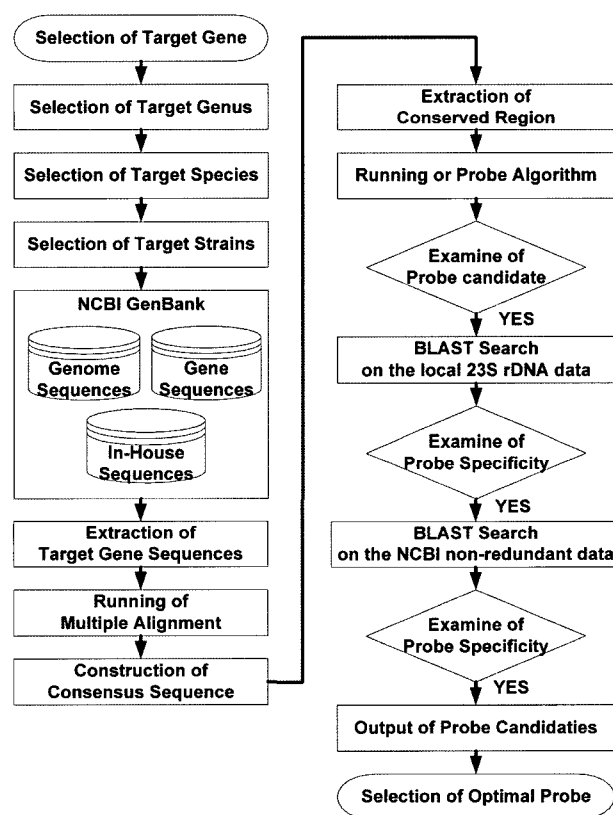
**Verification of Primers by Experiment**

The PCR-based technique has been exploited for the diagnosis of pathogenic bacteria. It is much less labor-intensive and tedious than the traditional methods. The PCR-based technique is also more cost-effective than the oligonucleotide array. Therefore, we tested the specificity of the representative primers by PCR.

The performance of the designed PCR primers was tested by using human pathogenic bacteria (reference strain), including 25 genera and 100 species. The genus included *Acinetobacter*, *Aeromonas*, *Bacillus*, *Bacteroides*,

*Campylobacter*, *Citrobacter*, *Clostridium*, *Corynebacterium*, *Enterobacter*, *Enterococcus*, *Escherichia*, *Haemophilus*, *Klebsiella*, *Listeria*, *Mycobacteria*, *Mycoplasma*, *Neisseria*, *Pseudomonas*, *Salmonella*, *Serratia*, *Shigella*, *Staphylococcus*, *Streptococcus*, *Vibrio*, and *Yersinia*.

The PCR reaction mixture was as follows: 100 mM KCl, 20 mM Tris HCl (pH 9.0), 1% Triton X-100, 10 mM deoxynucleotide triphosphates (dATP, dGTP, dTTP, and dCTP), 1.5 mM $MgCl_2$, 1 U Tag polymerase (all QIAGEN, U.S.A.), 0.1 μM each of the primers listed in Table 1. Five μl of extracted DNA was added to the reaction and the reaction mixture was heated to 94°C for 3 min, followed by 35 cycles of 94°C for 1 min, 50°C for 1 min, and 72°C for 1 min. A final extension was carried out at 72°C for 10 min.



**Fig. 9.** A simplified flow chart, describing the species-specific probes design of the system.

## Systems

We have designed a MySQL relational database to store the target gene sequences, consensus sequences, probe candidates, and BLAST results. We have also designed a Web interface and connected it to the database using CGI script. This allows one to query the database in a user-friendly manner. Using bioperl modules, various bioinformatics tools such as remote BLAST were performed.

## RESULTS AND DISCUSSION

The design of large-scale probes is laborious and tedious work if the analysis is not automated in some way. A couple of tools for this purpose already exist. However, neither of these tools supports extensive user interaction
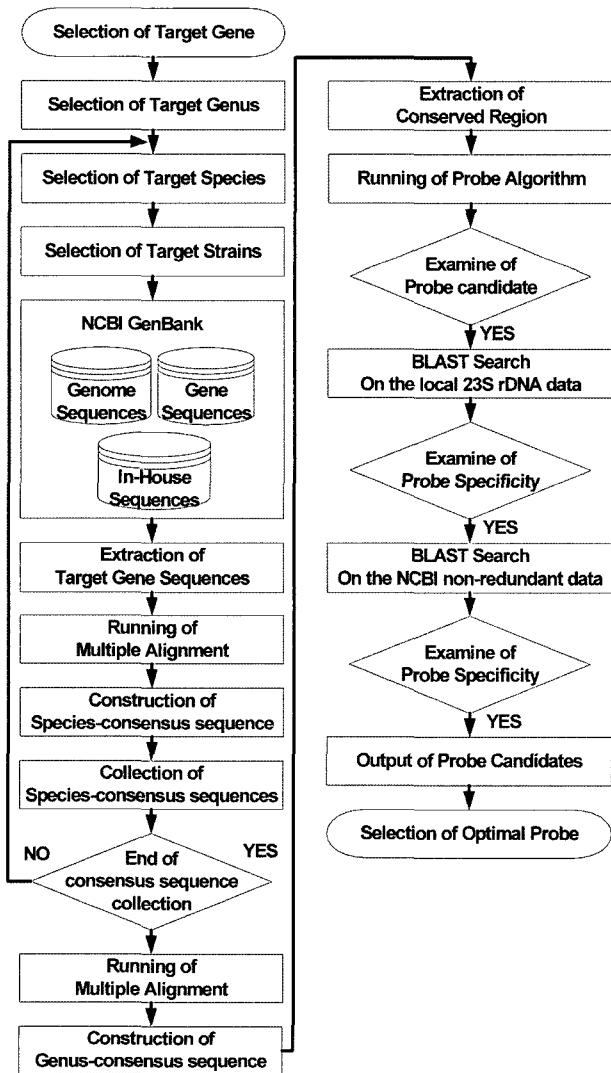


**A**

M N 1 2 3 4 5 6 7 8 9

**B**

M N 1 2 3 4 5 6 7 8 9 10
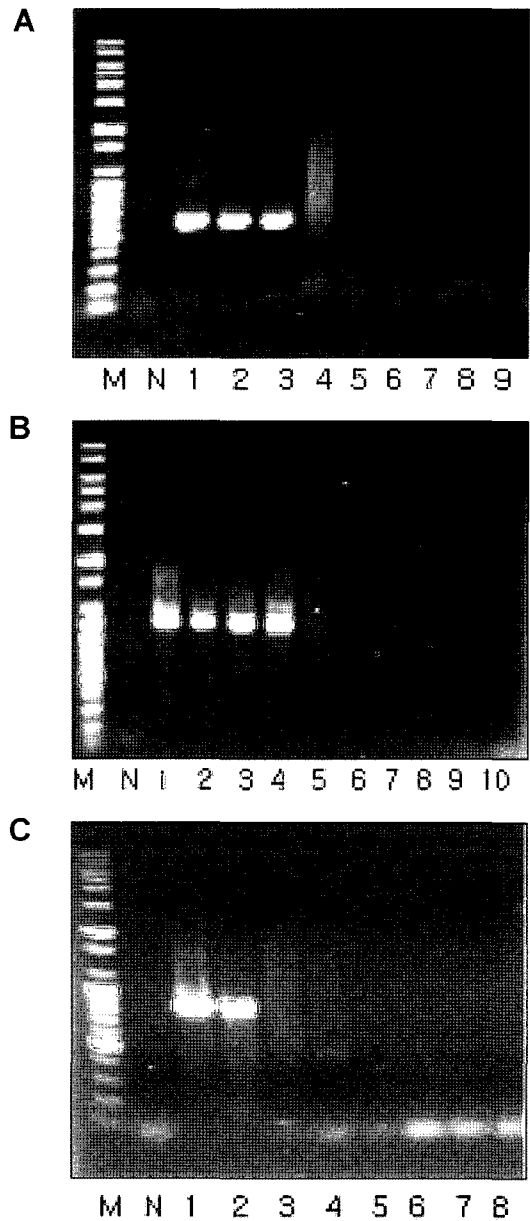
**C**

M N 1 2 3 4 5 6 7 8

**Fig. 11. A.** Gel electrophoresis of PCR products using *Enterococcus* genus-specific primers.
Lane M, 100-bp size markers; lane N, negative control; lane 1, *Enterococcus faecalis*; lane 2, *Enterococcus faedium*; lane 3, *Enterococcus hirae*; lane 4, *Aeromonas hydrophila*; lane 5, *Mycobacterium xenopii*; lane 6, *Mycobacterium falconis*; lane 7, *Streptococcus anginosus*; lane 8, Human blood DNA; lane 9, Hepetitis B virus DNA. **B.** Gel electrophoresis of PCR products using *Mycobacteria* genus-specific primers. Lane M, 100-bp size markers; lane N, negative control; lane 1, *Mycobacterium xenopi*; lane 2, *Mycobacterium flavescence*; lane 3, *Mycobacterium simiae*; lane 4, *Mycobacterium tuberculosis*; lane 5, *Aeromonas hydrophila*; lane 6, *Mycobacterium falconis*; lane 7, *Streptococcus anginosus*; lane 8, *Enterococcus faecalis*; lane 9, Human blood DNA; lane 10, Hepatitis B virus DNA. **C.** Gel electrophoresis of PCR products using *Streptococcus* genus-specific primers. Lane M, 100-bp size markers; lane N, negative control; lane 1, *Streptococcus anginosus*; lane 2, *Streptococcus bovis*; lane 3, *Aeromonas hydrophila*; lane 4, *Mycobacterium falconis*; lane 5, *Mycobacterium xenopi*; lane 6, *Enterococcus faecalis*; lane 7, Human blood DNA; lane 8, Hepatitis B virus DNA.



**Fig. 10.** A simplified flow chart, describing the genus-specific probes design of the system.

and they give only limited overview of the design process. Accordingly, an automated, Web-based, and integrated relational database design system has become essential.

We could get about 1,300 "23s rDNA" sequences from 41 genera, including 133 pathogenic bacteria species, within 5 h. In general, it requires more than 2 weeks manually. Multiple sequence alignment and probe detection algorithms are adopted for the target gene sequences, and BLAST system is used to evaluate the probe candidates.

In the present study, the performance of the designed PCR primers was demonstrated by using human pathogenic bacteria (reference strain), including 25 genera and 100 species. We selected the optimal genus-specific primers through PCR with the 3 genus (*Enterococcus*, *Mycobacteria*, and *Streptococcus*) as well as other pathogenic bacteria, human blood DNA, and hepatitis B virus DNA (Figs. 11A–11C).

PCR with the *Enterococcus* genus-specific primers, Entc-310F and Entc-909R, amplified an approximately 599-bp fragment only in the *Enterococcus*. No fragment for other pathogenic bacteria, human blood DNA, and hepatitis B virus DNA was found by PCR (Fig. 11A). PCR with the *Mycobacteria* genus-specific primers, MB-2089F and MB-3015R, amplified an approximately 962-bp fragment only in the *Mycobacteria*. No fragment for other pathogenic bacteria, human blood DNA, and hepatitis B virus DNA was found by PCR (Fig. 11B). PCR with the *Streptococcus* genus-specific primers, Str-791F and Str-1595R, amplified an approximately 804-bp fragment only in the *Streptococcus*. No fragment was found by PCR for other pathogenic bacteria, human blood DNA, and hepatitis B virus DNA (Fig. 11C).

Using the tested primers to design the system, we could obtain genus- and species-specific amplifications. A tutorial that demonstrates the collection of target gene sequences, construction of consensus sequences, collection of probe candidates, and *in-silico* evaluation of candidates is provided on the program Website.

## Acknowledgment

## REFERENCES

1. Altschul, S. F., T. L. Madden, A. A. Schaffer, J. Zhang, Z. Zhang, W. Miller, and D. J. Lipman. 1997. Gapped BLAST and PSI-BLAST: A new generation of protein database search programs. *Nucleic Acids Res.* **25**: 3389–3402.

2. Dennis A. B., I. Karsch-Mizrachi, D. J. Lipman, J. Ostel, and D. L. Wheeler. 2005. GenBank. *Nucleic Acids Res.* **33**: D34–D38.

3. Dong, X., L. Guangshan, W. Liyou, Z. Jizhong, and X. Ying. 2002. Primegens: Robust and efficient design of gene-specific probes for microarray analysis. *Bioinformatics* **18**: 1432–1437.

4. Fukushima, M., K. Kakinuma, H. Hayashi, H. Nagai, K. Ito, and R. Kawaguchi. 2003. Detection and identification of *Mycobacterium* species isolates by DNA microarray. *J. Clin. Microbiol.* **41**: 2605–2615.

5. Gunter, R., D. Michaela, F. M. Thomas, and D. Christoph. 2001. PrimeArray: Genome-scale primer design for DNA-microarray construction. *Bioinform. Applic. Note* **17**: 98–99.

6. Haas, S., M. Vingron, A. Poustka, and S. Wiemann. 1998. Primer design for large scale sequencing. *Nucleic Acids Res.* **26**: 2006–2012.

7. Park, H. K., H. J. Jang, E. S. Song, C. H. Chang, M. K. Lee, S. K. Jeong, J. H. Park, B. C. Kang, and C. M. Kim. 2005. Detection and genotyping of mycobacterium species from clinical isolates and specimens by oligonucleotide array. *J. Clin. Microbiol.* **43**: 1782–1788.

8. Patel, J. B. 2004. 16S rRNA *Gene Sequencing for Bacterial Pathogen Identification in the Clinical Laboratory. Molecular Diagnosis.* Vol. 6. No. 4.

9. Rouillard, J.-M., C. J. Herbert, and M. Zuker. 2002. OligoArray: Genome-scale oligonucleotide design for microarrays. *Bioinform. Applic. Note* **18**: 486–487.

10. Choi, J.-H., H.-Y. Jung, H.-S. Kim, and H.-G. Cho. 2000. PhyloDraw: A phylogenetic tree drawing system. *Bioinform. Applic. Note* **16**: 1056–1058.

11. Kampke, T., M. Kieninger, and M. Mecklenburg. 2001. Efficient primer design algorithm. *Bioinformatics* **17**: 214–225.

12. McCabe, K. M., Y.-H. Zhang, B.-L. Huang, E. A. Wager, and E. R. B. McCabe. 1999. Bacterial species identification after DNA amplification with a universal primer pair. *Molec. Genet. Metab.* **66**: 205–211.

13. Kim, B. S., H. M. Oh, H. J. Kang, S. S. Park, and J. S. Chun. 2004. Remarkable bacterial diversity in the tidal flat sediment as revealed by 16S rDNA analysis. *J. Microbiol. Biotechnol.* **14**: 205–211.

14. Kim, M. K., H. S. Kim, B. O. Kim, S. Y. Yoo, J. H. Seong, D. K. Kim, S. E. Lee, S. J. Choe, J. C. Park, B. M. Min, M. J. Jeong, D. K. Kim, Y. K. Shin, and J. K. Kook. 2004. Multiplex PCR using conserved and species-specific 16S rDNA primers for simultaneous detection of *Fusobacterium nucleatum* and *Actinobacillus actinomycetemcomitans*. *J. Microbiol. Biotechnol.* **14**: 110–115.

15. Lee, J. W., I. J. Jun, H. J. Kwun, K. L. Jang, and J. H. Cho. 2004. Direct identification of *Vibrio vulnificus* by PCR targeting elastase gene. *J. Microbiol. Biotechnol.* **14**: 284–289.

16. Park, H., H. Jang, C. Kim, B. Chung, C. L. Chang, S. K. Park, and S. Song. 2000. Detection and identification of mycobacteria by amplification of the internal transcribed

spacer regions with genus- and species-specific PCR primers. *J. Clin. Microbiol.* **38**: 4080–4085.

17. Proutski, V. and E. C. Holmes. 1996. Primer Master: A new program for the design and analysis of PCR primers. *Comput. Appl. Biosci.* **12**: 253–255.

18. Anthony, R. M. 2000. Rapid diagnosis of bacteremia by universal amplification of 23S ribosomal DNA followed by hybridization to an oligonucleotide array. *J. Clin. Microbiol.* **38**: 781–788.

19. Ryou, S. M., J. M. Kim, J. H. Yeon, H. L. Kim, H. Y. Go, E. K. Shin, and K. S. Lee. 2005. Species-specific cleavage by RNase E-like enzymes in 5S rRNA maturation. *J. Microbiol. Biotechnol.* **15**: 1100–1105.

20. Soini, H. and J. M. Musser. 2001. Molecular diagnosis of mycobacteria. *Clin. Chem.* **47**: 809–814.

21. Weckx, S., P. De Rijk, C. Van Broeckhoven, and J. Del-Favero. 2004. SNPbox: Web-based high-throughput primer design from gene to genome. *Nucleic Acids Res.* **32**: 170–172.

22. Thompson, J. D., D. G. Higgins, and T. J. Gibson. 1994. CLUSTAL W: Improving the sensitivity of progressive multiple sequence alignment through sequence weighting, positions-specific gap penalties and weight matrix choice. *Nucleic Acids Res.* **22**: 4673–4680.