

Automated Identification of Medically Important Bacteria by 16S rRNA Gene Sequencing Using a Novel Comprehensive Database, 16SpathDB[∇]

Patrick C. Y. Woo,^{1,2,3,4†*} Jade L. L. Teng,^{3†} Juilian M. Y. Yeung,^{3†} Herman Tse,^{1,2,3,4}
Susanna K. P. Lau,^{1,2,3,4*} and Kwok-Yung Yuen^{1,2,3,4}

State Key Laboratory of Emerging Infectious Diseases,¹ Research Centre of Infection and Immunology,² Department of Microbiology,³ and Carol Yu Centre for Infection,⁴ The University of Hong Kong, Hong Kong

Received 22 November 2010/Returned for modification 3 January 2011/Accepted 27 February 2011

Despite the increasing use of 16S rRNA gene sequencing, interpretation of 16S rRNA gene sequence results is one of the most difficult problems faced by clinical microbiologists and technicians. To overcome the problems we encountered in the existing databases during 16S rRNA gene sequence interpretation, we built a comprehensive database, 16SpathDB (<http://147.8.74.24/16SpathDB>) based on the 16S rRNA gene sequences of all medically important bacteria listed in the *Manual of Clinical Microbiology* and evaluated its use for automated identification of these bacteria. Among 91 nonduplicated bacterial isolates collected in our clinical microbiology laboratory, 71 (78%) were reported by 16SpathDB as a single bacterial species having >98.0% nucleotide identity with the query sequence, 19 (20.9%) were reported as more than one bacterial species having >98.0% nucleotide identity with the query sequence, and 1 (1.1%) was reported as no match. For the 71 bacterial isolates reported as a single bacterial species, all results were identical to their true identities as determined by a polyphasic approach. For the 19 bacterial isolates reported as more than one bacterial species, all results contained their true identities as determined by a polyphasic approach and all of them had their true identities as the “best match in 16SpathDB.” For the isolate (*Gordonibacter pamelaee*) reported as no match, the bacterium has never been reported to be associated with human disease and was not included in the *Manual of Clinical Microbiology*. 16SpathDB is an automated, user-friendly, efficient, accurate, and regularly updated database for 16S rRNA gene sequence interpretation in clinical microbiology laboratories.

In the last decade, as a result of the widespread use of PCR and DNA sequencing, 16S rRNA gene sequencing has played a pivotal role in accurate identification of medically important bacteria in both clinical microbiology and research laboratories (26). For the management of individual patients, accurate and objective identification of clinical isolates has assisted clinicians in the choice and duration of antibiotic therapy, as well as infection control measures. On the population scale, accurate identification has greatly improved our understanding of the epidemiology and, hence, the empirical treatment of infectious disease syndromes. In the past 10 years, our group and others have used this technology for the identification of a large number of medically important bacteria, resulting in major impacts on infectious diseases.

Interpretation of 16S rRNA gene sequence results is one of the most difficult problems faced by inexperienced clinical microbiologists and technical staff, despite the wide range of software and databases available. The best-known software and databases include GenBank (1), the Ribosomal Data-

base Project (RDP-II) (2, 3, 13), MicroSeq (15, 17), Ribosomal Differentiation of Medical Microorganisms (RIDOM) (4–6), and the SmartGene Integrated Database Network System (SmartGene IDNS) (16). The databases of RDP-II and SmartGene IDNS contain sequences downloaded from GenBank, whereas all sequences in the databases of RIDOM and MicroSeq were obtained by sequencing the 16S rRNA genes of bacterial strains from culture collections. Due to the large number of unvalidated 16S rRNA gene sequences in GenBank, it is often not easy for inexperienced users to decide whether the “first hit” or the “closest match” is the real identity of a bacterial isolate. As for the other software and databases, the usefulness is further limited by the choice of bacterial species in the database. If a bacterial species is not included in the database, it would never be the identity of an isolate. If the database includes bacterial species with minimal differences in their 16S rRNA gene sequences, and hence, that cannot be identified with confidence by 16S rRNA gene sequencing, they may also give rise to wrong identification if the software reports only that the “first hit” or “closest match” is the identity of the bacterium.

In view of these problems, in 2005 we started developing our own database, which includes the most representative 16S rRNA gene sequences of all medically important bacteria listed in the most current edition of the *Manual of Clinical Microbiology* (14), for identification of medically

* Corresponding author. Mailing address: Department of Microbiology, The University of Hong Kong, University Pathology Building, Queen Mary Hospital, Hong Kong. Phone: (852) 22554892. Fax: (852) 28551241. E-mail for P. C. Y. Woo: pcywoo@hkucc.hku.hk. E-mail for S. K. P. Lau: skplau@hkucc.hku.hk.

† P.C.Y.W., J.L.L.T., and J.M.Y.Y. contributed equally to the study.

[∇] Published ahead of print on 9 March 2011.

TABLE 1. Identification of clinical bacterial isolates using 16SpathDB

Strain no.	Identification by polyphasic approach	Reference	Identification results reported by 16SpathDB		Best match in 16SpathDB	Nucleotide identity (%)	Second-best match in 16SpathDB	Nucleotide identity (%)	Third-best match in 16SpathDB	Nucleotide identity (%)
			Bacterial species	Nucleotide identity (%)						
1	<i>Abiotrophia defectiva</i>	21	<i>Abiotrophia defectiva</i>	99.54	<i>Abiotrophia defectiva</i>	99.54	<i>Granulicatella adiacens</i>	92.77	<i>Granulicatella elegans</i>	92.67
2	<i>Actinobaculum schaalii</i>		<i>Actinobaculum schaalii</i>	99.35	<i>Actinobaculum schaalii</i>	99.35	<i>Actinobaculum massiliense</i>	95.69	<i>Actinobaculum suis</i>	94.52
3	<i>Actinobaculum urinale</i>		<i>Actinobaculum urinale</i>	100.00	<i>Actinobaculum urinale</i>	100.00	<i>Actinomyces massiliense</i>	91.57	<i>Actinobaculum schaalii</i>	91.01
4	<i>Actinomyces meyeri</i>		<i>Actinomyces meyeri</i>	100.00	<i>Actinomyces meyeri</i>	100.00	<i>Actinomyces odontolyticus</i>	97.59	<i>Actinomyces georgiae</i>	97.16
5	<i>Actinomyces odontolyticus</i>	22	<i>Actinomyces odontolyticus</i>	98.92	<i>Actinomyces odontolyticus</i>	98.92	<i>Actinomyces meyeri</i>	97.78	<i>Actinomyces turicensis</i>	96.07
6	<i>Actinomyces urogenitalis</i>		<i>Actinomyces urogenitalis</i>	100.00	<i>Actinomyces urogenitalis</i>	100.00	<i>Actinomyces stackii</i>	97.37	<i>Actinomyces radicans</i>	97.24
7	<i>Actinomyces turicensis</i>		<i>Actinomyces turicensis</i>	99.44	<i>Actinomyces turicensis</i>	99.44	<i>Actinomyces odontolyticus</i>	97.04	<i>Actinomyces meyeri</i>	96.65
8	<i>Aggregatibacter actinomycetemcomitans</i>		<i>Aggregatibacter actinomycetemcomitans</i>	98.62	<i>Aggregatibacter actinomycetemcomitans</i>	98.62	<i>Aggregatibacter aphrophilus</i>	94.12	<i>Haemophilus aegyptius</i>	93.73
9	<i>Aggregatibacter aphrophilus</i>	10	<i>Aggregatibacter aphrophilus</i>	99.59	<i>Aggregatibacter aphrophilus</i>	99.59	<i>Aggregatibacter segnis</i>	97.29	<i>Haemophilus parahaemolyticus</i>	95.31
10	<i>Aggregatibacter segnis</i>	10	<i>Aggregatibacter segnis</i>	99.59	<i>Aggregatibacter segnis</i>	99.59	<i>Aggregatibacter aphrophilus</i>	97.54	<i>Haemophilus influenzae</i>	95.73
11	<i>Agrobacterium tumefaciens</i>		<i>Agrobacterium tumefaciens</i>	98.86	<i>Agrobacterium tumefaciens</i>	98.86	<i>Ochrobactrum anthropti</i>	94.25	<i>Brucella melitensis</i>	93.95
12	<i>Arcanobacterium haemolyticum</i>		<i>Arcanobacterium haemolyticum</i>	99.80	<i>Arcanobacterium haemolyticum</i>	99.80	<i>Arcanobacterium phocae</i>	97.27	<i>Arcanobacterium pluranimalium</i>	94.39
13	<i>Arcobacter butzleri</i>		<i>Arcobacter butzleri</i>	99.64	<i>Arcobacter butzleri</i>	99.64	<i>Arcobacter cryaerophilus</i>	96.76	<i>Arcobacter skirrowii</i>	96.61
14	<i>Arcobacter cryaerophilus</i>	7	<i>Arcobacter cryaerophilus</i>	100.00	<i>Arcobacter cryaerophilus</i>	100.00	<i>Arcobacter skirrowii</i>	98.42	<i>Arcobacter butzleri</i>	97.18
15	<i>Averyella dalhousiensis</i>		<i>Averyella dalhousiensis</i>	99.18	<i>Averyella dalhousiensis</i>	99.18	<i>Enterobacter cancerogenus</i>	99.02	<i>Enterobacter aerogenes</i>	99.02
16	<i>Bifidobacterium breve</i>		<i>Bifidobacterium breve</i>	99.68	<i>Bifidobacterium breve</i>	99.68	<i>Bifidobacterium longum</i>	97.77	<i>Bifidobacterium bifidum</i>	95.96
17	<i>Burkholderia pseudomallei</i>	35	<i>Burkholderia pseudomallei</i>	100.00	<i>Burkholderia pseudomallei</i>	100.00	<i>Burkholderia mallei</i>	99.93	<i>Burkholderia thailandensis</i>	98.99
18	<i>Campylobacter coli</i>	7	<i>Campylobacter coli</i>	99.93	<i>Campylobacter coli</i>	99.80	<i>Campylobacter jejuni</i>	99.79	<i>Campylobacter lari</i>	98.33
19	<i>Campylobacter fetus</i>	7	<i>Campylobacter fetus</i>	99.79	<i>Campylobacter fetus</i>	99.85	<i>Campylobacter hyointestinalis</i>	97.98	<i>Campylobacter mucosalis</i>	96.40
20	<i>Campylobacter hyointestinalis</i>		<i>Campylobacter hyointestinalis</i>	100.00	<i>Campylobacter hyointestinalis</i>	100.00	<i>Campylobacter fetus</i>	99.01	<i>Campylobacter mucosalis</i>	96.06
21	<i>Campylobacter rectus</i>		<i>Campylobacter rectus</i>	99.34	<i>Campylobacter rectus</i>	99.34	<i>Campylobacter showae</i>	97.94	<i>Campylobacter concisus</i>	96.11
22	<i>Capnocytophaga sputigena</i>		<i>Capnocytophaga sputigena</i>	99.24	<i>Capnocytophaga sputigena</i>	99.24	<i>Capnocytophaga ochracea</i>	95.67	<i>Capnocytophaga cynodegmi</i>	92.24
23	<i>Citrobacter koseri</i>		<i>Citrobacter koseri</i>	99.92	<i>Citrobacter koseri</i>	99.92	<i>Citrobacter farmeri</i>	98.74	<i>Salmonella enterica</i>	98.74
24	<i>Clostridium baratii</i>	24	<i>Clostridium baratii</i>	98.74	<i>Clostridium baratii</i>	99.83	<i>Eubacterium budayi</i>	99.41	<i>Eubacterium nitritogenes</i>	99.16
25	<i>Clostridium difficile</i>	24	<i>Clostridium difficile</i>	99.16	<i>Clostridium difficile</i>	99.51	<i>Clostridium sporadicum</i>	98.28	<i>Clostridium indolis</i>	95.80
26	<i>Clostridium sporadicum</i>	24	<i>Clostridium sporadicum</i>	99.51	<i>Clostridium sporadicum</i>	98.99	<i>Clostridium buthefyi</i>	98.99	<i>Clostridium tenue</i>	96.91
27	<i>Clostridium buthefyi</i>	28	<i>Clostridium buthefyi</i>	98.28	<i>Clostridium buthefyi</i>	98.28	<i>Clostridium innocuum</i>	99.19	<i>Clostridium carnis</i>	93.86
28	<i>Clostridium innocuum</i>	24	<i>Clostridium innocuum</i>	99.19	<i>Clostridium innocuum</i>	99.19	<i>Holdemannia filiformis</i>	85.33	<i>Erysipelothrix inopinata</i>	84.82
29	<i>Clostridium orbiscindens</i>	24	<i>Clostridium orbiscindens</i>	100.00	<i>Clostridium orbiscindens</i>	100.00	<i>Bacteroides capillosus</i>	97.44	<i>Clostridium sporosphaeroideis</i>	86.62
30	<i>Clostridium paraputrificum</i>	24	<i>Clostridium paraputrificum</i>	99.42	<i>Clostridium paraputrificum</i>	99.42	<i>Clostridium carnis</i>	96.81	<i>Clostridium dispersicum</i>	96.53

TABLE 1—Continued

Strain no.	Identification by polyphasic approach	Reference	Identification results reported by 16SpathDB		Nucleotide identity (%)	Best match in 16SpathDB	Nucleotide identity (%)	Second-best match in 16SpathDB	Nucleotide identity (%)	Third-best match in 16SpathDB	Nucleotide identity (%)
			Bacterial species	Nucleotide identity (%)							
64	<i>Micrococcus luteus</i>	7	<i>Micrococcus luteus</i> <i>Micrococcus lysae</i> <i>Micrococcus antarcticus</i>	98.11 97.78 97.52	<i>Micrococcus luteus</i>	98.11	<i>Micrococcus lysae</i>	97.78	<i>Micrococcus antarcticus</i>	97.32	
65	<i>Moraxella osloensis</i>	7	<i>Moraxella osloensis</i>	99.53	<i>Moraxella osloensis</i>	99.53	<i>Moraxella lincolnii</i>	92.49	<i>Moraxella lacunata</i>	92.41	
66	<i>Mycobacterium chelonae</i>		<i>Mycobacterium chelonae</i> <i>Mycobacterium abscessus</i> <i>Mycobacterium immunogenium</i>	99.79 99.79 98.97	<i>Mycobacterium chelonae</i>	99.79	<i>Mycobacterium abscessus</i>	99.79	<i>Mycobacterium immunogenium</i>	98.97	
67	<i>Mycobacterium marinum</i>		<i>Mycobacterium marinum</i> <i>Mycobacterium ulcerans</i>	100.00 100.00	<i>Mycobacterium marinum</i>	100.00	<i>Mycobacterium marinum</i>	100.00	<i>Mycobacterium terrae</i>	98.71	
68	<i>Mycobacterium nonchromogenicum</i>	7	<i>Mycobacterium nonchromogenicum</i> <i>Mycoplasma hominis</i> <i>Nocardia cytaeigeorgica</i>	98.13 100.00 100.00	<i>Mycobacterium nonchromogenicum</i>	98.13	<i>Mycobacterium terrae</i>	96.69	<i>Mycobacterium cookii</i>	95.56	
69	<i>Mycoplasma hominis</i>	7	<i>Mycoplasma hominis</i>	100.00	<i>Mycoplasma hominis</i>	100.00	<i>Brevibacillus parabravis</i>	76.75	<i>Clostridium ramosum</i>	76.07	
70	<i>Nocardia cytaeigeorgica</i>	7	<i>Nocardia cytaeigeorgica</i>	100.00	<i>Nocardia cytaeigeorgica</i>	100.00	<i>Nocardia abscessus</i>	98.73	<i>Nocardia paucivonans</i>	98.33	
71	<i>Olsenella uli</i>	7	<i>Olsenella uli</i>	100.00	<i>Olsenella uli</i>	100.00	<i>Atopobium vaginae</i>	92.15	<i>Atopobium parvulum</i>	91.97	
72	<i>Pantoea dispersa</i>	7	<i>Pantoea dispersa</i>	99.92	<i>Pantoea dispersa</i>	99.92	<i>Enterobacter cancerogenus</i>	97.49	<i>Kluyvera cryocrescens</i>	97.49	
73	<i>Propionibacterium acnes</i>	7	<i>Propionibacterium acnes</i>	99.80	<i>Propionibacterium acnes</i>	99.80	<i>Propionibacterium avidum</i>	93.46	<i>Propionibacterium</i>	93.27	
74	<i>Propionibacterium avidum</i>	7	<i>Propionibacterium avidum</i> <i>Propionibacterium propionicum</i>	99.45 99.30	<i>Propionibacterium avidum</i>	99.45	<i>Propionibacterium propionicum</i>	99.30	<i>Propionibacterium acnes</i>	96.10	
75	<i>Providencia stuartii</i>	7	<i>Providencia stuartii</i>	99.77	<i>Providencia stuartii</i>	99.77	<i>Providencia rettgeri</i>	98.75	<i>Providencia rustigianii</i>	98.13	
76	<i>Pseudomonas oryzihabitans</i>	7	<i>Pseudomonas oryzihabitans</i>	99.79	<i>Pseudomonas oryzihabitans</i>	99.79	<i>Pseudomonas pseudocataligenes</i>	96.84	<i>Pseudomonas aeruginosa</i>	96.21	
77	<i>Salmonella enterica</i>	23	<i>Salmonella enterica</i> <i>Shewanella algae</i>	99.93 99.92	<i>Salmonella enterica</i>	99.93	<i>Citrobacter jammeri</i>	98.47	<i>Enterobacter cloacae</i>	98.46	
78	<i>Shewanella algae</i>	7	<i>Shewanella algae</i>	99.92	<i>Shewanella algae</i>	99.92	<i>Solobacterium moorei</i>	95.02	<i>Aeromonas hydrophila</i>	90.90	
79	<i>Solobacterium moorei</i>	8	<i>Solobacterium moorei</i>	98.33	<i>Solobacterium moorei</i>	98.33	<i>Bulleidia extracta</i>	93.20	<i>Holdemannia filiformis</i>	88.86	
80	<i>Staphylococcus aureus</i>	7	<i>Staphylococcus aureus</i>	99.63	<i>Staphylococcus aureus</i>	99.63	<i>Staphylococcus epidermidis</i>	97.66	<i>Staphylococcus caprae</i>	97.64	
81	<i>Staphylococcus epidermidis</i>	7	<i>Staphylococcus epidermidis</i> <i>Staphylococcus caprae</i>	100.00 99.21	<i>Staphylococcus epidermidis</i>	100.00	<i>Staphylococcus caprae</i>	99.21	<i>Staphylococcus capitis</i>	99.02	
82	<i>Staphylococcus lugdunensis</i>	18	<i>Staphylococcus lugdunensis</i> <i>Staphylococcus haemolyticus</i>	99.02 99.93 99.05	<i>Staphylococcus lugdunensis</i>	99.93	<i>Staphylococcus haemolyticus</i>	99.05	<i>Staphylococcus hominis</i>	98.84	
83	<i>Streptococcus anginosus</i>	33	<i>Streptococcus anginosus</i>	99.77	<i>Streptococcus anginosus</i>	99.77	<i>Streptococcus intermedius</i>	97.37	<i>Streptococcus constellatus</i>	96.62	
84	<i>Streptococcus dysgalactiae</i>	31	<i>Streptococcus dysgalactiae</i>	99.62	<i>Streptococcus dysgalactiae</i>	99.62	<i>Streptococcus iniae</i>	97.47	<i>Streptococcus agalactiae</i>	97.33	
85	<i>Streptococcus iniae</i>	7	<i>Streptococcus iniae</i>	99.62	<i>Streptococcus iniae</i>	99.62	<i>Streptococcus porcinus</i>	95.21	<i>Streptococcus dysgalactiae</i>	95.11	
86	<i>Streptococcus porcinus</i>	7	<i>Streptococcus porcinus</i>	99.76	<i>Streptococcus porcinus</i>	99.76	<i>Streptococcus agalactiae</i>	96.79	<i>Streptococcus iniae</i>	96.79	
87	<i>Streptococcus pyogenes</i>	12	<i>Streptococcus pyogenes</i>	100.00	<i>Streptococcus pyogenes</i>	100.00	<i>Streptococcus canis</i>	98.24	<i>Streptococcus agalactiae</i>	96.93	
88	<i>Tsukamurella pulmonis</i>	20	<i>Tsukamurella pulmonis</i> <i>Tsukamurella tyrosinosolvens</i> <i>Tsukamurella strandjordii</i> <i>Tsukamurella incheonensis</i>	99.13 99.05 98.90 98.74	<i>Tsukamurella pulmonis</i>	99.13	<i>Tsukamurella tyrosinosolvens</i>	99.05	<i>Tsukamurella strandjordii</i>	98.90	
89	<i>Tsukamurella tyrosinosolvens</i>	20	<i>Tsukamurella paurometabola</i> <i>Tsukamurella tyrosinosolvens</i> <i>Tsukamurella pulmonis</i> <i>Tsukamurella strandjordii</i> <i>Tsukamurella incheonensis</i>	98.66 100.00 99.69 99.39 99.24	<i>Tsukamurella tyrosinosolvens</i>	100.00	<i>Tsukamurella pulmonis</i>	99.69	<i>Tsukamurella strandjordii</i>	99.69	
90	<i>Veillonella atypica</i>	7	<i>Veillonella atypica</i> <i>Veillonella dispar</i>	98.95 98.42	<i>Veillonella atypica</i>	98.95	<i>Veillonella dispar</i>	98.42	<i>Veillonella parvula</i>	97.71	
91	<i>Vibrio furnissii</i>	7	<i>Vibrio furnissii</i> <i>Vibrio fluvialis</i>	99.62 98.85	<i>Vibrio furnissii</i>	99.62	<i>Vibrio fluvialis</i>	98.85	<i>Vibrio vulnificus</i>	97.02	

important bacteria using 16S rRNA gene sequencing. In this database, 16SpathDB, we sought to create an automated user-friendly platform that indicated the most likely identity of the 16S rRNA gene sequence of a medically important bacterium, as well as other medically important bacteria with similar 16S rRNA gene sequences that may be alternative identities, which the user should be aware of. In this article, we describe this comprehensive database of 16S rRNA gene sequences of medically important bacteria and its use for automated identification of these bacteria.

MATERIALS AND METHODS

Database design. 16SpathDB is a Web-based 16S rRNA gene sequence database for identification of medically important bacteria. MySQL was employed as the database back end to store the 16S rRNA gene sequence information, and PHP was used to generate HTML web pages for the user interface. The most representative 16S rRNA gene sequence of each medically important bacterial species listed in the most current edition of the *Manual of Clinical Microbiology* (14) was retrieved from GenBank by manual inspection according to the following criteria. First, strains with good phenotypic and/or genotypic characterization (e.g., type strains and strains with complete genomes sequenced) were preferred. Second, strains isolated from humans were preferred. Third, sequences with fewer undetermined bases were preferred. Fourth, longer sequences, especially those with better coverage of the 5' end, were preferred. For bacterial species with >2% intra-genomic difference in their 16S rRNA gene sequences and those with intervening sequences in their 16S rRNA genes, more than one 16S rRNA gene sequence for each species was included.

Identification of clinical bacterial isolates using 16SpathDB. To evaluate the usefulness of 16SpathDB, the 16S rRNA gene sequences of 91 nonduplicated medically important bacterial isolates we collected in our clinical microbiology laboratory in the past 10 years were input to the database for analysis (Table 1) (7–12, 18, 20–25, 27–31, 33–35). The exact identities of these isolates were determined by a polyphasic approach using a combination of phenotypic tests and 16S rRNA gene sequencing as described in our publications.

Availability and update of 16SpathDB. 16SpathDB is available at no charge at <http://147.8.74.24/16SpathDB> and will be updated periodically for every new edition of the *Manual of Clinical Microbiology*.

RESULTS

16SpathDB and functionality of the database. As of November 2010, 16SpathDB contained 1,014 16S rRNA gene sequences from 1,010 unique bacterial species.

Identification of medically important bacteria. The main goal for setting up 16SpathDB was to provide a convenient and efficient platform for identification of medically important bacterial isolates using their 16S rRNA gene sequences. The interfaces of the database are simple and user friendly. From the home page, one can enter the “query page” by clicking the “identify bacteria by 16S rRNA gene sequence” hyperlink (Fig. 1a). From this page, users can input one or more query 16S rRNA gene sequences of 50 bp to 1,800 bp by pasting the sequences in FASTA format in the textbox or uploading a file that contains the sequences in FASTA format by clicking the “browse” button. By clicking the “begin identification” button, the query sequence(s) will be aligned with each of the 16S rRNA gene sequences in 16SpathDB using the pairwise global alignment algorithm with free end gap penalty. The percent nucleotide identity calculated from the alignment between the query sequence and each of the sequences in 16SpathDB is then used for the evaluation of the identity of the query sequence, using the algorithm depicted in Fig. 2.

The results of the comparison will be shown on the results page. If there is one species in 16SpathDB with >98.0% nucleotide identity to the query sequence, this bacterial species, as well as the percent nucleotide identity between the query sequence and the sequence of the most likely bacterial species, will be reported (category 1) (Fig. 1b and 2). If there is more than one species in 16SpathDB with >98.0% nucleotide identity to the query sequence, the species that showed the highest nucleotide identity to the query sequence (“best match in 16SpathDB”), as well as those with 16S rRNA gene sequences having <1% difference from the “best match in 16SpathDB,” will be reported to alert the user that further tests, such as biochemical tests or sequencing additional genes, may be necessary to distinguish between the most probable identities (category 2) (Fig. 1c and 2). If there are no species in 16SpathDB with >98.0% nucleotide identity to the query sequence but there are one or more species in 16SpathDB with >96.0% nucleotide identity to the query sequence, only the genus will be reported (category 3) (Fig. 1d and 2). The user will also be reminded that further tests are necessary for definite species identification. If there is no species in 16SpathDB with >96.0% nucleotide identity to the query sequence, the results page will show “no species in 16SpathDB was found to be sharing high nucleotide identity to your query sequence” (category 4) (Fig. 1e and 2). This indicates that the query sequence may represent a bacterial species not included in the *Manual of Clinical Microbiology* or a novel bacterial species. Users are advised to perform a BLAST search against the GenBank nr database to differentiate between the two possibilities. By clicking the “run BLAST” hyperlink, users can enter the “query page” of the NCBI BLAST website. The cutoffs of 98.0% for reporting the identity at the species level and 96.0% for reporting the identity at the genus level were found to be the optimal cutoffs by using 400 16S rRNA gene sequences of bacteria isolated from patients found in the GenBank database as the query sequences (data not shown).

From the results page, users can click the bacterial species name and go to the page in GenBank that contains detailed information on the representative 16S rRNA gene sequence selected for that bacterial species in 16SpathDB. From the results page, users can also click “show top 10 matches” to show the 10 most closely matched sequences from 16SpathDB compared to the query sequence (Fig. 1f).

Browsing 16S rRNA gene sequences of medically important bacteria. In addition to identification of 16S rRNA gene sequences, users can also inspect the detailed contents of the database, as well as the information on the individual sequences. From the home page, one can enter the “sequence information” page by clicking the “browse bacteria 16S rRNA gene information” hyperlink (Fig. 3). From this page, the user can go to the page in GenBank that contains the detailed information on the 16S rRNA gene sequence selected for the corresponding bacterial species by entering the name of the bacterial species in the text box and click “view sequence information” and then the GenBank accession number. In addition, users can visualize the 10 bacterial species with 16S rRNA gene sequences having the highest nucleotide identities with the input bacterial species and

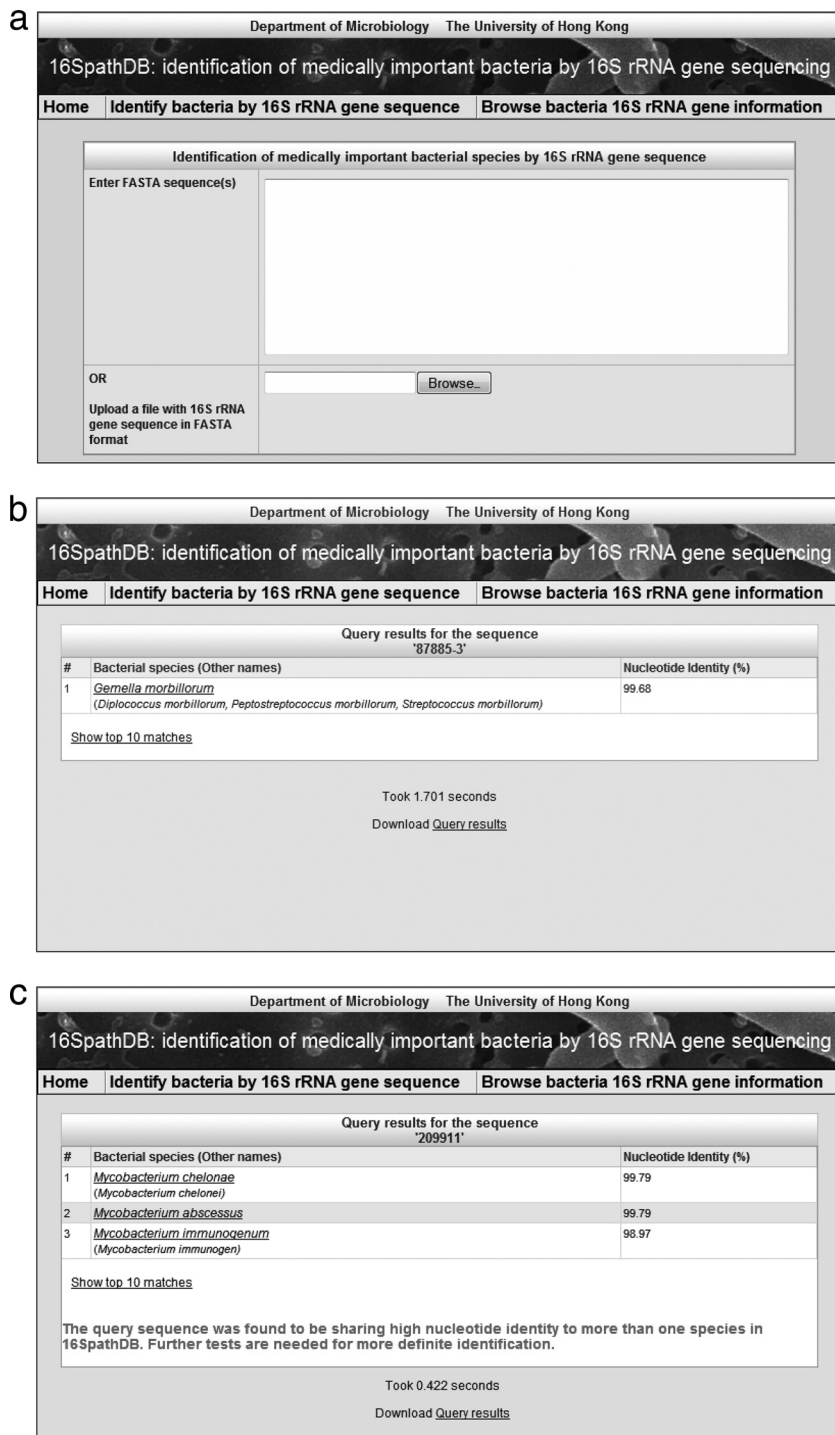


FIG. 1. (a) “Query page” of 16SpathDB. (b) Example of query results showing one species (*Gemella morbillorum*) in 16SpathDB with >98.0% nucleotide identity to the query sequence. (c) Example of query results showing more than one species (*Mycobacterium chelonae*, *Mycobacterium abscessus*, and *Mycobacterium immunogenum*) in 16SpathDB with >98.0% nucleotide identity to the query sequence. (d) Example of query results showing no species in 16SpathDB with >98.0% nucleotide identity to the query sequence but one or more species in 16SpathDB with >96.0% nucleotide identity to the query sequence. (e) Example of query results showing no species in 16SpathDB with >96.0% nucleotide identity to the query sequence. (f) Example of query results showing the identity of the query sequence (*Clostridium paraputrificum*) and the 10 most closely matched sequences from 16SpathDB compared to the query sequence.

d

Department of Microbiology The University of Hong Kong

16SpathDB: identification of medically important bacteria by 16S rRNA gene sequencing

Home Identify bacteria by 16S rRNA gene sequence Browse bacteria 16S rRNA gene information

Query results for the sequence '31121'

#	Bacterial species (Other names)	Nucleotide Identity (%)
1	<i>Clostridium</i> sp.	96.38

[Show top 10 matches](#)

There is one or more species with 16S rRNA gene sequence having nucleotide identity >96% to the query sequence. Further tests are needed for more definite identification.

Took 0.769 seconds

[Download Query results](#)

e

Department of Microbiology The University of Hong Kong

16SpathDB: identification of medically important bacteria by 16S rRNA gene sequencing

Home Identify bacteria by 16S rRNA gene sequence Browse bacteria 16S rRNA gene information

Query results for the sequence 'PW1492'

No species in 16SpathDB was found to be sharing high nucleotide identity to your query sequence

[Show top 10 matches](#)

Took 0.775 seconds

[Download Query results](#)

f

Department of Microbiology The University of Hong Kong

16SpathDB: identification of medically important bacteria by 16S rRNA gene sequencing

Home Identify bacteria by 16S rRNA gene sequence Browse bacteria 16S rRNA gene information

Query results for the sequence '99-27153'

#	Bacterial species (Other names)	Nucleotide Identity (%)
1	<i>Clostridium paraputrificum</i> (<i>Bacillus diaphthirus</i> , <i>Bacillus paraputrificus</i>)	99.15

[Show top 10 matches](#)

The top 10 closely matched bacterial species		
#	Bacterial species (Other names)	Nucleotide Identity (%)
1	<i>Clostridium paraputrificum</i> (<i>Bacillus diaphthirus</i> , <i>Bacillus paraputrificus</i>)	99.15
2	<i>Clostridium carnis</i> (<i>Bacillus carnis</i> , <i>Plectridium carnis</i>)	96.56
3	<i>Clostridium disporicum</i>	96.37
4	<i>Clostridium butyricum</i> (<i>Amylobacter navicula</i> , <i>Bacillus amylobacter</i> , <i>Bacillus butyricus</i> , <i>Bacillus navicula</i> , <i>Bacterium navicula</i> , <i>Clostridium naviculum</i> , <i>Clostridium pseudotetanicum</i> , <i>Metallacter amylobacter</i>)	95.95
5	<i>Clostridium beijerinckii</i> (<i>Clostridium rubrum</i>)	95.95
6	<i>Clostridium tertium</i> (<i>Bacillus tertius</i> , <i>Henrillus tertius</i> , <i>Plectridium tertium</i>)	95.94
7	<i>Eubacterium moniliforme</i> (<i>Bacillus moniliforme</i> , <i>Bacillus repazii</i> , <i>Cillobacterium moniliforme</i>)	95.69
8	<i>Clostridium baratii</i> (<i>Acuformis perennis</i> , <i>Clostridium baratii</i> , <i>Clostridium paraperfringens</i> , <i>Clostridium perenne</i>)	95.52
9	<i>Clostridium septicum</i> (<i>Bacillus septicus</i> , <i>Vibrio septicus</i>)	95.18
10	<i>Clostridium chauvoei</i> (<i>Bacillus chauvoei</i> , <i>Bacterium chauvoei</i> , <i>Clostridium chauvoei</i>)	95.10

FIG. 1—Continued.

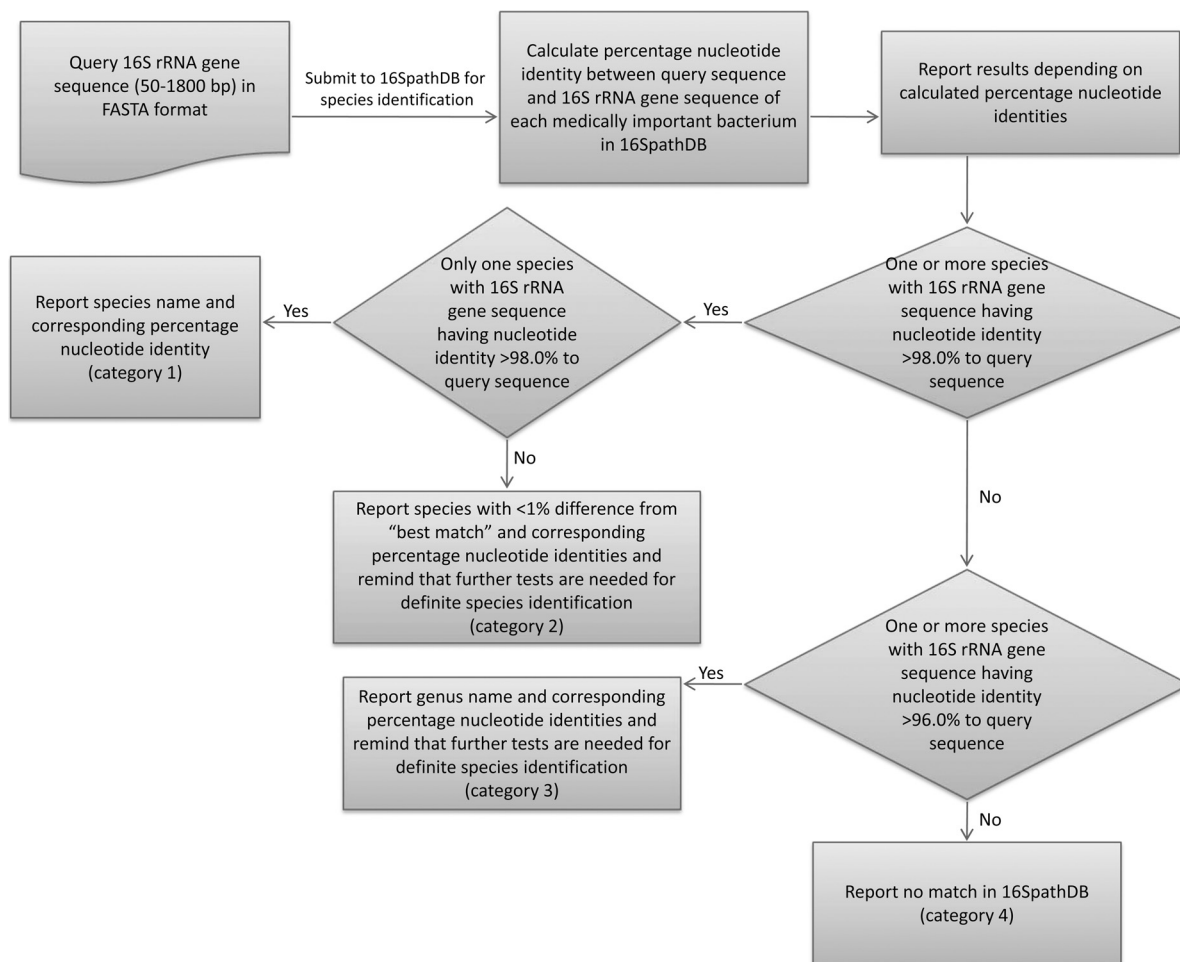


FIG. 2. Algorithm for reporting results in 16SpathDB.

inspect the detailed information on their 16S rRNA gene sequences in GenBank accordingly (Fig. 3).

Identification of clinical bacterial isolates using 16SpathDB. Among the 91 nonduplicated medically important bacterial isolates we collected in our clinical microbiology laboratory in the past 10 years (Table 1), 71 (78%) were reported by 16SpathDB as a single bacterial species having >98.0% nucleotide identity with the query sequence (category 1), 19 (20.9%) were reported by 16SpathDB as more than one bacterial species having >98.0% nucleotide identity with the query sequence (category 2), none was reported by 16SpathDB to the genus level (category 3), and 1 (1.1%) was reported by 16SpathDB as “no species in 16SpathDB was found to be sharing high nucleotide identity to your query sequence” (category 4). For the 71 bacterial isolates reported by 16SpathDB as a single bacterial species, all results were identical to the true identities of the isolates as determined by the polyphasic approach. For the 19 bacterial isolates reported by 16SpathDB as more than one bacterial species, all results contained the true identities of the isolates as determined by the polyphasic approach. In fact, all 19 of these isolates had their true identities as the “best match in 16SpathDB.”

DISCUSSION

Rapid and accurate interpretation of 16S rRNA gene sequence results is the cornerstone for identification of medically important bacteria by 16S rRNA gene sequencing. During the process of using 16S rRNA gene sequencing for identification of a large number of medically important bacteria in the past 5 years, we have developed an automated, user-friendly, and comprehensive database, 16SpathDB, of the most representative 16S rRNA gene sequences of all medically important bacteria listed in the most current edition of the *Manual of Clinical Microbiology* (14). In contrast to RDP-II and SmartGene IDNS software (Table 2), 16SpathDB includes only 16S rRNA gene sequences of medically important bacteria. Since more than 99.9% of the bacterial strains recovered from patients were included in the *Manual of Clinical Microbiology* (unpublished data), the inclusion of other bacterial species that have never been isolated from patients would create ambiguity during data interpretation, as the target users of 16SpathDB are technicians and clinical microbiologists who work on 16S rRNA gene sequencing for clinical isolates. As for the MicroSeq databases (Table 2), one of their main drawbacks for iden-

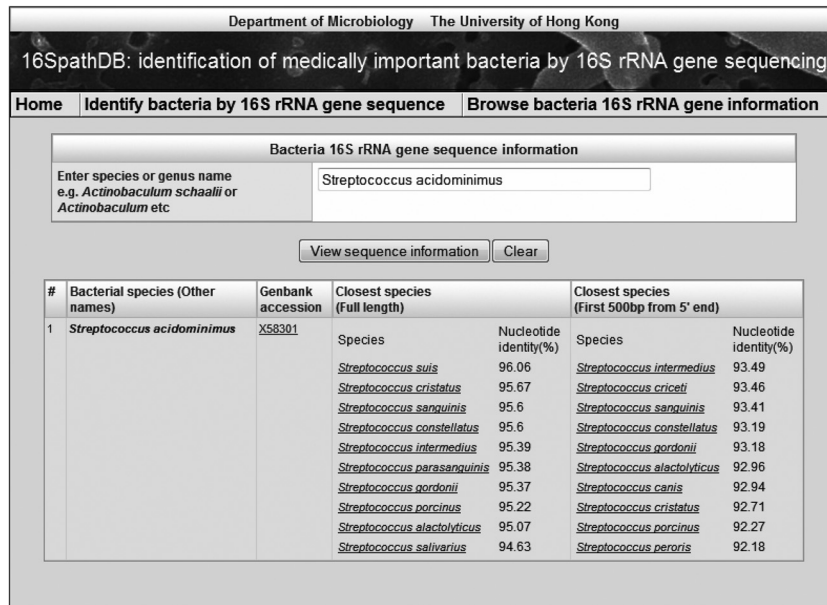


FIG. 3. Example showing the 10 bacterial species with 16S rRNA gene sequences having the highest nucleotide identities to the input bacterial species (*Streptococcus acidominimus*).

tification of medically important bacteria is that the databases do not include a significant number of medically important bacteria that 16S rRNA gene sequencing is able to identify. For example, 98 to 108 (53.3 to 67.1%), 38 to 39 (22.7 to 37.3%), and 23 to 39 (19.8 to 41.9%) medically important anaerobic, aerobic Gram-positive, and aerobic Gram-negative bacteria that should be confidently identified by 16S rRNA gene sequencing are not included in the full- and 500-MicroSeq databases, respectively (19, 32). If the query sequence is from one of these bacteria, the MicroSeq databases would automatically give wrong results. Furthermore, the results from the MicroSeq databases show only a single “identity” of the query sequence. Other bacterial species with similar 16S rRNA gene sequences, which may also be the identity of the isolate, are ignored. For example, it is well known that 16S rRNA gene sequencing is not useful for distinguishing some bacterial species, such as *Streptococcus pneumoniae*, *Streptococcus pseudopneumoniae*, *Streptococcus mitis*, and *Streptococcus oralis*, as their 16S rRNA gene sequences show more than 99% identity. In 16SpathDB, in addition to the species that shows the highest nucleotide identity to the query sequence, those species with 16S rRNA gene sequences having less than 1% difference from the species that showed the highest nucleotide identity to the query sequence will also be reported (Fig. 1c). This helps to alert the user that further tests have to be carried out in order to distinguish between these probable identities.

16SpathDB offers efficient and accurate analysis of 16S rRNA gene sequences from medically important bacteria. In the present study, all 91 bacteria recovered in our clinical microbiology laboratory in the past 10 years were successfully identified using 16SpathDB (Table 1). These 91 isolates included 55 aerobic and 36 anaerobic bacteria, and 62 Gram-positive bacteria, 28 Gram-negative bacteria, and one *Mycoplasma hominis* isolate. Among these 91 bacteria,

20.9% showed multiple possible identities, which reflected an inherent limitation of using 16S rRNA gene sequencing for bacterial identification. For this 20.9% of the strains, phenotypic tests or sequencing of additional gene loci should be performed to differentiate among the reported bacterial species. For example, when *Tsukamurella pulmonis* or *Tsukamurella tyrosinosolvans* (strain no. 88 and 89) (Table 1) were the exact identities of the query isolates, 16S rRNA gene sequencing was not able to distinguish them from the other medically important *Tsukamurella* spp. The API 20C AUX, and API 50 CH systems (bioMérieux, Lyon, France) have to be used to distinguish among these *Tsukamurella* spp. (20).

One limitation of 16SpathDB is that our database includes only the bacterial species that are known to be associated with infections, as described in the *Manual of Clinical Microbiology*. This was deliberately done because if those bacterial species that have never been reported to cause infections were also included, as in the RDP-II and Smart-Gen IDNS software, the proportion of results that showed multiple possible identities would be markedly increased. This would defeat the original purpose of designing the database, which is used for identification of medically important bacteria in clinical microbiology laboratories. For example, the 16S rRNA gene sequence of strain no. 50 (Table 1) was reported as “no species in 16SpathDB was found to be sharing high nucleotide identity to your query sequence.” This is because *Gordonibacter pamelaee* has never been reported to be associated with human disease and was not included in the *Manual of Clinical Microbiology*. In such a situation, users should perform a BLAST search against the GenBank nr database for the identification of the 16S rRNA gene sequence (30). Although the database will be updated regularly whenever there is a new edition of the *Manual of Clinical Microbiology*, users should bear in

TABLE 2. Comparison of 16SpathDB and other commonly used software for bacterial identification using 16S rRNA gene sequencing

Software	Yr of first description	Company/organization	Website	Partial/full 16S rRNA gene sequence included	Database size	Source of sequences	Cost	Quality control	Updates
Ribosomal Database Project (RDP)	1992	Michigan State University, East Lansing, MI	http://rdp.cmc.msu.edu/	Partial and full	1,418,497 (release 10.22)	GenBank	No	Partial	Periodically
MicroSeq Microbial Identification System ^a	1998	Applied Biosystems, Foster City, California	http://www.microseq.com	Partial and full ^b	1,834 and 1,261 in the two databases, respectively ^c	Sequence of 16S rRNA gene of one strain from each species	Yes	All type strains from culture collections	Periodically
Ribosomal Differentiation of Medical Microorganisms (RIDOM)	1999	Ridom GmbH, Würzburg, Germany	http://rdna2.ridom.de/	Partial	236 ^b	Sequences of 16S rRNA genes of medically relevant bacteria, mainly belonging to the <i>Neisseriaceae</i> , <i>Moraxellaceae</i> , and genus <i>Mycobacterium</i>	No	All strains from culture collections	Periodically
SmartGene IDNS software	2006	SmartGene GmbH, Switzerland	http://www.smartgene.com/mod_bacteria.html	Partial and full	243,000	GenBank	Yes	Partial	Daily
16SpathDB	2010	Department of Microbiology, University of Hong Kong, Hong Kong	http://147.8.74.24/16SpathDB	Partial and full	1,010	GenBank	No	All sequences manually selected from GenBank	Periodically

^a Contains two databases, MicroSeq ID 16S rDNA 500 Library v2.2, which contains sequences from the 5'-end 527 bp of 16S rRNA genes, and MicroSeq ID 16S rDNA Full Gene Library v2.0, which contains full 16S rRNA gene sequences.

^b Number from from <http://rdna2.ridom.de/ridom2/servelet/link?page=list> (8 Nov 2010).

mind that those bacterial species that have never been reported to be associated with infections may still have the potential to be so associated.

ACKNOWLEDGMENTS

This work is partly supported by the HKSAR Research Fund for the Control of Infectious Diseases of the Health, Welfare and Food Bureau, Research Grants Council Grant, and University Development Fund, The University of Hong Kong.

We thank Kit-Wah Leung and Ami M. Y. Fung for their critical comments on the database.

REFERENCES

- Benson, D. A., I. Karsch-Mizrachi, D. J. Lipman, J. Ostell, and D. L. Wheeler. 2008. GenBank. *Nucleic Acids Res.* **36**:D25–D30.
- Cole, J. R., et al. 2005. The Ribosomal Database Project (RDP-II): sequences and tools for high-throughput rRNA analysis. *Nucleic Acids Res.* **33**:D294–D296.
- Cole, J. R., et al. 2009. The Ribosomal Database Project: improved alignments and new tools for rRNA analysis. *Nucleic Acids Res.* **37**:D141–D145.
- Harmsen, D., J. Rothganger, M. Frosch, and J. Albert. 2002. RIDOM: Ribosomal Differentiation of Medical Micro-organisms Database. *Nucleic Acids Res.* **30**:416–417.
- Harmsen, D., J. Rothganger, C. Singer, J. Albert, and M. Frosch. 1999. Intuitive hypertext-based molecular identification of micro-organisms. *Lancet* **353**:291.
- Harmsen, D., et al. 2001. Diagnostics of neisseriaceae and moraxellaceae by ribosomal DNA sequencing: ribosomal differentiation of medical microorganisms. *J. Clin. Microbiol.* **39**:936–942.
- Lau, S. K., et al. 2006. Usefulness of the MicroSeq 500 16S rDNA bacterial identification system for identification of anaerobic Gram positive bacilli isolated from blood cultures. *J. Clin. Pathol.* **59**:219–222.
- Lau, S. K., et al. 2006. Bacteremia caused by *Solobacterium moorei* in a patient with acute proctitis and carcinoma of the cervix. *J. Clin. Microbiol.* **44**:3031–3034.
- Lau, S. K., et al. 2004. Anaerobic, non-sporulating, Gram-positive bacilli bacteraemia characterized by 16S rRNA gene sequencing. *J. Med. Microbiol.* **53**:1247–1253.
- Lau, S. K., et al. 2004. Characterization of *Haemophilus segnis*, an important cause of bacteremia, by 16S rRNA gene sequencing. *J. Clin. Microbiol.* **42**:877–880.
- Lau, S. K., et al. 2004. *Eggerthella hongkongensis* sp. nov. and *Eggerthella sinensis* sp. nov., two novel *Eggerthella* species, account for half of the cases of *Eggerthella* bacteremia. *Diagn. Microbiol. Infect. Dis.* **49**:255–263.
- Lau, S. K., P. C. Woo, T. C. Yim, A. P. To, and K. Y. Yuen. 2003. Molecular characterization of a strain of group A streptococcus isolated from a patient with a psoas abscess. *J. Clin. Microbiol.* **41**:4888–4891.
- Maidak, B. L., et al. 1999. A new version of the RDP (Ribosomal Database Project). *Nucleic Acids Res.* **27**:171–173.
- Murray, P. R., E. J. Baron, J. H. Tenover, M. R. Tenover, and M. A. Tenover (ed.). 2007. Manual of clinical microbiology, 9th ed. American Society for Microbiology, Washington, DC.
- Patel, J. B., et al. 2000. Sequence-based identification of *Mycobacterium* species using the MicroSeq 500 16S rDNA bacterial identification system. *J. Clin. Microbiol.* **38**:246–251.
- Simmon, K. E., A. C. Croft, and C. A. Petti. 2006. Application of SmartGene IDNS software to partial 16S rRNA gene sequences for a diverse group of bacteria in a clinical laboratory. *J. Clin. Microbiol.* **44**:4400–4406.
- Tang, Y. W., et al. 2000. Identification of coryneform bacterial isolates by ribosomal DNA sequence analysis. *J. Clin. Microbiol.* **38**:1676–1678.
- Tse, H., et al. 2010. Complete genome sequence of *Staphylococcus lugdunensis* strain HKU09-01. *J. Bacteriol.* **192**:1471–1472.
- Woo, P. C., et al. 2007. *In silico* analysis of 16S ribosomal RNA gene sequencing-based methods for identification of medically important anaerobic bacteria. *J. Clin. Pathol.* **60**:576–579.
- Woo, P. C., et al. 2009. First report of *Tsukamurella* keratitis: association between *T. tyrosinosolvens* and *T. pulmonis* and ophthalmologic infections. *J. Clin. Microbiol.* **47**:1953–1956.
- Woo, P. C., et al. 2003. *Granulicatella adiacens* and *Abiotrophia defectiva* bacteraemia characterized by 16S rRNA gene sequencing. *J. Med. Microbiol.* **52**:137–140.
- Woo, P. C., A. M. Fung, S. K. Lau, E. Hon, and K. Y. Yuen. 2002. Diagnosis of pelvic actinomycosis by 16S ribosomal RNA gene sequencing and its clinical significance. *Diagn. Microbiol. Infect. Dis.* **43**:113–118.
- Woo, P. C., A. M. Fung, S. S. Wong, H. W. Tsoi, and K. Y. Yuen. 2001. Isolation and characterization of a *Salmonella enterica* serotype Typhi variant and its clinical and public health implications. *J. Clin. Microbiol.* **39**:1190–1194.
- Woo, P. C., et al. 2005. *Clostridium* bacteraemia characterised by 16S ribosomal RNA gene sequencing. *J. Clin. Pathol.* **58**:301–307.

25. **Woo, P. C., et al.** 2007. Surgical site abscess caused by *Lactobacillus fermentum* identified by 16S ribosomal RNA gene sequencing. *Diagn. Microbiol. Infect. Dis.* **58**:251–254.
26. **Woo, P. C., S. K. Lau, J. L. Teng, H. Tse, and K. Y. Yuen.** 2008. Then and now: use of 16S rDNA gene sequencing for bacterial identification and discovery of novel bacteria in clinical microbiology laboratories. *Clin. Microbiol. Infect.* **14**:908–934.
27. **Woo, P. C., et al.** 2009. The complete genome and proteome of *Laribacter hongkongensis* reveal potential mechanisms for adaptations to different temperatures and habitats. *PLoS Genet.* **5**:e1000416.
28. **Woo, P. C., et al.** 2004. Bacteremia due to *Clostridium hathewayi* in a patient with acute appendicitis. *J. Clin. Microbiol.* **42**:5947–5949.
29. **Woo, P. C., D. M. Tam, S. K. Lau, A. M. Fung, and K. Y. Yuen.** 2004. *Enterococcus cecorum* empyema thoracis successfully treated with cefotaxime. *J. Clin. Microbiol.* **42**:919–922.
30. **Woo, P. C., et al.** 2010. First report of *Gordonibacter pamelaiae* bacteremia. *J. Clin. Microbiol.* **48**:319–322.
31. **Woo, P. C., et al.** 2003. Analysis of a viridans group strain reveals a case of bacteremia due to Lancefield group G alpha-hemolytic *Streptococcus dysgalactiae* subsp. *equisimilis* in a patient with pyomyositis and reactive arthritis. *J. Clin. Microbiol.* **41**:613–618.
32. **Woo, P. C., et al.** 2009. Guidelines for interpretation of 16S rRNA gene sequence-based results for identification of medically important aerobic Gram-positive bacteria. *J. Med. Microbiol.* **58**:1030–1036.
33. **Woo, P. C., et al.** 2004. “*Streptococcus milleri*” endocarditis caused by *Streptococcus anginosus*. *Diagn. Microbiol. Infect. Dis.* **48**:81–88.
34. **Woo, P. C., et al.** 2005. Life-threatening invasive *Helcococcus kunzii* infections in intravenous-drug users and *ermA*-mediated erythromycin resistance. *J. Clin. Microbiol.* **43**:6205–6208.
35. **Woo, P. C., G. K. Woo, S. K. Lau, S. S. Wong, and K. Yuen.** 2002. Single gene target bacterial identification. *groEL* gene sequencing for discriminating clinical isolates of *Burkholderia pseudomallei* and *Burkholderia thailandensis*. *Diagn. Microbiol. Infect. Dis.* **44**:143–149.