

Metagenomic Pyrosequencing and Microbial Identification

Joseph F. Petrosino,^{1,2} Sarah Highlander,^{1,2} Ruth Ann Luna,^{3,4} Richard A. Gibbs,^{2,5}
and James Versalovic^{1,3,4,5,6*}

BACKGROUND: The Human Microbiome Project has ushered in a new era for human metagenomics and high-throughput next-generation sequencing strategies.

CONTENT: This review describes evolving strategies in metagenomics, with a special emphasis on the core technology of DNA pyrosequencing. The challenges of microbial identification in the context of microbial populations are discussed. The development of next-generation pyrosequencing strategies and the technical hurdles confronting these methodologies are addressed. Bioinformatics-related topics include taxonomic systems, sequence databases, sequence-alignment tools, and classifiers. DNA sequencing based on 16S rRNA genes or entire genomes is summarized with respect to potential pyrosequencing applications.

SUMMARY: Both the approach of 16S rDNA amplicon sequencing and the whole-genome sequencing approach may be useful for human metagenomics, and numerous bioinformatics tools are being deployed to tackle such vast amounts of microbiological sequence diversity. Metagenomics, or genetic studies of microbial communities, may ultimately contribute to a more comprehensive understanding of human health, disease susceptibilities, and the pathophysiology of infectious and immune-mediated diseases.

© 2009 American Association for Clinical Chemistry

Metagenomics and the Human Microbiome

Metagenomics refers to culture-independent studies of the collective set of genomes of mixed microbial communities and applies to explorations of all microbial genomes in consortia that reside in environmental

niches, in plants, or in animal hosts. Examples in mammalian biology include studies of microbial communities on various mucosal surfaces and on the human skin. This review focuses on analytical strategies for identifying pathogens in mixed microbial communities via metagenomics.

Metagenomics and associated metastrategies have arrived at the forefront of biology primarily because of 2 major developments. The deployment of next-generation DNA-sequencing technologies in many centers has greatly enhanced capabilities for sequencing large meta-data sets. Technologic advances have created new opportunities for the pursuit of large-scale sequencing projects that were difficult to imagine just several years ago. The second key development is an emerging appreciation for the importance of complex microbial communities in mammalian biology and in human health and disease. The Human Microbiome Project was approved in May 2007 as one of 2 major components (in addition to the human epigenomics program) of RoadMap version 1.5 of the US NIH (1). The demands of this project have produced intense interest and a focus in genome centers to apply parallel DNA-sequencing technologies to human biology on a scale not previously witnessed.

The human microbiome is the entire population of microbes that colonize the human body, including the gastrointestinal tract, the genitourinary tract, the oral cavity, the nasopharynx, the respiratory tract, and the skin. The different microorganisms constituting the microbiome include bacteria, fungi (mostly yeasts), and viruses. Depending on the context, parasites may also be considered to compose part of the indigenous microbiota. The “metagenome” of microbial communities that occupy various sites in the body is estimated to be approximately 100-fold greater in terms of gene content than the human genome. These diverse and complex collections of genes encode a wide array of biochemical and physiological functions that may benefit the host, as well as neighboring microbes (1). We focus on bacterial populations because bacteria form a predominant group of the microbiome and have the most comprehensively documented phylogenetic data sets and classification systems. Most of the data gathered to date have been compiled with Sanger (dideoxy) sequencing platforms, but this review focuses on emerging parallel DNA-sequencing technologies based

¹ Department of Molecular Virology and Microbiology, Baylor College of Medicine, Houston, TX; ² Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX; ³ Department of Pathology, Baylor College of Medicine, Houston, TX; ⁴ Department of Pathology, Texas Children's Hospital, Houston, TX; ⁵ Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, TX; ⁶ Department of Pediatrics, Baylor College of Medicine, Houston, TX.

* Address correspondence to this author at: Department of Pathology, Texas Children's Hospital, 6621 Fannin MC 1-2261, Houston, TX 77030. Fax 832-825-0164; e-mail jamesv@bcm.edu.

Received July 31, 2008; accepted January 28, 2009.

Previously published online at DOI: 10.1373/clinchem.2008.107565

on pyrosequencing. Such next-generation sequencing systems introduce possibilities for deeply sequenced data collections and strategies aimed at microbial identification via single genetic targets or whole-genome methodologies.

Several important issues have recently emerged with respect to metagenomics and microbes. One issue is that the science of metagenomics, in contrast to individual microbial or animal genomes, is ultracomplex and challenged by the existence of vast unknown or knowledge “deserts.” Of the immense microbial taxonomic “space” in nature, only a restricted set of bacterial populations have been identified in the human body. As an example, the colonic microbiota is a vast ecosystem with approximately 800–1000 species per individual, but these estimates are in flux because the science of metagenomics and microbial pan-arrays is so new. Approximately 62% of the bacteria identified from the human intestine were previously unknown, and 80% of the bacteria identified by metagenomic sequencing were considered noncultivable (2). Only 9 of 70 bacterial phyla (divisions that vary in number depending on the taxonomic scheme) have been found in the human intestine, and 2 bacterial phyla, the *Firmicutes* and the *Bacteroidetes*, predominate in numbers (1, 3). As an example within the *Firmicutes*, the genus *Lactobacillus* includes >100 different species (<http://www.bacterio.cict.fr/l/Lactobacillus.htm>). To date, fewer than 20 *Lactobacillus* species have been found consistently in the mammalian gastrointestinal tract. These findings indicate that membership in indigenous communities is restricted to a limited subset of all bacteria, and bacterial populations are not randomly distributed in and on the human body. Preliminary studies suggest that the predominant species in the genitourinary tract and on skin sites are fundamentally different from the populations predominant in the gastrointestinal tract (4, 5).

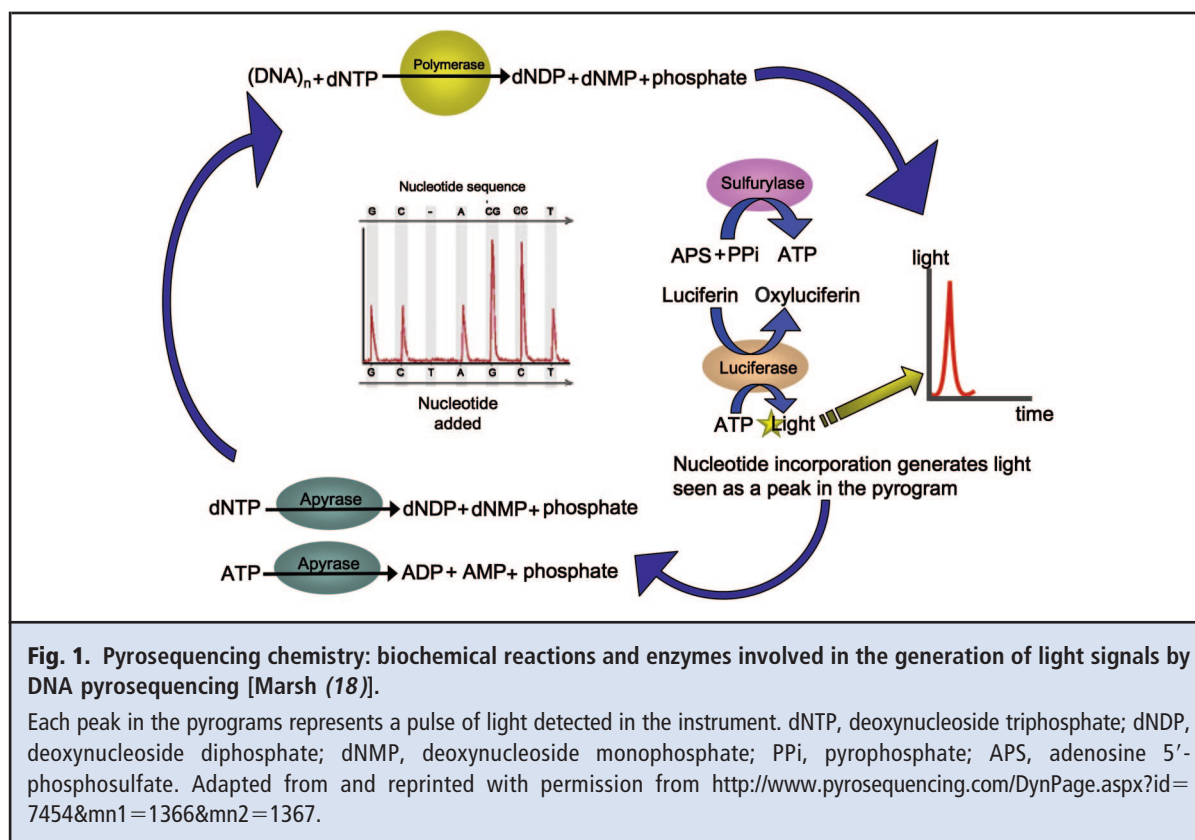
In contrast to the human genome, the human metagenomes may differ, depending on location (site) within the human body, age, and such environmental factors as diet. A key remaining question is whether a core human microbiome is definable (1). Different regions of the body, even in related and contiguous sites, may differ with respect to bacterial quantities and species composition. Bacterial species may not be randomly distributed in space or time. The large intestine contains highly complex microbial populations, and the relative proportions of different bacteria may vary in different regions of the large intestine. Culture-independent studies of cecal bacteria indicate that facultative anaerobes constitute a greater proportion of luminal bacteria, in contrast to the distal colon, where obligate anaerobes predominate. Metagenomics studies have highlighted differences between colonic

mucosa-associated populations and fecal bacterial populations in humans and nonhuman primates (6–12). In addition to differences with respect to sample type or body location, differences have also been noted with respect to sex in nonhuman primates, implying sexual dimorphism with respect to the human microbiome (12). To address the central challenge of microbial identification in the context of mixed species communities requires refining the primary strategies for DNA sequencing-based bacterial identification.

DNA Sequencing and Bacterial Identification

Pathogen identification in infectious diseases relies mostly on routine cultures and biochemical testing by means of semiautomated platforms in the clinical laboratory. The shift toward widespread adoption of nucleic acid sequencing for identification of microbial pathogens has been slowed by the user-intensive, highly technical nature of Sanger DNA sequencing. Nevertheless, several studies published in the 1990s indicated that sequencing of 16S rRNA genes could be useful for pathogen discovery and identification (13, 14). Prior studies of bacterial evolution and phylogenetics provided the foundation for subsequent applications of sequencing based on 16S rRNA genes (or 16S rDNA) for microbial identification (15). The initial studies were based on Sanger-sequencing strategies that included targeted sequencing of 16S rRNA genes (approximately 1.5 kb of target sequence). Such “long read” approaches enabled investigators and medical laboratory scientists to identify many individual genera and species that could not be identified with biochemical methods. Sequence-based identification could be established with a reasonable amount of confidence from relatively long reads and with the aid of sequence-classifier algorithms that included most of the 16S rDNA coding sequence; however, less than half of the coding sequence (approximately 500 bp), including several hypervariable regions, may be sufficient for genus- and species-level pathogen identification via Sanger sequencing (14, 16). As sequence targets for microbial identification have become more precisely defined, the introduction of pyrosequencing has provided a user-friendly approach for the clinical laboratory that has enabled more extensive sampling of microbial diversity with improved labor efficiencies (17).

Although a large body of phylogenetic data for microbial identification has been gathered via Sanger sequencing, new sequencing technologies have emerged that offer particular attractions for research and diagnostic laboratories. Specific genetic targets, such as hypervariable regions within bacterial 16S rRNA genes, may be amplified by the PCR and subjected to DNA pyrosequencing. DNA pyrosequencing, or sequencing



by synthesis, was developed in the mid 1990s as a fundamentally different approach to DNA sequencing (18). Sequencing by synthesis occurs by a DNA polymerase-driven generation of inorganic pyrophosphate, with the formation of ATP and ATP-dependent conversion of luciferin to oxyluciferin (Fig. 1). The generation of oxyluciferin causes the emission of light pulses, and the amplitude of each signal is directly related to the presence of one or more nucleosides. One important limitation of pyrosequencing is its relative inability to sequence longer stretches of DNA (sequences rarely exceed 100–200 bases with first- and second-generation high-throughput pyrosequencing chemistries). For the purposes of this review, *pyrosequencing* refers to the core chemistry, core technology, and low-throughput sequencing platforms (e.g., the Biotage PSQ 96 System) currently implemented in clinical laboratories. The term *454 sequencing* refers to high-throughput sequencing platforms (e.g., Roche/454 Life Sciences) for metagenomics that are based on pyrosequencing chemistry.

DNA pyrosequencing has been successfully applied in a variety of applications, including genotyping, single-nucleotide polymorphism detection, and microorganism identification (19). Pyrosequencing has been used to detect point mutations in antimicrobial or

antiviral resistance genes to explore the presence of drug-resistant microbes (20–22). The relatively short read lengths of DNA pyrosequencing have placed a premium on careful target selection and oligonucleotide primer placement. Pyrosequencing has been successfully applied to microbial identification by combining informative target selection (e.g., hypervariable regions within the 16S rRNA gene) and signature-sequence matching (23, 24). Despite the fact that DNA pyrosequencing yields relatively short read lengths and limited amounts of sequence data per pathogen or microbe, this strategy has been useful for microbial identification in different settings. As one example, careful selection of highly informative hypervariable regions within the 16S rRNA genes facilitated the implementation of routine pyrosequencing strategies for pathogen identification in a hospital setting (17).

Because of the relatively short read lengths, DNA-pyrosequencing applications for microbial identification have focused attention on hypervariable regions within small ribosomal-subunit RNA genes, especially 16S rRNA genes. Specific hypervariable regions have preferentially been used to identify different classes of bacteria via pyrosequencing (24, 25). Once DNA sequence data are generated, sequences must be analyzed with special considerations in mind to facilitate accu-

rate bacterial identification. First, different taxonomic classifications can be used for identification, and different species identifications may be generated, depending on the taxonomic scheme. The oldest and most traditional bacterial classification system is based on Bergey's taxonomy, which in recent years has attempted to merge phenotypic (e.g., biochemical) and molecular data to create a higher-order taxonomy (26). More recently developed taxonomic schemes include the systems proposed by Pace (27), Ludwig et al. (28), Hugenholtz (29), and the National Center for Biotechnology Information (NCBI).⁷ Multiple online databases have been developed on the basis of these different taxonomic schemes and currently provide convenient access to large rRNA sequence databases for clinical laboratories and research teams. The most prominent databases include the Ribosomal Database Project II (RDP II) (<http://rdp.cme.msu.edu/>) (30), Greengenes (<http://greengenes.lbl.gov>) (31), and ARB-SILVA (32). RDP II is based on Bergey's taxonomy, which contains a relatively small number of phyla (divisions). Greengenes includes multiple taxonomic schemes, allowing the results of queries made with different classification schemes to be compared. The ARB-SILVA database also offers a choice of microbial taxonomies, although it is more limited in its flexibility than Greengenes. Microbial identification depends on the taxonomic curation. As a case in point, the Pace and Hugenholtz lineages separately named 12 phylum-level lineages, and RDP II had not named any of these lineages (31). The taxonomic schemes varied with respect to the number of phyla; for example, there are a maximum of 88 phyla for the Pace and Hugenholtz curations and 31 phyla for the RDP classification system (based on Bergey's) (31). Therefore, in addition to the routine issues of "splitting" and "lumping" taxa in the different schemes, one is confronted with different phyla (divisions) and different corresponding subgroupings (e.g., class, order, family).

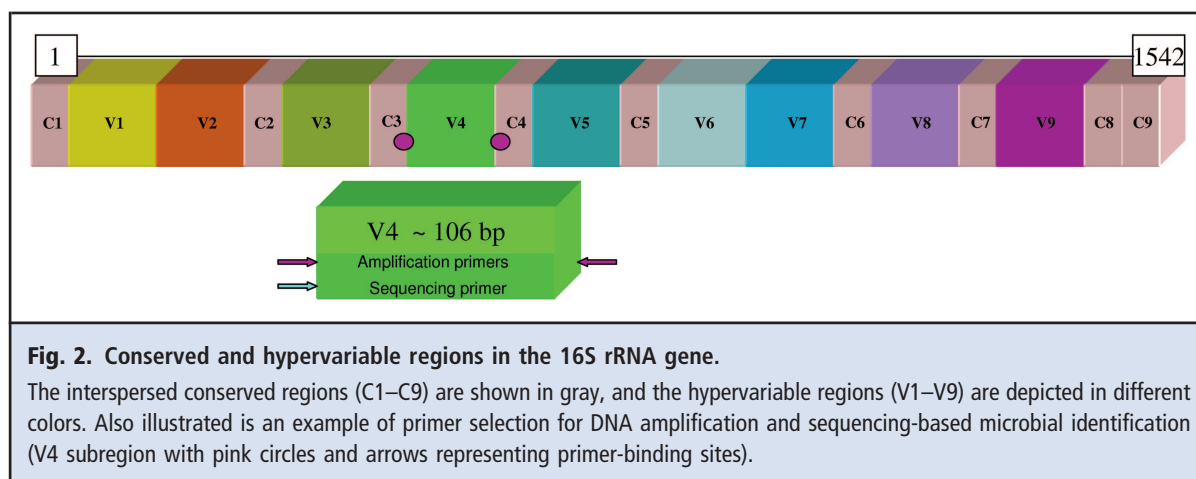
Online rRNA databases include a variety of software tools for sequence classification and multiple sequence alignments for facilitating microbial identification. The ARB software package is a widely used program suite that includes open-source, directly interacting software tools that are linked to an integrated microbial-sequence database (ARB-SILVA) (28). These software environments (Greengenes, RDP, ARB-SILVA) contain

sequence-query tools, sequence-alignment programs, and sequence editors. Greengenes provides a 16S rRNA workbench for sequence-based microbial identification with different query and sequence-alignment tools (31). Greengenes uses the NAST aligner tool (33) and generates output that is compatible with ARB software tools so that different open-source environments may be linked via the Internet for comprehensive studies of microbial populations. Different supervised sequence-classifier tools are available for matching test sequences with queried sequences. Compared with the Basic Local Alignment Search Tool (BLAST), supervised classifiers such as RDP SeqMatch demonstrated greater accuracy in finding the most similar rDNA sequences (34). The RDP-based SeqMatch *k*-nearest-neighbor (*k*-NN) classifier is effective at determining probable sequence identities on the basis of pairwise aligned distances. Alternatively, the RDP II group has developed its own naive Bayesian classifier that can be easily retrained as new sequences are incorporated into the rapidly expanding microbial-sequence databases (35). The Bayesian classifier uses information averaged within the entire genus and is less influenced by individual misplaced sequences. Sequence-query tools such as SeqMatch in RDP II enable relatively short query sequences ≥ 50 bases in length to yield accurate microbial identifications. Despite the use of 2 supervised classifiers with the same database, different results can be generated for particular sequences, particularly with phylogenetically broad genera such as *Clostridium* (35).

Next-Generation DNA-Sequencing Technologies – Pyrosequencing and 454

Until recently, Sanger-sequencing methods were primarily used to generate data in most microbial-genome and metagenomics sequencing projects. The rapid development of parallel, high-throughput sequencing technologies during the current decade has led to commercialization and widespread adoption of next-generation sequencing technologies. In contrast to a relatively homogeneous DNA-sequencing enterprise in the 1990s, current large-scale genome and metagenome sequencing projects are deploying multiple platforms and different sequencing chemistries in parallel. As of June 2008, 3 vendors of next-generation platforms commercially distribute machines for high-throughput sequencing: Roche/454 Life Sciences [Genome Sequencer 20 (GS20), FLX, LXR], Illumina/Solixa (Illumina G2), and Applied Biosystems (SOLiD). Different generations of the machines have been created, with different levels of performance (36, 37). 454 Life Sciences (now a subsidiary of Roche Diagnostics) was the one company that commercially developed py-

⁷ Nonstandard abbreviations: NCBI, National Center for Biotechnology Information; RDP II, Ribosomal Database Project II; BLAST, Basic Local Alignment Search Tool; *k*-NN, *k*-nearest-neighbor algorithm; OTU, operational taxonomic unit; GS20, Genome Sequencer 20; WGA, whole-genome amplification; prok-MSA, prokaryotic multiple sequence alignment; ASAP, automated simultaneous analysis phylogenetics; DNAML, DNA maximum likelihood.



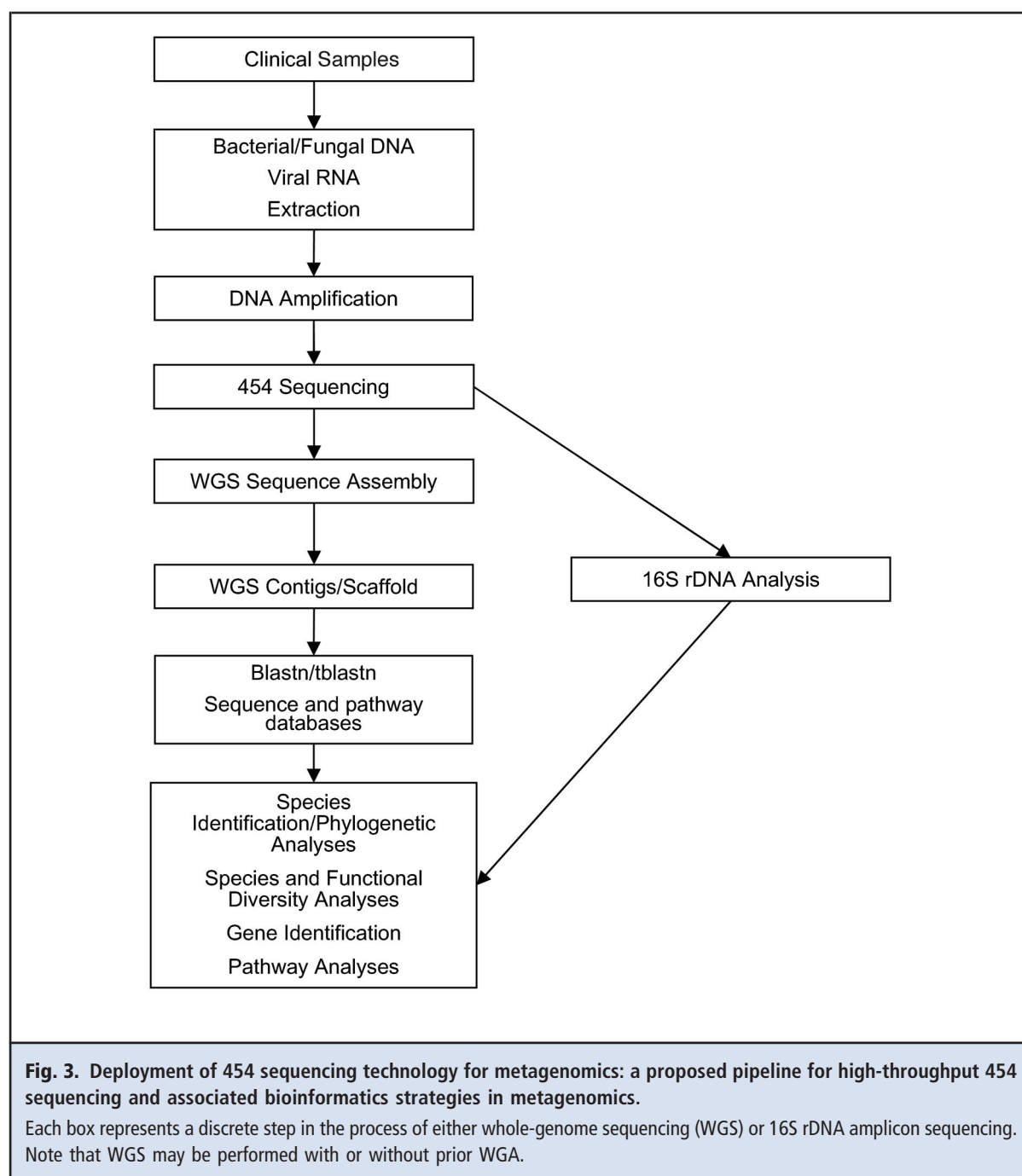
pyrosequencing for metagenomics; thus, we have used *454 sequencing* in this review to refer to high-throughput pyrosequencing.

With respect to 454 sequencing, third-generation platforms that provide longer read lengths are now emerging. The first-generation instrument, the GS20, yielded 100-bp reads and 30–60 Mb per run. The second- and third-generation instruments include the 454-FLX and 454-LXR platforms, respectively. The 454-FLX was released in 2006 and yielded 250-bp reads and approximately 150 Mb/run. The 454-LXR, released in 2008, yields demonstrably higher read lengths, exceeding 350 bp and approximately 400 Mb/run. The 454 instruments are the most widely deployed next-generation sequencing systems currently in the scientific community, and these pyrosequencing-based platforms preceded other high-throughput platforms, such as the Illumina/Solexa and SOLiD technologies mentioned above. Each 454 platform uses a modern adaptation of DNA-pyrosequencing chemistry (36, 37). Cited on the vendor's Web site are >100 publications, including 15 bacterial and 13 metagenomics reports (<http://www.454.com/news-events/publications.asp>). The first project for sequencing the human genome based on 454 sequencing was recently published (38). Generally, the sequencing community regards the 454 technology as advantageous because of the technical robustness of the chemistry. The relatively long reads generated by 454 sequencing allow more frequent unambiguous mapping to complex targets than the products of the other next-generation technologies, which feature shorter reads. During the past decade, sequencing read lengths have improved because of refinements in pyrosequencing biochemistry, such as the addition of recombinant enzymes including single-stranded binding protein (39, 40). Advances in microfluidics technologies within instruments have increased the speed of sequencing reaction cycles so that more cycles

can be performed per unit time in second- and third-generation sequencers. Additionally, the large numbers of reads per run that are possible with 454 technology deliver much greater depth of coverage for metagenomic sequencing than Sanger sequencing.

Metagenomics: Sequencing of 16S rDNA Amplicons

Metagenomics strategies may be directed at examining microbial composition or the broader issue of tackling phylogenomic diversity of highly complex microbial populations. One basic approach is to identify microbes in a complex community by exploiting universal and conserved targets, such as rRNA genes. By amplifying selected target regions within 16S rRNA genes (Fig. 2), microbes (specifically bacteria and archaea) can be identified by the effective combination of conserved primer-binding sites and intervening variable sequences that facilitate genus and species identification (Fig. 3). The 16S rRNA gene in bacteria consists of conserved sequences interspersed with variable sequences that include 9 hypervariable regions (V1–V9, Fig. 2). The lengths of these hypervariable regions range from approximately 50 bases to 100 bases, and the sequences differ with respect to variation and in their corresponding utility for universal microbial identification. Reads obtained by 454 sequencing encompass multiple hypervariable regions with the second-generation platforms such as the FLX. Third-generation 454 sequencing platforms such as the LXR will generate reads exceeding 350 bp and further facilitate the sequencing of multiple hypervariable regions. A recent study documented that the longest stretch of totally conserved bases in 16S rDNA was only 11 bases but that the longest strings of absolutely conserved bases were only 1–4 bases in most areas of this gene (41). This stark reality for a highly conserved gene highlights the enormous challenge with any met-



agenomics strategy. Different hypervariable regions demonstrated different efficacies with respect to species calls in different genera, and the V2 and V3 regions were most effective for universal genus identification (42). In a separate study, parallel analysis of 3 different hypervariable regions of 16S rDNA sequence (V2–V3, V4–V5, and V6–V8 regions) was effective in determining the composition of bacterial consortia in maize rhi-

zospheres (43). As a universal approach to the identification of bacterial pathogens, a 2-region approach yielded bacterial-genus identifications in approximately 90% of isolates not amenable to biochemical identification (17). These studies highlight the degree of variability in the representation of operational taxonomic units (OTUs), which depends on the hypervariable region used for the analysis.

To obtain a medically meaningful microbial identification, genus- or species-level classification is important. With 16S rRNA gene sequence data, genera and species are typically distinguished at levels of 95% and 97% pairwise sequence identities, respectively (44). Strains may be distinguished at the level of 99% pairwise sequence identity, although alternative molecular methods provide strain-typing approaches that are more feasible in today's clinical laboratory than DNA sequencing. Ultimately, strain-level resolution will depend on whole-genome sequencing strategies, and these methods may eventually supplant established molecular typing methods. Sequencing accuracy becomes mission-critical when metagenomics is combined with the need for identifications at the genus/species levels that may affect medical management in the future. Accurate, proofreading, and thermostable DNA polymerases and the application of temperature gradients during PCR amplification represent key considerations for maximizing the specificity of DNA amplification prior to 454 sequencing. Improving the accuracy of target amplification can minimize subsequent errors produced in high-throughput pyrosequencing.

Next-generation pyrosequencing of individual genomes and the assembly of many overlapping reads appear to yield a sequencing accuracy comparable to the current gold standard of Sanger sequencing, with error rates of 0.03%–0.07%, depending on the study (45–48). Generating consensus sequences from the assembly of overlapping reads of a single genome is not an option available for metagenomics studies, and newly developed strategies are required to minimize error rates for such community sequencing endeavors. One recent study with first-generation parallel pyrosequencing (GS20) reported the quantification of per-base error rates and error-reduction strategies, such as the removal of all reads containing any sequence ambiguities, inexact matches to the primer sequences, and read-length anomalies (49). The final conclusion of this report was that parallel pyrosequencing could surpass the accuracy of Sanger capillary sequencing in metagenomics applications. High-throughput next-generation sequencing greatly increases the depth of coverage in sequencing projects so that rich amounts of microbial diversity can be analyzed. Current 454 sequencing runs typically generate 300 000–400 000 reads per run. Our own investigations suggest that 454 sequencing can detect rare minority organisms in a microbial community, whereas approaches with relatively low depth of coverage, such as Sanger sequencing, miss these microbes entirely (S.H. and J.F.P., unpublished data).

Several metagenomics studies of the human gastrointestinal tract based on Sanger sequencing of 16S rDNA amplicons that have been published in the past several years indicate differences in composition and

the relative predominance of a few bacterial phyla. One metagenomic study described the complex gastrointestinal microbiota as spanning only 9 of 55 bacterial phylogenetic groups, with 2 predominant phyla, *Bacteroidetes* and *Firmicutes* (50). In these studies, it was clear that 7 other bacterial phyla were represented less well (*Actinobacteria*, *Fusobacteria*, *Proteobacteria*, *Verucomicrobia*, *Cyanobacteria*, *Spirochaetes*, and *VadinBE97*) (2, 51, 52). Thus far, a single archaeon, *Methanobrevibacter smithii*, has been identified as a member of the gut microbiota. Recent 454 sequencing–based metagenomics studies based on 16S rDNA amplicons have provided glimpses into the relative power of such investigations. A survey of the gut microbiome of a nonhuman primate (macaque) by 454 sequencing yielded 141 000 sequences from 100 uncultured samples obtained from 12 macaques and demonstrated clear differences, depending on anatomic location, age, and sex of the animals (12). Comparative metagenomics studies of the gut microbiomes of humans, mice, and macaques have shown clearly defined clusters, depending on the mammalian species (12). The extension of 16S rDNA amplicon sequencing strategies to whole-genome sequencing strategies potentially expands the abilities of high-throughput sequencing systems to comprehensively assess microbial diversity and to identify pathogens or “pathogenic communities.”

Whole-Genome Shotgun Sequencing

Microbial 16S rDNA sequencing is considered the gold standard for characterization of microbial communities, but 16S rDNA sequencing may not be sufficiently sensitive for comprehensive microbiome studies. rRNA gene–based sequencing can detect the predominant members of the community, but these approaches may not detect the rare members of a community with divergent target sequences. Primer bias and the low depth of sampling account for some of these limitations, which could be improved with 454 sequencing of entire microbial genomes. To overcome the limitations of single gene–based amplicon sequencing by pyrosequencing, whole-genome shotgun sequencing has emerged as an attractive strategy for assessing complex microbial diversity in mixed populations. Whole genome–based approaches offer the promise of more comprehensive coverage by high-throughput, parallel DNA-sequencing platforms, because they are not limited by sequence conservation or primer-binding site variation within a specific target (Fig. 3). Whole-genome approaches enable scientists to identify and annotate diverse arrays of microbial genes that encode many different biochemical or metabolic functions. Novel genes and functions are being discovered because of the massive data sets obtained in

whole-genome shotgun sequencing of marine samples (53). Ultimately, the assessment of aggregate biological functions or community phenotypes based on functional metagenomics may depend on whole-genome metagenomic sequencing strategies. Arguably, whole-genome approaches provide the only bona fide strategies for true metagenomics studies. The challenges and limitations of whole-genome strategies include the relatively large amounts of starting material required, potential contamination of metagenomic samples with host genetic material, and high numbers of genes of unknown function or lacking quality annotation.

One key aspect of whole-genome sequencing strategies is the requirement for greater amounts of input genomic DNA for comprehensive metagenomics studies. Whole-genome amplification (WGA) may be deployed to generate ample amounts of DNA for whole-genome shotgun sequencing. WGA represents an effective technology for enabling whole-genome shotgun sequencing, because precious human samples (e.g., skin) may yield limited amounts of total DNA after extraction; however, this strategy may introduce amplification bias prior to high-throughput sequencing. Edwards et al. used WGA on samples from a deep mine and concluded that sequence bias can be minimized (54). In our estimation, WGA represents a viable approach for studies of metagenomic DNA samples and can be performed with commercially available polymerases, such as the Phi 29 DNA polymerase (REPLI-g; Qiagen). In addition to the possibility of amplification bias, the potential of WGA to coamplify contaminating human (host) DNA poses a significant challenge, and such host DNA coamplification may overwhelm the bacterial DNA sequence data in the sample. Different subtraction strategies for human DNA sequences are being developed to minimize this possible barrier.

Next-Generation Microbial-Identification Strategies: Metagenomics and Informatics

The primary challenge for metagenomics studies at the analytical end is how to obtain accurate microbial identification for hundreds or thousands of species in a reasonable time and for a reasonable cost. Current bioinformatics throughput is too slow and not sufficiently automated for large-scale projects such as the Human Microbiome Project. High-throughput methods of metagenomics rDNA analyses are needed by the scientific community and are currently in development. Clearly, sufficient computational power is necessary, although distributed computing networks and robust server technology may eventually meet current metagenomics data-analysis demands in research settings. The clinical laboratory will need to greatly enhance its

computing infrastructure and pipelines in the near future to accommodate this demand, especially in academic centers and universities.

Beginning with sequence collection and verification, algorithms must be in place to trim sequences and to vet the quality of individual reads via various strategies (Fig. 3) (49). Sequence trimming uses various algorithms to remove primer and low-quality sequence data before sequence assembly. Huse et al. have performed error analyses of V6 sequences generated by 454 sequencing and have described methods for filtering low-quality sequence data to produce robust data sets for 16S rDNA analyses (49). Once the sequence reads have been trimmed of primer and low-quality sequences, sequences can be aligned with sequence-alignment programs such as NAST (33) or MUSCLE (55, 56). Another problem is that the PCR may generate sequence chimeras because of errors that couple disparate DNA sequences during the amplification process. Chimera-checking software has been developed so that amplicons can be vetted for the presence of “sequence hybrids” in software environments such as Greengenes (31) and RDP (30), and with tools such as Bellerophon (57) or Pintail (58).

Once high-quality sequences have been obtained from mixed species communities, the next challenge is to accurately identify many microbes in parallel. Sequences can be identified with facile classifiers such as the Bayesian Classifier in the RDP system (35) and can be compared with robust multisequence alignment programs such as NAST (33) or MUSCLE (55, 56). Existing software environments such as RDP, Greengenes, or ARB-SILVA include multisequence alignment programs that can be effectively coupled with sequence editors in an integrated fashion. For large data sets, 16S rDNA sequences may be binned with programs such as FastGroupII (59), and tallies of OTUs may be generated from these bins. Aligned sequences may also be classified against databases such as prokMSA (prokaryotic multiple sequence alignment) in Greengenes, and tallies of phyla may be examined and ultimately displayed as relative-abundance histograms so that differences in proportions of different bacterial groups can be compared.

Novel informatics approaches such as CARMA enable sequences encoding protein segments as short as 27 amino acid residues from whole-genome sequencing projects to be applied in microbial-identification strategies for comparative metagenomics (60). High-throughput informatics approaches must be developed to cope with the demands of next-generation DNA sequencing. One new strategy, automated simultaneous analysis phylogenetics (ASAP) (61), offers an automated strategy for phylogenomics that may facilitate analyses of high volumes of sequence data, especially in future whole genome-based microbi-

ome explorations. In addition to accurate microbial identification, indices and algorithms have been developed to assess microbial diversity in the context of the microbiome. Phylogenetic distance matrices may be constructed in programs such as DNAML (62). Distance matrices may be transferred to DOTUR (63) for construction of collector's curves, rarefaction curves, calculations of Chao and ACE richness estimates, and computations of Simpson and Shannon indices of diversity. Reductions in microbial diversity have been associated with human disease phenotypes (64), and such diversity indices may be relevant to medical diagnostics in the future.

Special Challenges: Fungal and Viral Metagenomics

Whole-genome sequencing of metagenomic samples is likely to reveal many bacteriophage, prophage, and eukaryotic viral sequences, but viral metagenomics analyses have shown that 60% of the sequences in a viral preparation are unique, thus representing unknown viral species. As such, viral sequences may be missed by whole-genome sequencing. New viral sequences have been isolated from clinical respiratory samples (65–67). Phages have been identified in the oral cavity, urine, sputum, and serum (68). Finkbeiner et al. have done a similar study with fecal samples from pediatric patients with diarrhea (69). The studies of Allander et al. and Finkbeiner et al. involved the use of a random PCR technique to amplify nucleic acid for cloning and subsequent Sanger sequencing. Preparation of viral nucleic acids may include the filtration of samples to remove host and bacterial cells, followed by treatment of the filtrate with nucleases to remove host nucleic acids (65, 67). Such virome sequencing strategies could easily be adapted to high-throughput 454 sequencing platforms. In the area of eukaryotic metagenomics, limited studies have been performed on fungal diversity in soil (70–72) and fungi associated with plants (73). The internal transcribed spacer regions downstream of 18S rRNA genes may be useful for fungal identification (74).

Future Directions and Deeper Considerations

The science of metagenomics is currently in its pioneering stages of development as a field, and many tools and technologies are undergoing rapid evolution. In addition to paradigmatic shifts toward next-generation DNA-sequencing technology based on novel chemistries, bioinformatics tools are also being redefined in fundamental ways to accommodate the large volumes of data. In addition to massive data sets, new questions are being posited that challenge the abilities of current algorithms to deliver meaningful answers in the context of biology and medicine. The

open-source software movement and “wikinomics,” or mass collaboration approaches in biology, have already established a foundation for the metagenomics arena with software environments such as ARB (28). Complementary strategies for microbial identification that depend on pan-microbial microarrays with known sequences, such as the Phylochip (75) or the Virochip (76), may provide practical approaches for metagenomics in the clinical laboratory. Although metagenomics is not yet ripe for routine application in the clinical laboratory setting, the rapid progress with the human microbiome in recent years means that the clinical laboratory community must consider how meta-approaches in biology may be relevant to disease risk assessment, diagnosis, and management in the future world of personalized medicine. In addition to the phenotypic dimension of human biology, such as gene expression profiling, proteomics, and metabolomics, perhaps we need to extend our concept of the human genome to include the more comprehensive and plastic human metagenome in laboratory medicine. Future diagnostic tests may consider sequence polymorphisms and implied biological functions in our microbial communities as part of the dynamic assessment of health status and disease management.

Author Contributions: All authors confirmed they have contributed to the intellectual content of this paper and have met the following 3 requirements: (a) significant contributions to the conception and design, acquisition of data, or analysis and interpretation of data; (b) drafting or revising the article for intellectual content; and (c) final approval of the published article.

Authors' Disclosures of Potential Conflicts of Interest: Upon manuscript submission, all authors completed the Disclosures of Potential Conflict of Interest form. Potential conflicts of interest:

Employment or Leadership: R.A. Gibbs is the Director of the Baylor College of Medicine Human Genome Sequencing Center.

Consultant or Advisory Role: None declared.

Stock Ownership: None declared.

Honoraria: None declared.

Research Funding: J.F. Petrosino, S. Highlander, R.A. Gibbs, and J. Versalovic have active funding from the NIH (NHGRI), including support for the Jumpstart phase of the Human Microbiome Project. J.F. Petrosino has current support from the NIH/NIAID-sponsored Western Regional Center of Excellence for Biodefense and Emerging Infectious Disease (WRCE, RCE VI). S. Highlander has current support from the US Department of Agriculture. J. Versalovic currently receives support from the NIH (NIDDK R01 DK065075; NCCAM R01 AT004326; NCCAM R21 AT003482), the Office of Naval Research, and the Defense Advanced Research Projects Agency (DARPA).

Expert Testimony: None declared.

Role of Sponsor: The funding organizations played no role in the design of study, choice of enrolled patients, review and interpretation of data, or preparation or approval of manuscript.

Acknowledgments: The authors acknowledge Tiffany Morgan for her assistance with manuscript preparation.

References

- Turnbaugh PJ, Ley RE, Hamady M, Fraser-Liggett CM, Knight R, Gordon JI. The Human Microbiome Project. *Nature* 2007;449:804–10.
- Eckburg PB, Bik EM, Bernstein CN, Purdom E, Dethlefsen L, Sargent M, et al. Diversity of the human intestinal microbial flora. *Science* 2005;308:1635–8.
- Wilson M. *Bacteriology of humans: an ecological perspective*. Malden, MA: Blackwell Publishing; 2008. p 266–326.
- Thies FL, König W, König B. Rapid characterization of the normal and disturbed vaginal microbiota by application of 16S rRNA gene terminal RFLP fingerprinting. *J Med Microbiol* 2007;56:755–61.
- Grice EA, Kong HH, Renaud G, Young AC, Bouffard GG, Blakesley RW, et al. A diversity profile of the human skin microbiota. *Genome Res* 2008;18:1043–50.
- Delgado S, Suarez A, Mayo B. Identification of dominant bacteria in feces and colonic mucosa from healthy Spanish adults by culturing and by 16S rDNA sequence analysis. *Dig Dis Sci* 2006;51:744–51.
- Lucke K, Miehlke S, Jacobs E, Schuppler M. Prevalence of *Bacteroides* and *Prevotella* spp. in ulcerative colitis. *J Med Microbiol* 2006;55:617–24.
- Prindiville T, Cantrell M, Wilson KH. Ribosomal DNA sequence analysis of mucosa-associated bacteria in Crohn's disease. *Inflamm Bowel Dis* 2004;10:824–33.
- Hold GL, Pryde SE, Russell VJ, Furrie E, Flint HJ. Assessment of microbial diversity in human colonic samples by 16S rDNA sequence analysis. *FEMS Microbiol Ecol* 2002;39:33–9.
- Hayashi H, Sakamoto M, Benno Y. Phylogenetic analysis of the human gut microbiota using 16S rDNA clone libraries and strictly anaerobic culture-based methods. *Microbiol Immunol* 2002;46:535–48.
- Suares A, Bonnet R, Sutren M, Godon JJ, Gibson GR, Collins MD, Dore J. Direct analysis of genes encoding 16S rRNA from complex communities reveals many novel molecular species within the human gut. *Appl Environ Microbiol* 1999;65:4799–807.
- McKenna P, Hoffmann C, Minkah N, Aye PP, Lackner A, Liu Z, et al. The macaque gut microbiome in health, lentiviral infection, and chronic enterocolitis. *PLoS Pathog* 2008;4:e20.
- Relman DA, Falkow S, LeBoit PE, Perkocha LA, Min KW, Welch DF, Slater LN. The organism causing bacillary angiomatosis, peliosis hepatis, and fever and bacteremia in immunocompromised patients. *N Engl J Med* 1991;324:1514.
- Kolbert CP, Rys PN, Hopkins M, Lynch DT, Germer JJ, O'Sullivan CE, et al. 16S ribosomal DNA sequence analysis for identification of bacteria in a clinical microbiology laboratory. In: Persing DH, Tenover FD, Versalovic J, Tang Y-W, Unger ER, Relman DA, White TJ, eds. *Molecular microbiology: diagnostic principles and practice*. Washington, DC: ASM Press 2004. p 361–77.
- Winker S, Woese CR. A definition of the domains Archaea, Bacteria and Eucarya in terms of small subunit ribosomal RNA characteristics. *Syst Appl Microbiol* 1991;14:305–10.
- Kolbert CP, Persing DH. Ribosomal DNA sequencing as a tool for identification of bacterial pathogens. *Curr Opin Microbiol* 1999;2:299–305.
- Luna RA, Fasciano LR, Jones SC, Boyanton BL Jr, Ton TT, Versalovic J. DNA pyrosequencing-based bacterial pathogen identification in a pediatric hospital setting. *J Clin Microbiol* 2007;45:2985–92.
- Ronaghi M, Karamohamed S, Pettersson B, Uhlen M, Nyren P. Real-time DNA sequencing using detection of pyrophosphate release. *Anal Biochem* 1996;242:84–9.
- Marsh S. Pyrosequencing applications. *Methods Mol Biol* 2007;373:15–24.
- Hopkins KL, Arnold C, Threlfall EJ. Rapid detection of *gyrA* and *parC* mutations in quinolone-resistant *Salmonella enterica* using pyrosequencing technology. *J Microbiol Methods* 2007;68:163–71.
- Lindbäck E, Unemo M, Akhras M, Gharizadeh B, Fredlund H, Pourmand N, Wretling B. Pyrosequencing of the DNA gyrase gene in *Neisseria* species: effective indicator of ciprofloxacin resistance in *Neisseria gonorrhoeae*. *APMIS* 2006;114:837–41.
- Yang ZJ, Tu MZ, Liu J, Wang XL, Jin HZ. Comparison of amplicon-sequencing, pyrosequencing and real-time PCR for detection of YMDD mutants in patients with chronic hepatitis B. *World J Gastroenterol* 2006;12:7192–6.
- Jonasson J, Olofsson M, Monstein HJ. Classification, identification and subtyping of bacteria based on pyrosequencing and signature matching of 16S rDNA fragments. *APMIS* 2002;110:263–72.
- Tarnberg M, Jakobsson J, Jonasson J, Forsum U. Identification of randomly selected colonies of lactobacilli from normal vaginal fluid by pyrosequencing of the 16S rDNA variable V1 and V3 regions. *APMIS* 2002;110:802–10.
- Monstein H, Nikpour-Badr S, Jonasson J. Rapid molecular identification and subtyping of *Helicobacter pylori* by pyrosequencing of the 16S rDNA variable V1 and V3 regions. *FEMS Microbiol Lett* 2001;199:103–7.
- Garrity GM, Brenner DJ, Krieg NR, Staley JT. *Bergey's manual of systematic bacteriology*. 2nd ed. New York: Springer; 2005.
- Pace NR. A molecular view of microbial diversity and the biosphere. *Science* 1997;276:734–40.
- Ludwig W, Strunk O, Westram R, Richter L, Meier H, Yadhukumar, et al. ARB: a software environment for sequence data. *Nucleic Acids Res* 2004;32:1363–71.
- Hugenholtz P. Exploring prokaryotic diversity in the genomic era. *Genome Biol* 2002;3:REVIEWS0003.
- Cole JR, Chai B, Farris RJ, Wang Q, Kulam-Syed-Mohideen AS, McGarrell DM, et al. The Ribosomal Database Project (RDP-II): introducing myRDP space and quality controlled public data. *Nucleic Acids Res* 2007;35:D169–72.
- DeSantis TZ, Hugenholtz P, Larsen N, Rojas M, Brodie EL, Keller K, et al. Greengenes, a chimera-checked 16S rRNA gene database and workbench compatible with ARB. *Appl Environ Microbiol* 2006;72:5069–72.
- Pruesse E, Quast C, Knittel K, Fuchs BM, Ludwig W, Peplies J, Glockner FO. SILVA: a comprehensive online resource for quality checked and aligned ribosomal RNA sequence data compatible with ARB. *Nucleic Acids Res* 2007;35:7188–96.
- DeSantis TZ, Dubosarskiy I, Murray SR, Andersen GL. Comprehensive aligned sequence construction for automated design of effective probes (CASCADE-P) using 16S rDNA. *Bioinformatics* 2003;19:1461–8.
- Cole JR, Chai B, Farris RJ, Wang Q, Kulam SA, McGarrell DM, et al. The Ribosomal Database Project (RDP-II): sequences and tools for high-throughput rRNA analysis. *Nucleic Acids Res* 2005;33:D294–6.
- Wang Q, Garrity GM, Tiedje JM, Cole JR. Naive Bayesian classifier for rapid assignment of rRNA sequences into the new bacterial taxonomy. *Appl Environ Microbiol* 2007;73:5261–7.
- Margulies M, Egholm M, Altman WE, Attiya S, Bader JS, Bemben LA, et al. Genome sequencing in microfabricated high-density picolitre reactors. *Nature* 2005;437:376–80.
- Bentley DR. Whole-genome re-sequencing. *Curr Opin Genet Dev* 2006;16:545–52.
- Wheeler DA, Srinivasan M, Egholm M, Shen Y, Chen L, McGuire A, et al. The complete genome of an individual by massively parallel DNA sequencing. *Nature* 2008;452:872–6.
- Ronaghi M. Improved performance of pyrosequencing using single-stranded DNA-binding protein. *Anal Biochem* 2000;286:282–8.
- Mashayekhi F, Ronaghi M. Analysis of read length limiting factors in pyrosequencing chemistry. *Anal Biochem* 2007;363:275–87.
- Baker GC, Smith JJ, Cowan DA. Review and re-analysis of domain-specific 16S primers. *J Microbiol Methods* 2003;55:541–55.
- Chakravorty S, Helb D, Burday M, Connell N, Alland D. A detailed analysis of 16S ribosomal RNA gene segments for the diagnosis of pathogenic bacteria. *J Microbiol Methods* 2007;69:330–9.
- Schmalenberger A, Schwiager F, Tebbe CC. Effect of primers hybridizing to different evolutionarily conserved regions of the small-subunit rRNA gene in PCR-based microbial community analyses and genetic profiling. *Appl Environ Microbiol* 2001;67:3557–63.
- Peterson DA, Frank DN, Pace NR, Gordon JI. Metagenomic approaches for defining the pathogenesis of inflammatory bowel diseases. *Cell Host Microbe* 2008;3:417–27.
- Goldberg SM, Johnson J, Busam D, Feldblyum T, Ferreira S, Friedman R, et al. A Sanger/pyrosequencing hybrid approach for the generation of high-quality draft assemblies of marine microbial genomes. *Proc Natl Acad Sci U S A* 2006;103:11240–5.
- Moore MJ, Dhangra A, Soltis PS, Shaw R, Farmerie WG, Foltz KM, Soltis DE. Rapid and accurate pyrosequencing of angiosperm plastid genomes. *BMC Plant Biol* 2006;6:17.
- Wicker T, Schlagenhauf E, Graner A, Close TJ, Keller B, Stein N. 454 sequencing put to the test using the complex genome of barley. *BMC Genomics* 2006;7:275.

48. Gharizadeh B, Herman ZS, Eason RG, Jejelowo O, Pourmand N. Large-scale pyrosequencing of synthetic DNA: a comparison with results from Sanger dideoxy sequencing. *Electrophoresis* 2006;27:3042–7.
49. Huse SM, Huber JA, Morrison HG, Sogin ML, Welch DM. Accuracy and quality of massively parallel DNA pyrosequencing. *Genome Biol* 2007;8:R143.
50. Ley RE, Peterson DA, Gordon JL. Ecological and evolutionary forces shaping microbial diversity in the human intestine. *Cell* 2006;124:837–48.
51. Backhed F, Ley RE, Sonnenburg JL, Peterson DA, Gordon JL. Host-bacterial mutualism in the human intestine. *Science* 2005;307:1915–20.
52. Zoetendal EG, Vaughan EE, de Vos WM. A microbial world within us. *Mol Microbiol* 2006;59:1639–50.
53. Yooshef S, Sutton G, Rusch DB, Halpern AL, Williamson SJ, Remington K, et al. The Sorcerer II Global Ocean Sampling expedition: expanding the universe of protein families. *PLoS Biol* 2007;5:e16.
54. Edwards RA, Rodriguez-Brito B, Wegley L, Haynes M, Breitbart M, Peterson DM, et al. Using pyrosequencing to shed light on deep mine microbial ecology. *BMC Genomics* 2006;7:57.
55. Edgar RC. MUSCLE: a multiple sequence alignment method with reduced time and space complexity. *BMC Bioinformatics* 2004;5:113.
56. Edgar RC. MUSCLE: Multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res* 2004;32:1792–7.
57. Huber T, Faulkner G, Hugenholtz P. Bellerophon: a program to detect chimeric sequences in multiple sequence alignments. *Bioinformatics* 2004;20:2317–9.
58. Ashelford KE, Chuzhanova NA, Fry JC, Jones AJ, Weightman AJ. At least 1 in 20 16S rRNA sequence records currently held in public repositories is estimated to contain substantial anomalies. *Appl Environ Microbiol* 2005;71:7724–36.
59. Yu Y, Breitbart M, McNairnie P, Rohwer F. FastGroupII: a web-based bioinformatics platform for analyses of large 16S rDNA libraries. *BMC Bioinformatics* 2006;7:57.
60. Krause L, Diaz NN, Goesmann A, Kelley S, Nattkemper TW, Rohwer F, et al. Phylogenetic classification of short environmental DNA fragments. *Nucleic Acids Res* 2008;36:2230–9.
61. Sarkar IN, Egan MG, Coruzzi G, Lee EK, DeSalle R. Automated simultaneous analysis phylogenetics (ASAP): an enabling tool for phylogenomics. *BMC Bioinformatics* 2008;9:103.
62. Olsen GJ, Matsuda H, Hagstrom R, Overbeek R. fastDNAML: a tool for construction of phylogenetic trees of DNA sequences using maximum likelihood. *Comput Appl Biosci* 1994;10:41–8.
63. Schloss PD, Handelsman J. Introducing DOTUR, a computer program for defining operational taxonomic units and estimating species richness. *Appl Environ Microbiol* 2005;71:1501–6.
64. Manichanh C, Rigottier-Gois L, Bonnaud E, Gloux K, Pelletier E, Frangeul L, et al. Reduced diversity of faecal microbiota in Crohn's disease revealed by a metagenomic approach. *Gut* 2006;55:205–11.
65. Allander T, Andreasson K, Gupta S, Bjerkner A, Bogdanovic G, Persson MA, et al. Identification of a third human polyomavirus. *J Virol* 2007;81:4130–6.
66. Allander T, Emerson SU, Engle RE, Purcell RH, Bukh J. A virus discovery method incorporating DNase treatment and its application to the identification of two bovine parvovirus species. *Proc Natl Acad Sci U S A* 2001;98:11609–14.
67. Allander T, Tammi MT, Eriksson M, Bjerkner A, Tiveljung-Lindell A, Andersson B. Cloning of a human parvovirus by molecular screening of respiratory tract samples. *Proc Natl Acad Sci U S A* 2005;102:12891–6.
68. Gorski A, Weber-Dabrowska B. The potential role of endogenous bacteriophages in controlling invading pathogens. *Cell Mol Life Sci* 2005;62:511–9.
69. Finkbeiner SR, Allred AF, Tarr PI, Klein EJ, Kirkwood CD, Wang D. Metagenomic analysis of human diarrhea: viral detection and discovery. *PLoS Pathog* 2008;4:e1000011.
70. Anderson IC, Campbell CD, Prosser JL. Diversity of fungi in organic soils under a moorland–Scots pine (*Pinus sylvestris* L.) gradient. *Environ Microbiol* 2003;5:1121–32.
71. Anderson IC, Campbell CD, Prosser JL. Potential bias of fungal 18S rDNA and internal transcribed spacer polymerase chain reaction primers for estimating fungal biodiversity in soil. *Environ Microbiol* 2003;5:36–47.
72. Hunt J, Boddy L, Randerson PF, Rogers HJ. An evaluation of 18S rDNA approaches for the study of fungal diversity in grassland soils. *Microb Ecol* 2004;47:385–95.
73. Smit E, Leeflang P, Glandorf B, van Elsland JD, Wernars K. Analysis of fungal diversity in the wheat rhizosphere by sequencing of cloned PCR-amplified genes encoding 18S rRNA and temperature gradient gel electrophoresis. *Appl Environ Microbiol* 1999;65:2614–21.
74. White T, Bruns T, Lee S, Taylor J. Amplification and direct sequencing of fungal ribosomal RNA genes for phylogenetics. In: Innis MA, Gelfand D, Sninsky J, White T, eds. *PCR protocols: a guide to methods and applications*. New York: Academic Press; 1990. p 315–22.
75. Flanagan JL, Brodie EL, Weng L, Lynch SV, Garcia O, Brown R, et al. Loss of bacterial diversity during antibiotic treatment of intubated patients colonized with *Pseudomonas aeruginosa*. *J Clin Microbiol* 2007;45:1954–62.
76. Wang D, Coscoy L, Zylberberg M, Avila PC, Boushey HA, Ganem D, DeRisi JL. Microarray-based detection and genotyping of viral pathogens. *Proc Natl Acad Sci U S A* 2002;99:15687–92.