



















Diversity of uncultured organisms explored by rRNA sequencing

David A. Stahl, David J. Lane, Gary J. Olsen and Norman R. Pace Science, New Series, Vol. 224, No. 4647 (Apr. 27, 1984), pp. 409-411 Published by: American Association for the Advancement of Science

Analysis of Hydrothermal Vent–Associated Symbionts by Ribosomal RNA Sequences

Abstract. Elibiosomal RNA 1000A1 sequences were auch on establish he physifreie effluitations of produces in the discussion of the authority of the order malgochemonatorophy on their inversebrate hosts. Two submarine hydrothemol andials, the verture filter and average marking heaptophysical and the elion elliphysic estimated from rubinoshearing tasses assessment into author form, and indecidate sequences determined and related to one dis SRMS in a phylicing receivable task and the sequence of the sequence of the sequence of the intervention of the sequence of the sequence of the sequence of the intervention of the sequence of the sequence of the sequence of the intervention of the sequence of the sequence of the sequence of the intervention of the sequence of the sequence of the sequence of the intervention of the sequence of the sequence of the sequence of the intervention of the sequence of the sequence of the sequence of the intervention of the sequence of the sequence of the sequence of the intervention of the sequence of the sequence of the sequence of the intervention of the sequence of the sequence of the sequence of the intervention of the sequence of the sequence of the sequence of the intervention of the sequence of the sequence of the sequence of the intervention of the sequence of the sequence of the sequence of the intervention of the sequence of the sequence of the sequence of the intervention of the sequence of the sequence of the sequence of the intervention of the sequence of the sequence of the sequence of the intervention of the sequence of the sequence of the sequence of the sequence of the intervention of the sequence of the sequence of the sequence of the intervention of the sequence of the sequence of the sequence of the intervention of the sequence of the sequence

One approach to characterizing and
tivable organisms is to establish th
phylogenetic relationships to bett
known organisms by appropriate maci
molecular sequence comparisons (5). I
bosomal RNA's (rRNA) seem well-su
ed among cellular macromolecules
such analyses because of their ubig
tous distribution, functional constant
high conservation of primary structu
and apparent freedom from artifacts

	(trunk wall and tropnosome) and Calyp-	
	togena magnifica (gill tissue) and live	
	specimens of Solemya velum were ob-	
	tained (7); gill and foot tissues were	
86-	excised and frozen immediately upon	
sed	receipt. Total RNA was isolated from	
ent	homogenized tissues extracted with hot	
:na	phenol and sodium dodecyl sulfate and	
ere	fractionated by polyacrylamide gel elec-	
eir	trophoresis (Fig. 1A). After elution, the	
tic	mixtures of 5S rRNA's (host and symbi-	
ted	ont) were labeled at their 5' termini with	
085	[y-12P]ATP (adenosine triphosphate) and	
10	polynucleotide kinase or at their 3' ter-	
	mini with [5'-32P]pCp (C, cytosine) and	
	RNA ligase and were resolved by elec-	
ul-	trophoresis on 8 percent polyacrylamide	
teir	sequencing gels (Fig. 1B). All 5S rRNA's	
er-	were sequenced from both termini by	
ro-	enzymatic and chemical partial diges-	
Ri-	tions (Fig. 1C). The derived sequences	
uit-	and the alignments used for phylogenetic	
for	analysis are shown in Fig. 2.	
ui-	The relation of the symbiont 55	
cy,	rRNA's to those of better-known orga-	
re,	nisms is best understood as a phyloge-	
of	netic tree (Fig. 3). The branch lengths	



sdp OOO **Prokaryotic Taxonomy**

- Current taxonomy (nearly) coherent with phylogeneny.
- Taxonomy is informed by phenotype.
- Taxon boundaries are circumscribed by experts to attempt to give groups meaningful to the practitioners
- Uncultivated organisms not included

Phylogenetic Analysis vs. Classification

- Classification is conceptually easier to interpret.
- Often preferred when the groups are well understood.
- Phylogenetic methods are preferred for new groups or when the placement is not clear.



Cole et al. (2005) Nucleic Acids Rese 3, Database issue: D294–D296 0.1186/s40168-015-0093-6

Nucleic Acids Research

The Ribosomal Database Project (RDP-II): sequences and tools for high-throughput rRNA analysis

ABSTRACT

://rdp.cme.msu.edu

e Ribo

The RDP ar

NBSTRACT in the Ribosonal Database Project (RDP-II) provide the research committee of the research research committee of the research committee o	The Ribosom high-through	al Database Project (RDP-II): sequences and tools for put rRNA analysis
The reasoning is the second in the second second is the second second is the second se	ABSTRACT	
DESCRIPTION view to a detailed results view for any query sequence. In this	the research commu rRNA gene sequenci and a phylogenetica work for these data. are made available th rdp.orme.msu.edu/). contains 101 632 bas sequences in align throughput loois fit identification of relat testing, data navigat are provided. The Rt or comments is rdpt	Sequence Match is a complete re-implementation of the original Sequence Match method (1). Sequence Match finds sequences similar to a user's query sequences using a vord matching strategy not requiring prior alignment. Sequence Match is more accurate than BLAST (6) at finding colosyl related rRNA sequences (Table 1). The related sequences returned by Sequence Match serve as a good starting point for more detailed examination of relatedness by classical phylogenetic or other methods. The initial result page presents a <i>k</i> -nearest neighbor (<i>k</i> -NN) classifier assignment of the query value of <i>k</i> , as well as the three data filters can be changed at will in this view. The user can switch from the summary <i>k</i> -NN
then to a demined results then for any query sequences in ans	DESCRIPTION	view to a detailed results view for any query sequence. In this







60 0

SeqMatch Math

Given query sequence *A* and training sequence *B*, the k-mer similarity between *A* and *B* is defined as:

 $S_{AB} \equiv \frac{|k-mers \text{ in common}|}{\min(|k-mers \text{ in } A|, |k-mers \text{ in } B|)}$

The training sequences with the highest S_{AB} are the *nearest-neighbors* of query A

SeqMatch classifies query A into the same taxa as its nearest-neighbors



Wang et al. (2007) AEM 73(16): 5261-5267



Naive Bayesian classifier for rapid assignment of rRNA sequences into the new bacterial taxonomy.

Inco the new bacterial taxonomy. The Ribosomal Database Project (RDP) Classifier, a naive Bayesian classifier, can rapidy and accurately classify haterial 168 rRAN sequences in the new higher-order taxonomy proposed in Bergey's *taxonomic Outline of the Prokaryotas* (and ed., release 50, Springer-Verlag, New York, NY, 2004). It provides taxonomic assignments from domain to genus, with confluence estimates for each assignment. The majority of classifications (98%) were of high estimated confidence (245%) and high accursy (98%), addition to being tested with the corpus of 5.20,95 rRAN sequences from Bergey's outline, the RDP Classifier was tested with a corpus of 2.3,095 rRNA sequences as assigned by the NCBI late their thermative higher-order taxonomy. The results from leave-on-eout testing on both corpors show that the down to the genus level, and the majority of the classification errors appear to be due to anomalie in the down to the genus level, and the majority of the classification errors appear to be due to anomalie in the variant classing and the analysis of libraries of thousands of sequences. Another related tool, propresented beyers terror rates. Net RDP Classifier is suitable both for the analysis of single rRNA sequences and for the analysis of libraries of thousands of sequences. Another related tool, prevent libraries. It combines the RDP Classifier with a statistical test to flag taxa differentially prevented between samples. The RDP Classifier and RDP Library Compare are available induced http://tycc.ms.ms.edu.

http://dp.cme.msu.edu/. Starting in the mid-1980s, Carl Wosen revolutionized the field of microbiology with his rRNA-based phylogenetic comparisons delinearing the three main branches of life (28). Today, rRNA-based analysis remains a central method in microbiology, used not only to explore microbid idversity but also as a day-to-day method for bacterial identification. Identification methods are conceptually easier to interpret than molecular phylogenetic analyses and are often preferred when the groups are veal understood. Most rRNA identification (dissification)

showed that the Bayesian method can still be optimal even when this attribute independency is violated. The method has also been reported to perform well on problems issuing to the dasafication of sequence data, such as the dasafication of test documents, that more a high-dimensional feature space (utility) and (b), the set of the dasafication of test documents of the took and services related to rRNA sequences to the research community. As of Jamay 2007, the RPP maintains occes 700,000 bacterial sequences and averages over 5,000 new sequences





Estimating P(V G)								
	k-mer Set V							
		CGGCUAA	GUAAUAC	ACGGAGG	GCCGCGG			
B. su	btilis	+	+	-	-			
B. cla	usii	+	-	+	-			
B. sm	ithii	+	+	-	-			
B. fle	xus	+	+	-				
B. rui	B. ruris + +							
Prob	Probability: 99% 80% 40% 01%							
$P(V Bacillus) = (0.99 \times 0.80 \times 0.40 \times 0.01) = 0.003168 \text{ or } \approx 1 \text{ in } 316$ $P(V Sinobaca) = (0.99 \times 0.01 \times 0.20 \times 0.01) = 0.000019 \text{ or } \approx 1 \text{ in } 50505$								
When you have eliminated the impossible, whatever remains, however improbable, must be the truth.								
Arthur Conan Doyle								
http://rdp.cme	http://rdp.cme.msu.edu							







Bacterial and Archaeal Genomes Available at NCBI (9/3/2018) (Compare to >3 million rRNA genes)

11403 Complete Genomes
2025 Complete Chromosome
70950 Scaffolds
72311 Contigs only
156689 Total Genomes

Majority from a small number of species



Microbial Genomes from Uncultured Organisms

- Single Cell Genomes: Single microbial cells are separated before sequencing

 Issues: Incomplete genomes, enzymatic DNA amplification causes artifacts
- Metagenomic Binning: Grouped from metagenomic assemblies
 - **Issues:** Incomplete genomes, may mix allelic variants, contamination an issue



















- 100 130 genes commonly used
- Two main methods
 - Concatenate multiple gene sequences
 - Supertrees (combine multiple trees)























Average Identity Methodology

- Determine all orthologous pairs
- Measure the percent identity for each pair
- Take the average

://rdp.cme.msu.edu

Need to find comparable genes for Average Identity methods

- Homologous:
 - The existence of shared ancestry between a pair of genes.
- Orthologous:
 - Inherited by two organisms from the same ancestral sequence. (Usually same function.)
- Paralogous:
 - Originally created by a duplication event within a single genome. (May have different functions.)





ANI & AAI: Pros and Cons

- Advantages of Average Identity:
 - Takes into account all related data
 - Can be used to classify organisms into existing or new clades
- Disadvantages of Average Identity:
 - Can not be used directly for phylogenetic analysis
 - Collapses multidimensional into a single distance measure

http://rdp.cme.msu.edu



	W32.4W38 Nucleis Acide Research, 2018, Vol. 46, Web Server issue Published online 14 June doi: 10.1093/marifity467	2018
	The Microbial Genomes Atlas (MiGA) webserver: taxonomic and gene diversity analysis of Archaea an Bacteria at the whole genome level	nd
	Luis M. Rodriguez-R ¹ , Santosh Guntu ⁷ , William T. Harvey ¹ , Ramon Rosselló-Mora ³ , James M. Tiedje ^{2,4} , James R. Cole ² and Konstantinos T. Konstantinidis ^{1,5,*}	
	ABSTRACT	
	The small subunit ribosomal RNA gene (165 rRNA) has been successfully used to catalogue and study the diversity of prokaryotic species and communi- ties but it offers limited resolution at the species and finer levels, and cannot represent the whole- genome diversity and fluidity. To overcome these lim- ies (MIGA), a webserver that allows the classification of an unknown query genomic sequence, complete or partial, against all taxonomically classified taxa with available genome sequences, as well as com- parisons to other related genomes including uncuti- vated ones, based on the genome-saggregate Average Nucleotide and Amino Acid Identity (ANI)AAI) con- cepts. MIGA Integrates best practices in sequence be raw reads or assemblies from isolate genomes, single-cell sequences, and metagenome.	
http://rdp.cme.m	msu.edu genomes (MAGs). Further, MiGA can take as input hundreds of closely related genomes of the same or	1















