

---

# Database on the structure of large ribosomal subunit RNA

---

Peter De Rijk, Yves Van de Peer, Sabine Chapelle and Rupert De Wachter\*

Departement Biochemie, Universiteit Antwerpen (UIA), Universiteitsplein 1, B-2610 Antwerpen, Belgium

---

## ABSTRACT

**A database on large ribosomal subunit RNA is made available. It contains 258 sequences. It provides sequence, alignment and secondary structure information in computer-readable formats. Files can be obtained using ftp.**

## INTRODUCTION

This paper presents a comprehensive database of large ribosomal subunit RNA (further abbreviated as LSU rRNA) structures. Our goal is to offer researchers on-line access to LSU rRNA sequences in the form of an alignment containing secondary structure information, in a format suitable for use in computer programs. Literature references and accession numbers in sequence databases are included as well as taxonomic information.

This database will conceivably be used to perform phylogenetic analysis, to find primers or probes, and to elicit the secondary structure of newly determined sequences. It can also be an invaluable tool to find sequence errors introduced during gel reading or typing. These errors cause anomalies in the sequence or structure which can often be easily detected by alignment to a set of known sequences and comparison of their possible secondary structure.

New entries or updates in the EMBL sequence database (1) are continuously scanned for new LSU rRNA sequences using the Current Sequence Awareness program (a service of the Belgian EMBnet Node). These sequences are used to update older entries, or added to the database as new entries. They are then aligned, and their secondary structure is investigated and incorporated into the alignment using the program DCSE (2). When anomalies or errors are found in the annotations in the sequence libraries, these are corrected when possible. A note indicating the changes made is added.

## CONTENTS OF THE DATABASE

Only complete or reasonably complete sequences are being incorporated into the database. Partial sequences are excluded when the combined length of the sequence segments in *Escherichia coli* LSU rRNA homologous to the sequenced segments, amounts to less than 70% of the total *Escherichia coli* sequence. The database currently contains 258 sequences, viz. 42 eukaryal, 16 archaeal, 81 bacterial, 36 plastidial and 83 mitochondrial sequences.

Table 1 shows a list of species for which the LSU rRNA structure is recorded in the database. The same taxonomic

classification is used as in the small ribosomal subunit rRNA database (3). For the domain Eukarya, the taxonomic classification of the species is according to Brusca and Brusca (4) for the Animalia, according to Cronquist (5) for the higher plants, according to Ainsworth *et al.* (6) for the zygomycetes and ascomycetes, according to Moore (7) for the basidiomycetes, and according to Margulis *et al.* (8) for the remaining eukaryotes, viz. the Protoctista.

For the Bacteria and the Archaea, the classification is based on the construction of evolutionary trees, explained in more detail in a previous compilation of small ribosomal subunit RNA sequences (9). In short, evolutionary trees are constructed by the neighbor-joining method (10) for all new sequences retrieved from the EMBL (1) and/or GenBank (11) nucleotide libraries. According to the phylogenetic position of the sequences, they are assigned to one of the taxa listed in Fig. 1 of a previous compilation (9) and described essentially by Woese and coworkers (12, 13).

## HETEROGENEITY IN SEQUENCE AND CHAIN LENGTH

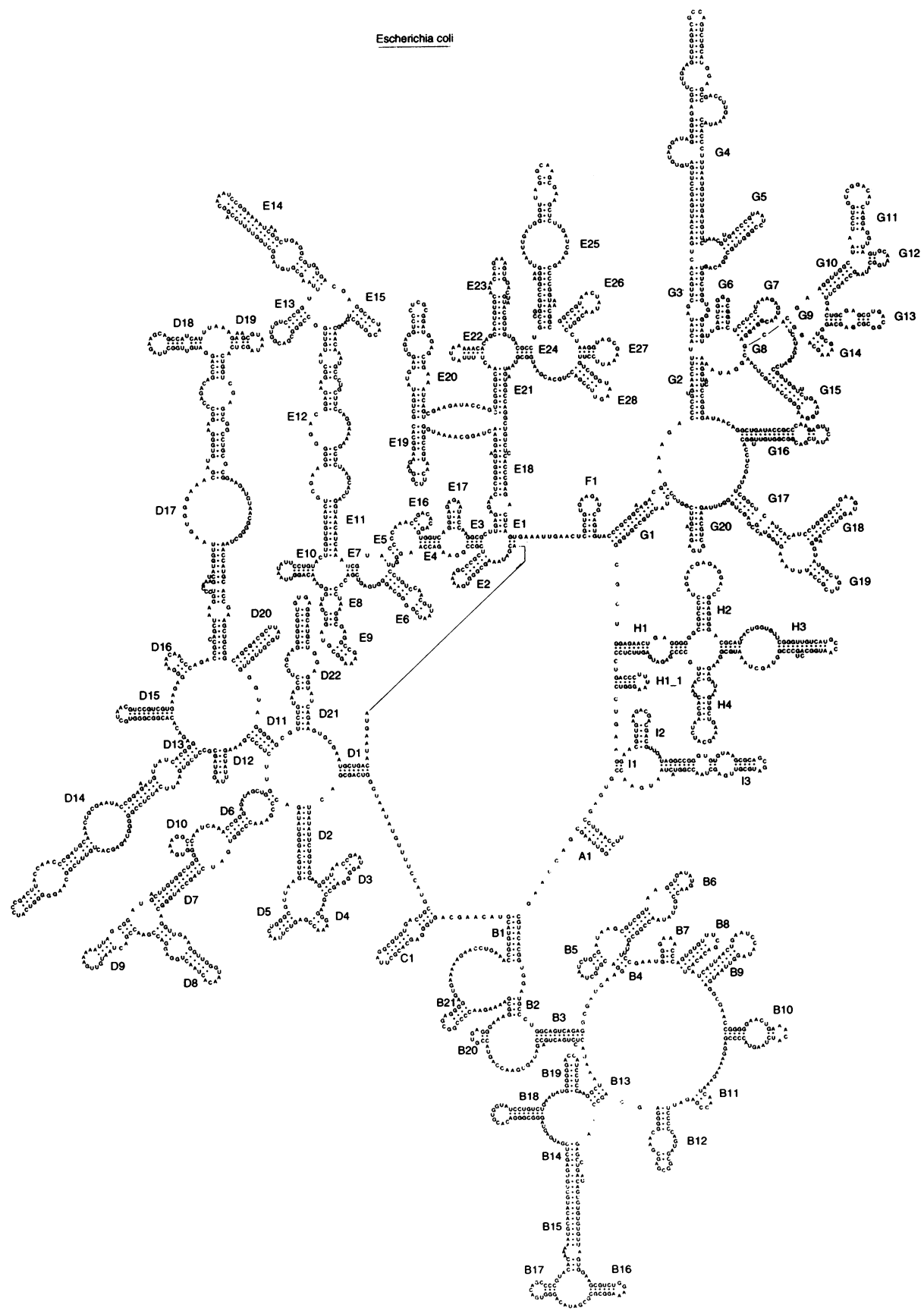
Bacterial, archaeal and plastidial LSU sequences have a relatively constant length of approximately 2900 nucleotides. However, eukaryotic sequences show a much greater diversity, ranging in length from sizes comparable to those of the bacteria to over 5000 bases in the *Homo sapiens* sequence. The presence of extra nucleotides seems to be restricted mainly to several extremely variable areas, which occupy a constant position relative to the more conserved parts of the sequences (14,15). Sequence variation is even larger in mitochondria. The molecules found in animal and kinetoplastid mitochondria even miss large parts of the sequence conserved in other LSU rRNAs, and can be under 1000 nucleotides in size. Plant and fungi mitochondrial LSU rRNAs have chain lengths comparable to or larger than those found in bacteria.

## SECONDARY STRUCTURE MODEL

Figures 1 and 2 show secondary structure models for a procaryotic (*Escherichia coli*) and a eukaryotic (*Saccharomyces cerevisiae*) LSU rRNA. A core structure is conserved in the majority of eukaryotic and bacterial LSU rRNAs. In the mitochondria of kinetoplastids and animals several helices of this core are absent. Other mitochondria have most of the helices of the core, although the structural variability is higher than among bacteria. The variable insertion regions in Eukaryotes can have

---

\*To whom correspondence should be addressed



**Figure 1.** Secondary structure model for *Escherichia coli* LSU rRNA. The sequence is written clockwise from 5' to 3' terminus.

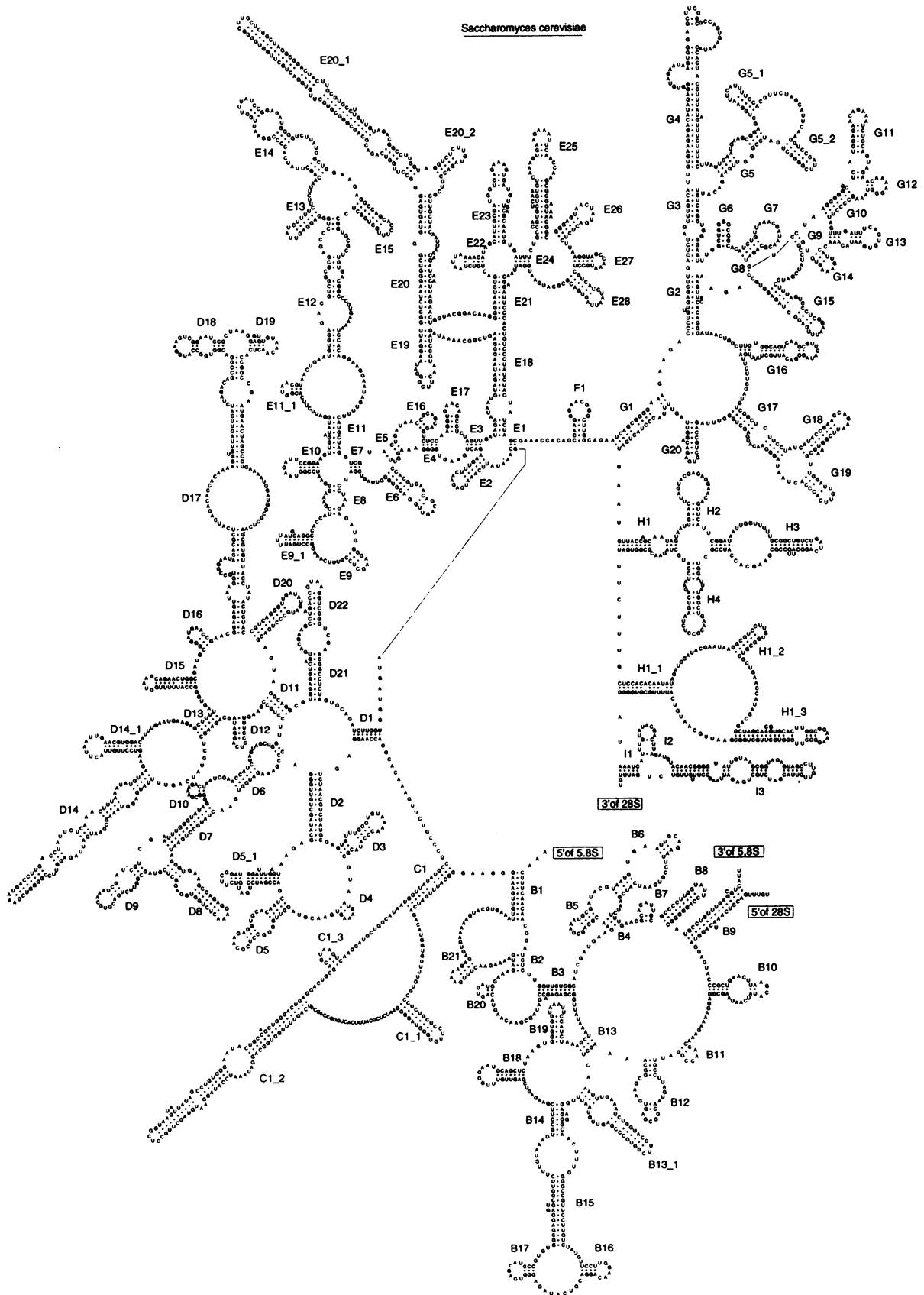


Figure 2. Secondary structure model for *Saccharomyces cerevisiae* LSU rRNA.

**Table 1.** List of species for which LSU rRNA structure is recorded in the database<sup>a</sup>

<b>ARCHAEA</b>	
<b>CRENARCHAEOTA</b>	<b>METHANOBACTERIALES</b>
Desulfurococcus mobilis	Methanobacterium thermoautotrophicum
Pyrobaculum islandicum	<b>METHANOCOCCALES</b>
Sulfolobus acidocaldarius	Methanococcus vannielii
Sulfolobus solfataricus	<b>METHANOMICROBIUM GROUP</b>
Thermofilum pendens DSM 2475	Methanospirillum hungatei
Thermoproteus tenax	<b>THERMOCOCCALES</b>
<b>EURYARCHAEOTA</b>	Thermococcus celer
<b>HALOBACTERIA</b>	<b>THERMOPLASMA</b>
Halobacterium halobium	Thermoplasma acidophilum
Halobacterium maris-mortui	
Halococcus morrhuae ATCC 17082	
Haloflex volcanii	
<b>ARCHAEOGLOBALES</b>	
Archaeoglobus fulgidus	
<b>BACTERIA</b>	
<b>PROTEOBACTERIA ALPHA</b>	Clostridium botulinum 2
Bradyrhizobium japonicum DSM 30131	Clostridium botulinum 3 ATCC 25765
Rhodobacter capsulatus DSM 938	Clostridium botulinum 4
Rhodobacter sphaeroides 1	Clostridium botulinum 5 NCTC 7272
Rhodobacter sphaeroides 2	Clostridium tyrobutyricum
Rhodobacter sphaeroides 3	Lactobacillus confusus NCDO 1586
Rhodospseudomonas palustris DSM 126	Lactobacillus delbrueckii subsp. bulgari
<b>PROTEOBACTERIA BETA</b>	Lactococcus lactis subsp. cremoris DSM 20069
Bordetella avium	Lactococcus lactis subsp. lactis 1 DSM 20481
Bordetella bronchiseptica	Lactococcus lactis subsp. lactis 2
Bordetella parapertussis	Leuconostoc carnosum
Bordetella pertussis	Leuconostoc mesenteroides
Neisseria gonorrhoeae	Leuconostoc oenos
Neisseria meningitidis	Leuconostoc paramesenteroides
Pseudomonas cepacia DSM 50181	Listeria monocytogenes 1 ATCC 19115
<b>PROTEOBACTERIA GAMMA</b>	Listeria monocytogenes 2
Aeromonas hydrophila	Mycoplasma flocculare
Escherichia coli 1	Mycoplasma hyopneumoniae ATCC 27719
Escherichia coli 2	Mycoplasma pneumoniae
Escherichia coli 3	Pectinatus frisingensis DSM 20465
Escherichia coli 4	Peptococcus niger DSM 20475
Plesiomonas shigelloides NCIMB 9242	Staphylococcus aureus ATCC 12600
Ruminobacter amylophilus	Staphylococcus carnosus DSM 20501
Pseudomonas aeruginosa ATCC 10145	Streptococcus oralis DSM 20066
Pseudomonas perfectomarina	Streptococcus parauberis NCDO 2020
<b>PROTEOBACTERIA EPSILON</b>	Streptococcus thermophilus DSM 20617
Campylobacter jejuni ATCC 43431	Streptococcus uberis NCDO 2038
<b>GRAM POSITIVES AND RELATIVES, HIGH G+C</b>	<b>CYANOBACTERIA</b>
Frankia sp. 1	Anacystis nidulans
Frankia sp. 2	<b>FLAVOBACTERIA AND RELATIVES</b>
Micrococcus luteus	Flavobacterium odoratum ATCC 4651
Mycobacterium kansasii ATCC 12478	Flexibacter flexilis ATCC 23079
Mycobacterium leprae 1	<b>GREEN SULFUR</b>
Mycobacterium leprae 2	Chlorobium limicola ATCC 8327
Streptomyces ambifaciens ATCC 23877	<b>PLANCTOMYCES AND RELATIVES</b>
Streptomyces griseus KCTC 9080	Pirellula marina
<b>GRAM POSITIVES AND RELATIVES, LOW G+C</b>	<b>SPIROCHETES</b>
Bacillus alcalophilus	Borrelia burgdorferi 1
Bacillus anthracis	Borrelia burgdorferi 2
Bacillus cereus NCTC 11143	Borrelia burgdorferi 3
Bacillus globisporus DSM 4	Leptospira interrogans
Bacillus licheniformis DSM 13	<b>RADIORESISTANT MICROCOCCI AND RELATIVES</b>
Bacillus sp.	Thermus thermophilus
Bacillus stearothermophilus 1	<b>THERMOTOGALES</b>
Bacillus stearothermophilus 2	Thermotoga maritima
Bacillus subtilis 1	
Bacillus subtilis 2	
Clostridium botulinum 1	
<b>EUKARYA</b>	
<b>ANIMALIA</b>	<b>MAGNOLIOPSIDA</b>
<b>CHORDATA</b>	Arabidopsis thaliana
<b>VERTEBRATA</b>	Brassica napus

Table 1. (cont.)

<p><b>AMPHIBIA</b> Xenopus borealis Xenopus laevis 1 Xenopus laevis 2</p> <p><b>MAMMALIA</b> Homo sapiens Mus musculus Rattus norvegicus</p> <p><b>UROCHORDATA</b> Herdmania momus</p> <p><b>ARTHROPODA</b> <b>INSECTA</b> Aedes albopictus Drosophila melanogaster</p> <p><b>NEMATODA</b> <b>SECERNENTEA</b> Caenorhabditis elegans</p>	<p>Citrus limon Fragaria ananassa Lycopersicon esculentum Sinapis alba</p>
<p><b>FUNGI</b> <b>ZYGYMYCOTINA</b> <b>ZYGYMYCETES</b> Mucor racemosus</p> <p><b>ASCOMYCOTINA</b> <b>HEMIASCOMYCETES</b> Candida albicans Saccharomyces cerevisiae Schizosaccharomyces pombe 1 Schizosaccharomyces pombe 2</p> <p><b>UNCERTAIN AFFILIATION</b> Pneumocystis carinii</p> <p><b>BASIDIOMYCOTINA</b> <b>HETEROBASIDIOMYCETES</b> Cryptococcus neoformans 1 Cryptococcus neoformans 2</p>	<p><b>PROTOCTISTA</b> <b>APICOMPLEXA</b> <b>COCCIDIA</b> Toxoplasma gondii 1 Toxoplasma gondii 2 Toxoplasma gondii 3</p> <p><b>CHYTRIDIOMYCOTA</b> <b>OOMYCOTA</b> Phytophthora megasperma</p> <p><b>CILIOPHORA</b> Tetrahymena pyriformis Tetrahymena thermophila</p> <p><b>DICTYOSTELIDA</b> Dictyostelium discoideum</p> <p><b>DINOFLLAGELLATA</b> Prorocentrum micans</p> <p><b>EUGLENIDA</b> Euglena gracilis</p> <p><b>PLASMODIAL SLIME MOLDS</b> <b>MYXOMYCOTA</b> Didymium iridis Physarum polycephalum</p> <p><b>RHIZOPODA</b> <b>LOBOSEA</b> Entamoeba histolytica</p> <p><b>ZOOMASTIGINA</b> <b>DIPLOMONADIDA</b> Giardia ardeae Giardia intestinalis Giardia muris</p> <p><b>KINETOPLASTIDA</b> Crithidia fasciculata Trypanosoma brucei</p>
<p><b>PLANTAE</b> <b>MAGNOLIOPHYTA</b> <b>LILIOPSIDA</b> Oryza sativa</p>	
<b>PLASTIDS</b>	
<p><b>PLANTAE</b> <b>BRYOPHYTA</b> <b>MARCHANTIOPSIDA</b> Marchantia polymorpha</p> <p><b>MAGNOLIOPHYTA</b> <b>LILIOPSIDA</b> Oryza sativa Zea mays</p> <p><b>MAGNOLIOPSIDA</b> Alnus incana Conopholis americana Epifagus virginiana 1 Epifagus virginiana 2 Nicotiana tabacum 1 Nicotiana tabacum 2 Pisum sativum</p>	<p>Chlamydomonas humicola Chlamydomonas indica Chlamydomonas iyengarii Chlamydomonas komma Chlamydomonas mexicana Chlamydomonas moewusii Chlamydomonas pallidostigmatica Chlamydomonas peterii Chlamydomonas pischmannii Chlamydomonas reinhardtii Chlamydomonas sp. Chlamydomonas starrii Chlamydomonas zebra Chlorella ellipsoidea</p> <p><b>PHAEOPHYTA</b> Pylaiella littoralis</p> <p><b>EUGLENIDA</b> Astasia longa Euglena gracilis 1 Euglena gracilis 2 Euglena gracilis 3 Euglena gracilis 4</p> <p><b>RHODOPHYTA</b> Palmaria palmata</p>
<p><b>PROTOCTISTA</b> <b>CHLOROPHYTA</b> <b>CHLOROPHYCEAE</b> Nanochlorum eucaryotum Chlamydomonas eugametos Chlamydomonas frankii Chlamydomonas geitleri Chlamydomonas gelatinosa</p>	
<b>MITOCHONDRIA</b>	
<p><b>ANIMALIA</b> <b>CHORDATA</b> <b>VERTEBRATA</b></p>	<p><b>MOLLUSCA</b> <b>BIVALVIA</b></p>

Table 1. (cont.)

<b>MAMMALIA</b>	<i>Mytilus edulis</i>
<i>Aepyceros melampus</i>	
<i>Antilocapra americana</i>	<b>NEMATODA</b>
<i>Balaenoptera musculus</i>	<b>SECERNENTEA</b>
<i>Balaenoptera physalus</i>	<i>Ascaris suum</i>
<i>Bos taurus</i>	<i>Caenorhabditis elegans</i>
<i>Boselaphus tragocamelus</i>	
<i>Capra hircus</i>	<b>FUNGI</b>
<i>Cephalophus maxwelli</i>	<b>ASCOMYCOTINA</b>
<i>Cervus unicolor</i>	<b>HEMIASCOMYCETES</b>
<i>Damaliscus dorcas</i>	<i>Saccharomyces cerevisiae</i> 1
<i>Didelphis virginiana</i>	<i>Saccharomyces cerevisiae</i> 2
<i>Gazella thomsoni</i>	<i>Schizosaccharomyces pombe</i>
<i>Halichoerus grypus</i>	
<i>Homo sapiens</i> 1	<b>PLECTOMYCETES</b>
<i>Homo sapiens</i> 2	<i>Aspergillus nidulans</i>
<i>Homo sapiens</i> 3	<i>Penicillium chrysogenum</i>
<i>Hydropotes inermis</i>	
<i>Kobus ellipsiprymnus</i>	<b>PYRENOMYCETES</b>
<i>Madoqua kirkii</i>	<i>Neurospora crassa</i>
<i>Muntiacus reevesi</i>	<i>Podospora anserina</i>
<i>Mus musculus</i>	
<i>Odocoileus virginianus</i>	<b>PLANTAE</b>
<i>Phoca vitulina</i>	<b>BRYOPHYTA</b>
<i>Rattus norvegicus</i> 1	<b>MARCHANTIOPSIDA</b>
<i>Rattus norvegicus</i> 2	<i>Marchantia polymorpha</i>
<i>Tragelaphus imberbis</i>	
<i>Tragulus napu</i>	<b>MAGNOLIOPHYTA</b>
	<b>LILIOPSIDA</b>
<b>AVES</b>	<i>Triticum aestivum</i>
<i>Anas platyrhynchos</i>	<i>Zea mays</i>
<i>Cairina moschata</i>	
<i>Gallus gallus</i>	<b>MAGNOLIOPSIDA</b>
	<i>Oenothera berteriana</i>
<b>AMPHIBIA</b>	
<i>Rana catesbeiana</i>	<b>PROTOCTISTA</b>
<i>Xenopus laevis</i>	<b>APICOMPLEXA</b>
	<b>HEMATOZOA</b>
<b>OSTEICHTHYES</b>	<i>Plasmodium falciparum</i>
<i>Crossostoma lacustre</i>	
<i>Cyprinus carpio</i>	<b>CHLOROPHYTA</b>
<i>Neoceratodus foresteri</i> 1	<b>CHLOROPHYCEAE</b>
<i>Neoceratodus foresteri</i> 2	<i>Chlamydomonas eugametos</i>
<i>Protopterus</i> sp.	<i>Chlamydomonas reinhardtii</i>
<i>Latimeria chalumnae</i>	<i>Prototheca wickerhamii</i>
	<i>Scenedesmus obliquus</i>
<b>ECHINODERMATA</b>	
<b>ECHINOIDEA</b>	<b>CILIOPHORA</b>
<i>Paracentrotus lividus</i>	<i>Paramecium aurelia</i>
<i>Strongylocentrotus purpuratus</i>	<i>Paramecium primaurelia</i> 1
	<i>Paramecium tetraurelia</i> 2
<b>ARTHROPODA</b>	<i>Tetrahymena pyriformis</i> 1
<b>MALACOSTRACA</b>	<i>Tetrahymena pyriformis</i> 2
<i>Artemia franciscana</i>	
<i>Artemia salina</i>	<b>ZOOMASTIGINA</b>
	<b>KINETOPLASTIDA</b>
<b>INSECTA</b>	<i>Critidia fasciculata</i>
<i>Aedes albopictus</i>	<i>Critidia oncopelti</i>
<i>Apis mellifera</i>	<i>Herpetomonas mariadeanei</i>
<i>Apis mellifera ligustica</i>	<i>Herpetomonas megaseiiae</i>
<i>Drosophila melanogaster</i>	<i>Herpetomonas muscarum</i>
<i>Drosophila yakuba</i>	<i>Herpetomonas samuelpessoai</i>
<i>Locusta migratoria</i>	<i>Leishmania tarentolae</i>
<i>Spodoptera frugiperda</i>	<i>Leptomonas</i> sp.
	<i>Trypanosoma brucei</i> 1
	<i>Trypanosoma brucei</i> 2

<sup>a</sup>In some cases, species names are listed several times followed by a sequential number, because multiple LSU rRNA sequences have been determined, usually by different authors. These sequences are not necessarily the same because they may originate from different varieties or strains, or from different genes, of the same species. The taxonomy followed for the three domains Eukarya, Archaea, and Bacteria, is as explained in the text. Plastidial and mitochondrial structures are listed according to the systematics followed for the host organism. In the case of Archaea and Bacteria, the species name is followed by the culture collection name and number if specified by the author.

differences in length of up to 900 bases. The structure of some of these regions has not yet been conclusively determined. The alignment and proposed secondary structure of the mitochondrial LSU rRNAs is less dependable because of the larger variability in both length and sequence.

The secondary structure of the molecule is treelike, with the helices forming branches which end either in a hairpin or in a multibranching loop. The stem of the tree joins the 5' and 3' end of bacterial LSU rRNAs. From this stem emanates a central multibranching loop. In Eukarya, and probably in Archaea the

stem helix is not present, but the central loop is. The following provisional helix numbering system is used in Figs 1 and 2. Structures branching from the central loop are labeled A to I, starting with the stem helix. Within each of these structures, helices bear a different number when they are separated by a multibranching loop. All numbering is sequentially from 5' to 3'. Structural elements specific to certain taxa are named after the preceding core helix followed by an underscore and number. The helix numbering may have to be revised if additional structural elements are identified in the future.

#### AVAILABILITY AND FORMAT OF THE DATABASE

The LSU rRNA database will be made available through anonymous ftp on the server uiam3.uia.ac.be (143.169.8.1). The files will also be made available to the EMBL nucleotide library for distribution. Researchers who cannot obtain the database through these channels, can request the database or parts thereof on magnetic media from the authors. The authors can be contacted by electronic mail to dwachter@reks.uia.ac.be or rrna@reks.uia.ac.be. On the server, a file called 'readme' will be present which describes the latest state of the database, giving the contents of the files and directories, and a description of the programs available for format conversion, alignment editing (2) and phylogenetic tree construction (16).

In order to simplify access to the database, each sequence is stored in a separate file, together with information about this sequence. The names of these files are produced from the species name by taking characters of the genus and species names. These are preceded by a code describing the phylogenetic group to which the species belongs. This makes it possible to either retrieve specific sequences using the full file name, or to retrieve a set of sequences belonging to a phylogenetic group using wild cards. Several sequence files can be integrated into one alignment using a program available on the server.

The format of the files is very simple, so that the files can be used readily by computer programs, or can easily be converted to formats used by specific programs. The files start with a few header lines which contain data about the sequence such as the accession number and literature reference. These are followed by the organism name. The sequence comes next. It consists of a range of nucleotide symbols interspersed with gap symbols necessary for alignment. The sequence end is indicated by an asterisk. The beginning and end of secondary structure elements are indicated by insertion of special symbols. Special 'helix numbering' files are present for researchers who wish to use the secondary structure information. When these are incorporated into an alignment, they indicate the name of each different helix segment.

When a sequence consists of several fragments resulting from processing, or of several exons, the sequence of each part ends with an asterisk, and has its own header containing the accession number, literature reference and a description of the sequence segment. However, the segments are stored in the same file and have the same organism name.

Users of the database are requested to cite this paper.

#### ACKNOWLEDGEMENTS

Peter De Rijk is research assistant of the National Fund for Scientific Research. Our research was supported by the Programme on Interuniversity Poles of Attraction (contract 23)

of the Office for Science Policy Programming of the Belgian State, and by the Fund for Collective Fundamental Research.

#### REFERENCES

1. Rice, C.M., Fuchs, R., Higgins, D.G., Stoehr, P.J. and Cameron, G.N. (1993) *Nucleic Acids Res.* **21**, 2967–2971.
2. De Rijk P. and De Wachter, R. (1993) *Comput. Applic. Biosci.*, **9**, 735–740
3. Van de Peer, Y., Van den Broeck, I., De Rijk, P. and De Wachter, R. *Nucleic Acids Res.*, this issue
4. Brusca, R.C. and Brusca, G.J. (1990) *Invertebrates*, Sinauer Associates, Inc. Sunderland.
5. Cronquist, A. (1971) *Introductory Botany*, Harper & Row, New York.
6. Ainsworth, G.C., Sparrow, F.K. and Sussman, A.S. (1973), *The Fungi: an Advanced Treatise*, Academic Press, New York, Vol. 4A.
7. Moore, R.T. (1988) in Moriarty, Ch. (ed.), *Taxonomy putting plants and animals in their place*. Royal Irish Academy, Dublin, pp. 61–88.
8. Margulis, L., Corliss, J.O., Melkonian, M. and Chapman, D.J. (eds.) (1990) *Handbook of Protozoists*, Jones and Bartlett Publishers, Boston.
9. Neefs, J.-M., Van de Peer, Y., De Rijk, P., Chapelle, S. and De Wachter, R. (1993) *Nucleic Acids Res.* **21**, 3025–3049.
10. Saitou, N. and Nei, M. (1987) *Mol. Biol. Evol.* **4**, 406–425.
11. Benson, D., Lipman, D.J. and Ostell, J. (1993) *Nucleic Acids Res.* **21**, 2963–2965.
12. Woese, C.R. (1987) *Microbiol. Rev.* **51**, 221–271.
13. Olsen, G.J., Woese, C.R. and Overbeek, R. (1994) *J. Bacteriol.* **176**, 1–6.
14. Veldman, G.M., Klootwijk, J., de Regt, V.C.H.F., Planta, R.J., Brantlant, C., Krol, A. and Ebel, J.-P. (1981) *Nucleic Acids Res.*, **9**, 6935–6952
15. Michot, B., Hassouna, N. and Bachelier, J.-P. (1984) *Nucleic Acids Res.*, **12**, 4259–4279
16. Van de Peer, Y. and De Wachter, R. (1993) *Comput. Applic. Biosci.*, **9**, 177–182