

<b>Exercise 13</b>	<b>Phylogeny of Prokaryotes based on 16S rRNA genes</b>
<b>Advisor</b>	Thomas Horath, <a href="mailto:horath@botinst.unizh.ch">horath@botinst.unizh.ch</a> , 01 634 82 41
<b>Reading</b>	Chapters in BBOM 9 <sup>th</sup> : 12.4 - 12.7 Chapters in BBOM 10 <sup>th</sup> : 11.4-11.8 BBOM: Madigan M.T., J.M. Martinko and J. Parker: "Brock - Biology of Microorganisms", 9 <sup>th</sup> edition, 1999. 10 <sup>th</sup> edition, 2003. Prentice Hall.
<b>Objectives</b>	<ul style="list-style-type: none"> <li>• Get to know some useful sequence analysis databases available on the internet.</li> <li>• Get to know a computer-based resource for sequence analysis.</li> <li>• Why did the 16S rRNA molecule become so famous?</li> <li>• Will sequence analyses conducted on different databases yield different results ?</li> </ul>
<b>Background</b>	<p>Originally, the taxonomic <b>classification of prokaryotes</b> depended exclusively on <b>phenotypic characteristics</b> such as</p> <ul style="list-style-type: none"> <li>• <b>shape</b>: coccus, rod, spirillum, vibrio, etc. (BBOM 9<sup>th</sup>, 3.4; BBOM 10<sup>th</sup>, 4.4),</li> <li>• <b>motility</b>: movement by gliding, with flagellum or flagella (monotrich, polytrich, peritrich), etc. (BBOM 9<sup>th</sup>, 3.11, BBOM 10<sup>th</sup>, 4.10-4.11)</li> <li>• <b>behavior</b>: chemotactic, phototactic, magnetotactic (BBOM 9<sup>th</sup>, 3.12; BBOM 10<sup>th</sup>, 4.12)</li> <li>• <b>membrane structure</b>: e.g. ester-lipids vs. ether-lipids (BBOM 9<sup>th</sup>, 3.5; BBOM 10<sup>th</sup>, 4.5),</li> <li>• <b>cell inclusions and surface structures</b>: slime layers, capsules, glycogen, sulfur, magnetosomes, spores etc. (BBOM 9<sup>th</sup>, 3.13-3.15; BBOM 10<sup>th</sup>, 4.13-4.15)</li> <li>• <b>metabolism</b>: phototrophic, chemotrophic, lithotrophic, organotrophic, autotrophic, heterotrophic (BBOM 9<sup>th</sup>, 4.1; BBOM 10<sup>th</sup>, 5.14)</li> <li>• <b>resistance</b> to antibiotics (BBOM 9<sup>th</sup>, 18.12; BBOM 10<sup>th</sup>, 20.12)</li> <li>• <b>cell wall</b>: Gram-positive or Gram-negative, LPS (BBOM 9<sup>th</sup>, 3.7, 3.8; BBOM 10<sup>th</sup>, 4.8-4.9)</li> <li>• <b>pathogenicity, virulence</b> etc. (BBOM 9<sup>th</sup>, 1.7, 1.8; BBOM 10<sup>th</sup>, 1.4, 1.5)</li> <li>• <b>and other characters</b>: Pigments, temperature tolerance, ecotype, etc.</li> </ul> <p>Some of these properties turned out to be good distinguishing characteristics (e.g., the Gram stain), while others are not (e.g., cell shape).</p> <p><b>Genotypic classification</b> based on nucleotide sequence comparison of <b>16S ribosomal RiboNucleicAcid (16S rRNA)</b> genes is now available as an additional taxonomic tool (BBOM 9<sup>th</sup> chapters 12.4 &amp; 12.5; BBOM 10<sup>th</sup>, 11.4, 11.5). 16S rRNA, along with the 23S rRNA, has properties which predestine it as a <b>universal phylogenetic marker</b>. All living organisms, prokaryotes as well as eukaryotes, contain the small (16S or 18S) and the large (23S or 28S) subunit ribosomal RNA where they always have the same function (BBOM Fig. 12.7 BBOM 10<sup>th</sup>, Fig. 11.8). Ribosomal RNA must have been present since very early in the development of life, because it is essential for protein synthesis. Any Mutation in the 16S rRNA gene directly can affect the proper functioning of the ribosome and lead to the elimination of less efficient organisms through selection. One may assume, therefore, that the 16S rRNA genes contain a large number of <b>highly conserved sequence patterns</b>. There are a number of sequence differences which did not impair on the</p>

	<p>functioning of the ribosome, however, and which were maintained over evolutionary times. These can be used to distinguish <b>phylogenetically different organisms</b>.</p> <p>There are regions on the 16S rRNA which are quite conserved and others which are variable. Comparing the differences in the <b>base sequence of 16S rRNA genes</b> is, therefore, a good tool to detect evolutionary changes and phylogenetic relatedness of organisms. The question remains, however, how many point mutations might have occurred at one base position up to now, from A to G to U and back to A for instance. This uncertainty and the fact that we do not know when the particular mutations took place make it difficult to use the 16S rRNA as an exact <b>evolutionary clock</b>.</p> <p>One branch of <b>bio-informatics</b> is studying the relatedness between organisms based on sequence comparison. The molecular sequence data obtained worldwide from sequencing projects are collected in public <b>databases</b>. With this information, one can test novelty and possibly function of new sequences, search for homologous patterns and regulatory domains and create hypothetical phylogenetic relationships called evolutionary trees. The information which is encoded in sequences can be analyzed by comparing against existing sequences whose functions are already known. The computer uses algorithms to find similarities between the query sequence and every sequence in the database and scores them according to the degree of relative similarity. Different databases might lead to differences in scoring depending on the algorithm and the dataset used. It is reasonable, therefore, to conduct searches on at least two different database libraries in order to obtain the best possible homologous sequences.</p> <p>As a course exercise, we will use the database provided by the National Center for Biotechnology Information (NCBI) and apply the <b>BLAST algorithm</b> to the nucleotide sequences given below (BLAST = Basic Local Alignment Search Tool). The BLAST algorithm searches for patches of similarity, and it is based on ungapped sequence alignments, i.e. the alignments created by BLAST do not allow for gaps, but BLAST does allow multiple hits on the same sequence. With this type of statistical model one increases search speed but one reduces sensitivity, i.e. BLAST might miss certain matches.</p> <p>The basic concept underlying all <b>comparing and tree building methods</b> is: Take two sequences, make them the same length and divide the number of identical nucleotides by the number of all nucleotides in one sequence (and gaps if there are any). What you get is the basic similarity percentage between two sequences (BBOM 9<sup>th</sup> Fig. 12.9; BBOM 10<sup>th</sup>, 11.10). Presently, similarities of less than 95% define different genera, less than 97% different species.</p>
<b>Literature</b>	<ul style="list-style-type: none"> <li>• Stackebrandt, Erko and Michael Goodfellow (eds.) 1991. "Nucleic Acid Techniques in Bacterial Systematics", John Wiley &amp; Sons; Chichester, New York, Brisbane, Toronto, Singapore. ISBN: 0-471-92906-9</li> <li>• Peruski Jr., Leonard F. and Anne Harwood Peruski. 1997. "The Internet and the New Biology - Tools for Genomic and Molecular Research". American Society for Microbiology, Washington DC. ISBN 1-55581-119-1</li> </ul>
<b>www. Links</b>	<ul style="list-style-type: none"> <li>• U.S. National Library of Medicine and National Institutes of Health; <b>National Center for Biotechnology Information (NCBI)</b>: <a href="http://www.ncbi.nlm.nih.gov/">http://www.ncbi.nlm.nih.gov/</a></li> <li>• The <b>Ribosomal Database Project (RDP)</b>: <a href="http://rdp.cme.msu.edu/html/index.html">http://rdp.cme.msu.edu/html/index.html</a></li> </ul>

	<ul style="list-style-type: none"> <li>• The <b>European Molecular Biology Laboratory</b>: <a href="http://www.embl-heidelberg.de/">http://www.embl-heidelberg.de/</a></li> <li>• EMBL's <b>European Bioinformatics Institute (EBI)</b>: <a href="http://www.ebi.ac.uk/index.html">http://www.ebi.ac.uk/index.html</a></li> <li>• The <b>Oligonucleotide Probe Base</b>: <a href="http://www.microbial-ecology.net/probebase/default.asp">http://www.microbial-ecology.net/probebase/default.asp</a></li> </ul>
<b>Exercises</b>	<p>Below you will find two 16S rDNA nucleotide sequences. Try to find out the following, applying the NCBI internet resources:</p> <ol style="list-style-type: none"> <li>1. What are the sequences available in BLAST that share the highest relative level of similarity with the sequences given?</li> <li>2. To which organisms do sequences 1 and 2 belong?</li> <li>3. What is the minimal length of the nucleotide sequence to be typed into the BLAST window to still get the correct nearest organism?</li> <li>4. Is this length always the same or does it depend on the fragment location which you choose from the entire sequence?</li> <li>5. Take a small fragment of your favorite sequence and exchange one or more nucleotides in it. Does this new sequence still have the same nearest relatives?</li> <li>6. On the Ribosomal Database Project it is possible to create sequence based trees. Can we fit the inhabitant of lake Cadagno into its adequate group?</li> </ol> <p><b>Genomic sequences to be analyzed with BLAST:</b></p> <ol style="list-style-type: none"> <li>1. <b>A famous inhabitant of lake Cadagno:</b>  CGTGGCGGTATGCTTAACACATGCAAGTCAAGCTCAAGGTCTTCGGATTGAGTAG  CGTGGCGGACGGGTGAGTAAAGCGTGGGAATCTGCCTTGCAAGTGGGGGATAACCCG  GGGAAACTCGGGCTAATACCGCATACGCCCTACGGGGGAAAGGGGGCTTTGGCTCT  CGTTGCAAGATGAGCCACGTCCGATTAGCTAGTTGGTAGGGTAAAGGCCTACCAAG  GCGACGATSGGTGCTGCTGCTGAGAGGATGACCAGCCACACTGGGACTGAGACACG  GCCCAGACTCCTACGGGAGGCAGCAGTGGGAATATTGGACAATGGGGGAAACCTG  ATCCAGCAATACCGCGTGTGTGAAGAAGGCCTGCGGGTTGTAAGCACTTTTCAGTGG  GAAAGAAAACCTGGTGGTTAATACCCATCGGCTTTGACGTTACTCACAAAAGAAGCAC  CGGCTAACTCCGTGCCAGCAGCCGCGGTAATACGGAGGGTGCAAGCGTTAATCGGA  ATTACTGGGCGTAAAGCGCACGTAGGCGGCGCCGTCAGTCCGATGTGAAAGCCCTG  GGCTTAACCTGGGAAGTGCATTGGATACTGCGGCGCTAGAATGTGAAAGAGGGGAGT  GGAATTCAGGTGTAGCGGTGAAATGCGTAGAGATCTGGAGGAACACCAAGTGGCGAA  GGCGGCTCCCTGGTTCAACATTGACGCTGAGGTGCGAAAGNGTGGGTAGCAAACAG  GNTTAGATACCCTGGTAGTCCACGCNGTAAACGATGTGCACTAGCCGTTGGGTCCAT  TTAAGGGCTTAGTGGCGCATAAACGCGATAAGTCGACCGCCCTGGGGAGTACGGCCG  CAAGGTTAAACTCAAAGGAATTGACGGGGGCCCCGACAAAGCGGTGGAGCATGTGTT  TTAATTCGATGCAACGCGAAAAACCTTACCAGCCCTTGACATCCTCGGAATCTTGCA  AGATGTGAGAGTGCCTTCGGGAACCGAGAGACAGGTGCTGCATGGCTGTCGTCAGC  TCGTGTCGTGAGATGTTGGGTAAAGTCCCGTAACGACGCGCAACCCTTGTCCTTAGTT  GCCAGCGCGTCAAGGCGGGAACTCTAAGGAGACTGCCGGTGATAAACCGGAGGAAG  GTGGGGATGACGTCAAGTCATCATGGCCCTTATGGGCTGGGCTACACACGTGCTACA  ATGGCCGGTACAGAGCGTTGCGACCCCGCAGGGGTGAGCCAATCGCAGAAAACCGG  TCGTAGTCCGGATCGCAGTCTGCAACTCGACTGCGTGAAGTCGGAATCGCTAGTAAT  CGCGAATCAGCATGTCGSGGTGAATACGTTCCCGGGCCTTGACACACCGCCSGTC  ACACCATGGGAGTTGGTTGCACCAGAAGTAGATCGCTTAACCGCAAGAMGGGCGTTT  ACCACGGTGTGTACTACTGACTGGGGTGAAGTCGTACAAGG </li> </ol>

	<p>2. <b>Something present in humans?:</b></p> <p>GCTAAACCTAGCCCCAAACCCACTCCACCTTACTACCAGACAACCTTAGCCAAACCA  TTTACCCAAATAAAGTATAGGCGATAGAAATTGAAACCTGGCGCAATAGATATAGTACC  GCAAGGGAAAGATGAAAAATTATAACCAAGCATAATATAGCAAGGACTAACCCCTATAC  CTTCTGCATAATGAATTAAGTAAATAAAGTAAAGGAGAGAGCCAAAGCTAAGACCC  CCGAAACCAGACGAGCTACCTAAGAAGAGCTAAAAGAGCACACCCGTCTATGTAGCA  AAATAGTGGGAAGATTTATAGGTAGAGGCGACAAACCTACCGAGCCTGGTGATAGCTG  GTTGTCCAAGATAGAATCTTAGTTCAACTTTAAATTTGCCACAGAACCCTCTAAATCC  CCTTGTAATTTAACTGTTAGTCCAAAGAGGAACAGCTCTTTGGACACTAGGAAAAAAC  CTTGAGAGAGAGTAAAAATTTAACACCCATAGTAGGCCTAAAAGCAGCCACCAATTA  AGAAAGCGTTCAAGCTCAACACCCACTACCTAAAAAATCCCAAACATATAACTGAACT  CCTCACACCCAATTGGACCAATCTATCACCTATAGAAGAACTAATGTTAGTATAAGTA  ACATGAAAAACATTCTCCTCCGCATAAGCCTGCGTCAGATTAACAACTGAACTGACAA  TTAACAGCCCAATATCTACAATCAACCAACAAGTCATTATTACCCTCACTGTCAACCC  AACACAGGCATGCTCATAAGGAAAGGTTAAAAAAGTAAAGGAAGCTCGGCAATCTTA  CCCCGCCTGTTTACCAAAAACATCACCTCTAGCATCACCGATTAGAGGCACCGCC  TGCCCGATGACACATGTTTAAACGGCCGCGGTACCCTAACCGTGCAAGGTAGACATA  TCACTTGTTCTTAAATAGGGACCTGTATGAATGGCTCCACGAGGGTTCACTGTCT  CTTACTTTTAAACAGTGAAATTGACCTGCCCGTGAAGAGGCGGGCATAACACAGCAA  GACGAGAAGACCCTATGGAGCTTTAATTTATTAATGCAACAGTACCTAACAAACCCA  CAGGTCTAAACTACCAAAACCTGCATTAAAAATTTTCGGTTGGGGCAGCTCGGAGCA  GAACCCAACCTCCGAGCAGTACATGCTAAGACTTCACCAGTCAAAGCGAACTACTATA  CTCAATTGATCCAATAACTTGACCAACGGAACAAGTTACCCTAGGGATAACAGCGCAA  TCCTATTCTAGAGTCCATATCAACAATAGGGTTTACGACCTCGATGTTGGATCAGGAC  ATCCCGATGGTGCAGCCGCTATTAAAGGTTTCGTTTGTTCACGATTAAGTCTACGT  GATCTGAGTTCAGACCGGAGTAATCCAGGTCGGTTTCTATCTACCTTCAAATTCCTCC  CTGTACGAAAGGACAAGAGAAATAAGGCCTACTTCAAAAGCGCCTTCCCCGTAAAT  GATATCATCTCAACTTAGTATTATACCCACACCCACCCAAGAACAGGGTTT</p>
<b>Equipment</b>	Internet work stations
<b>Rules &amp; Precautions</b>	No microbial risks, only computer viruses....
<b>Experiences gained</b>	<ul style="list-style-type: none"> <li>• Familiarize yourself with a few theoretical and practical aspects of computer-aided molecular sequence analysis techniques</li> <li>• Experience how the internet can be used to get phylogenetic information from a nucleotide sequence</li> <li>• Perform sequence database searches</li> <li>• Practice to download and understand sequence records from international data bases</li> <li>• Learn about the usefulness and limitations of some computer algorithms for sequence analysis</li> </ul>
<b>Timing</b>	90 minutes
<b>Reporting</b>	Take notes on the exercise and present the results in your report and to the class.
<b>Questions to be answered</b>	<ol style="list-style-type: none"> <li>1. What is the RDP?</li> <li>2. What is an evolutionary distance (ED)?</li> <li>3. How old is Earth?</li> <li>4. How old are the oldest known microfossils?</li> <li>5. Why are ribosomal RNA genes good evolutionary chronometers? Why are they not?</li> <li>6. How does the universal tree support the idea that early Earth was very hot?</li> <li>7. Describe why the results of phylogenetic community analyses of microbial habitats have been rather surprising.</li> <li>8. Why are we using rRNA genes for studying phylogenetic relationships?</li> <li>9. What effect do slightly different sequences have on the positioning in a phylogenetic group?</li> </ol>