# Environmental Microbiology: Bioinformatic exercises II

## Identification, classification and phylogenetic tree construction of 16S rRNA sequences

26.01.2015

# II Molecular (Culture-Independent) Analyses of Microbial Communities
## "Molecular Microbial Ecology"

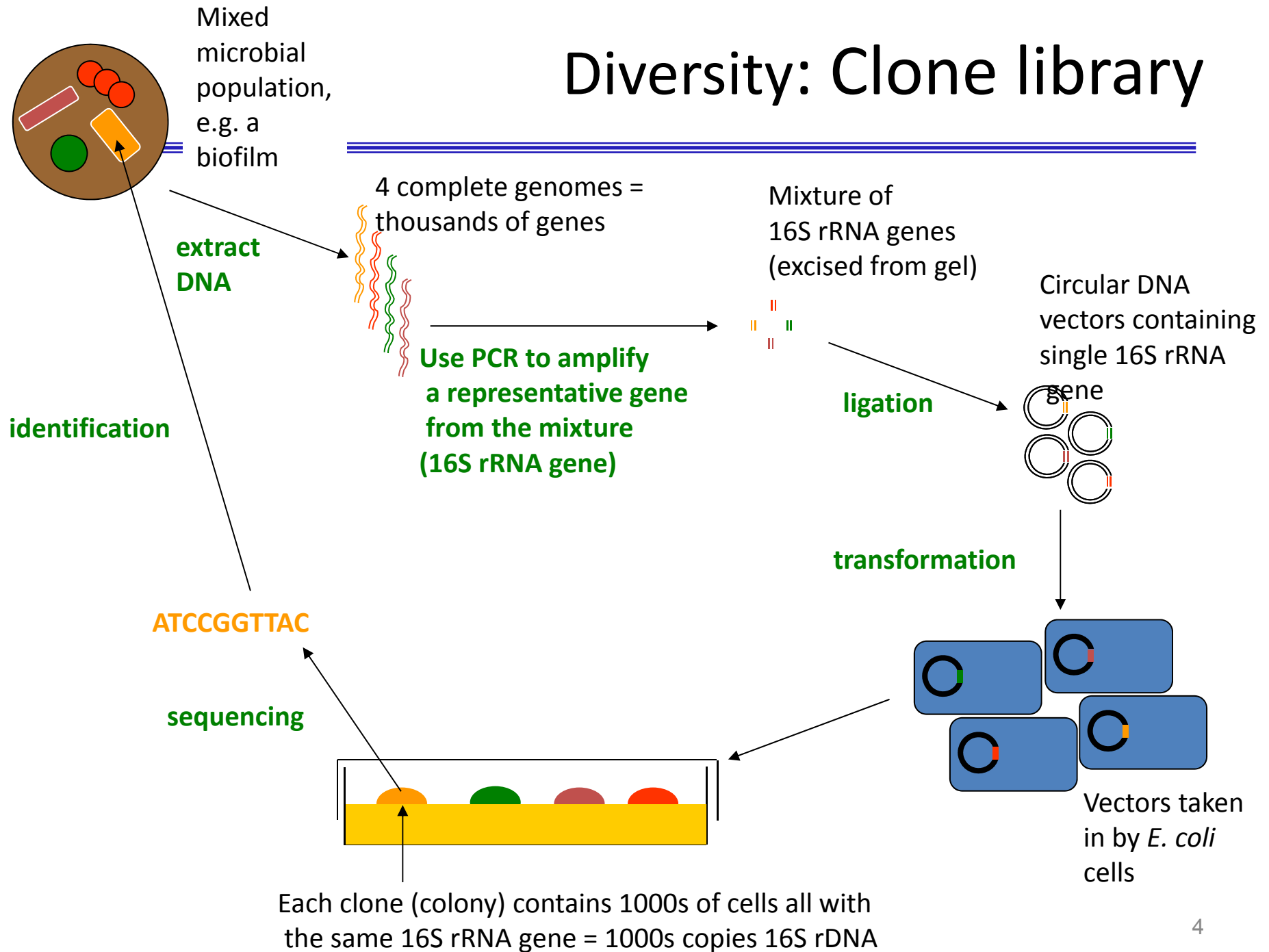## IIC Linking specific genes to specific organisms using PCR
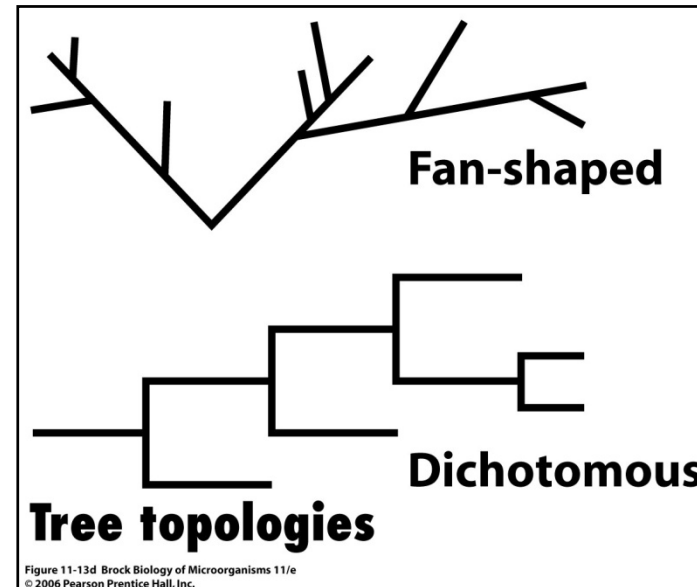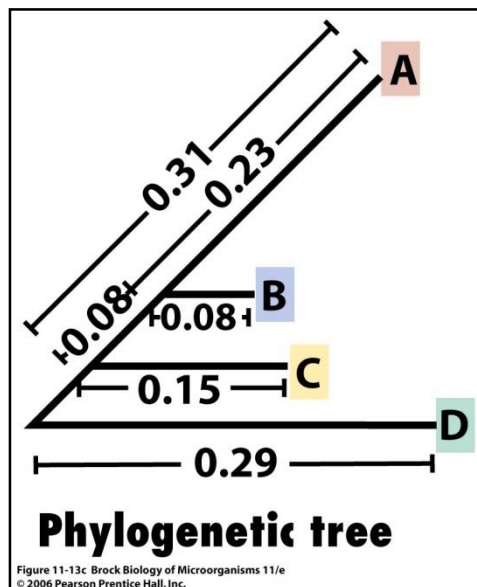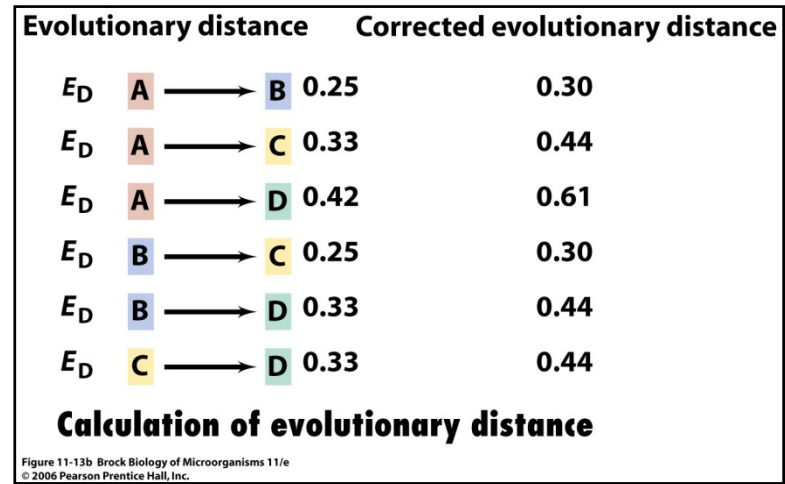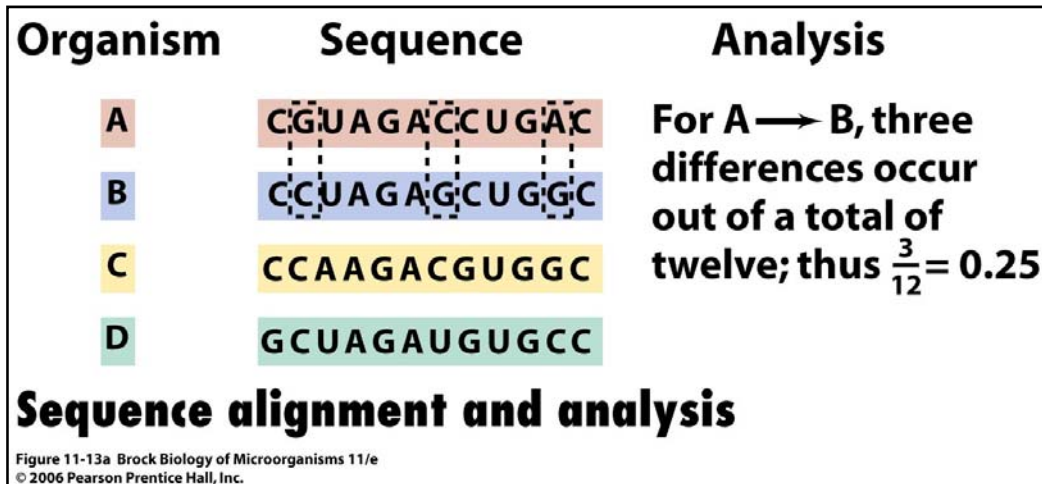
# Molecular analysis of diversity based on 16S rDNA

- **Cloning** – generation of a clone library.
  - Enables the study of 16S genes in isolation.
- **Denaturing gradient gel electrophoresis** (DGGE).
  - Separates fragments of the 16S DNA, which have a different sequence.
- **Terminal restriction fragment length polymorphism** (T-RFLP)
  - Digestion of PCR products (16S DNA) resulting in different Fragments
- **All** of the above methods depend on **PCR amplification of the target genes** from the environment.

# Diversity: Clone library

Mixed microbial population, e.g. a biofilm

**extract DNA**

4 complete genomes = thousands of genes

**Use PCR to amplify a representative gene from the mixture (16S rRNA gene)**

Mixture of 16S rRNA genes (excised from gel)

Circular DNA vectors containing single 16S rRNA gene

**ligation**

**transformation**

Vectors taken in by *E. coli* cells

**identification**

ATCCGGTTAC

**sequencing**

Each clone (colony) contains 1000s of cells all with the same 16S rRNA gene = 1000s copies 16S rDNA

4

# Sequence Analyses & Phylogenetic Tree Construction

## Sequence alignment and analysis

| Organism | Sequence | Analysis |
|----------|----------|----------|
| A | CGUAGACCUGAC | For A → B, three differences occur out of a total of twelve; thus $\frac{3}{12} = 0.25$ |
| B | CCUAGAGCUGGC | |
| C | CCAAGACGUGGC | |
| D | GCUAGAUGUGCC | |

Figure 11-13a Brock Biology of Microorganisms 11/e
© 2006 Pearson Prentice Hall, Inc.

## Calculation of evolutionary distance

| Evolutionary distance | | | | Corrected evolutionary distance |
|---|---|---|---|---|
| $E_D$ | A → B | 0.25 | | 0.30 |
| $E_D$ | A → C | 0.33 | | 0.44 |
| $E_D$ | A → D | 0.42 | | 0.61 |
| $E_D$ | B → C | 0.25 | | 0.30 |
| $E_D$ | B → D | 0.33 | | 0.44 |
| $E_D$ | C → D | 0.33 | | 0.44 |

Figure 11-13b Brock Biology of Microorganisms 11/e
© 2006 Pearson Prentice Hall, Inc.

## Phylogenetic tree

A — 0.23 — 0.31
B — 0.08 — 0.08
C — 0.15
D — 0.29

Figure 11-13c Brock Biology of Microorganisms 11/e
© 2006 Pearson Prentice Hall, Inc.

## Tree topologies

Fan-shaped

Dichotomous

Figure 11-13d Brock Biology of Microorganisms 11/e
© 2006 Pearson Prentice Hall, Inc.

# Sequencing results

You find the sequence file (MyExerziseSeq) on your desktop

```
>myExercizeSeq
AGAGTTTGATCATGGCTCAGATTGAACGCTGGCGGCAGGCCTAACACATGCAAGTCGAACGGTAACAGGAAGAAGCTTGCTTCTTTGCTGACGAG
TGGCGGACGGGTGAGTAATGTCTGGGAAACTGCCTGATGGAGGGGGATAACTACTGGAAACGGTAGCTAATACCGCATAACGTCGCAAGACCAAA
GAGGGGGACCTTCGGGCCTCTTGCCATCGGATGTGCCCAGATGGGATTAGCTAGTAGGTGGGGTAACGGCTCACCTAGGCGACGATCCCTAGCTG
GTCTGAGAGGATGACCAGCCACACTGGAACTGAGACACGGTCCAGACTCCTACGGGAGGCAGCAGTGGGGAATATTGCACAATGGGCGCAAGCCT
GATGCAGCCATGCCGCGTGTATGAAGAAGGCCTTCGGGTTGTAAAGTACTTTCAGCGGGGAGGAAGGGAGTAAAGTTAATACCTTTGCTCATTGA
CGTTACCCGCAGAAGAAGCACCGGCTAACTCCGTGCCAGCAGCCGCGGTAATACGGAGGGTGCAAGCGTTAATCGGAATTACTGGGCGTAAAGCG
CACGCAGGCGGTTTGTTAAGTCAGATGTGAAATCCCCGGGCTCAACCTGGGAACTGCATCTGATACTGGCAAGCTTGAGTCTCGTAGAGGGGGGT
AGAATTCCAGGTGTAGCGGTGAAATGCGTAGAGATCTGGAGGAATACCGGTGGCGAAGGCGGCCCCCTGGACGAAGACTGACGCTCAGGTGCGAA
AGCGTGGGGAGCAAACAGGATTAGATACCCTGGTAGTCCACGCCGTAAACGATGTCGACTTGGAGGTTGTGCCCTTGAGGCGTGGCTTCCGGAGC
TAACGCGTTAAGTCGACCGCCTGGGGAGTACGGCCGCAAGGTTAAAACTCAAATGAATTGACGGGGGCCCGCACAAGCGGTGGAGCATGTGGTTT
AATTCGATGCAACGCGAAGAACCTTACCTGGTCTTGACATCCACAGAACTTTCCAGAGATGGATTGGTGCCTTCGGGAACTGTGAGACAGGTGCT
GCATGGCTGTCGTCAGCTCGTGTTGTGAAATGTTGGGTTAAGTCCCGCAACGAGCGCAACCCTTATCTTTTGTTGCCAGCGGTCCGGCCGGGAAC
TCAAAGGAGACTGCCAGTGATAAACTGGAGGAAGGTGGGGATGACGTCAAGTCATCATGGCCCTTACGACCAGGGCTACACACGTGCTACAATGG
CGCATACAAAGAGAAGCGACCTCGCGAGAGCAAGCGGACCTCATAAAGTGCGTCGTAGTCCGGATTGGAGTCTGCAACTCGACTCCATGAAGTCG
GAATCGCTAGTAATCGTGGATCAGAATGCCACGGTGAATACGTTCCCGGGCCTTGTACACACCGCCCGTCACACCATGGGAGTGGGTTGCAAAAG
AAGTAGGTAGCTTAACCTTCGGGAGGGCGCTTACCACTTTGTGATTCATGACTGGGGTGAAGTCGTAACAAGGTAACCGTAGGGGAACC
```

# What does this sequence encode? Really an rRNA sequence? Is it the entire coding region?

# Topics

- Identification of the sequence using BLAST against GenBank (Nucleotide BLAST at NCBI) and against the 16S rRNA sequence database at NCBI
- Classification of the sequence, identification of closest relatives and phylogenetic tree construction using the the Ribosomal Database Project (RDP) (a specialized database rRNA analyses)
- Aligning sequences using Clustal Omega it the EMBL-EBI server

# BLAST - **B**asic **L**ocal **A**lignment **S**earch **T**ool

The program

- compares nucleotide or protein sequences to sequence databases

- finds regions of local similarity between sequences

- calculates the statistical significance of matches

- BLAST can be used to infer functional and evolutionary relationships between sequences as well as help identify members of gene families.

# Sequence identification

http://www.ncbi.nlm.nih.gov/ → BLAST

# Nucleotide BLAST
# (searching the GenBank database)

# Nucleotide BLAST

Copy/paste query sequence (>mySeq) into the query sequence entry field → BLAST

# Nucleotide BLAST results



- The graphical display shows the sequences that BLAST was able to align
- Alignment scores are represented on the color bar at the top of the figure, with scores going from low (black) to high (red).
- The numbered line below the color bar represents the amino acid sequence of query sequence
- Below it are various sequences from several databases that were found to align to the query. The precise position of each sequence relative to the query sequence indicates the areas of sequence similarity.

# Nucleotide BLAST results

# Nucleotide BLAST results



- E.g. the first K-12 substr. MG1655 hit in the list is only a certain region of the genome (with 100% identity to the query)
- the Genbank link in the alignment section suggests that the sequence encodes a 16S rRNA
- But it is not clear if it encodes an entire gene

# Nucleotide BLAST results



- The first hit in the list is only a certain region of the genome (with 100% identity to the query)
- the Genbank link in the alignment section suggests that the sequence encodes a 16S rRNA
- But it is not clear if it encodes an entire gene → click on graphics

# Nucleotide BLAST results



Inspect the sequence, compare your query with the „hit"
(puprle), zoom into the sequence etc.
Is the Query sequence  a full length 16S rDNA gene, which
parts are missing?

# Limitations of GenBank

- GenBank entries can contain
  - Entire genes
  - Portions of genes
  - Many genes
- GenBank entries can be of uneven quality
  - Can be duplicates and/or inaccurate
  - The database is not a selection center
  - All data is treated equally
- GenBank entries are not the final word on particular genes
  - They have no authoritative biological meaning
  - They merely keep track of what was done
- Gene-centric databases are needed to compile everything that is known on a given gene and to correct potential errors

# Nucleotide BLAST

Choose the 16S rRNA sequence database in the nucleotide blast databases menu to ensure that the query sequence encodes a 16S rRNA and to find the entire gene

# BLAST results

# BLAST results



One specific hit for the E. coli strain K-12 MG1655 in the gene centric 16S rRNA database at NCBI (all hits are reference sequences)

# BLAST results



The pairwise alignment shows that the sequence obtain experimentally does not comprise the full 16SrRNA encoding gene

# BLAST results

# Nucleotide BLAST results

- The obtained sequence indeed encodes a 16S rRNA
- Most likely from *E. coli* K-12 MG1655
  - Full sequence coverage
  - 100% identical

# rRNA sequence databases

For more specialized and accurate classification and phylogenies of rRNA (rDNA) sequences

- http://rdp.cme.msu.edu
- http://greengenes.lbl.gov/
- http://www.arb-silva.de/

# Ribosomal database project (rdp) database

# rdp database project

# Upload sequence



Choose the file
MyExercize to
upload

# Classifying the query sequence

# Classification results

# Classification results



As already indicated by the BLAST searches, you can now be quite sure that:
- The obtained sequence indeed encodes a 16S rRNA
- Most likely from the genus *Escherichia/Shigella*

# Finding the closest relatives to your sequence - Seqmatch



Choose sequence match on the rdp database start page

# Seqmatch

# Seqmatch

**Seqmatch :: Query Sequences Status**

**Running Jobs:** 2
**Pending Jobs:** 0

**Status:** running

**Current Time:** Fri Jan 23 06:49:39 EST 2015

**Progress:** 0% completed

refresh     cancel

# Seqmatch - Results

# Seqmatch - Results



The most similar sequences are found in *E. coli* strains (K-12 derivatives) with 100% identity

# Seqmatch - Results

# Four steps in building a phylogenetic tree

1. Choosing the sequence type and set
2. Alignment of sequence data (in rdp this is done automatically)
3. Search for the best tree (in rdp a distance based method i.e. special form of the neighbor joining method is used)
4. Evaluation of tree reproducibility (bootstrapping)

# Building a phylogenetic tree

# Building a phylogenetic tree

# Building a phylogenetic tree



There are more then 300,000 prokaryotic sequences available in the database
Which sequences should be used to studie phylogenetic relationships?

If you want to show broader relationships in the bacterial domain or a phylum, you can select few sequences representative for higher taxonomic ranks → choose genome browser

# Sequence selection

# Sequence selection



Then select Proteobacteria…

# Sequence selection



And finally check the „+" left to Gammaproteobacteria, thereby you select all gammaproteobacterial representative sequences to infer phylogeny (remember: *E. coli* is a gammaproteobacterium)

# Sequence selection

# Sequence selection



Nitrosomonadales → Nitrosospira → klick the „+" left to  Nitrosospira multiformis

# Sequence selection

Then click on tree builder



Nitrosomonadales → Nitrosospira → klick the „+" left to Nitrosospira multiformis

# Building the phylogenetic tree – Tree builder

# Building the phylogenetic tree – Tree builder



Then create the tree (confirm all requests from Java etc.), this might take a few minutes

# The final tree



This is the final result of the tree construction
You could now refine the study further by choosing a sequence set of lower taxonomic rank, e.g. Enterobacteriaceae

# The final tree

# Some further informations:
# Reading Your Tree

- There's a lot of vocabulary in a tree

- **Nodes** correspond to common ancestors

- The **root** is the oldest ancestor
  - Often artificial
  - Only meaningful with a good outgroup

- Trees can be un-rooted

- Branch lengths are only meaningful when the tree is scaled and refer to the degree of differences

# Building a phylogenetic Tree

- There are two types of tree-reconstruction methods
  - Distance-based methods
  - Statistical methods

- Statistical methods are the most accurate
  - Maximum likelihood of success
  - Parsimony

- Statistical methods take more time
  - Limited to small datasets

# Distance-based Methods for Tree Reconstruction

- Distance-based methods are the most popular

  – Neighbor Joining (NJ)

  – UPGMA

- Distance-based methods involve 2 steps:

  – Measure the distances between pairs of sequences in the MSA

  – Transform the distance matrix into a tree

# Bootstrapping

- Use bootstrapping to verify the solidity of each node

- ClustalW and Phylip do bootstrap operations automatically

- Bootstrapping involves these steps:
  - Select a subset of your MSA
  - Redo the tree
  - Repeat this operation N times (100 or 1000 times if you can)
  - Compute a consensus tree of the N trees
  - Measure how many of the N trees agree with the consensus tree on each node

- Each node gets a bootstrap figure between 0 and N

- High bootstrap ⇔ good node

# Doing a sequence alignment

- Different alignment programms available, e.g. Clustal, Muscel, Kalign etc.

- Can be downloaded as stand alone software (expasy.org)

- Or run on servers, e.g. ebi-embl (http://www.ebi.ac.uk/)

# Doing a sequence alignment

- Go to http://www.ebi.ac.uk/



Choose services and then DNA & RNA

# Doing a sequence alignment

- Scroll down the page and select Clustal Omega



Clustal Omega

Multiple sequence alignment of DNA or protein sequences. Clustal Omega replaces the older ClustalW alignment tools.

# Doing a sequence alignment

## Multiple Sequence Alignment

Clustal Omega is a new multiple sequence alignment program that uses seeded guide trees and HMM profile-profile techniques to generate alignments between **three or more** sequences. For the alignment of two sequences please instead use our pairwise sequence alignment tools.

### STEP 1 - Enter your input sequences

Enter or paste a set of PROTEIN ▾ sequences in any supported format:

Or, upload a file: [ Durchsuchen... ] Keine Datei ausgewählt.

Upload the file myrdp_download_26_seqs.fas, this is the same sequence set used in the rdp database analysis

### STEP 2 - Set your parameters

OUTPUT FORMAT [ Clustal w/o numbers ▾ ]

The default settings will fulfill the needs of most users and, for that reason, are not visible.

[ More options... ] (Click here, if you want to view or change the default settings.)

### STEP 3 - Submit your job

☐ Be notified by email (Tick this box if you want to be notified by email when the results are available)

[ Submit ]   ...and submit the job

# Doing a sequence alignment



- Scroll down to inspect the alignment
- Try out the options (e.g. coloring)
- A tree building option is also implemented try it out and compare the tree to that obtained from rdp

# Doing a sequence alignment

- You can also download the alignment file which can than be loaded to alignment viewer and editing software available at e.g. www.**expasy**.org/

- For example bioedit

# Most important databases

- www.ncbi.nlm.nih.gov  (The US site of the joint international DNA sequence repository (GenBank))

- www.ddbj.nig.ac.jp (Its counterpart in Japan)

- www.ebi.ac.uk/embl/ (Its counterpart in Europe (EMBL)); with links to

- www.expasy.org/sprot/ this is a very good starting point when analyzing proteins