# Algal Diversity and Phylogeny Part II

# What about Algal Phylogeny?

-How algal groups (orders, families or species) have evolved?

-Is it possible to look back into the past?

-Is it possible to discover how the species or families are related to each other?

-Is it possible to uncover when species have appeared and how they become distributed in our planet after millions and millions of years?

Yes, it can be done!  just as it is possible to determine if you are related to your family members (paternity tests)
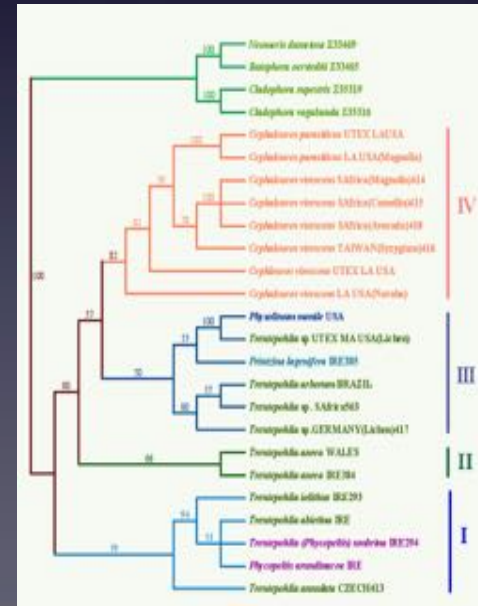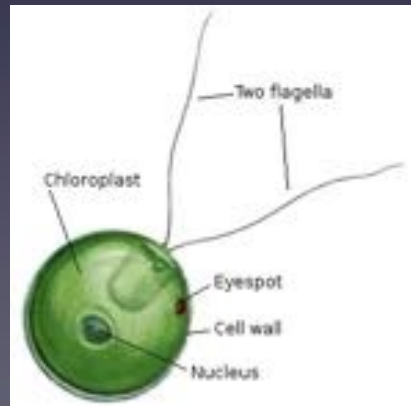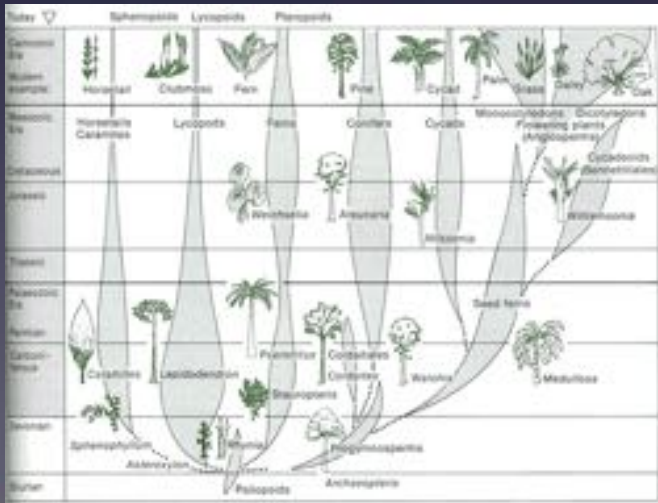
# Phylogeny reconstruction

The goal of phylogeny reconstruction is to understand patterns and processes of evolution (to explain the occurrence of particular characteristics in specific groups of organisms)

This is particularly important in the absence of fossil record

These characteristics have been typically morphological (for centuries!) but with molecular methods a new field has emerged!
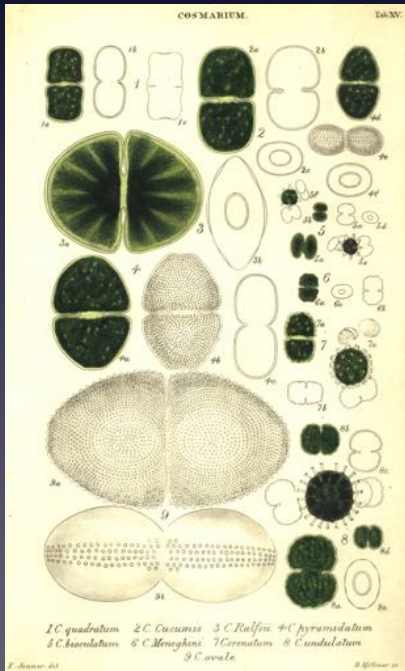
There are a variety of inference methods (phenetics, parsimony, maximum likelihood, and bayesian). These methods result in tree diagrams

# Molecular sequences vs. Structural data

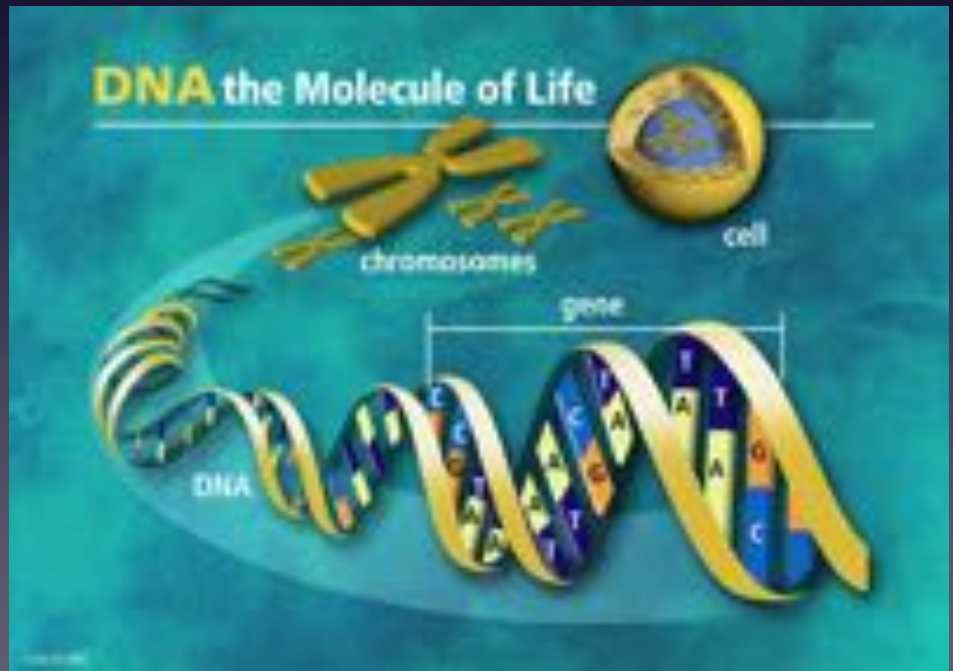**Molecular data are highly numerous (base pairs)**

**Although molecular data undergo parallel or convergent evolution, just as morphological data, its use often results in greater phylogenetic resolution than does use of morphological data alone**



*Cosmarium*
Num. of recognized species: over 500
Few morphological characters

*Versus*

Potentially over 3 billions bp available for comparison

# Choosing molecular phylogenetic approaches
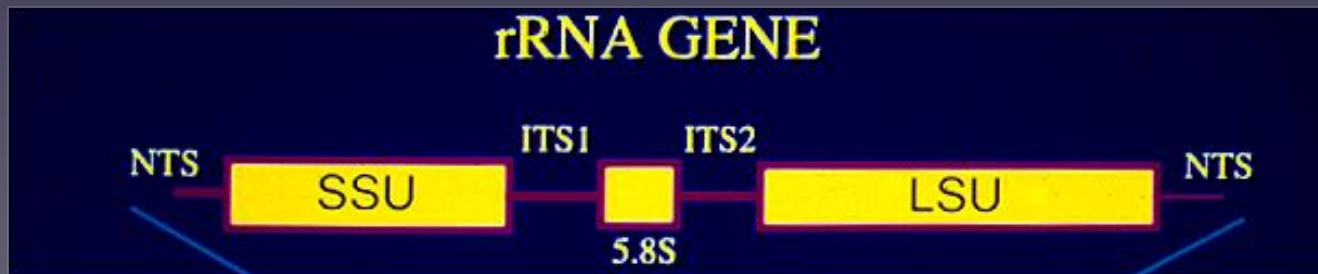
**Many different molecules are available for comparison**

**The choose of the molecule should be adequate for the specific goals**
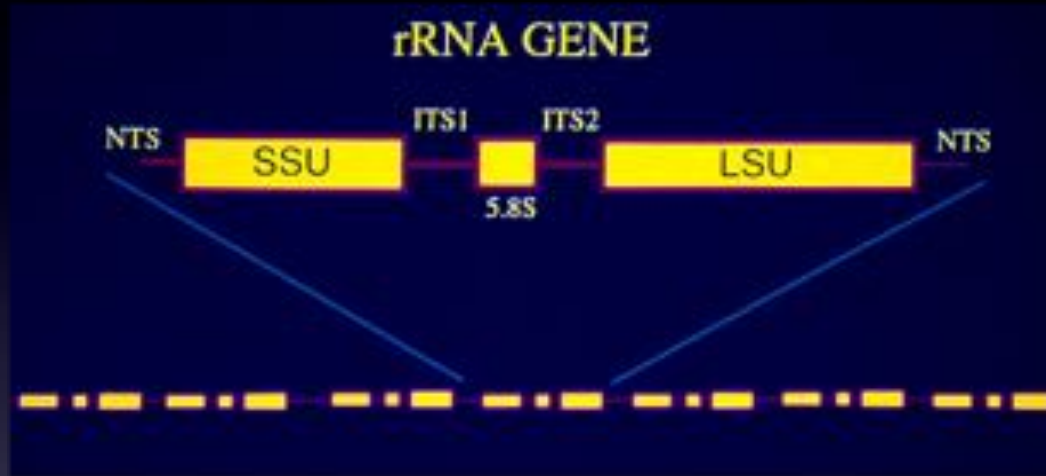
**Level of variation is important**

Ribosomal RNA sequences (rDNA) are widely used:

• Large subunit (LSU or 28S) LESS VARIABLE

• Small subunit (SSU or 18S) LESS VARIABLE with regions of high variability

• Internal transcribed spacer (ITS) HIGHLY VARIABLE

# Ribosomal RNA gene (rDNA sequences, nuclear-encoded):

Divergences greater than 500 million years ago (major algal divisions) requires sequencing of very slowly evolving genes like the Small and Large (SSU and LSU) subunits ribosomal RNA



Conservation is derived from their essential role in cell function

rDNA is present in all eukaryotes; with a high number of characters (SSU @ 1800 bp, ITS @600 bp, LSU @ 2800 bp)

It occurs in <u>tandem repeats</u>, so it can be easily amplified

# Rubisco gene (chloroplast-encoded):

Large subunit rubisco or *rbc*L is commonly used

*rbc*L has a higher base substitution rate and is much more useful for more "recent" divergence events inside of orders, families and particularly genera and species level

No insertions or deletions: 1467 bp

Easy to align

# Generating, identifying, and evaluating optimal phylogenetic trees

Trees can be generated by

a) Distance methods

b) Maximum parsimony (cladistic methods)



Samples originally corresponding to *Klebsormidium flaccidum* are spread along the *rbc*L tree

In Distance methods the number of differences between all pairs of taxa are determined

Then these numbers are are used to group taxa to form a Dendrogram (tree diagram) which attempts to accommodate all of the pairwise distances

| Species | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---------|---|---|---|---|---|---|---|---|---|----|
| I | A | T | A | T | - | C | G | T | A | T |
| II | G | T | A | T | A | C | G | T | A | T |
| III | A | T | G | G | A | C | G | T | G | C |
| Outgroup | G | C | G | T | A | T | G | C | A | C |

|  | I | II | III | Outgroup |
|---|---|---|---|---|
| I | - | 2/10 = 0.20 | 5/10 = 0.50 | 0.70 |
| II | | - | 0.50 | 0.50 |
| III | | | - | 0.60 |
| Outgroup | | | | - |

P= # differences/10 sites

The table above summarizes the calculation of pairwise distances between the gene sequences for four hypothetical species

The coefficients provide a simple summary of how similar (or different) each sequence is from the other
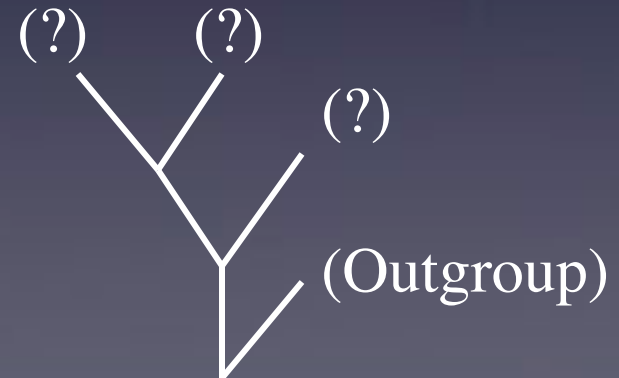Sequence I and II are more alike to each other than either is to III

<u>Parsimony methods</u> operate under the assumption (which is not always correct) that evolution operates in the most efficient (parsimonious) manner, i.e., the most accurate phylogeny is the one requiring the fewest number of changes

| Species | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---------|---|---|---|---|---|---|---|---|---|----|
| I | A | T | A | T | - | C | G | T | A | T |
| II | G | T | A | T | A | C | G | T | A | T |
| III | A | T | G | G | A | C | G | T | G | C |
| Outgroup | G | C | G | T | A | T | G | C | A | C |

Given 4 taxa (I, II, III, and outgroup) how many different ways these branches can be arranged in a tree?

(?)   (?)

(?)

(Outgroup)

There are only 3 possible branch arrangments/combinations

When we superimpose in these three possible trees the changes (-) necessary to explain each tree from the ancestral state (Outgroup) then we get different values (changes) for each tree:

I (T)   II (T)

III (C)

Outgroup(C)

1 change

I (T)   III (C)

II (T)

Outgroup(C)

2 changes

II (T)   III (C)

I (T)

Outgroup(C)

2 changes

—   = C to T change

**The left tree is the most parsimonious because it requires only 1 change!**

The entire sequence must be considered in the analysis; in this case, that means all 10 positions, not just the 10th

Once all 10 positions are considered, the tree, called Cladogram, on the left is still considered the most parsimonious, as it takes only 10 changes (rather than 11 or 12)



**Cladograms**

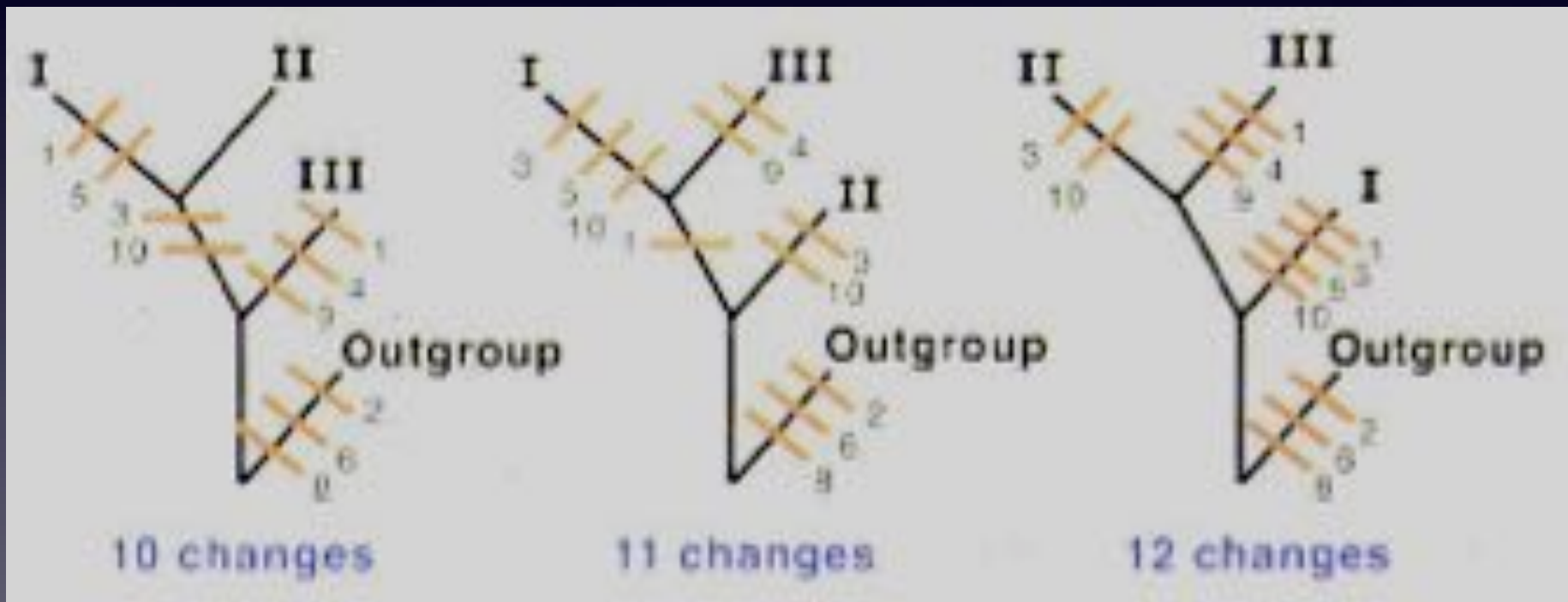The number of possible phylogenies (trees) increases with the number of taxa examined. There are 2 million possible trees for 10 taxa; for 50 taxa the number is $3 \times 10^{74}$

Therefore an <u>Exhaustive</u> search (when every possible tree is analyzed) is performed when the number of taxa is less than 12.

For studies with 12 or more taxa, <u>Heuristic</u> searches (trial and error approaches) are performed.

How reliable a tree can be?

Each individual node in the tree is evaluated by calculation of a <u>Bootstrap</u> value: it simulates the collection of replicate data sets through repeated re-sampling of the data, including some portions and leaving out others each time. The number of times a particular branch is recovered is determined as a percent value. Branches with 50% or less are considered poorly supported

# Example of a cladogram, note the similarity among branch lengths
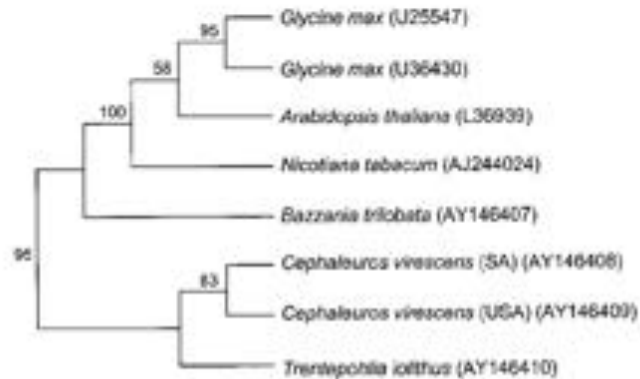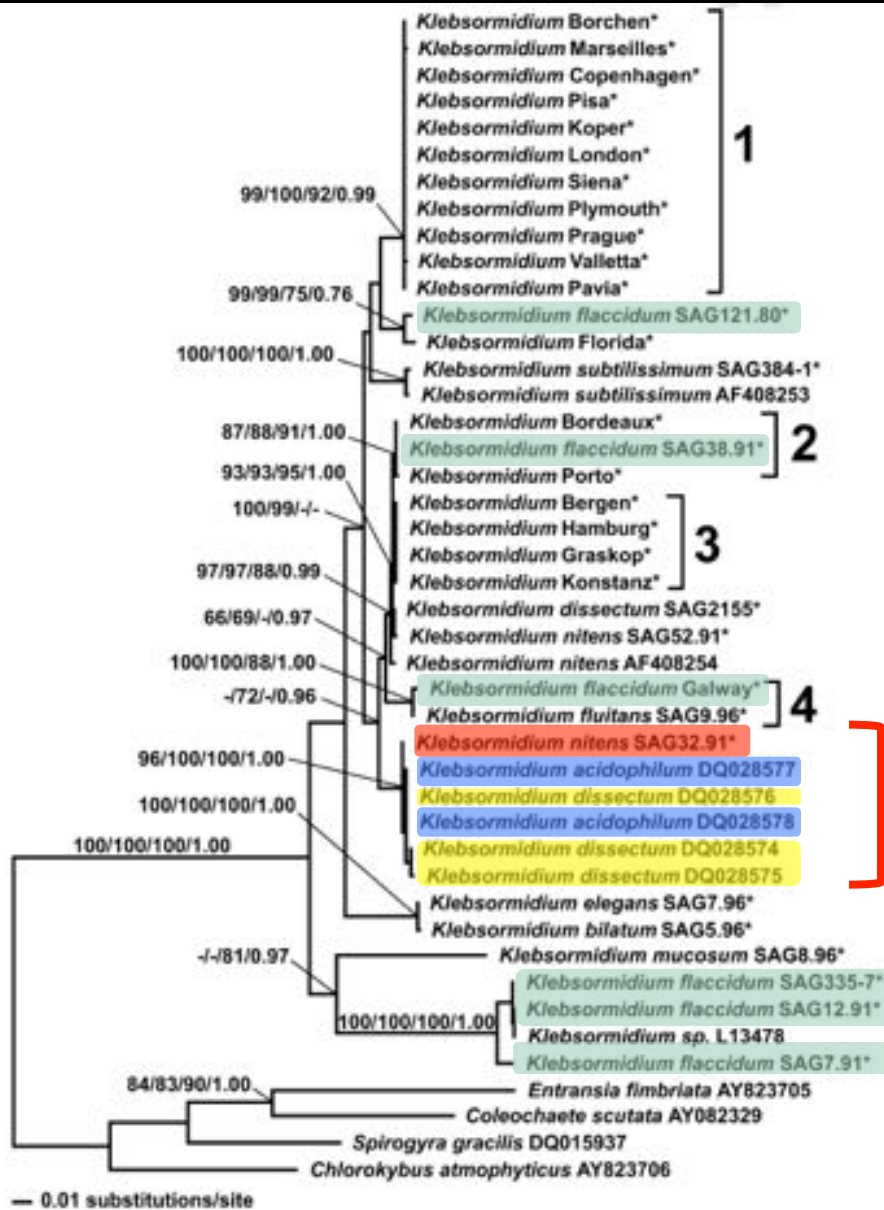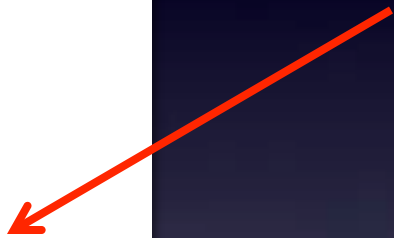


**Fig. 1.** Maximum-parsimony tree based on partial phragmo-plastin sequences, estimated from exhaustive search. A total of 584 positions are represented in the phylogram. GenBank accession numbers are given in parentheses. Tree length = 1090; CI = 0·868; RI = 0·636. Bootstrap values represent percentages derived from 10 000 replicates.

FIG. 2. Phylogram inferred from maximum-likelihood analysis of the *rbc*L gene in *Klebsormidium* and outgroup taxa, with bootstrap support (BP) and Bayesian posterior probabilities (PP) indicated at the nodes. From left to right, support values at nodes correspond to neighbor-joining BP, maximum-parsimony BP, maximum-likelihood BP, and Bayesian PP. New sequences produced in this study are marked with an asterisk. The four clades containing strains from urban environments are numbered as reported in the text.

Samples originally corresponding to *Klebsormidium flaccidum* are spread along the *rbc*L tree

On the other hand, samples corresponding to different morphological species are grouped together

This is a <u>phylogram</u>, an evolutionary tree where branch lengths are proportional to the amount of character change

# The application of Phylogeny

**C**an be used to define monophyletic groups of algal species with some confidence

**W**hen other data (biochemical, ultrastructural, etc) are mapped onto phylogenies, clade-specific attributes are often revealed!

**Phylogenetic information can form a robust basis for predicting the extent to which physiologies or ecological behavior can be extrapolated from taxa that have been studied to related form that have not**

<u>**Evolutionary utility**</u>

**-Biodiversity**

**-Evolution of life history and sexual reproduction**

<u>**-Ecological utility**</u>

**-Detection of microorganisms in natural environments**

**-Evaluation of gene expression of primary producers in horizontal, vertical or geographic gradients**

Reading: Assembling the Tree of Life

**The biological revolution caused by electron microscopy in the 1970's is repeated today by the Gene analysis!**



**What is next?**

# Maybe Phylogenomics?

Complete genome sequences is allowing to compare organisms in a new way to understand the origin and evolution of the genome

By combining the use of two disciplines, <u>Molecular Phylogenetics</u> (*the study of the evolutionary relationships among organisms and genes*) and <u>Genomics</u> (*the study of the organization and evolution of genes and genomes*)



The green algal ancestry of land plants as revealed by the chloroplast genome, by Turmel *et al., in press*