# Sequence manipulation. Retrieving sequences from GenBank

**Rafael Medina (rafael.medina.bry@gmail.com)**
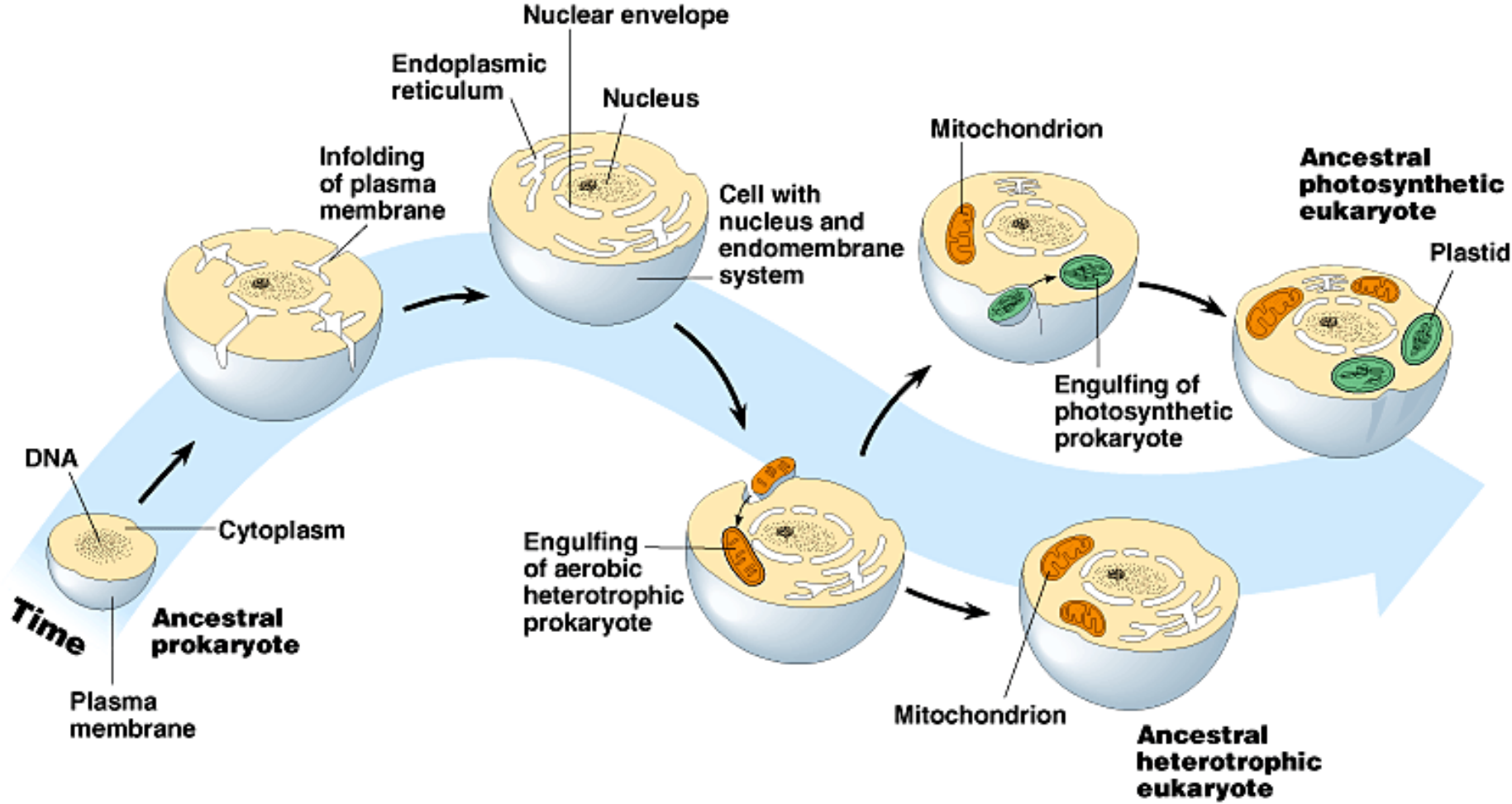**Yang Liu (yang.liu@uconn.edu)**

*atp*B-*rbc*L

# Molecular evolution and phylogeny of the *atpB–rbcL* spacer of chloroplast DNA in the true mosses

**Tzen-Yuh Chiang and Barbara A. Schaal**

Two universal primers, *rbcL-1* (5′-AACACCAGCTTTRAATC-CAA-3′) and *atpB-1* (5′-ACATCKARTACKGGACCAATAA-3′), were developed for amplifying and sequencing the *rbcL-atpB* spacers (Chiang et al. 1998) from the sequences of *Marchantia* (Umesono et al. 1988), tobacco (Shinozaki et al. 1986), and rice (Nishizawa and Hirai 1987). The PCR amplification protocol utilized two units of Taq polymerase (New England BioLab), the Taq buffer (500 mM KCl, 100 mM Tris–HCl, pH 9.0, and 1.0% Triton X–100), 2.5 mM $MgCl_2$, 10 pmol of each primer, and 8 mM dNTP in 100 µL reaction. PCR amplification was carried out in 30 cycles of 94°C denaturing for 45 s, 57°C annealing for 1 min 15 s, and 72°C extension for 1 min 15 s, followed by 72°C extension for 10 min and 4°C for storing. PCR products were polyacrylamide-
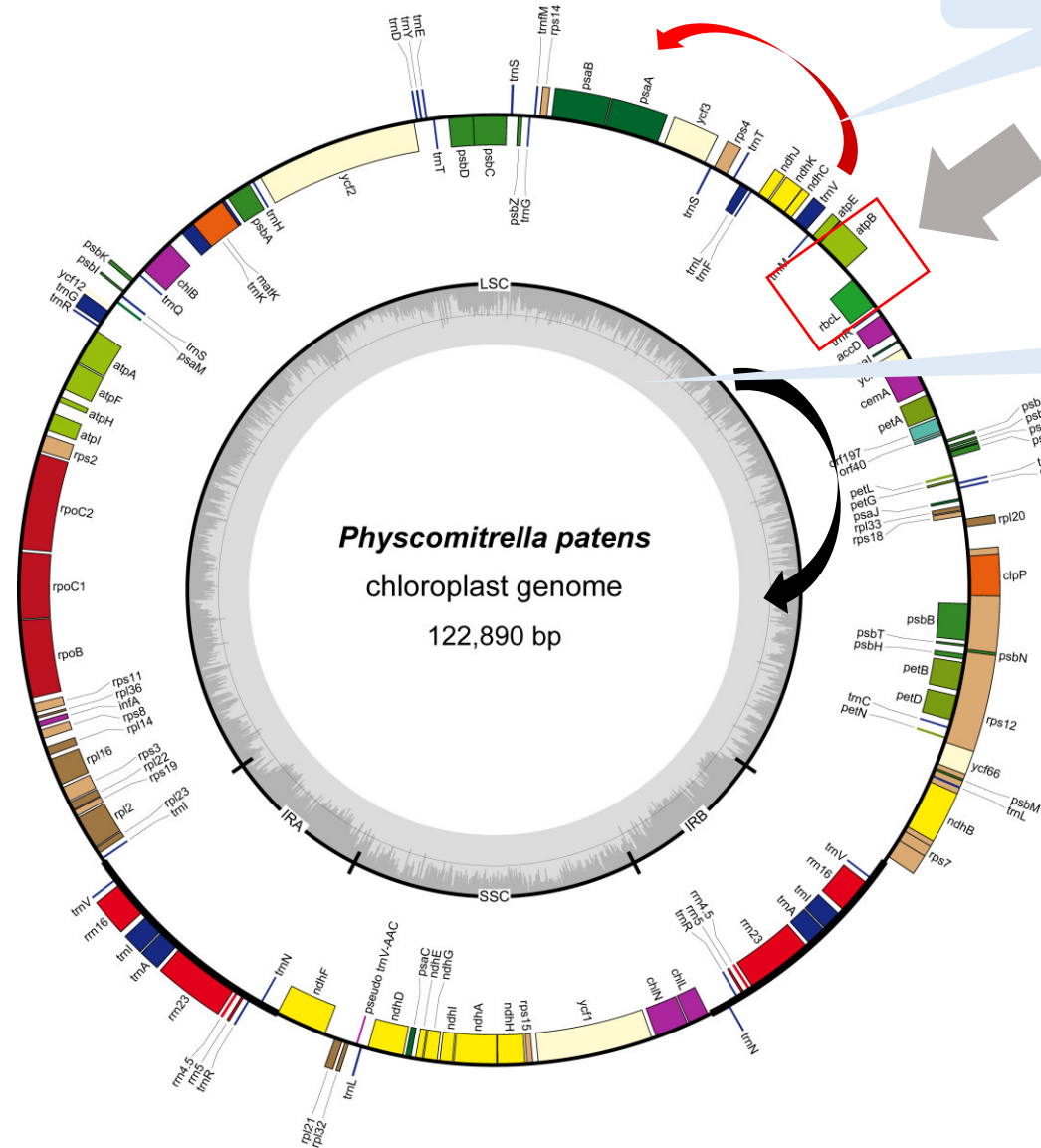
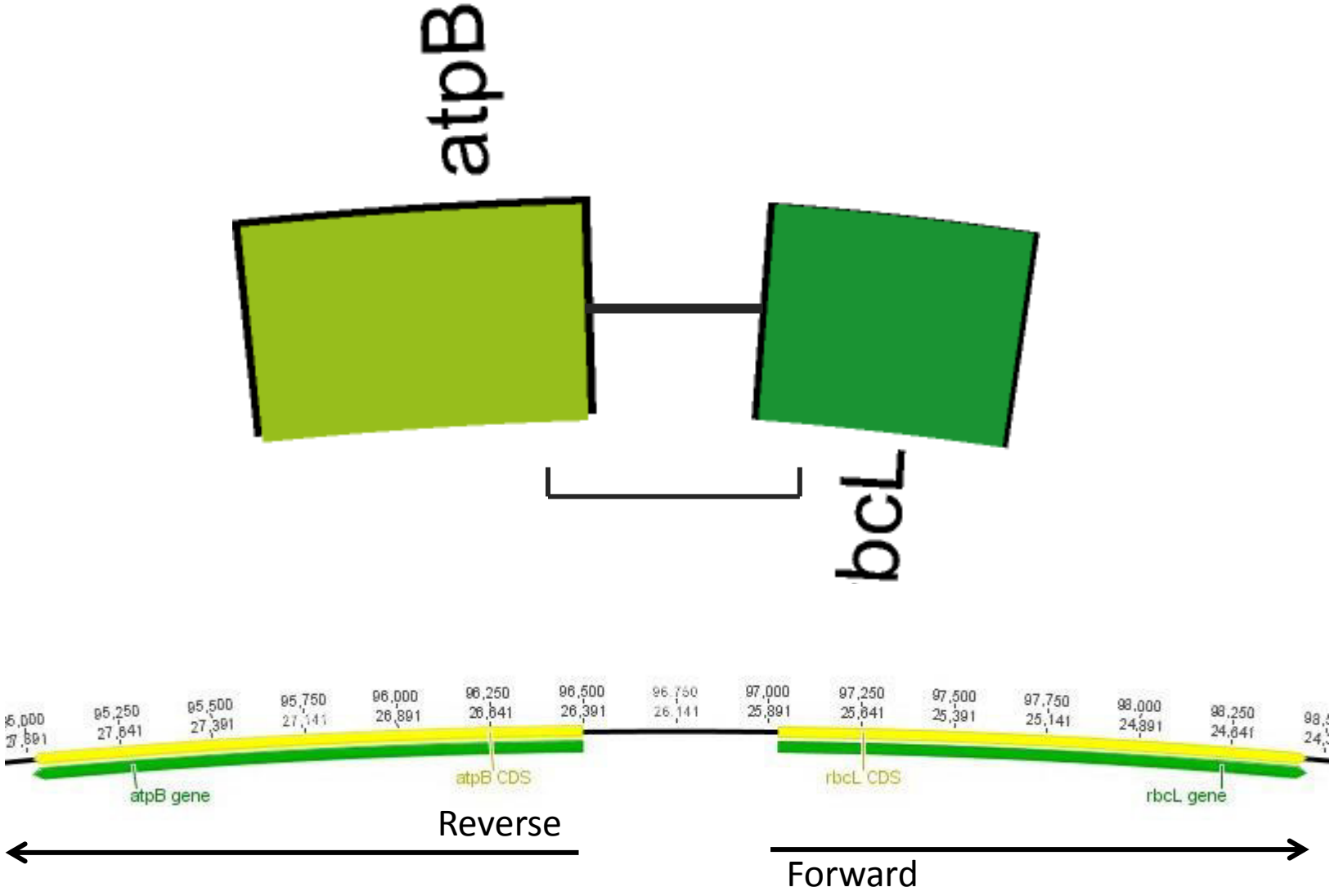# Origin of plant chloroplast genome

# *Physcomitrella* (moss) chloroplast genome



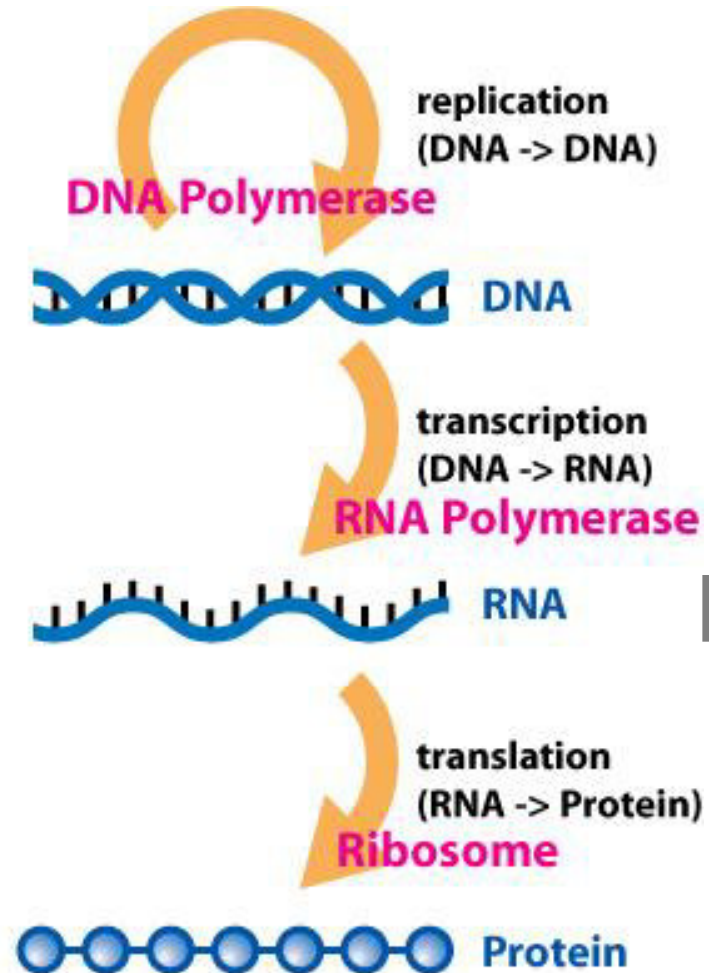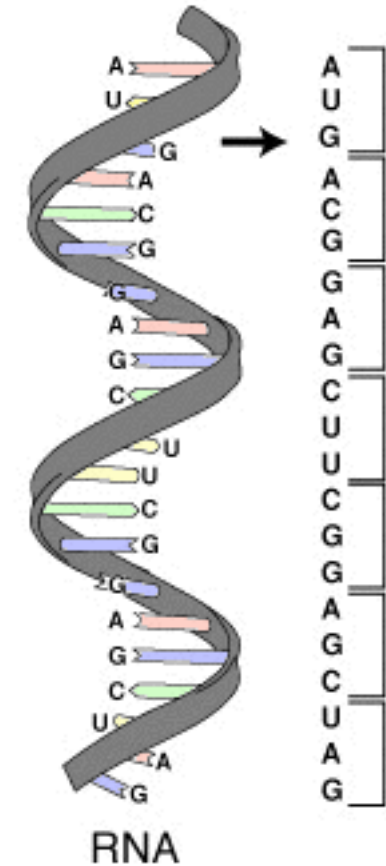Outer genes: reverse

Inner genes: direct

*Physcomitrella patens*
chloroplast genome
122,890 bp

Terasawa et al. *MBE* 2007

# *atp*B-*rbc*L

# Central dogma

Start codons: ATG
Stop codons: TAA, TGA, TAG



replication
(DNA -> DNA)
**DNA Polymerase**
DNA

transcription
(DNA -> RNA)
**RNA Polymerase**
RNA

translation
(RNA -> Protein)
**Ribosome**
Protein

Dhorspool

RNA

Ribonucleic acid

TransControl

# *Plagiothecium latebricola atp*B-*rbc*L

# Transcription and translation a gene with intron

# Organize sequences



➢ Use a text editor: WordPad, Notepad ++ (windows); TextEdit, TextWrangler (Mac)

➢ One file for one gene, always include GenBank **accession number** for a sequence

➢ Molecular programs: BioEdit v7.2, MEGA v5, PhyDe

# BioEdit v7.2

# BioEdit v7.2

# MEGA v5

# GenBank

➢ GenBank: an open access sequence database collecting nucleotide sequences and their protein translations

➢ GenBank is found in 1982 by the National Center for Biotechnology Information (NCBI), which belongs to the National Institutes of Health (NIH)

➢ By August 2014, GenBank has 174 million loci, 165 billion bases, and for more than 300,000 organisms

➢ Doubling every 18 months

# *Plagiothecium latebricola atp*B-*rbc*L

# **Retrieving sequences** from GenBank

➢ Search by gene and organism name

➢ Search by a query sequence
- Using BLAST (Basic Local Alignment Search Tool)

# Acquiring sequences from GenBank: **Search by name**



(http://www.ncbi.nlm.nih.gov/)

# Acquiring sequences from GenBank: **Search by name**

# Acquiring sequences from GenBank: **Search by name**

# Acquiring sequences from GenBank: **Search by name**

# **Retrieving** sequences from GenBank

➢ Search by gene and organism name

➢ Search by a query sequence
  • Using BLAST (Basic Local Alignment Search Tool)

# Acquiring sequences from GenBank: **BLAST**



http://blast.ncbi.nlm.nih.gov/

# Acquiring sequences from GenBank: **BLAST**



- **blastn**
  Search nucleotide database by nucleotide query
- **blastp**
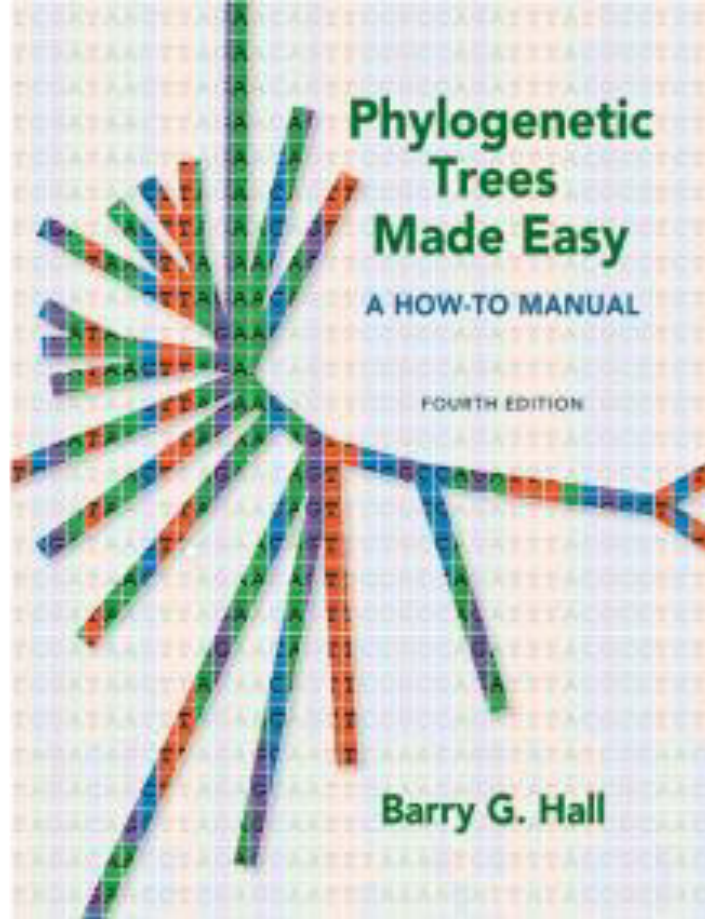  Search protein database by protein query
- **blastx**
  Search protein database by translated nucleotide query
- **tblastn**
  Search translated nucleotide database by protein query
- **tblastx**
  Search translated nucleotide database by translated nucleotide query

Barry G. Hall
**Phylogenetic Trees Made Easy: A How-To Manual**
Sinauer Associates, Inc.; Fourth edition (April 30, 2011)

**Table of Contents:**