

Joint Bayesian Estimation of Alignment and Phylogeny

Benjamin D. Redelings and Marc A. Suchard

Outline

- Introduction
- Existing Methods
- Bayesian Coestimation and Inference
- Results

Introduction

- There are problems with existing alignment and phylogeny methods that can produce inaccurate results
- Introduce a new method based on Bayesian inference to simultaneously estimate the alignment and phylogeny
- Results show that trees remain consistent when more taxa are added

Existing Methods

- Two stages:
 - Sequence alignment
 - Phylogeny reconstruction

Existing Methods

- Works well when alignments are well resolved
- Spurious or no results when alignments are uncertain
- Exaggerated support for phylogenies
- Biased towards the alignment guide tree

Existing Methods

- Solutions:
 - Remove ambiguous regions from alignments
 - Code ambiguity information into the alignment
 - Optimization alignment: coestimation of alignments and phylogenies in a parsimony framework

Bayesian Coestimation

- Joint estimation of alignments and phylogenies
- Weighted by posterior probabilities
- No external guide tree—eliminates bias
- More accurate substitution and indel models

Bayesian Coestimation

- Markov chain Monte Carlo (MCMC) sampling
- Simplifying assumption: indels occur independently of branch length
 - Can be modeled using a symmetrical pair-hidden Markov model (pair-HMM)

$$P(\mathbf{H} \mid \mathbf{E}) = \frac{P(\mathbf{E} \mid \mathbf{H}) \cdot P(\mathbf{H})}{P(\mathbf{E})}$$

Bayes' Theorem

\mathbf{Y} = set of sequences

\mathbf{A} = multiple alignment

τ = unrooted tree topology

\mathbf{T} = branch lengths

Θ = substitution parameters

Λ = indel parameters

Variables

$$P(\mathbf{Y}, \tau, \mathbf{T}, \Theta \mid \mathbf{A}) = P(\mathbf{Y} \mid \tau, \mathbf{T}, \Theta, \mathbf{A}) \\ \times P(\tau, \mathbf{T}) \times P(\Theta)$$

Traditional Probabilistic Model

$$P(\mathbf{Y}, \mathbf{A}, \tau, \mathbf{T}, \Theta, \Lambda) = P(\mathbf{Y} \mid \mathbf{A}, \tau, \mathbf{T}, \Theta) \times P(\mathbf{A} \mid \tau, \Lambda) \\ \times P(\tau, \mathbf{T}) \times P(\Theta) \times P(\Lambda)$$

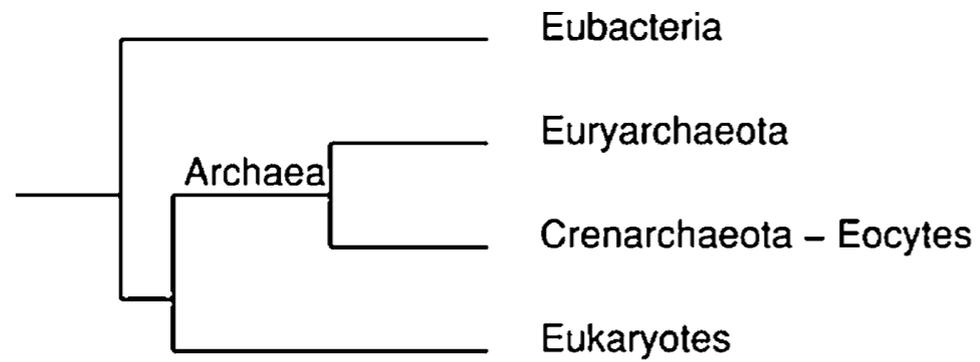
Coestimation Probabilistic Model

Posterior Sampling

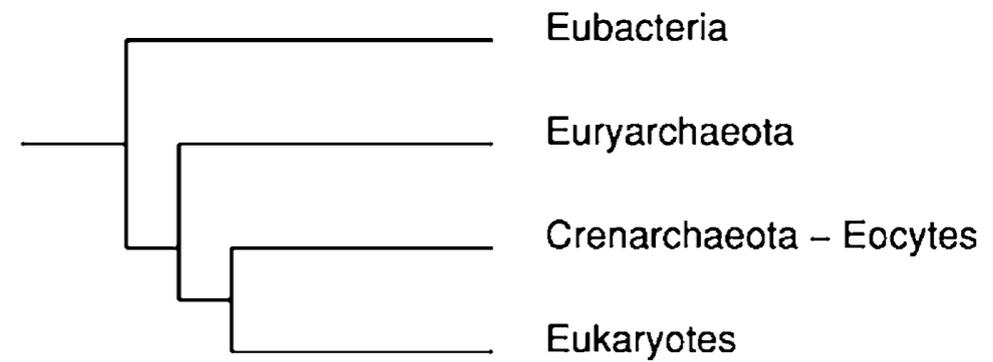
- Sample from the model posterior $P(A, \tau, T, \Theta, \Lambda | Y)$ using Markov Chain Monte Carlo (MCMC)
- Update the topology and each branch length, pairwise branch alignment, and internal node sequence length and alignment at least once each iteration
- Topology is updated with nearest-neighbor interchange
- Parameters sampled a Poisson number of times, with mean $1 + B/3$, where B is the number of branches

Results

- Early Branching in the Tree of Life: Resolving the Archaea
- Archaeal hypothesis: Archaea form a monophyletic group
- Eocyte hypothesis: Archaea are paraphyletic and that the eocyte Archaea are more closely related to Eukaryotes than to other Archaea



(a)



(b)

Archaeal vs Eocyte Hypothesis

(a): Archaeal

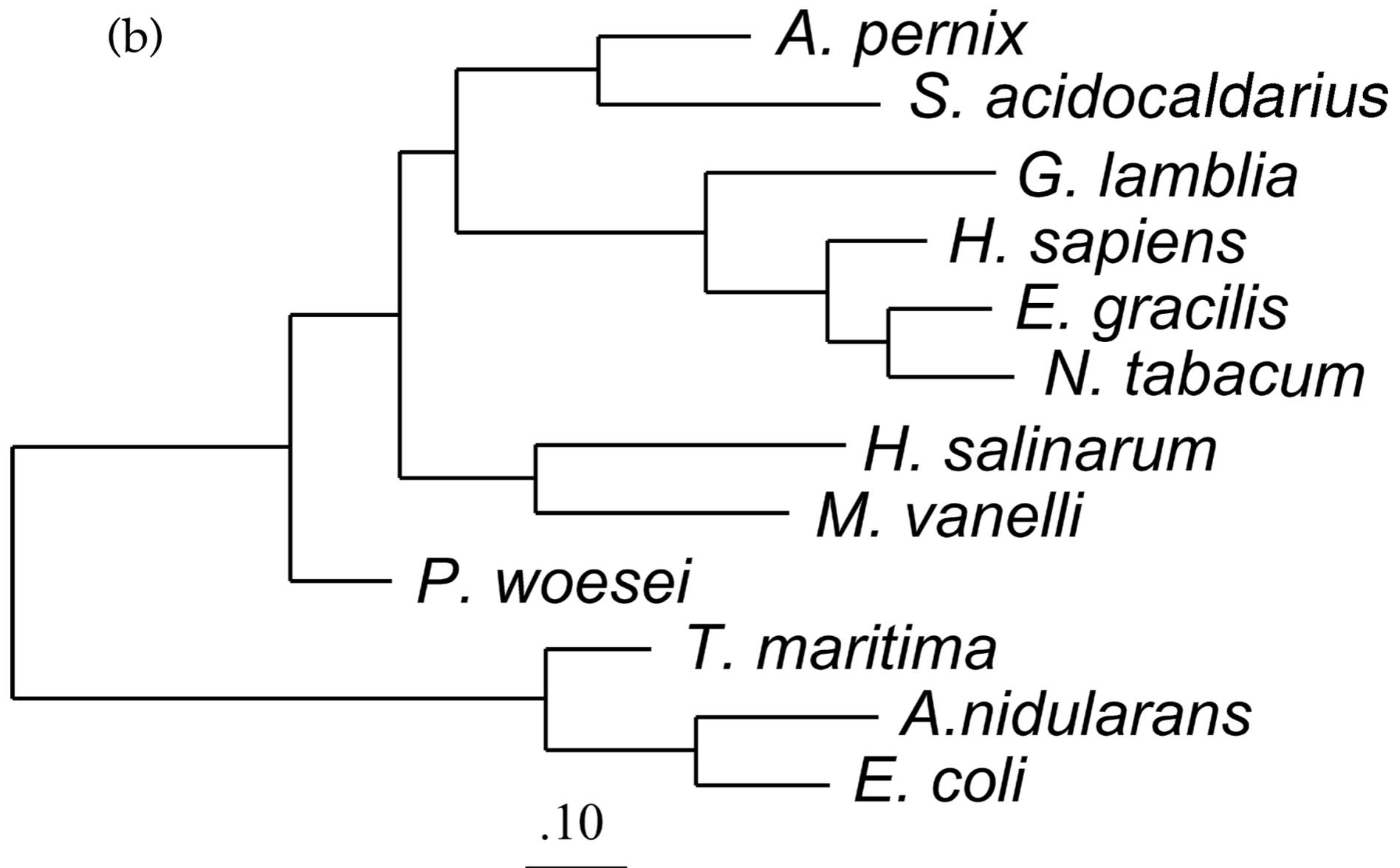
(b): Eocyte

Species Examined

Taxa	Domain	Order	Note
<i>Homo sapiens</i>	Eukaryotes	Metazoa	Human beings
<i>Nicotiana tabacum</i>	Eukaryotes	Plantae	Tobacco plant
<i>Euglena gracilis</i>	Eukaryotes	Protista	Photosynthetic single cell
<i>Giardia lamblia</i>	Eukaryotes	Diplomonadida	Intestinal protist
<i>Sulfolobus acidocaldarius</i>	Archaea	Crenarchaeota	High pH thermophile
<i>Aeropyrum pernix</i>	Archaea	Crenarchaeota	Anaerobic thermophile
<i>Pyrococcus woesei</i>	Archaea	Euryarchaeota	Thermophile
<i>Halobacterium salinarum</i>	Archaea	Euryarchaeota	Methanogen
<i>Methanococcus vannelli</i>	Bacteria	Aquificae	Thermophile
<i>Thermotoga maritima</i>	Bacteria	Thermotogales	Thermophile
<i>Anacystis nidularans</i>	Bacteria	Cyanobacteria	Photosynthetic, aquatic
<i>Escherichia coli</i>	Bacteria	Proteobacteria	Intestinal symbiont

Results

- 5– and 12–taxon datasets
- EF-1 α /Tu protein: at least 26% conserved
- Near threshold of 20~25% where homology becomes difficult to detect



12-Taxon Tree