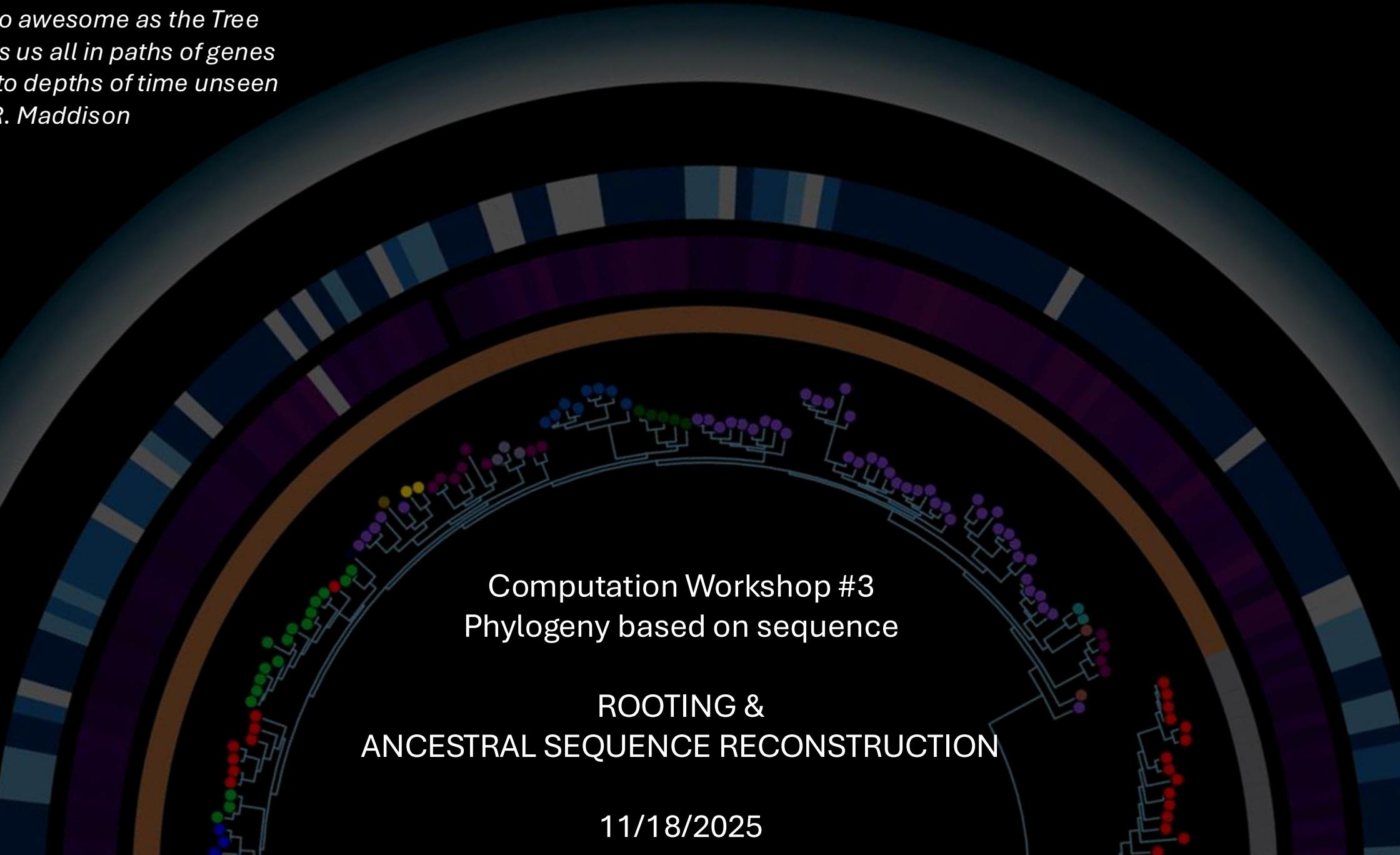


*I think that I shall never see
A thing so awesome as the Tree
That links us all in paths of genes
Down into depths of time unseen*
- David R. Maddison

A circular phylogenetic tree diagram is centered in the background. The tree branches out from a central point at the bottom, with various tips colored in shades of red, green, blue, purple, and yellow. The background of the slide features concentric circles in dark blue, purple, and orange.

Computation Workshop #3
Phylogeny based on sequence

ROOTING &
ANCESTRAL SEQUENCE RECONSTRUCTION

ASR Workflow

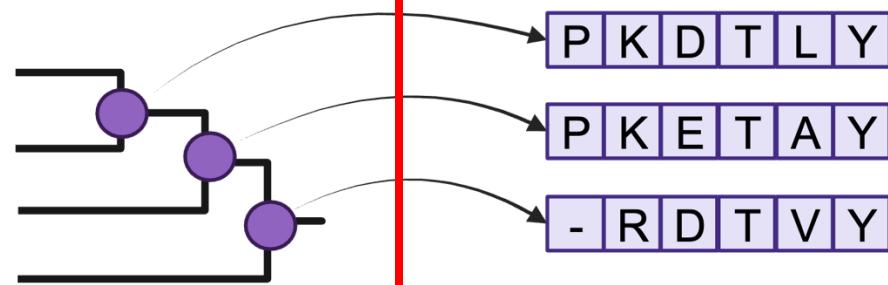
Sequence Curation

P	K	D	T	L	Y
P	R	D	T	A	Y
	K	D	T	V	Y
P	K	E	T	-	Y

Multiple Sequence
Alignment

P	K	D	T	L	Y
P	K	D	T	A	Y
P	K	E	T	-	Y
-	R	E	T	V	Y

Phylogenetic
Reconstruction



Ancestral Sequence
Inference

ASR Workflow

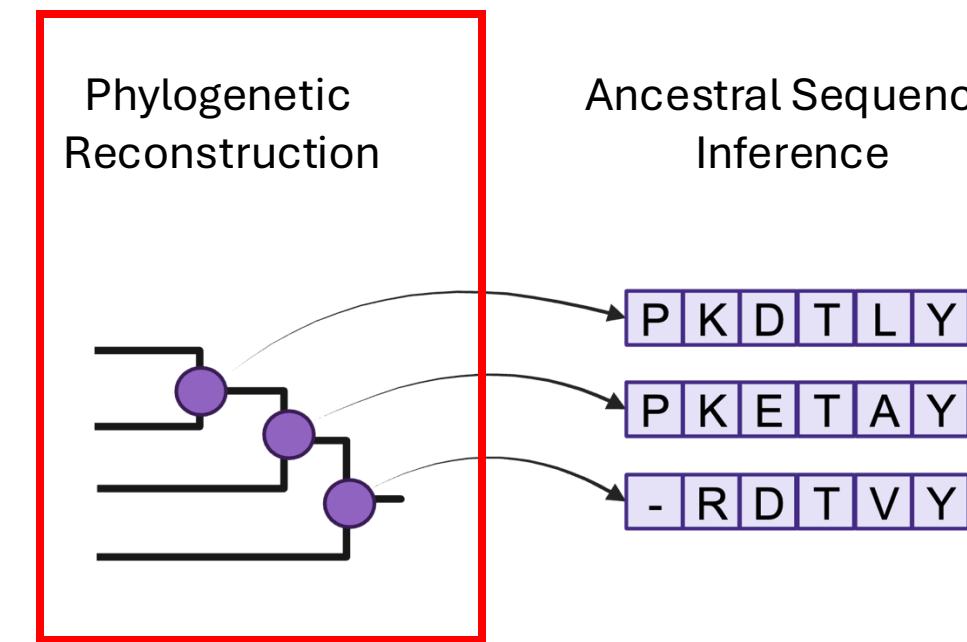
Sequence Curation

P	K	D	T	L	Y
P	R	D	T	A	Y
	K	D	T	V	Y
P	K	E	T	-	Y

Multiple Sequence
Alignment

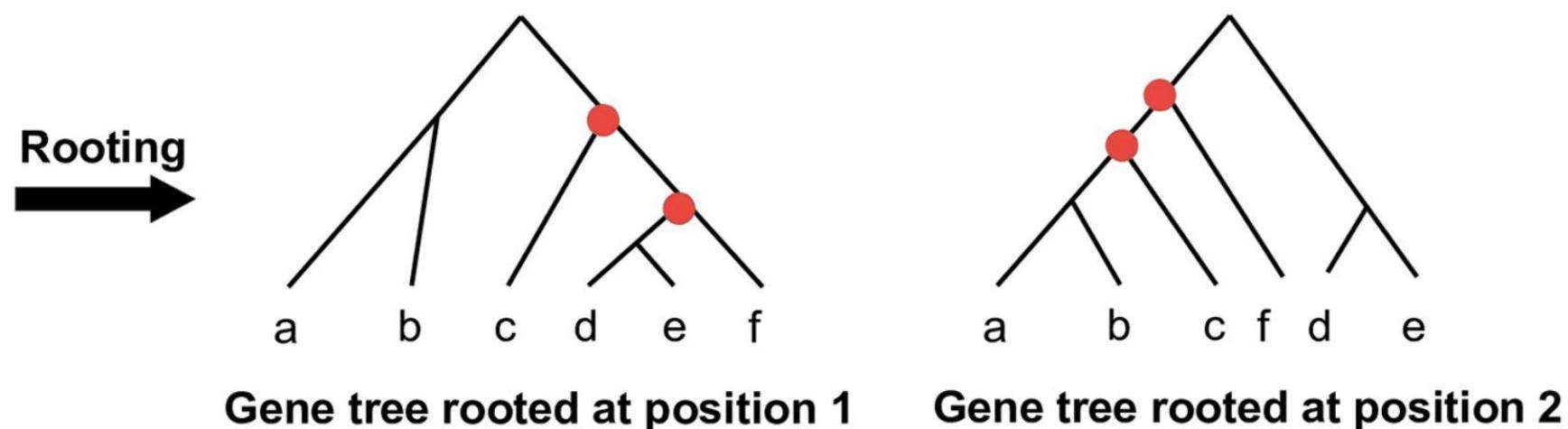
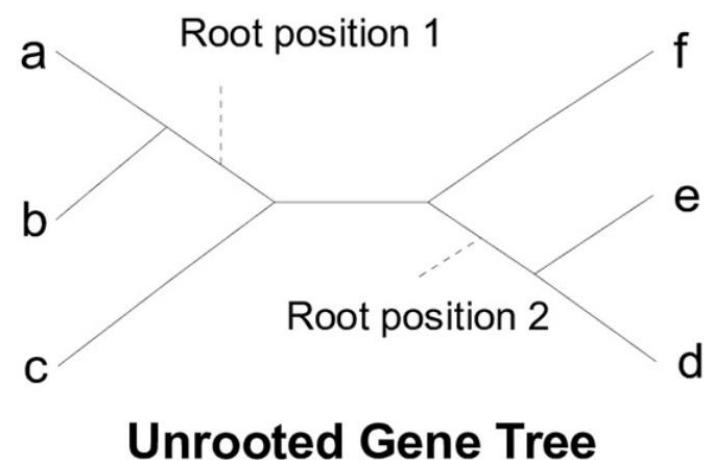
P	K	D	T	L	Y
P	K	D	T	A	Y
P	K	E	T	-	Y
-	R	E	T	V	Y

Phylogenetic
Reconstruction



Ancestral Sequence
Inference

Unrooted vs Rooted Tree



No ancestor-descendant relationship

Ancestral relationships are determined

How to decide where to place the root? 🧐

Phylogenetic Tree Rooting Methods



Species Tree Rooting

Outgroup rooting

Phylogenetic reconciliation

Based on Branch Lengths

Midpoint rooting

Minimal variance

Minimal Ancestral Deviation



Using Molecular Clock

Molecular clock rooting

Bayesian molecular clock
rooting

Non-reversible evolutionary
models

Relaxed clock models

Phylogenetic Reconciliation for Gene Trees

Duplication-Transfer-Loss
(DTL) reconciliation

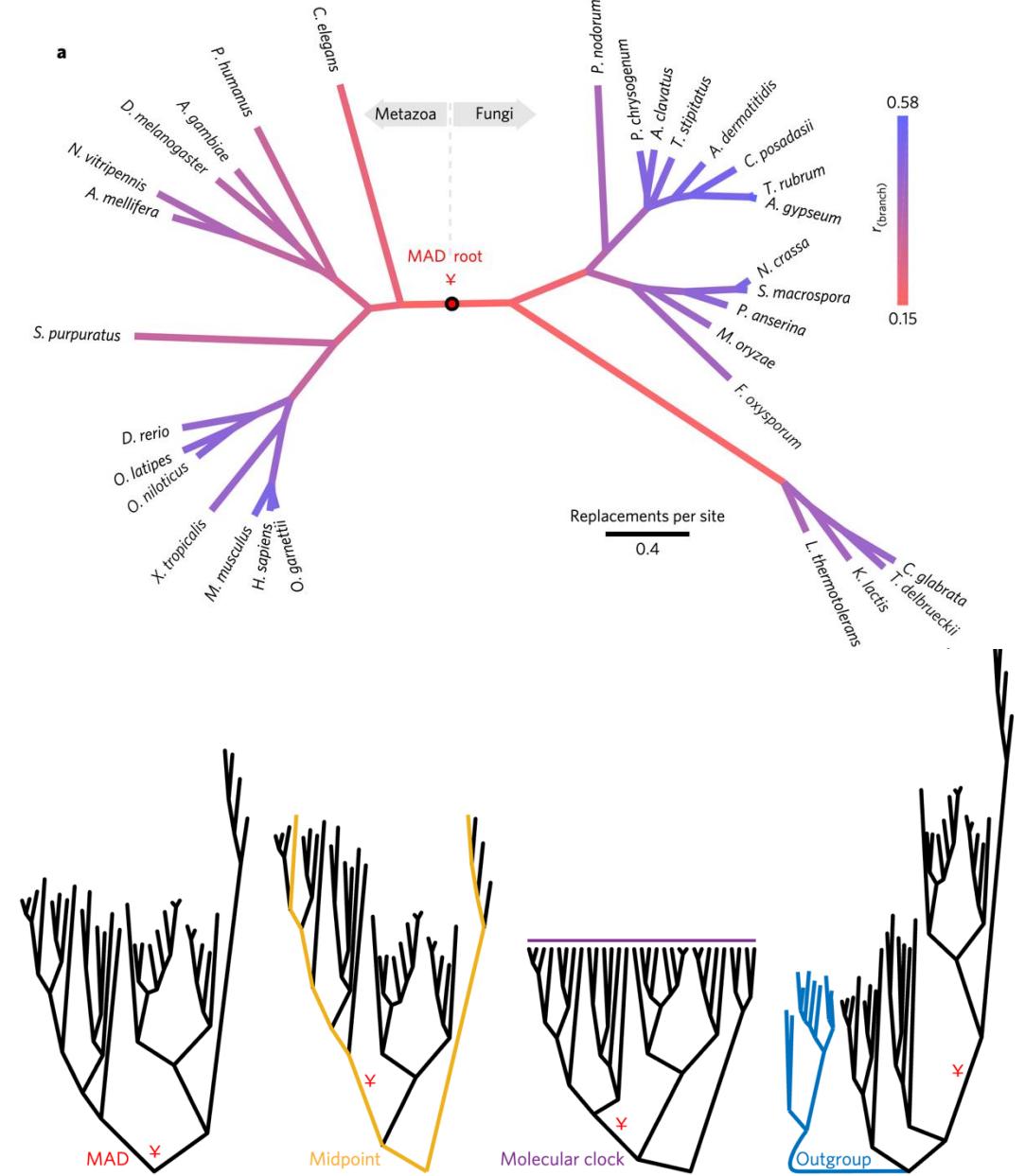
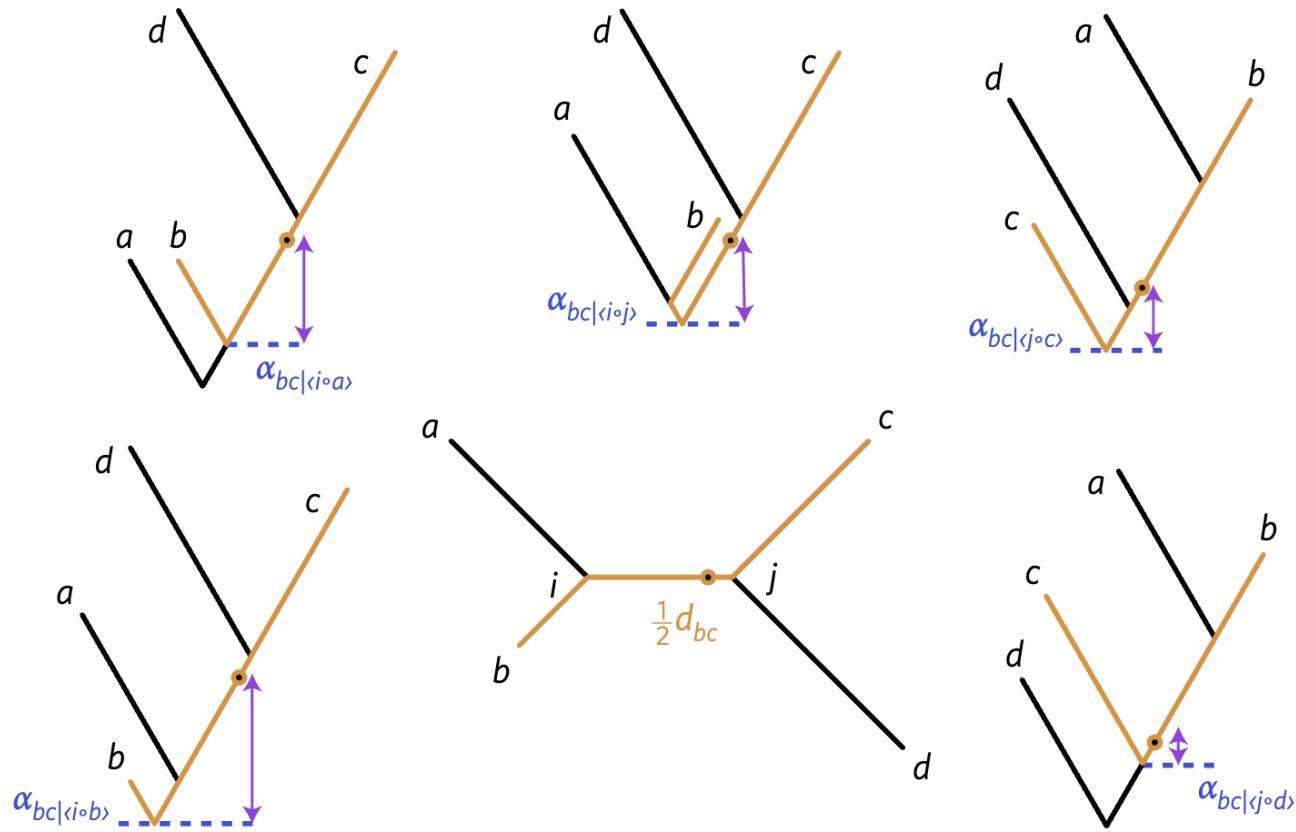
ALE rooting



Rooting based on branch lengths

- **Midpoint:** unrooted tree is rooted at the mid-point of the longest path in the tree (Farris, *American Naturalist*, 1972)
- **Minimum Variance:** improve upon midpoint rooting, a root that minimizes root-to-leaf distance variance (Mai et al., *PLOS One*, 2017)
- **Minimal Ancestor Deviation (MAD):** deviations of the midpoint criterion for all possible root positions (Tria et al., *Nat Eco Evo*, 2017)

How does MAD work?



Which one to use?

RESEARCH ARTICLE

Assessing the accuracy of phylogenetic rooting methods on prokaryotic gene families

Taylor Wade¹, L. Thiberio Rangel², Soumya Kundu^{1✉}, Gregory P. Fournier², Mukul S. Bansal^{1,3*}

1 Department of Computer Science & Engineering, University of Connecticut, Storrs, CT, United States of America, **2** Department of Earth, Atmospheric & Planetary Sciences, Massachusetts Institute of Technology, Cambridge, MA, United States of America, **3** Institute for Systems Genomics, University of Connecticut, Storrs, CT, United States of America

Abstract

Almost all standard phylogenetic methods for reconstructing gene trees result in *unrooted* trees; yet, many of the most useful applications of gene trees require that the gene trees be correctly rooted. As a result, several computational methods have been developed for inferring the root of unrooted gene trees. However, the accuracy of such methods has never been systematically evaluated on prokaryotic gene families, where horizontal gene transfer is often one of the dominant evolutionary events driving gene family evolution. In this work, we address this gap by conducting a thorough comparative evaluation of five different rooting methods using large collections of both simulated and empirical prokaryotic gene trees. Our simulation study is based on 6000 true and reconstructed gene trees on 100 species and characterizes the rooting accuracy of the four methods under 36 different evolutionary conditions and 3 levels of gene tree reconstruction error. The empirical study is based on a large, carefully designed data set of 3098 gene trees from 504 bacterial species (406 Alphaproteobacteria and 98 Cyanobacteria) and reveals insights that supplement those gleaned from the simulation study. Overall, this work provides several valuable insights into the accuracy of the considered methods that will help inform the choice of rooting methods to use when studying microbial gene family evolution. Among other findings, this study identifies parsimonious Duplication-Transfer-Loss (DTL) rooting and Minimal Ancestor Deviation (MAD) rooting as two of the most accurate gene tree rooting methods for prokaryotes and specifies the evolutionary conditions under which these methods are most accurate, demonstrates that DTL rooting is highly sensitive to high evolutionary rates and gene tree error, and that rooting methods based on branch-lengths are generally robust to gene tree reconstruction error.

Tutorial time

- Running MAFFT



ASR Workflow

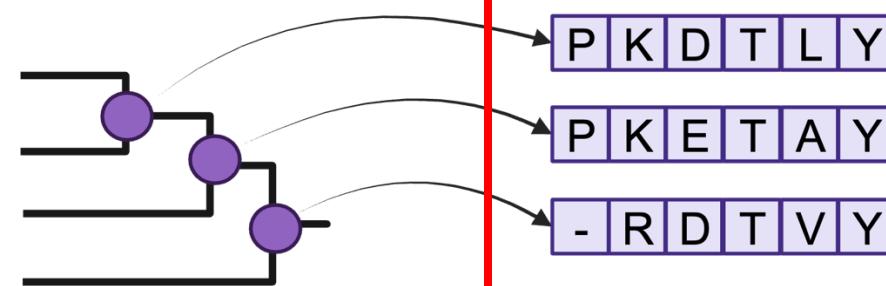
Sequence Curation

P	K	D	T	L	Y
P	R	D	T	A	Y
	K	D	T	V	Y
P	K	E	T	-	Y

Multiple Sequence
Alignment

P	K	D	T	L	Y
P	K	D	T	A	Y
P	K	E	T	-	Y
-	R	E	T	V	Y

Phylogenetic
Reconstruction



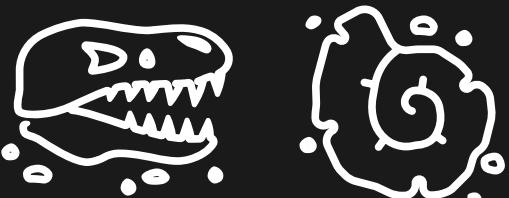
Ancestral Sequence
Inference

P	K	D	T	L	Y
P	K	E	T	A	Y
-	R	D	T	V	Y

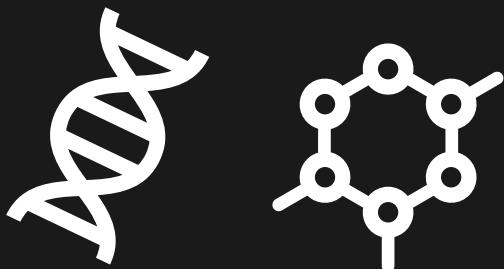
First time proposing the idea of reconstructing the ancestral genes/proteins

Two records of life:

1) Fossil Records



2) Molecular Records



Chemical Paleogenetics
Molecular "Restoration Studies" of Extinct Forms of Life

LINUS PAULING and EMILE ZUCKERKANDL*

*Division of Chemistry and Chemical Engineering, California Institute
of Technology, Pasadena, California, USA***

Attention is attracted to the possibility of reconstructing the amino-acid sequence of ancestral polypeptide chains by virtue of a comparison between the amino-acid sequences of related polypeptide chains found in contemporary organisms. A tentative partial structure is proposed for two ancestral hemoglobin polypeptide chains. Some perspectives of paleobiochemistry are outlined.

What is Ancestral Sequence Reconstruction (ASR)?

ASR is a computational method that reconstructs ancestral sequences based on the molecular data of present-day (extant) organisms using a model of molecular evolution

Alignment of
extant molecular
sequences



A model with rates
of different
substitutions

Phylogenetic
tree

What ASR is **NOT**:



A time machine

- Represent hypotheses
- Uncertainty coming from statistical inferences
- May not be the “true” ancestor



A molecular clock

- ASR does not estimate divergence time
- But related tools can be used



A machine learning

- ASR uses classical statistical models
- No training data, no pattern learning
- ML = maximum likelihood (not machine learning)

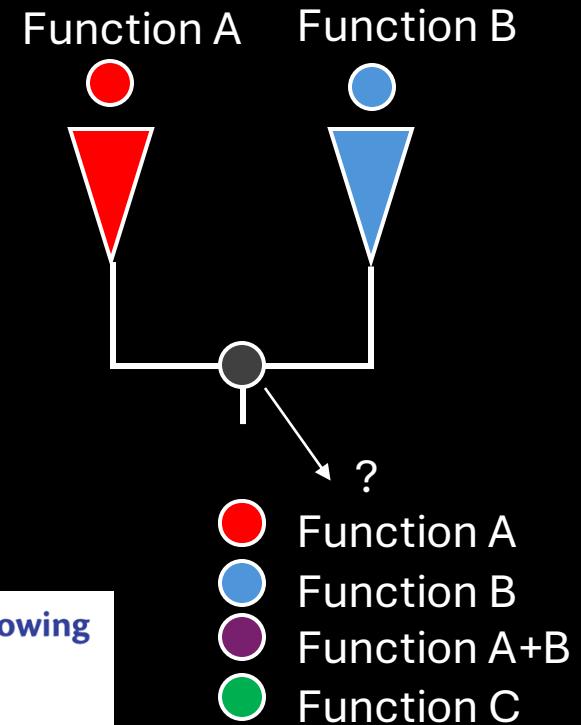


Why do we reconstruct ancestral sequences?

- What were the functions of ancestral proteins?
- Which property/function evolved first?
- Was the ancestral protein a specialist/generalist?
- Which sequence and structural changes gave rise to a new function?
- Can activities of ancestral proteins be aligned with environmental conditions?

We can get some insights into:

- Mechanisms of protein evolution
- Sequence and structural foundations of functional adaptation
- Environmental conditions in the past



Reconstruction of Nitrogenase Predecessors Suggests Origin from Maturase-Like Proteins

Amanda K. Garcia ,¹ Bryan Kolaczkowski ,² and Betül Kaçar^{1,*}

Emergence of an Orphan Nitrogenase Protein Following Atmospheric Oxygenation

Bruno Cuevas-Zuviría,^{1,2,†} Amanda K. Garcia ,^{1,†} Alex J. Rivier,¹ Holly R. Rucker,¹ Brooke M. Carruthers,¹ and Betül Kaçar ,^{1,*}

Kinetic Analysis Suggests Evolution of Ribosome Specificity in Modern Elongation Factor-Tus from “Generalist” Ancestors

Arindam De Tarafder ,¹ Narayan Prasad Parajuli ,¹ Soneya Majumdar ,¹ Betül Kaçar ,^{2,3} and Suparna Sanyal ,^{1,*}

Ancient nitrogenases are ATP dependent

Derek F. Harris ,¹ Holly R. Rucker ,² Amanda K. Garcia ,² Zhi-Yong Yang ,¹ Scott D. Chang ,² Hannah Feinsilber ,¹ Betül Kaçar ,² Lance C. Seefeldt ,¹

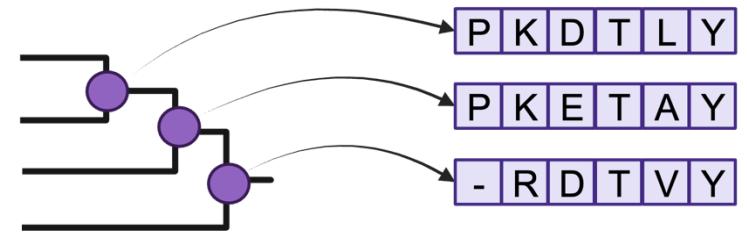
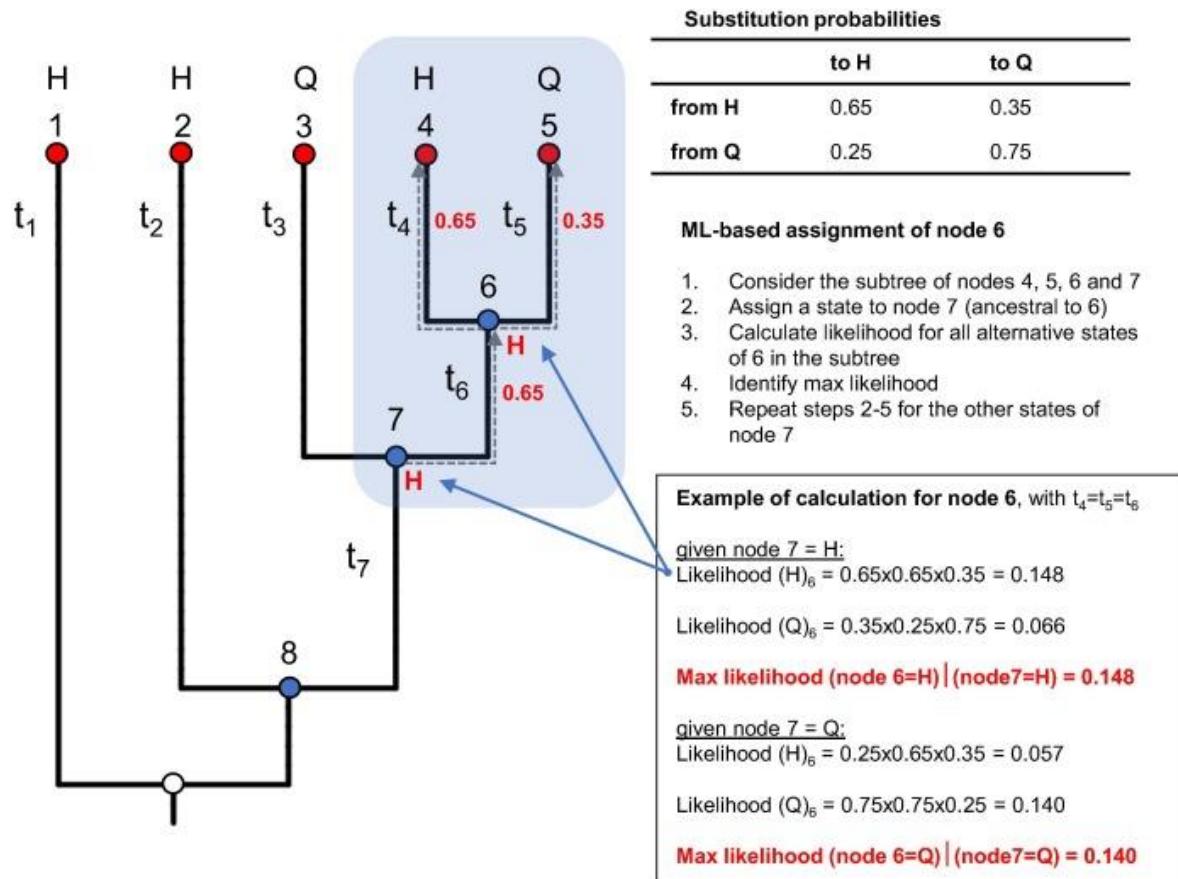
Earliest Photic Zone Niches Probed by Ancestral Microbial Rhodopsins

Cathryn D. Sephus,^{1,2,3} Eryim Fer,^{1,2,3} Amanda K. Garcia ,^{1,2} Zachary R. Adam,^{4,5} Edward W. Schwieterman ,^{5,6} and Betül Kaçar ,^{1,2,*}

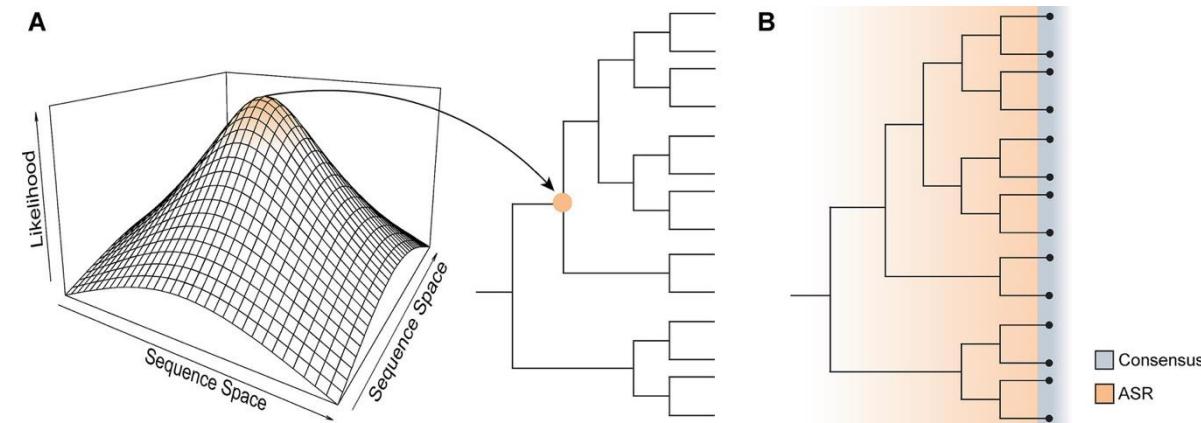
Evolutionary Dynamics of RuBisCO: Emergence of the Small Subunit and its Impact Through Time

Kaustubh Amritkar, Bruno Cuevas-Zuviría, Betül Kaçar Author Notes

4. Ancestral Sequence Inference



ASR is different than the consensus sequence!



(b)

Example ancestral posterior probability matrix

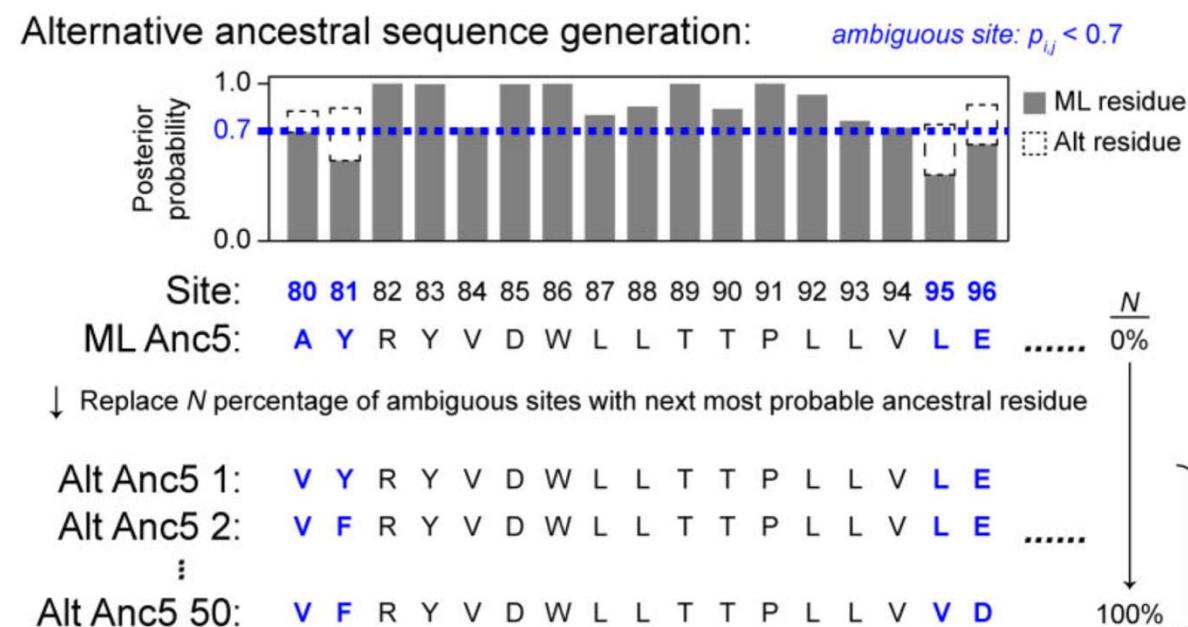
Node	Site	-	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
Anc5	1	$p_{1,1}$	$p_{1,2}$	$p_{1,3}$	$p_{1,4}$	$p_{1,5}$	$p_{1,6}$	$p_{1,7}$	$p_{1,8}$	$p_{1,9}$	$p_{1,10}$										$p_{1,21}$	
:	:	:																				
Anc5	X	$p_{X,1}$	$p_{X,2}$	$p_{X,3}$	$p_{X,4}$	$p_{X,5}$	$p_{X,6}$	$p_{X,7}$	$p_{X,8}$	$p_{X,9}$	$p_{X,10}$										$p_{X,21}$	

p_{ij} = posterior probability of ancestral residue, j , at site, i

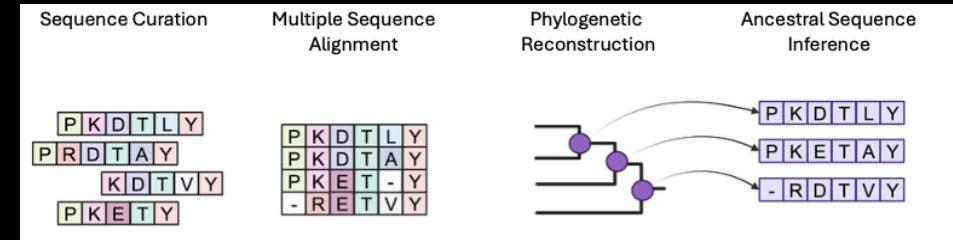
Robustness of Reconstructed Ancestral Protein Functions to Statistical Uncertainty

Geeta N. Eick,^{1,2} Jamie T. Bridgman,¹ Douglas P. Anderson,^{1,3} Michael J. Harms,^{1,3} and Joseph W. Thornton*,⁴

certainty about their primary sequence. Most studies to date that have addressed this question have done so by generating variants of the ML ancestral sequence, each of which contains a plausible alternate amino acid at one of the ambiguously reconstructed sites (typically defined as amino acids with a posterior probability above some arbitrary but reasonable cutoff, such as 0.2). The experimental characterization is



Dealing with uncertainty



- Reminder: ASR is not a time machine! Most likely ancestral sequences may not be true ancestral sequence
- Uncertainties come from each step of ASR workflow: sequence dataset, alignment, evolutionary model, tree, etc...
- Carefully articulate scientific questions, conclusions given the uncertainty in ASR

Are your conclusions sensitive to ASR parameters or robust to the uncertainties?

The robustness can be tested at each step by using different methods and compare the results.

E.g., Maybe alignment generates significant uncertainty:

But we can test different alignment methods and compare the ancestral sequences resulted from each
If they all yield proteins with similar experimental properties, we can be confident in our conclusions

E.g., Maybe evolutionary model generates uncertainty:

We can test second best evolutionary model and, again, compare the ancestral sequences resulted from each
If both result ancestral proteins with properties, we can be confident in our conclusions

We have the ancestral sequences...

What is next?

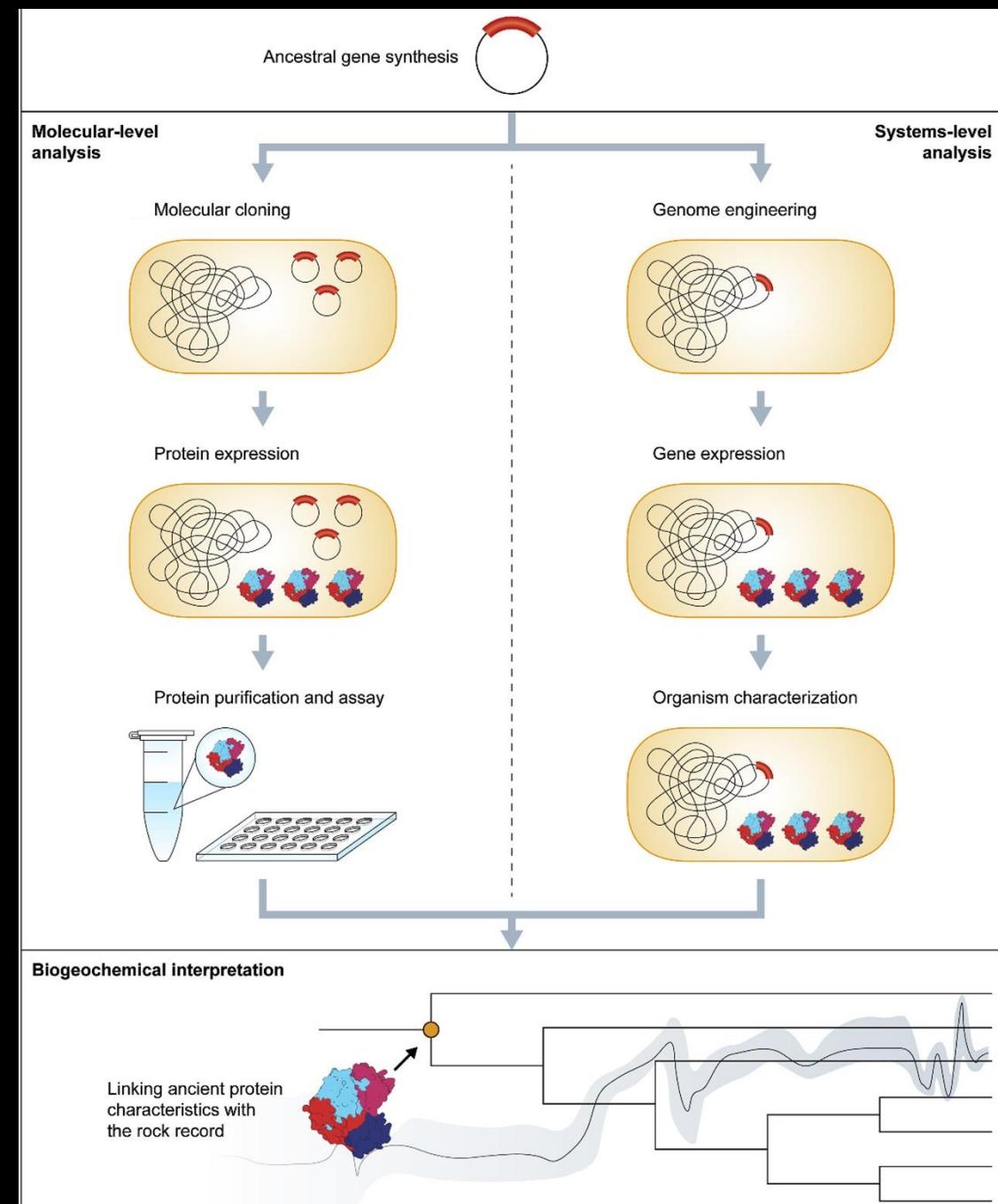
Important considerations moving toward next steps:

- Which model organism to use to capture property of interest?
- Codon optimization to improve translation efficiency in expression system
- Which genes to replace? Location of genetic insertion?
- *In vivo* or *in vitro* characterization?

Impact of cellular interactions and environments...

Ancient proteins presumably adapted toward environments, interactions partners, and substrates.

How to account for mismatches between ancient proteins and modern environments?





Any question?

Acknowledgement:

- Some slides are modified from Dr. Amanda Garcia's ASR Workshop presentation

Tutorial time

- ASR

