**Tree of Life** 

# EVE 161: Microbial Phylogenomics

# Class #6: The Tree of Life, Example Paper

UC Davis, Winter 2018 Instructor: Jonathan Eisen Teaching Assistant: Cassie Ettinger

### Where we are going and where we have been

**Tree of Life** 

- <u>Previous Class</u>:
  5. Phylogeny
- <u>Current Class:</u>
  - 6. Phylogeny example
- <u>Next Class:</u>
  - rRNA Sequencing from uncultured

# **Eisen Office Hours**

#### Tree of Life

- Today 10:30-11:30 Storer 5331
- Monday 10:00-11:00 GBSF 5311

- Task: Summarize one figure from one of the papers for class, ~ 5 minutes.
  - There will be 2-3 presenters per class.
  - Each person needs to do a different figure so you need to contact the other person to coordinate.
  - Can meet with Eisen, Ettinger before hand to discuss.
- Sign up for date on Google Sheet <u>goo.gl/bmxCDw</u> include your name & email address.

 Select 1-2 papers on one of the topics of the course (approval needed)

- Review the paper and write up a summary of your assessment of the paper (more detail on this later)
- Present a short summary of what you did to the class
- Ask and answer questions about your and other people's reviews in the discussion forum on Canvas

# Palfrey et al.

Syst. Biol. 59(5):518–533, 2010 © The Author(s) 2010. Published by Oxford University Press, on behalf of the Society of Systematic Biologists. All rights reserved. For Permissions, please email: journals.permissions@oxfordjournals.org DOI:10.1093/sysbio/syq037 Advance Access publication on July 23, 2010

#### Broadly Sampled Multigene Analyses Yield a Well-Resolved Eukaryotic Tree of Life

LAURA WEGENER PARFREY<sup>1</sup>, JESSICA GRANT<sup>2</sup>, YONAS I. TEKLE<sup>2,6</sup>, ERICA LASEK-NESSELQUIST<sup>3,4</sup>, HILARY G. MORRISON<sup>3</sup>, MITCHELL L. SOGIN<sup>3</sup>, DAVID J. PATTERSON<sup>5</sup>, AND LAURA A. KATZ<sup>1,2,\*</sup>



Tree of Life

• Key Things You Want to Know More About?



#### Tree of Life

• Why did I select this paper?



# Palfrey et al.

Syst. Biol. 59(5):518–533, 2010 © The Author(s) 2010. Published by Oxford University Press, on behalf of the Society of Systematic Biologists. All rights reserved. For Permissions, please email: journals.permissions@oxfordjournals.org DOI:10.1093/sysbio/syq037 Advance Access publication on July 23, 2010

#### Broadly Sampled Multigene Analyses Yield a Well-Resolved Eukaryotic Tree of Life

LAURA WEGENER PARFREY<sup>1</sup>, JESSICA GRANT<sup>2</sup>, YONAS I. TEKLE<sup>2,6</sup>, ERICA LASEK-NESSELQUIST<sup>3,4</sup>, HILARY G. MORRISON<sup>3</sup>, MITCHELL L. SOGIN<sup>3</sup>, DAVID J. PATTERSON<sup>5</sup>, AND LAURA A. KATZ<sup>1,2,\*</sup>



### What Does 1st Authorship Mean?

**Tree of Life** 

### Broadly Sampled Multigene Analyses Yield a Well-Resolved Eukaryotic Tree of Life

LAURA WEGENER PARFREY<sup>1</sup>, JESSICA GRANT<sup>2</sup>, YONAS I. TEKLE<sup>2,6</sup>, ERICA LASEK-NESSELQUIST<sup>3,4</sup>, HILARY G. MORRISON<sup>3</sup>, MITCHELL L. SOGIN<sup>3</sup>, DAVID J. PATTERSON<sup>5</sup>, AND LAURA A. KATZ<sup>1,2,\*</sup>

<sup>1</sup>Program in Organismic and Evolutionary Biology, University of Massachusetts, 611 North Pleasant Street, Amherst, MA 01003, USA; <sup>2</sup>Department of Biological Sciences, Smith College, 44 College Lane, Northampton, MA 01063, USA; <sup>3</sup>Bay Paul Center for Comparative Molecular Biology and Evolution, Marine Biological Laboratory, 7 MBL Street, Woods Hole, MA 02543, USA; <sup>4</sup>Department of Ecology and Evolutionary Biology, Brown University, 80 Waterman Street, Providence, RI 02912, USA; <sup>5</sup>Biodiversity Informatics Group, Marine Biological Laboratory, 7 MBL Street, Woods Hole, MA 02543, USA; <sup>6</sup>Present address: Department of Epidemiology and Public Health, Yale University School of Medicine, New Haven, CT 06520, USA;

\* Correspondence to be sent to: Laura A. Katz, 44 College Lane, Northampton, MA 01003, USA; E-mail: lkatz@smith.edu.

Laura Wegener Parfrey and Jessica Grant have contributed equally to this work.

Received 30 September 2009; reviews returned 1 December 2009; accepted 25 May 2010 Associate Editor: Cécile Ané

Tree of Life

#### **Tree of Life**

# Introduction



Perspectives on the structure of the eukaryotic tree of life have shifted in the past decade as molecular analyses provide hypotheses for relationships among the approximately 75 robust lineages of eukaryotes. These lineages are de ned by ultrastructural identities (Patterson 1999) -patterns of cellular and subcellular organization revealed by electron microscopy - and are strongly supported in molecular analyses (Parfrey et al. 2006; Yoon et al. 2008). Most of these lineages now fall within a small number of higher level clades, the supergroups of eukaryotes (Simpson and Roger 2004; Adl et al. 2005; Keeling et al. 2005). Several of these clades— Opisthokonta, Rhizaria, and Amoebozoa— are increasingly well supported by phylogenomic (Rodr'ıguez-Ezpeleta et al. 2007a; Burki et al. 2008; Hampl et al. 2009) and phylogenetic (Parfrey et al. 2006; Pawlowski and Burki 2009), analyses, whereas support for "Archaeplastida" predominantly comes from some phylogenomic studies (Rodríguez-Ezpeleta et al. 2005; Burki et al. 2007) or analyses of plastid genes (Yoon et al. 2002; Parfrey et al. 2006). In con- trast, support for "Chromalveolata" and Excavata is mixed, often dependent on the selection of taxa in- cluded in analyses (Rodr'iguez-Ezpeleta et al. 2005; Parfrey et al. 2006; Rodríguez-Ezpeleta et al. 2007a; Burki et al. 2008; Hampl et al. 2009). We use quotation marks throughout to note groups where uncertainties remain. Moreover, it is dif cult to evaluate the overall stability of major clades of eukaryotes because phyloge- nomic analyses have 19 or fewer of the major lineages and hence do not suf ciently sample eukaryotic diver- sity(Rodríguez-Ezpeletaetal.2007b;Burkietal.2008; Hampl et al. 2009), whereas taxon-rich analyses with 4 or fewer genes yield topologies with poor support at deep nodes (Cavalier-Smith 2004; Parfrey et al. 2006; Yoon et al. 2008).

#### Tree of Life

Perspectives on the structure of the eukaryotic tree of life have shifted in the past decade as molecular analyses provide hypotheses for relationships among the approximately 75 robust lineages of eukaryotes. These lineages are de ned by ultrastructural identities (Patterson 1999) -patterns of cellular and subcellular organization revealed by electron microscopy - and are strongly supported in molecular analyses (Parfrey et al. 2006; Yoon et al. 2008). Most of these lineages now fall within a small number of higher level clades, the supergroups of eukaryotes (Simpson and Roger 2004; Adl et al. 2005; Keeling et al. 2005). Several of these clades— Opisthokonta, Rhizaria, and Amoebozoa — are increasingly well supported by phylogenomic (Rodríguez-Ezpeleta et al. 2007a; Burki et al. 2008; Hampl et al. 2009) and phylogenetic (Parfrey et al. 2006; Pawlowski and Burki 2009), analyses, whereas support for "Archaeplastida" predominantly comes from some phylogenomic studies (Rodríguez-Ezpeleta et al. 2005; Burki et al. 2007) or analyses of plastid genes (Yoon et al. 2002; Parfrey et al. 2006). In con-trast, support for "Chromalveolata" and Excavata is mixed, often dependent on the selection of taxa in- cluded in analyses (Rodríguez-Ezpeleta et al. 2005; Parfrey et al. 2006; Rodríguez-Ezpeleta et al. 2007a; Burki et al. 2008; Hampl et al. 2009). We use quotation marks throughout to note groups where uncertainties remain. Moreover, it is dif cult to evaluate the overall stability of major clades of eukaryotes because phyloge- nomic analyses have 19 or fewer of the major lineages and hence do not suf ciently sample eukaryotic diver- sity(Rodríguez-Ezpeletaetal.2007b;Burkietal.2008; Hampl et al. 2009), whereas taxon-rich analyses with 4 or fewer genes yield topologies with poor support at deep nodes (Cavalier-Smith 2004; Parfrey et al. 2006; Yoon et al. 2008).

#### Tree of Life

Estimating the relationships of the major lineages of eukaryotes is difficult because of both the ancient age of eukaryotes (1.2–1.8 billion years; Knoll et al. 2006) and complex gene histories that include heterogeneous rates of molecular evolution and paralogy (Maddison 1997; Gribaldo and Philippe 2002; Tekle et al. 2009). A further issue obscuring eukaryotic relationships is the chimeric nature of the eukaryotic genome—not all genes are vertically inherited due to lateral gene transfer (LGT) and endosymbiotic gene transfer (EGT)---that can also mislead efforts to re-construct phylogenetic relationships (Andersson 2005; Rannala and Yang 2008; Tekle et al. 2009). This is especially true among photosynthetic lineages that comprise "Chromalveolata" and "Archaeplastida" where a large portion of the host genome (approximately 8–18%) is

#### 10011 Ct al. 2000 J.

Estimating the relationships of the major lineages of eukaryotes is difficult because of both the ancient age of eukaryotes (1.2–1.8 billion years; Knoll et al. 2006) and complex gene histories that include heterogeneous rates of molecular evolution and paralogy ex(Maddison 1997; Gribaldo and Philippe 2002; Tekle et al. 2009). A further issue obscuring eukaryotic relationships is the chimeric nature of the eukaryotic genome—not all genes are vertically inherited due to lateral gene transfer (LGT) and endosymbiotic gene transfer (EGT)-that can also mislead efforts to reconstruct phylogenetic relationships (Andersson 2005; Rannala and Yang 2008; Tekle et al. 2009). This is especially true among photosynthetic lineages that comprise "Chromalveolata" and "Archaeplastida" where a large portion of the host genome (approximately 8-18%) is

#### Tree of Life

Estimating the relationships of the major lineages of eukaryotes is difficult because of both the ancient age of eukaryotes (1.2–1.8 billion years; Knoll et al. 2006) and complex gene histories that include heterogeneous rates of molecular evolution and paralogy x(Maddison 1997; Gribaldo and Philippe 2002; Tekle et al. 2009). A further issue obscuring eukaryotic relationships is the chimeric nature of the eukaryotic genome—not all genes are vertically inherited due to lateral gene transfer (LGT) and endosymbiotic gene transfer (EGT)-that can also mislead efforts to reconstruct phylogenetic relationships (Andersson 2005; Rannala and Yang 2008; Tekle et al. 2009). This is especially true among photosynthetic lineages that comprise "Chromalveolata" and "Archaeplastida" where a large portion of the host genome (approximately 8-18%) is

# LGT and Endosymbiosis



There is a long-standing debate among systematists as to the relative benefits of increasing gene or taxon sampling (Hillis et al. 2003; Cummings and Meyer 2005; Rokas and Carroll 2005). Both approaches improve phylogenetic reconstruction by alleviating either stochastic or systematic phylogenetic error (e.g., Rokas and Carroll 2005; Hedtke et al. 2006). Stochastic error results from too little signal in the data (e.g., single to few gene trees) to estimate relationships and results in poorly resolved trees with low support, especially at deep levels (Swofford et al. 1996; Rokas and Carroll 2005). The problems of stochastic error are amplified for deep relationships, such as relationships among major clades of eukaryotes (Roger and Hug 2006). Many researchers opt to increase the number of genes, exemplified by phylogenomic studies, which alleviates stochastic error and yields well-resolved trees that are highly supported (Rokas and Carroll 2005; Burki et al. 2007; Hampl et al. 2009). However, analyses of many genes are still vulnerable to systematic error and often include very few lineages.

22

There is a long-standing debate among systematists as to the relative benefits of increasing gene or taxon sampling (Hillis et al. 2003; Cummings and Meyer 2005; Rokas and Carroll 2005). Both approaches improve phylogenetic reconstruction by alleviating either stochastic or systematic phylogenetic error (e.g., Rokas and Carroll 2005; Hedtke et al. 2006). Stochastic error results from too little signal in the data (e.g., single to few gene trees) to estimate relationships and results in poorly resolved trees with low support, especially at deep levels (Swofford et al. 1996; Rokas and Carroll 2005). The problems of stochastic error are amplified for deep relationships, such as relationships among major clades of eukaryotes (Roger and Hug 2006). Many researchers opt to increase the number of genes, exemplified by phylogenomic studies, which alleviates stochastic error and yields well-resolved trees that are highly supported (Rokas and Carroll 2005; Burki et al. 2007; Hampl et al.

Tree of Life

Systematic error results from biases in the data that mislead phylogenetic reconstruction, yielding incorrect sister group relationships that do not reflect historical relationships; the most well known of these is longbranch attraction (Felsenstein 1978). Incongruence can also arise from conflicts between gene trees and species trees resulting from population genetic processes or the chimeric nature of eukaryotic genomes (Maddison 1997; Rannala and Yang 2008). Systematic errors can be detected and eliminated by several methods that are often combined, including using more realistic models of sequence evolution (e.g., Rodríguez-Ezpeleta et al. 2007b), removing rapidly evolving genes and/or taxa that cause errors (Brinkmann et al. 2005), and by increasing taxonomic sampling (Zwickl and Hillis 2002; Hedtke et al. 2006). Increased taxon sampling has been shown to improve phylogenetic accuracy even when the additional taxa contain large amounts of missing data (Philippe et al. 2004; Wiens 2005; Wiens and Moen 2008). In contrast, the abundance of data in phylogenomic studies can yield highly supported, but incorrect relationships caused by these systematic biases (Philippe et al. 2004; Hedtke et al. 2006; Jeffroy et al. 2006; Rokas and Chatzimanolis 2008). Taxon-rich analyses provide a method for testing the accuracy of relationships that receive high BS support in phylogenomic analyses (Zwickl and Hillis 2002; Heath et al. 2008).

Systematic error results from biases in the data that mislead phylogenetic reconstruction, yielding incorrect sister group relationships that do not reflect historical relationships; the most well known of these is longbranch attraction (Felsenstein 1978). Incongruence can also arise from conflicts between gene trees and species trees resulting from population genetic processes or the chimeric nature of eukaryotic genomes (Maddison 1997; Rannala and Yang 2008). Systematic errors can be detected and eliminated by several methods that are often combined, including using more realistic models of sequence evolution (e.g., Rodríguez-Ezpeleta et al. 2007b), removing rapidly evolving genes and /or taxa that cause errors (Brinkmann et al. 2005), and by increasing taxonomic sampling (Zwickl and Hillis 2002; Hedtke et al.



# Introduction

#### Tree of Life





# Methods

#### Tree of Life

- What data did they use?
- What genes?
- What taxa?
- How did they get the data?



### Taxa and Genes

#### Data set Assembly

Taxa and genes were selected to maximize taxonomic diversity and evenness given the availability of molecular data. This strategy was used to improve phylogenetic accuracy by breaking up long branches with dense sampling across the eukaryotic tree (Hillis 1998). The classifications systems of Patterson (1999) and Adl et al. (2005) were used as guides as we aimed to sample eukaryotic diversity by including representatives of as many lineages defined by ultrastructural identities as possible (Table S2). These lineages have generally proven to be robust as they are well supported in molecular analyses (e.g., Adl et al. 2005; Parfrey et al. 2006; Yoon et al. 2008), including the current study, and they represent monophyletic groups that serve as a proxy for taxonomic diversity. Our data set has representatives from 72 lineages, including 53 of the 71 lineages plus 7 of 200 unplaced genera as defined in Patterson (1999). Additionally, we include 3 unplaced lineages isolated more recently, Malawimonas jakobiformis (O'Kelly and Nerad 1999), Breviata anathema (Walker et al. 2006),

### Taxa and Genes

un minimito j.

To maximize gene sampling for diverse taxa, we include markers historically targeted by polymerase chain reaction-based analyses (e.g., SSU-rDNA, actin, elongation factor  $1\alpha$ ; Table S3) plus commonly sequenced ESTs (e.g., ribosomal proteins, 14-3-3; Table S3). The comprehensively sampled SSU-rDNA and the historical markers facilitate inclusion of many additional taxa for which only these genes have been characterized (Table S4). The minimum sequence data required for inclusion were nearly full-length SSU-rDNA, which provided the core of information necessary for phylogenetic placement with large amounts of missing data (Wiens and Moen 2008).





### Alignments

SSU-rDNA sequences were hand curated for target taxa by removing introns, unalignable regions, nonnuclear rDNAs, and misannotated sequences. This alignment was crucial to overall accuracy because nearly half of the target taxa are represented only by SSU-rDNA, thus several alignment and masking methods were assessed to ensure the robustness of the SSUrDNA alignment. SSU-rDNA sequences were aligned by HMMER (Eddy 2001), version 2.1.4 with default settings, taking secondary structure into account. HMMER used a set of previously aligned sequences to model the secondary structure of a sequence. The training alignment for building the model, consisting of all available SSU-rDNA eukaryote sequences (as of December 2008) aligned according to their secondary structure, was downloaded from the European Ribosomal Database (Wuyts et al. 2002). An additional SSU-rDNA alignment was constructed in MAFFT 6 implemented in

### Alignments

dowinoaded nom the European Nicosomar Database (Wuyts et al. 2002). An additional SSU-rDNA alignment was constructed in MAFFT 6 implemented in SeaView (Galtier et al. 1996) with the E-INS-i algorithm (Katoh and Toh 2008). Both alignments were further edited manually in MacClade v4.08 (Maddison D.R. and Maddison W.P. 2005). To assess the effect of rate heterogeneity on the SSU-rDNA topologies, we partitioned the data matrices into 8 rate classes using the general time-reversible (GTR) model with invariable sites and rate variation among sites following a discrete gamma distribution, as implemented in HyPhy version .99b package (Kosakovsky Pond et al. 2005). We then ran analyses without the fastest and two fastest rate classes, resulting in 1197 and 1019 characters, respectively. However, the reduced data sets resulted in less resolution in the backbone without improving apparent the long-branch attraction. Thus, we used the alignment generated in MAFFT and masked with GBlocks (Talavera and Castresana 2007) and by eye in MacClade, resulting in 867 unambiguously aligned characters.

### **Protein Database Searches**

#### Tree of Life

Assembly of the protein data set relied on a custombuilt pipeline and database that combined Perl and Python scripts to identify homologs from diverse eukaryotes. Our goal in developing this pipeline was to ensure that we captured the broadest possible set of sequences given the tremendous heterogeneity among microbial eukaryotes. All available protein and EST



data from our target taxa (Table S4) were downloaded from GenBank in January 2009 and ESTs were analyzed in all 6 translated frames to identify correct sequences for our alignment. A fasta file of 6 sequences representing the six "supergroups" was created for each target gene and used to query our database of target taxa by BLASTp. Results were limited by length, e-value, and identity, and all sequences with greater than 1% divergence within each taxon were retained for assessment of paralogy. The resulting sequences were aligned with ClustalW (Thompson et al. 1994) and the resulting single gene alignments were assessed by eye to remove nonhomologous sequences.

### Alignments and Paralogs

The inferred amino acid sequences for each of the protein genes from our data pipeline were combined with the new sequences generated for this study and again aligned in Clustal W (Thompson et al. 1994). The alignment was adjusted by eye in MacClade (Maddison D.R. and Maddison W.P. 2005). As these alignments included all paralogs extracted from the pipeline, individual gene trees were examined to choose appropriate orthologs. For example, in cases where paralogs formed a monophyletic group, the shortest branch sequence was retained. When paralogs fell into multiple locations on the tree, we aimed to maintain orthologous groups that included the greatest taxonomic representation. The individual gene alignments were then concatenated to build a 16 gene, 451-taxon matrix with 6578 unambiguously aligned characters, including SSU-rDNA. All other data sets were constructed by removing taxa and/or genes from this matrix. All data matrices are available at TreeBASE (submission ID S10562).
### Creation of Subdata Matrices

We created an array of data matrices by subsampling our full data matrix of 16 genes (15 protein-coding genes plus SSU-rDNA) and 451 taxa (denoted all:16) in order to assess the impact of taxon sampling, missing data, and gene sampling. First, seven data sets were created to assess the impact of missing data and taxon sampling (summarized in Table 1). The least inclusive of these contained 16 genes and all 88 taxa that had at least 10 of the 16 genes (10:16), which resulted in 17% missing data. Similarly, the 6:16 and 4:16 matrices include all taxa with at least 6 and 4 of the targeted 16 genes, respectively. SSU-rDNA is ubiquitously sampled in our data set and many phylogenetic hypotheses are based on SSU-rDNA genealogies. To address the concern that SSU-rDNA was driving our results, we deleted it from each of the 16 gene data sets resulting in 9:15, 5:15, 3:15, and all:15 matrices.

genes and tana.

Photosynthetic lineages have chimeric genomes that are composed of genes originating both from the host eukaryote, the endosymbiotic plastid (through EGT), and, in cases of secondary or greater endosymbiosis, from the symbiont nucleus. If genes of multiple origins were retained in our concatenated data set, the resulting conflicting signal between host, symbiont, and plastid could mislead phylogenetic reconstruction. This chimerism may contribute to the instability observed for photosynthetic lineages without clear sister groups (red algae, green algae, glaucocystophytes, cryptomonads, and haptophytes). Thus, we used 2 methods to detect discordance among loci that could indicate EGT. First, the 16 genes from representatives of each of these photosynthetic lineages were analyzed by top BLASTp hit. We

### *Phylogenetic Analyses*

Genealogies for this study were constructed almost exclusively in RaxML. The MPI version of RaxML 7.0.4 with rapid bootstrapping was used (Stamatakis et al. 2008). The SSU-rDNA partition was analyzed with GTR+gamma as this was the best fitting model available in RAxML, according to MrModelTest (Nylander 2004). ProtTest (Abascal et al. 2005) was used to select the appropriate model of sequence evolution for the amino acid data using the 9:15 data set. The WAG amino acid replacement matrix was found to be the best-fitting model for the concatenated data, but the rtREV amino acid replacement matrix was the best for some of the individual partitions and both WAG and rtREV were among the top 3 models for all but 1 gene (and with similar likelihood scores). We ran our data under both WAG and rtREV models and found consistent results, indicating that our interpretations are robust to at least

## Bootstrapping

#### **Tree of Life**

this level of model choice. The results presented are from the WAG analyses and the rtREV analyses differed only in level of BS for key nodes (usually  $\pm 5$ points). In initial analyses, the appropriate number of independent bootstrap replicates was determined for each data set using bootstopping criteria in RAxML 7.0.4 as implemented on Cyberinfrastructure for Phylogenetic Research (CIPRES) portal 2 (Miller et al. 2009). All analyses stopped after 200 or fewer replicates, except all:16, which stopped after 400 replicates. In later analyses, using the MPI version of RAxML, which does not implement a bootstopping criterion, 200 rapid bootstrap replicates followed by a full maximum-likelihood search was used for all analyses except all:16, for which 600 bootstrap replicates were run. Because of the computational cost of the all:16 analysis, this was run as 6 separate analyses: 100 bootstraps followed by a full maximum likelihood search and 5 other runs of 100 bootstraps each. These data were combined in RAxML to complete the analysis. We found no significant difference in comparisons between fast and slow RAxML bootstrap methods (Fig. S1i), which we tested because the fast bootstrapping method in RAxML can produce misleading results particularly for long-branch taxa (Leigh 2008). The results of rapid bootstrapping are shown.

## **Other Methods**

#### Tree of Life

#### **3110 1111**

To investigate the stability of our tree topology under different analytic methods, select data sets were analyzed with Bayesian approaches and Parsimony (Fig. S1s–v). Parsimony analysis of 10:16, implemented in Paup\* (Swofford 2002), yielded a less resolved version of the RAxML topology (i.e. Excavata as a polytomy) that is generally concordant with the more resolved tree obtained by maximum-likelihood methods. The

## Methods

#### Tree of Life

## Questions about Methods?



## Results and Discussion



## RESULTS AND DISCUSSION

## Robust Topology of the Eukaryotic Tree of Life

Many major clades were consistently recovered across our analyses (Fig. 1 and Table 1). These stable groups receive moderate to strong support in analyses with limited missing data (Fig. 2) and less support as missing data increases. The Opisthokonta, which includes animals and funginand the beterogeneous clade Rhizaria

## Fig 1: 451 Taxa and some of the 16 genes



FIGURE 1. Most likely eukaryotic tree of life reconstructed using all 451 taxa and all 16 genes (SSU-rDNA plus 15 protein genes). Major nodes in this topology are robust to analyses of subsets of taxa and genes, which include varying levels of missing data (Table 1). Clades in bold are monophyletic in analyses with 2 or more members except in all:15 in which taxa represented by a single gene were sometimes misplaced. Numbers in boxes represent support at key nodes in analyses with increasing amounts of missing data (10:16, 6:16, 4:16, and all:16 analyses; see Table 1 for more details). Given uncertainties around the root of the eukaryotic tree of life (see text), we have chosen to draw the tree rooted with the well-supported clade Opisthokonta. Dashed line indicates alternate branching pattern seen for Amoebozoa in other analyses. Long branches, indicated by //, have been reduced by half. The 6 lineages labeled by \* represent taxa that are misplaced, probably due to LBA, listed from top to bottom with expected clade in parentheses. These are Protopalina japonica (Stramenopiles), Aggregata actopians (Apicon lexa), Mikroytos mackini (Haplosporidia), Centropyxis laevigata (Tubulinea), Marteilioides changmaensis (unplaced), and Cochliopadium spiniferum (Amoebozoa).

#### **Tree of Life**





**Tree of Life** 









Euryarchaea

## Fig. 2: 88 Taxa each w/ 10 or more of the 16 genes

#### **Tree of Life**



FIGURE 2. Most likely eukaryotic tree of life reconstructed with 10:16, which includes 88 taxa (each with 10 or more of the genes analyzed in this study) and 16 genes (SSU-rDNA plus 15 protein genes). Thickened lines receive >95% bootstrap support. Other notes as in Methods section and Figure 1.







94 94 94 94 94 94 94 94 94 94	- Tetrahymena thermophila aramecium tetraurelia la uncinata lis	SAR: Rhizaria Stramenopiles Alveolates
100 99 94 82 Gromia, Paradinium Corallomyxa Theratromyxa Allantion Gromia, Paradinium Corallomyxa Halantion Metopion Metopion	Foraminifera Polycystina Acantharea Haplosporidia, Plasmodiophora Euglyphida, Cercomonadida Thaumatomonads Phaeodarea Desmothoracids, Gymnophrea	SAR: Rhizaria
Paramonas Filobiopsids	Chlorarachniophytes Diatoms Brown algae St Dinoflagellates Apicomplexa Ciliates	ramenopiles Alveolates

**Tree of Life** 



**Tree of Life** 





Eurvarchaea





Euglenozoa
Excavata

Heterolobosea
Jakobids

Malawimonas
Parabasalids

Preaxostyla
Fornicata



#### Dieviala analiterita Capitella capitata Aplysia californica Schistosoma mansoni Drosophila melanogaster Caenorhabditis elegans Homo sapiens Gallus gallus Branchiostoma floridae Strongylocentrotus purpuratus — Ciona intestinalis Oscarella carmela - Mnemiopsis leidyi Opisthokonta Nematostella vectensis Monosiga brevicollis Amoebidium parasiticum — Sphaeroforma arctica Capsaspora owczarzaki Candida albicans - Saccharomyces cerevisiae Schizosaccharomyces pombe Phanerochaete chrysosporium - Allomyces macrogynus Spizellomyces punctatus Encephalitozoon cuniculi



#### **Tree of Life**

- (reviewed in Roger and Simpson 2009; Tekle et al. 2009).
- . In our analyses, we find at best moderate support for
  - "Unikonta" (Table 1), but concatenated analyses such as
- these cannot resolve the root.
  - In evoluting the tradeoffe hetween increasing taxo-



#### Tree of Life

	10:16	6:16	4:16	all:16	9:15	5:15	3:15	all:15
Supported clades								
Opisthokonta	99 <sup>a</sup>	<b>97</b> <sup>a</sup>	<b>97</b> <sup>a</sup>	69	<b>100</b> <sup>a</sup>	99 <sup>a</sup>	85	19
Rĥizaria	<b>100</b> <sup>a</sup>	<b>99</b> <sup>a</sup>	<b>94</b> <sup>a</sup>	82	<b>100</b> <sup>a</sup>	<b>100</b> <sup>a</sup>	47	29
SAR	<b>97</b> <sup>a</sup>	<b>98</b> <sup>a</sup>	63	22	<b>100</b> <sup>a</sup>	<b>100</b> <sup>a</sup>	32	19
Rhizaria + stramenopiles	<b>94</b> <sup>a</sup>	<b>94</b> <sup>a</sup>	57	26	<b>92</b> <sup>a</sup>	<b>96</b> <sup>a</sup>	29	18
Excavata	83	77	65	6	84	76	44	19
Amoebozoa	59	46	49	nm	68	56	44	5
"Unikonta"	63	39	21	nm	54	50	15	3
Weak/unsupported hypotheses								
"Archaeplastida"	nm	nm	nm	nm	nm	nm	nm	nm
"Chromalveloata"	nm	nm	nm	nm	nm	nm	nm	nm
Cryptomonads + haptophytes	33	50	nm	29	38	56	22	25
Haptophytes + SAR	nm	nm	15	nm	nm	nm	nm	nm
Alveolates + stramenopiles	nm	nm	nm	nm	nm	nm	nm	nm
Red algae + green algae	nm	nm	nm	nm	nm	nm	nm	nm
Red, Green, Glauco, Hapto, Crypt	47	32	nm	9	39	27	16	8
Data set statistics								
Number of taxa	88	111	160	451	88	111	160	240
Number of lineages	26	30	45	72	26	30	45	54
% Missing data (characters)	17	25	38	69	19	28	43	59

TABLE 1. Support for major clades of eukaryotes in analyses containing varying levels of taxon inclusion and missing data

Note: Supported clades are stable across analyses, albeit with decreasing support as the percentage of missing data increases. Bootstrap support values from RAxML analyses. Support values greater than 75 are indicated by bold text and greater than 85 are indicated with a. nm = nonmonophyletic. Column headings describe the data sets. For example, "10:16" includes all taxa that have at least 10 of the 16 genes, with a total of 88 taxa representing 26 lineages and containing 17% missing data. The "all:15" includes the protein-coding genes from all taxa and contains 59% missing data. See Table S2 for lineagesand Figure S1a–h for individual trees. <sup>a</sup>Support values greater than 85.

### Just Rhizaria



FIGURE 3. Maximum likelihood tree of Rhizaria reconstructed with 103 Rhizaria taxa and 16 genes. The SSU-rDNA partition was analyzed with GTR+gamma and proteins with rtREV. Thickened lines receive >80% bootstrap support in all analyses. Node support in boxes from Rhizaria:4-gene, Rhizaria:16-gene, all:16 analyses. Taxa with new data are in bold. Dashed lines indicate nonmonophyly.

### Just Excavata



proteins with rtREV. See Figure 3 for other notes.

## Consensus



FIGURE 5. Summary of major findings—the evolutionary relationships among major lineages of eukaryotes. Clades have been collapsed into those that we view to be strongly supported. The many polytomies represent uncertainties that remain.

## **Results and Discussions**



moj ao a babar member or receible observe

### CONCLUSIONS

The robust tree of life emerging from this study demonstrates the benefits of improved taxon sampling for reconstructing deep phylogeny as our analyses produce stable topologies that include a broad representation of eukaryotes. The current study, combined with insights from other studies referenced herein, has refined the eukaryotic tree of life from over 70 major lineages (Patterson 1999) to ~16 major groups (Fig. 5, http://eutree.lifedesks.org/). Most significantly, we attribute the stability of major clades (e.g., Excavata, Amoebozoa, Opisthokonta, and SAR) to broader taxonomic sampling combined with analyses of sufficient characters (16 genes or 6578 characters). In our view, inclusion of more taxa coupled with carefully chosen genes is necessary to further resolve the 16 or so major lineages of microbial eukaryotes for which sister group relationships remain uncertain.

## UPGMA

# Unweighted Pair Group Method with Arithmetic mean (UPGMA) algorithm

## The True Tree



## Distance Matrix For True Tree

## TABLE 27.6. Distance matrix

OTUs	Α	В	С	D	Ε	F
А	0	2	4	6	6	8
В	2	0	4	6	6	8
С	4	4	0	6	6	8
D	6	6	6	0	4	8
E	6	6	6	4	0	8
F	8	8	8	8	8	0

## Collapse Diagonal

## TABLE 27.7. Diagonal matrix—Step 1

OTUs	Α	B	С	D	Ε
В	2				
С	4	4			
D	6	6	6		
E	6	6	6	4	
F	8	8	8	8	8

## Identify Lowest D

OTUs	A	B	С	D	E
B	<u>2</u>				
С	4	4			
D	6	6	6		
E	6	6	6	4	
F	8	8	8	8	8

## Join Those Two Taxa

• Create branch with length = D

• 2 • A-----B

## Make D from Node Equal



## Create New Distance Matrix

• Merge Two OTUs joined in previous step (AB)

• 
$$D_{x,AB} = 0.5 * (D_{x,A} + D_{x,B})$$
#### New Matrix

#### TABLE 27.8. Diagonal matrix—Step 2

OTUs	AB	С	D	Ε	
C	4				
	4	6			
F	6	6	1		
L	0	0	4	0	
Г	0	0	0	0	

### UPGMA

TABLE 27.9. Example of UPGMA tree construction						
Step	Cycle 1	Cycle 2	Cycle 3	Cycle 4	Cycle 5	Cycle 6
Distance matrix	OTUs     A     B     C     D     E       B     2     -     -     -       C     4     4     -     -       D     6     6     6     -       E     6     6     6     4       F     8     8     8     8     8	OTUs AB C D E C 4 D 6 6 E 6 6 4 F 8 8 8 8	OTUS AB C DE C 4 DE 6 6 DE 8 8 8	OTUS ABC DE DE <mark>6</mark> F 8 8	OTUs <u>ABCDE</u> F <mark>8</mark>	No new matrix
ldentify smallest D	$A \leftrightarrow B = 2$	$AB \leftrightarrow C = 4$ $D \leftrightarrow E= 4$	$AB \leftrightarrow DE = 6$ $C \leftrightarrow DE = 6$	ABC↔DE	ABCDE↔F	
Taxa joined	A and B	D and E	AB and C	ABC and DE	ABCDE and F	
Subtree	1 A 1 B	2 D 2 E	1 1 A 2 C	1 1 A 1 2 C 1 2 D E	1 1 A 1 2 C 1 2 D 4 F	Root 1 1 A 1 1 B 2 C 1 2 D 4 F
Comments on tree drawing	The distance between A and B is 2 units. A sub- tree is drawn with the branch point halfway between the two. Thus, each branch is 1 unit in length.	Branching done as in Step 1. Because the distance from AB to C is also 4, that pair could have been selected as well.	First a subtree is drawn with AB and C: 2 AB C The the AB subtree is attached to the AB branch at a point equal to the length of the A and B branches.	The tree is first done as in Step 3 with the ABC and DE subtrees replacing the branches.	The tree is now complete but unrooted.	The tree can then be rooted using midpoint rooting which tries to balance all the tips to reach the same end point. Note this is the tree that we started with to build the distance matrix.

#### Compare to True Tree



#### But ...

#### • What is evolutionary rates not equal



# UPGMA with Unequal rates

TABLE 27.10. UPGMA tree construction errors						
	Cycle 1	Cycle 2	Cycle 3	Cycle 4	Cycle 5	Cycle 6
Distance matrix	A       B       C       D       E         B       5       -       -       -       -         C       4       7       -       -       -         D       7       10       7       -       -         E       6       9       6       5       -         F       8       11       8       9       8	ACBDEB4D710-E695-F81189	ACBDEB4-DE6.59.5F811	ACB DE DE <mark>8</mark> F 9.5 9.5	ABCDE F 9	No new matrix
ldentify smallest D	$A \leftrightarrow C = 4$	$D \leftrightarrow E = 5$	$AC \leftrightarrow B = 4$	$ACB \leftrightarrow DF = 8$	$ABCDE \leftrightarrow F = 9$	
Taxa joined Subtree	A and C $ \begin{array}{c} 2 \\ 2 \\ C \end{array} $	D and E 2.5 2.5 E	A and C with B 1 $2$ A 2 C 3 B	ABC with DE	ABCDE with F	$\begin{bmatrix} 1 & 1 & A \\ 1 & 1 & B \\ 2 & C \\ 1 & 2 & C \\ 1 & 2 & D \\ 4 & F \end{bmatrix}$
Comments					F Note how this is not the same as the starting tree	

#### Compare to True Tree



# Neighbor Joining

### Start with Star Tree



# Calculate S

- For each OTU, calculate a measure (S) as follows.
- S is the sum of the distances (D) between that OTU and every other OTU, divided by N-2 where N is the total number of OTUs.
- This is a measure of the distance an OTU is from all other OTUs

# Calculate Pairwise Distances

• Calculate the distance D<sub>ij</sub> between each OTU pair (e.g., I and J).

# Identify Closest Pair

• Identify the pair of OTUs with the minimum value of  $D_{ij} - S_i - S_j$ 

# Joining Taxa

• As in UPGMA, join these two taxa at a node in a subtree.

# Calculate Branches

- Calculate branch lengths. Unlike UPGMA, neighbor joining does not force the branch lengths from node X to I (D<sub>xi</sub>) and to J (D<sub>xj</sub>) to be equal, i.e., does not force the rate of change in those branches to be equal. Instead, these distances are calculated according to the following formulas
- $D_{xi} = 1/2 D_{ij} + 1/2 (S_i S_j)$
- $D_{xj} = 1/2 D_{ij} + 1/2 (S_j S_i)$

## Calculate New Matrix

• Calculate a new distance matrix with I and J merged and replaced by the node (X) that joins them. Calculate the distances from this node to the other tips (K) by:

• 
$$D_{xk} = (D_{ik} + D_{jk} - D_{ij})/2$$

#### Distance Matrix



- $S_A = (5+4+7+6+8) / 4 = 7.5$
- $S_B = (5+7+10+9+11) / 4 = 10.5$
- $S_C = (4+7+7+6+8) / 4 = 8$
- $S_D = (7+10+7+5+9) / 4 = 9.5$
- $S_E = (6+9+6+5+8) / 4 = 8.5$
- $S_F = (8+11+8+9+8) / 4 = 11$

# M

•  $M_{ij} = D_{ij} - S_i - S_j$ 

- Smallest are
- $M_{AB} = 5 7.5 10.5 = -13$
- $M_{DE} = 5 9.5 8.5 = -13$
- Choose one of these (AB here)

# New Tree

- Create a node (U) that joins pair with lowest  $M_{ij}$  such that
- $S_{IU} = D_{IJ}/2 + (S_I S_J)/2$
- Merge U with other taxa
- Join I and J according to S above and make all other taxa in form of a star



### New Matrix

•  $D_{XU} = D_{IX} + D_{JX} - D_{IJ}$  where I and J are those selected from above.

TABLE 27.11. Neighbor-joining example						
	Cycle 1	Cycle 2	Cycle 3	Cycle 4	Cycle 5	
Distance matrix	A     B     C     D     E       B     5     -     -     -       C     4     7     -     -       D     7     10     7     -       E     6     9     6     5       F     8     11     8     9     8	U <sub>1</sub> C D E C 3 D 6 7 E 5 6 5 F 7 8 9 8	$\begin{array}{c cccc} U_1 & C & U_2 \\ C & 3 & & \\ U_2 & 3 & 4 & \\ F & 7 & 8 & 6 \end{array}$	$ \begin{array}{cccccccc} U_2 & U_3 \\ U_3 & 2 \\ F & 6 & 6 \end{array} $	U <sub>4</sub> F 5	
Step 1						
S calculations $S_x = (\text{sum all } D_x)/(N-2),$ where N is the # of OTUs in the set.	$\begin{split} S_{\rm A} &= (5{+}4{+}7{+}6{+}8)/4 = 7.5\\ S_{\rm B} &= (5{+}7{+}10{+}9{+}11)/4 = 10.5\\ S_{\rm C} &= (4{+}7{+}7{+}6{+}8)/4 = 8\\ S_{\rm D} &= (7{+}10{+}7{+}5{+}9)/4 = 9.5\\ S_{\rm E} &= (6{+}9{+}6{+}5{+}8)/4 = 8.5\\ S_{\rm F} &= (8{+}11{+}8{+}9{+}8)/4 = 11 \end{split}$	$S_{U_1} = (3+6+5+7)/3 = 7$ $S_C = (3+7+6=8)/3 = 8$ $S_D = (6+7+5+9)/3 = 9$ $S_E = (5+6+5+8)/3 = 8$ $S_F = (7+8+9+8)/3 = 10.6$	$S_{U_1} = (3+3+7)/2 = 6.5$ $S_C = (3+4+8)/2 = 7.5$ $S_{U_2} = (3+4+6)/2 = 6.5$ $S_F = (7+8+6)/2 = 10.5$	$S_{\cup 2} = (2+6)/1 = 8$ $S_{\cup 3} = (2+6)/1 = 8$ $S_{\rm F} = (6+6)/1 = 12$	Because N – 2 = 0, we cannot do this calculation.	
Step 2						
Calculate pair with smallest ( <i>M</i> ), where $M_{ij} = D_{ij} - S_i - S_j$ .	Smallest are $M_{AB} = 5 - 7.5 - 10.5 = -13$ $M_{DE} = 5 - 9.5 - 8.5 = -13$ Choose one of these (AB here).	Smallest is $M_{CU_1} = 3 - 7 - 8 = -12$ $M_{DE} = 5 - 9 - 8 = -12$ Choose one of these (DE here).	Smallest is M <sub>CU1</sub> = 3 – 6.5 – 7.5 = –11	Smallest is $M_{U_2F} = 6 - 8 - 12 = -14$ $M_{U_3F} = 6 - 8 - 12 = -14$ $M_{U_2U_3} = 2 - 8 - 8 = -14$ Choose one of these ( $M_{U_2U_3}$ here).		
Step 3				2 9		
Create a node (U) that joins pair with lowest $M_{ij}$ such that $S_{I\cup} = D_{ij}/2 + (S_i - S_j)/2$ .	U <sub>1</sub> joins A and B: $S_{AU1} = D_{AB}/2 + (S_A - S_B)/2 = 1$ $S_{BU1} = D_{AB}/2 + (S_B - S_A)/2 = 4$	U <sub>2</sub> joins D and E: $S_{DU2} = D_{DE}/2 + (S_D - S_E)2 = 3$ $S_{EU2} = D_{DE}/2 + (S_E - S_D)/2 = 2$	U <sub>3</sub> joins C and U <sub>1</sub> : $S_{CU_3} = D_{CU_1}/2 + (S_C - S_{U_1})/2 = 2$ $S_{U_1U_3} = D_{CU_1}/2 + (S_{U_1} - S_C)/2 = 1$	U <sub>4</sub> joins U <sub>2</sub> and U <sub>3</sub> : $S_{U_2U_4} = D_{U_2U_3}/2 + (S_{U_2} - S_{U_3})/2 = \frac{1}{2}$ $S_{U_3U_4} = D_{U_2U_3}/2 + (S_{U_3} - S_{U_2})/2 = \frac{1}{2}$	For last pair connect 1 with branch = D. 1 Here $D_{\cup 4\mathrm{F}}$ = 5.	
Step 4						
Join <i>i</i> and <i>j</i> according to <i>S</i> above and make all other taxa in form of a star. Branches in black are of unknown length. Branches in red are of known length.	C = C = B C = A $C = BC = BC = BC = BC = BC = AB = A$	$ \begin{array}{c} C \\ D \\ 3 \\ U_2 \\ U_1 \\ 1 \\ A \end{array} $	$ \begin{array}{cccc} D & & C & B \\  & & U_2 & 2 & U_1 & 4 \\  & & U_2 & 2 & U_1 & 1 \\  & & E & 2 & U_3 & 1 & A \\  & & & F & F \end{array} $	$ \begin{array}{ccccc} D & 3 & C & 4 \\ & & 2 & 2 & 0 \\ & & & 2 & 0 \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & $	$ \begin{array}{ccccc} D & 3 & C & B \\ & & & 2 & U_1 & 4 \\ & & & U_2 & U_3 & U_1 & 1 \\ & & & E & 2 & U_4 & A \\ & & & 5 & F \\ \end{array} $	
Step 5	F	F			Comments	
Calculate new distance					Note this is the same	

tree we started with

(drawn in unrooted

form here).

Calculate new distance matrix of all other taxa to U with  $D_{x\cup} = D_{ix} + D_{jx} - D_{ij}$ , where *i* and *j* are those selected from above.

#### Compare to True Tree

