Unravelling complexity: Exploring the evolution of microbial eukaryotes through phylogenomics and phylogenetics

Anna Karnkowska Institute of Evolutionary Biology University of Warsaw



Institute of Evolutionary Biology



Courtesy of Sebastien Colin

about me





about me

 Taxonomy & Phylogenetics of Euglenophyta



My first (published) tree ever in 2005!



FIG. 2. The phylogenetic tree of the 18S rDNA sequence obtained by Bayesian inference (model GTR + G + I). Numbers at the essential nodes show posterior probabilities of the tree bipartitions as well as the bootstrap values/ decay indices obtained for the main clades by MP analysis and bootstrap values obtained by NJ and ML analysis (model GTR + I + G). The support for the remaining nodes (numbered) is listed in Table 5. Branches leading to nodes with support of less than 50% are collapsed, and those without substantial support (pp<0.95, bs<0.75) are not listed.

about me

 Taxonomy & Phylogenetics of Euglenophyta

- Evolution of mitochondria
- Reductive evolution
- Eukaryotic cell evolution
- Phylogenomics & comparative genomics of microbial euks.

- Evolution of organelles
- Eukaryotic cell evolution
- Reductive evolution
- Eukaryotic rhodopsins diversity & evolution
- Freshwater protists diversity and interactions
- Phylogenomics & comparative genomics













Evolution & genomics of microbial eukaryotes





microbial eukaryotes (protists) constitute majority of lineages on the **eukaryotic tree of life**



By Tricholome - Own work, CC BY-SA 4.0, <u>https://commons.wikimedia.org/w/index.php?curid=91622328</u> Burki et al. 2020 morphological diversity of protist



https://oceans.taraexpeditions.org

functional diversity of microbial eukaryotes



How the eukaryotic **cell complexity** evolved? What drives **diversification of eukaryotes**?

Endosymbiosis

Origin, evolution and fate of endosymbiotic organelles



Interactions of microbial eukaryotes in the changing environment





EXPERIMENTS



BIOINFORMATICS







Diversity & interactions of microbial eukaryotes



Who is there? What are the interactions? Which environmental factors influence interactions? Do all cells of the same species have the same (endo)symbionts?

Overview

- Tree of eukaryotes history
- Phylogenomics and the tree of eukaryotes
- Why do we need to resolve the tree of euakryotes?
- Origin of eukaryotes
- Metagenomics in phylogenomics

۶⁵5

• Phylogenomics as a tool to address (my) evolutionary questions

First tree of eukaryotes with a kingdom Protista

Ernst Heckl - 1866



New kingdom of Fungi Whittaker 1967



Molecular phylogeny

Carl Woese 1988

- only one phylogenetic marker rDNA
- three domains of life
- new relationships have been detected



The Next Generation Sequencing Revolution



Aimin Yang, Wei Zhang, Jiahao Wang, Ke Yang, Yang Han and Limin Zhang, CC BY-SA 4.0 < https://creativecommons.org/licenses/by-sa/4.0>, via Wikimedia Commons

More molecular data

The tree shown is one of two shortest trees found by parsimony analysis of concatenated four protein-encoding genes EF-1 α , actin, α -tubulin, and β -tubulin amino acid sequence



Baldauf et al. 2000, Science

More molecular data

Eukaryotes subdivided into supergroups





Holy Grail of "Excavata"

- 143 gene sumermatrix
- maximum likelihood method, using RAxML under WAG+ Γ model
- Bayesian analyses, implemented in PhyloBayes version 2.1c employing the CAT+ Γ model











Removal of long-branch taxa





Removal of long-branch gene sequences





Hampl et al. 2009

- Many of the general biases that apply to single-gene phylogenetics are exacerbated in a phylogenomic context
- The goal of a phylogenetic inference is to minimize the noise and maximize the true phylogenetic signal



More data & errors

 Saturated positions in the data set (stochastic errors) - arise when the number of positions in an alignment is small and the random background noise will have a neutralizing effect on the positions that contain the phylogenetic signal.

More data should have more phylogenetic signal

 Model misspecifications (systematic errors) - heterogeneity of nucleotide or amino acid composition, which tends to incorrectly cluster together species sharing the same composition, or the heterogeneity of the evolutionary rates, which can result in the LBA artefact.

More data do not solve systematic errors, yet since more data are available strategies to diminish the known sources of systematic errors can be applied

With more data there is a risk of amplification of errors!

Burki, 2014 & 2020

Eukaryotic phylogeny based on genomic and transcriptomic data



Trends in Ecology & Evolution

Trends in Ecology & Evolution 2020 3543-55DOI: (10.1016/j.tree.2019.08.008)

Typical pipeline

1. Inferred amino acid sequences of proteins retrieved from genomic or transcriptomic data

- 2. Adding sequences from new taxa to the well-curated sets of orthologous genes
- 3. Single-gene trees analysis to identify lateral gene transfers, incorrect paralog selections, and various contaminants (still mainly manual).
- 4. A subset of taxa is selected for the actual analysis to evenly cover the relevant phylogenetic breadth while excluding problematic species (e.g., those with limited data, extreme evolutionary rates in many genes, etc.).
- 5. Genes are concatenated into a phylogenomic 'supermatrix'.
- 6. Usually, both ML and Bayesian analyses are conducted. Various evolutionary models are employed (siteheterogeneous models), with choice often constrained by computational logistics.
- 7. Testing whether results are robust to perturbations of the data, especially excluding data most likely to foster incorrect phylogenetic inference (e.g., the fastest-evolving species, sites, or genes).

Discoveries of new evolutionary lineages







Hemimastix kukwesjijk

EUKARYOTA SUPERGROUPS

Lax et al. Nature, 2018

Hemimastigophora is a novel supra-kingdom-level lineage of eukaryotes

Gordon Lax, Yana Eglit, Laura Eme, Erin M. Bertrand, Andrew J. Roger & Alastair G. B. Simpson

Nature 564, 410–414 (2018) Cite this article



It was a cold April day on a local trail in Nova Scotia during a hike. Sporadic sampling is a known professional hazard in protistology, since all we really need for an explosion of interesting things to look at is a few millilitres of water or dirt.



Yana soaked the soil in distilled water for about a month, periodically observing what was living there'. One evening, small, ~20 μ m long, **lanceolate cells showed up that behaved not quite like any 'common' flagellate**. Hemimastigophora is placed outside of all established eukaryote supergroups

Unrooted phylogeny inferred from **351** genes and 61 taxa. Most of data comes from transcriptomes, *Hemimastix* and *Spironema* data from single cells!



Lax et al. Nature, 2018

Microbial predators form a new supergroup of eukaryotes

Denis V. Tikhonenkov ⊠, Kirill V. Mikhailov, Ryan M. R. Gawryluk, Artem O. Belyaev, Varsha Mathur, Sergey A. Karpov, Dmitry G. Zagumyonnyi, Anastasia S. Borodina, Kristina I. Prokina, Alexander P. Mylnikov, Vladimir V. Aleoshin & Patrick J. Keeling

Nature 612, 714–719 (2022) Cite this article





"They have different behaviours. One group eats its prey. It literally just swallows cells whole. The other has a tooth — they'll gang up in a swarm and nibble it to death" – I prefer to call them nibblers and lions

Provora at the tree of eukaryotes





320 orthologous gene groups for 69 species

Bayesian inference approach implemented in PhyloBayes and the maximum-likelihood approach of IQ-TREE.

Anything new in 2024?

Meteora!

Maximum likelihood phylogeny 70,471 sites across <mark>254 genes over</mark> 66 taxa LG + MAM60 + Γ model

Bars on the right indicate % coverage by gene (above) and by site (below).




Anything new in 2024?

Meteora!

Maximum likelihood phylogeny 70,471 sites across <mark>254 genes over</mark> 66 taxa LG + MAM60 + Γ model

Bars on the right indicate % coverage by gene (above) and by site (below).



Validating robustness of the *Meteora* + Hemimastigophora clade

- 1. An analysis that excluded three taxa identified as long-branching outliers -> maximal support for the *Meteora* + Hemimastigophora clade
- 2. Recoding the amino acid data into a reduced alphabet of 4 classes -> robustly supported *Meteora* + Hemimastigophora
- 3. Removal of the fastest evolving sites in 10% increments (fast-site removal analysis) -> *Meteora* + Hemimastigophora as maximally supported until 30% sites remaining
- 4. Random subsampling of 50% of the genes -> robust support for *Meteora* + Hemimastigophora. The phylogenetic signal for the *Meteora* + Hemimastigophora relationship is broadly distributed across genes.

Why do we need well resolved tree of eukaryotes?

Evolution of cellular systems Intracellular transport



https://www.nature.com/scitable/topicpage/ how-do-proteins-move-through-the-golgi-14397318/



b



Flagellar aparatus





Transmission electron microscope micrograph

Yubuki & Leander, 2013

Reconstruction of the cytoskeleton supporting the flagellar apparatus



Evolution of orgenelles gene content of mitochondial genome



Presence or absence of genes in mitochondrial genomes of various eukaryotes is shown by closed and open boxes, respectively.

Roger et al. 2017

With more data and well resolved tree of eukaryotes we've learned that the last common ancestor was complex eukaryotic cell



The long march of eukaryotes





- Comparative genomic analyses have reconstructed a complex last eukaryotic common ancestor.
- But how and in which order these complex eukaryotic features evolved in an Asgard archaea?

Roger et al. 2021

Metagenomics has revealed the vast diversity of Archaea, including the recently described Asgard superphylum





Phylogenomic analyses have placed the Asgard archaea as the closest prokaryotic relatives of eukaryotes





Genomic investigation of Asgard archaea showed that they carry several genes believed to be **eukaryotic specific**

			Ribosomal proteins	Informational proteins	Trafficking machinery	Cytoskeleton	Ubiquitin system
5 4 4 3 3 2 2 1 5 1		———— Eukarya	 L41e L14e L13e S25e S30e L13e L13e L13e L28e 	 ELF1 RNA pol α* RPB8 or RpoG RPC34 TOPOIB DNA pol ε 	 ESCRT-I: VPS28 ESCRT-I: steadiness box domain⁴ ESCRT-II: EAP30 domain ESCRT-II: VPS22 ESCRT-II: VPS22/24/46⁵ ESCRT-III: VPS22/24/60⁶ ESCRT-III: VPS22/24/60	 Actin¹ ARP 2/3 complex, subunit 4 Profilin Gelsolin domain protein Tubulin[#] 	 Ubiquitin-related modifier 1 Ubiquitin domain protein Ub-activating protein E1 Ub-conjugating protein E2 Putative E3 UFM1 domain (Ub-like protein) Putative deubiquitin protein
		Lokiarchaeota	$\bigcirc \bullet \bullet \bullet \bullet \bullet \bullet \bullet \circ \bigcirc \bigcirc$	000000	00000000000	00000	0000000
		Odinarchaeota	$\bigcirc \bigcirc $	000000	000000000000000000000000000000000000000		0000000
		Thorarchaeota	$\bigcirc \bigcirc $	000000	\bigcirc	00000	0000000
		——— Heimdallarchaeota	$\bigcirc \bullet \bullet \bullet \bullet \bullet \bigcirc \bullet \bullet$	000000	000000 00000	00000	0000000
			0000000000	000000	0000 • 00000	00000	0000000
		- Crenarchaeota	000000000	000000	0000 • 00000	00000	0000000
			000000000	000000	0000 0 00000	00000	0000000
		Bathyarchaeota	$\bigcirc \bigcirc $	000000	0000 • 00000	00000	0000000
		- 'Aigarchaeota'	$\bigcirc \bigcirc $	000000	0000 • 00000	00000	0000000
		Thaumarchaeota	0000	\bigcirc	0000 • 00000	00000	0000000
		Euryarchaeota	•••••	000000	0000 0 00000	00000	0000000

4 3 2

Inference and reconstruction of the heimdallarchaeial ancestry of eukaryotes



Hod in Norse mythology, is a blind god, associated with night and darkness.

Received: 23 April 2021

Accepted: 10 May 2023

Eme et al. 2023

Expandind genomic diversity of Asgard archaea with metagenomics

63 new Asgard archaeal metagenome-assembled genomes (MAGs) from samples obtained from 11 locations around the world



Eme et al. 2023

Exploring new phylogenomic markers

- typically ribosomal proteins have been used in phylogenomic analyses of archaea and eukaryotes (RP56), however they might co-evolve and therefore might contribute to phylogenetic artefacts
- independent new marker dataset comprising 57 proteins of archaeal origin in eukaryotes
- the NM57 proteins are mostly involved in informational, metabolic and cellular processes, but do not include ribosomal proteins
- they are longer and therefore putatively more phylogenetically informative compared with the RP56 markers
- broader functional distribution of NM57 markers is less likely to cause phylogenetic reconstruction artefacts

Eukarya as sister to Hod





BI based on 278 archaeal taxa, 15,733 AA positions The concatenation was SR4-recoded and analysed using the CAT+GTR model Schematic representation of the shift in the position of eukaryotes (grey branches) in ML and BI analyses under different treatments Eme et al. 2023

Metagenomic data to the rescue

If MAGs can be recovered from metagenomic data, they can be used as a source of data for phylogenomics in a similar manner as genomes (with higher risk of contamination)





The genomic analysis of microbial communities by extraction and sequencing of their DNA, which allows studying organsims directly in their natural environment

Metagenome of lake

Protists 10%

Carter ...

*

Prokaryotes 90%

How to obtain eukaryotic genomes from metagenomes?





Karlicki et al. *Bioinformatics*, 2022

Other tools - Eukfinder

With sufficient reads, Eukfinder efficiently assembles high-quality near-complete nuclear and mitochondrial genomes from diverse Blastocystis subtypes from metagenomic data without the aid of a reference genome.



Zhao et al. 2023 BiorXiv



How the eukaryotic **cell complexity** evolved? What drives **diversification of eukaryotes**?

Endosymbiosis

plastids' evolution



De Vries et al. 2016

Trends in Plant Science



McGrath 2020



Plastids and cyanobacteria have been recovered repeatedly as a monophyletic group pointing to the cyanobacterial origin of primary plastids

0.10

Rodriguez-Ezpeleta et al. (2005) Curr Biol

Cyanobacteria which became a chloroplast

- *Gloeomargarita lithophora* is the closest extant cyanobacteria to plastids
- old cyanobacterial lineage
- freshwater cyanobacteria (!)





Ponce-Toledo et al., 2017

Secondary endosymbioses in several lineages of microbial eukaryotes



-----**** Green Algae Cand **Red Algae** Glaucophytes Paulinella Euglenids Cryptomonads Haptophytes Ciliates Chlorarachniophytes **/** () Heterokonts Apicomplexa Dinoflagellates Dinophysis

Lepidodinium

Karenia

Kryptoperidinium

Keeling (2004) Am J Bot

Phylogenomics of plastids



secondary "red" endosymbionts

Sevcikova et al. 2015



Nuclear-encoded genes vs plastid-encoded genes



Cryptic serial endosybiosis



Burki et al. 2016

Plastid endosymbiosis



What are the initial steps of endosymbiosis? What are the steps in the transition from endosymbiont to organelle? Which comes first, gene transfer or cellular fixation? A targeting-ratchet model for the endosymbiotic origin of plastids Keeling, 2013





kleptoplasty

transient association between host and endosymbiont, which might resemble the initial steps of the establishing endosymbiosis






Kleptoplasty is widespread among eukaryotes





Rapaza viridis – mixotrophic euglenid

feeds on a specific strain of green algae, Tetraselmis sp. anine. ARCHAEPLASTIDA Chloroplas

AMORPHEA OPISTHOKONTA CRUM Glaucophyta SITAAH AAST

Rapaza viridis

Yamaguchi, Yubuki & Leander (2013)



Karnkowska et al. 2023

Rapaza viridis phylogeny





Modified from Leliaert *et al., Crit. Rev. Plant Sci.* 31:1-46 (2012) updated 25 Oct 2013



Tetraselmis sp.

Rapaza viridis (Euglenophyceae) and its kleptoplast as a model for studying early stages of endosymbiosis

Rapaza viridis

Which comes first, gene transfer or cellular fixation?





Genes encoding chloroplast proteins are transferred to nuclear genome via endosymbiotic gene transfer (EGT).

Plastid proteins encoded in nuclear genome are **targeted** to the plastids.

Searching for plastid-targeted proteins in Rapaza transcriptome

Sequencing and RNAseq analysis



Plastid proteins identification



Searching for translocon components



Euglena translocon components Zahonova et al. 2018

274 proteins

Translocon components of R. viridis



TO(

64

34

159

TIC

21

20

62

110

40)

32

-55

ClpC ,

More components of translocon than in core Euglenophytes Diverse origin of genes encoding translocon components

Karnkowska et al. (2023)



Are there any plastid proteins encoded in nuclear genome of *Rapaza* targeted to kleptoplasts?

Searching for plastid-targeted proteins in Rapaza transcriptome

Sequencing and RNAseq analysis



Plastid proteins identification

274 proteins



Targeting signal prediction Yes!

signal peptide protein transit peptide

Phylogenetic analyses of candidate proteins



Searching for plastid-targeted proteins in Rapaza transcriptome

Phylogenetic analyses of candidate proteins



60% originated from green algae

40% from other evolutionary lineages

proteins derived from different lineages as the kleptoplast indicate that transfer of genes preceded establishment of the endosymbiont

origin of plastids via kleptoplasty gene transfer comes before cellular fixation

> transfer of genes encoding plastid proteins to the nuclear genome

targeting of plastid proteins to the kleptoplast



Karnkowska et al. (2023)

mitochondria



The genome sequence of *Rickettsia prowazekii* and the origin of mitochondria



Andersson (1998), Nature

Alphaproteobacterial phylogenomics and origin of mitochondria

- The tree of alphaproteobacteria is not resolved
- Number of genes of alphaproteobacterial origin in mitochondrial genomes is limited (not much phylogenetic signal)



Environmental MAGs allow to expand the taxa sampling





D

Munoz-Gomez et al. 2022

Expanded dataset favour mitochondria as sister to known Alphaproteobacteria



- (A) Phylogenetic tree for the *Alphaproteobacteria* and mitochondria derived from a site-heterogeneous analyses of an untreated dataset.
- (B) Phylogenetic tree for the *Alphaproteobacteria* and mitochondria derived from a site-heterogeneous analysis of a dataset from which 50% of the most compositionally heterogeneous sites had been removed. Munoz-Gomez et al. 2022

Which came **first** nucleus or mitochondrion?

Archezoa hypothesis

Thomas Cavalier-Smith, 1989





Embley and Martin (2006) Nature

Archezoa

- early branching eukaryotes
- lack of introns
- no sexual reproduction
- lack of peroxisomes
- lack of Golgi apparatus
- lack of mitochondria



Trichomonas

Giardia

Microsporidia

Archamoebae

Archezoa posses mitochondria related organelles (MROs)



mitochondrion

hydrogenosome

mitosomes

Common origin of mitochondria and MROs

Cpn60 chaperonin

- mitochondrial chaperonins
- ATP transporters
- mitochondrial membrane transport proteins
- Fe-S cluster assembly proteins



cpn60 localization Jerlström-Hultqvist et al. 2013



Evolution of mitochondria

MROs are widespread at the tree of eukaryotes



Archaeplastida

Archezoa hypothesis rejected

- All "archezoa" possess:
- mitochondrial genes in nuclear genomes
- degenerate derivatives of mitochondria
- they do not group together on the modern tree of life

Common ancestor of all eukaryotes possessed mitochondria



mitochondrion

hydrogenosome

mitosomes

Does amitochondriate eukaryote exist?





cryptic mitochondria

Keeling, 2007

A eukaryotic cell with typical mitochondria



Oxidatvie phosphorylation, and other essential process such as Fe-S cluster formation takes part in mitochondria

Fe-S clusters are cofactors of essential enzymes

An anaerobic eukaryote with mitochondrion related organelles (MROs)



0

hydrogenosome

mitosome

Loss of oxidative phosphorylation in anaerobic eykaryotes

mitochondria in 'Excavata'





- found in the intestinal tracts of animals
- sexual reproduction debatable
- no peroxisomes
- no Golgi apparatus
- no mitochondria



Monocercomonoides - microaerophilic, commensal of animals





lack of mitochondria in Monocercomonoides under the microscope

genome analysis of *Monocercomonoides*



Monocercomonoides is less divergent than Parabasalids and Diplomonads

Karnkowska et al. (2016)

Searching for mitochondrial proteins

Mitochondrial outer membrane targeted proteins (TA) Proteins with mitochondrial localization signal β-barrel mitochondrial outer membrane proteins



Karnkowska et al. (2016)

mitochondrial membrane transport proteins






Lateral gene transfer (LGT) of SUF system



Karnkowska et al. Curr Biol (2016)

SUF proteins in Preaxostyla and representatives of bacteria



Bacillus subtilis (Firmicutes) Thermotoga marina (Thermotogae) Spirochaeta caldaria (Spirochaetes) Chloroflexus sp. Y-400-fl (Chloroflexi) Thiomonas intermedia (Proteobacteria) Escherichia coli (Proteobacteria) Mycobacterium tubercolosis (Actinobacteria) Owenweeksia hongkongensis (Bacteroidetes) Sulfolobus tokodaii (Crenarcheota) Haloferax volcanii (Euryarcheota)



SUF system is widespread in Preaxostyla





Vacek et al. 2018

Monocercomonoides an amitochondriate eukaryote



Lateral gene transfer (LGT) of bacterial genes encoding sulfur mobilization pathway (SUF) and loss of reduced mitochondrion

Loss of a mitochondrial organelle

LGT of SUF system resulted in relocation of the pathway to the cytosol



Endosymbiosis can be undone!

Karnkowska et al. (2016)