# *Incongruence*

**Antonis Rokas**
***Department of Biological Sciences, Vanderbilt University***
**http://www.rokaslab.org**
**@RokasLab**

# *Brief Bio*

| | | | |
|---|---|---|---|
| **93-98** | B.Sc. | **Biology** | Univ. of Crete, Greece (Advisor: L. Zouros) |
| **98-01** | Ph.D. | **Evolutionary Ecology** | Edinburgh Univ., Scotland (Advisor: G. Stone) |
| **02-05** | PostDoc | **Evolutionary Genomics** | Univ. Wisconsin-Madison (Advisor: S. Carroll) |
| **05-07** | Res. Scientist | **Fungal Genomics** | Broad Institute |
| **07-now:** Faculty | | **Evolutionary Biology** | Vanderbilt University |

*http://www.rokaslab.org*
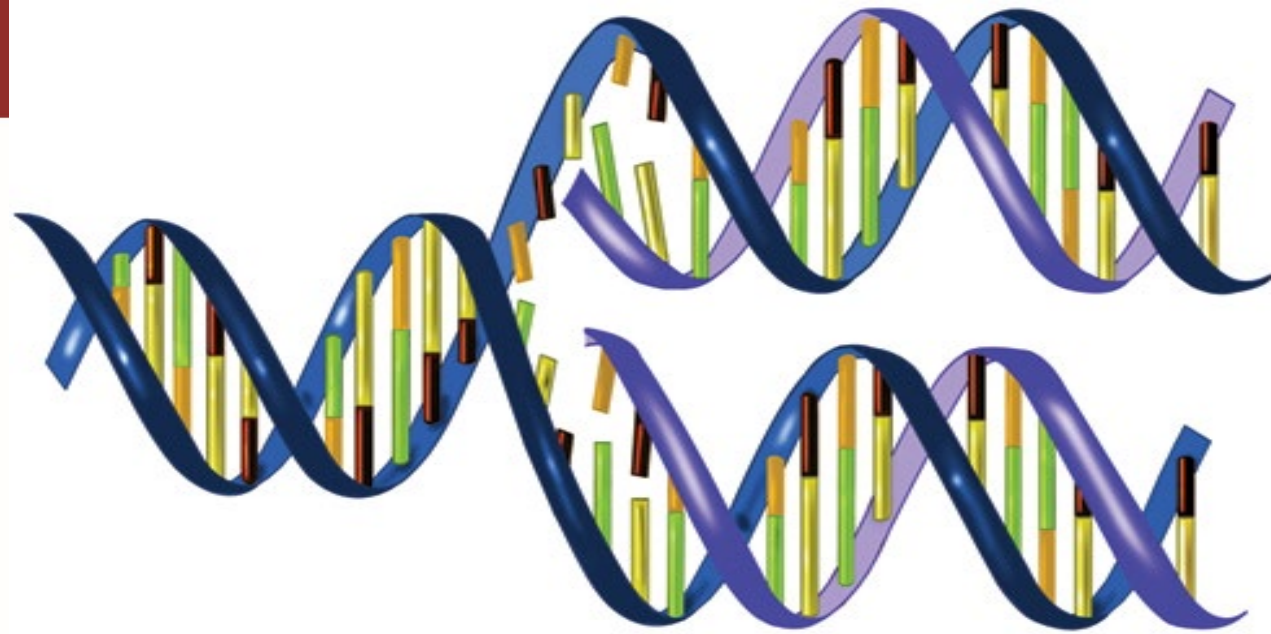
Antonis (me!) & Scott Handley

"The genome is, it's a fossil record; the genome is a landscape; the genome is a whole geography of distributions. […] you might think the genome's just a boring string of letters [...] The genome is a storybook that's been edited for a couple of billion years, and you could take it to bed, like *A Thousand and One Arabian Nights*, and read a different story, in the genome, every night."
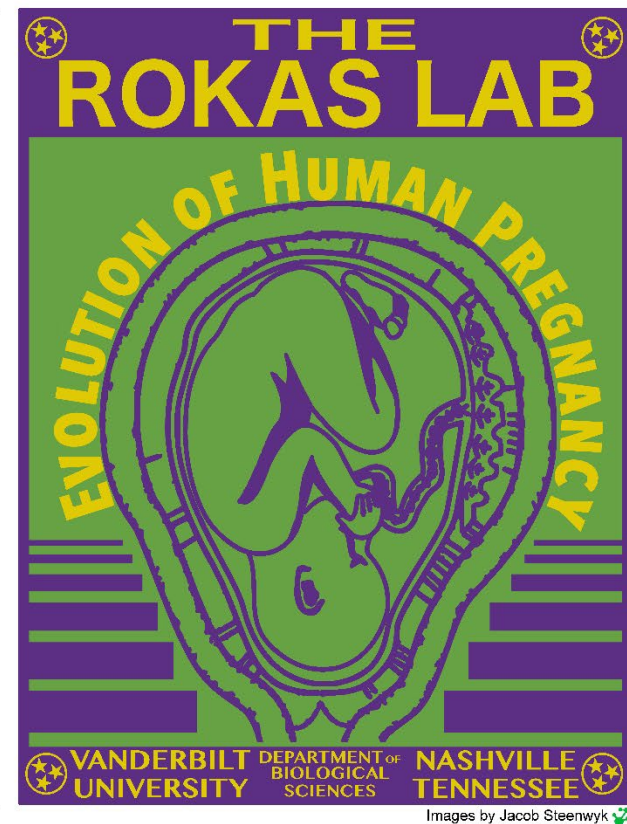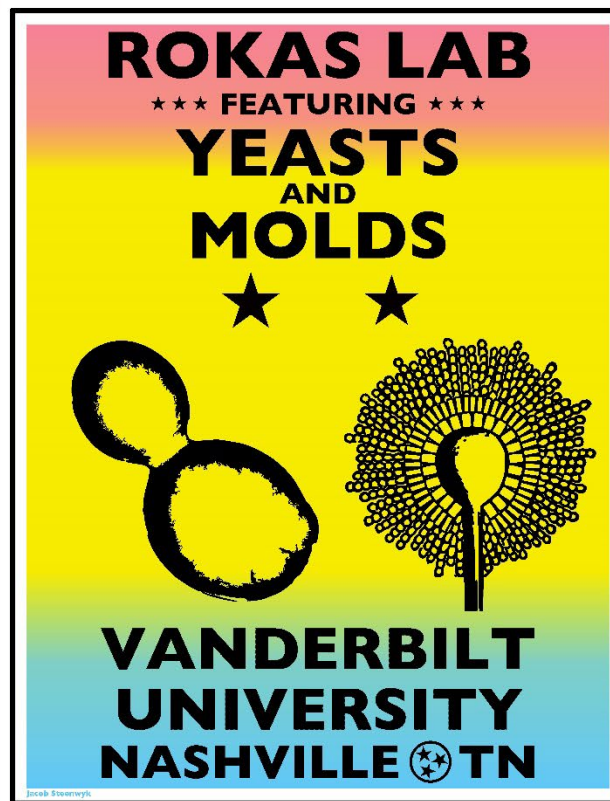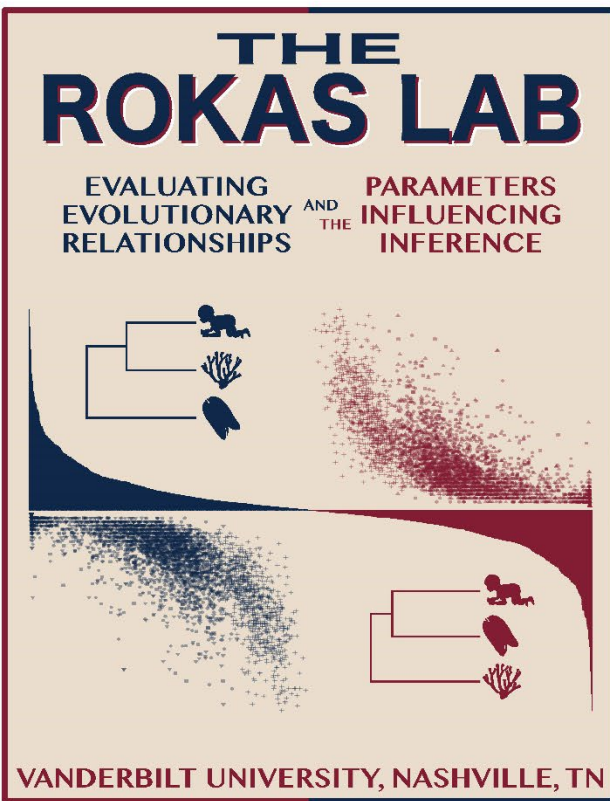
Eric Lander

**We study the DNA record to gain insight into evolutionary patterns and processes using computational and experimental approaches**
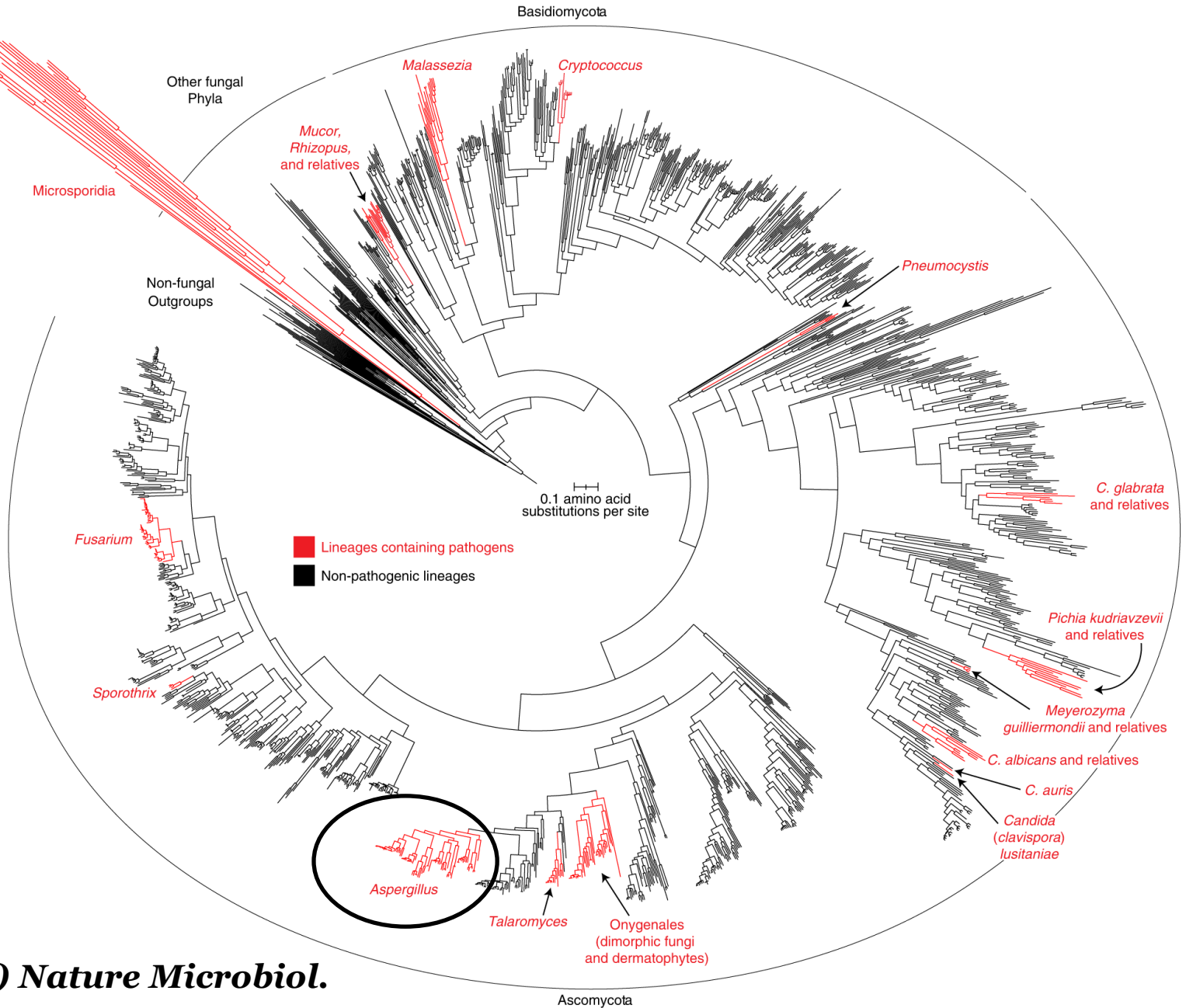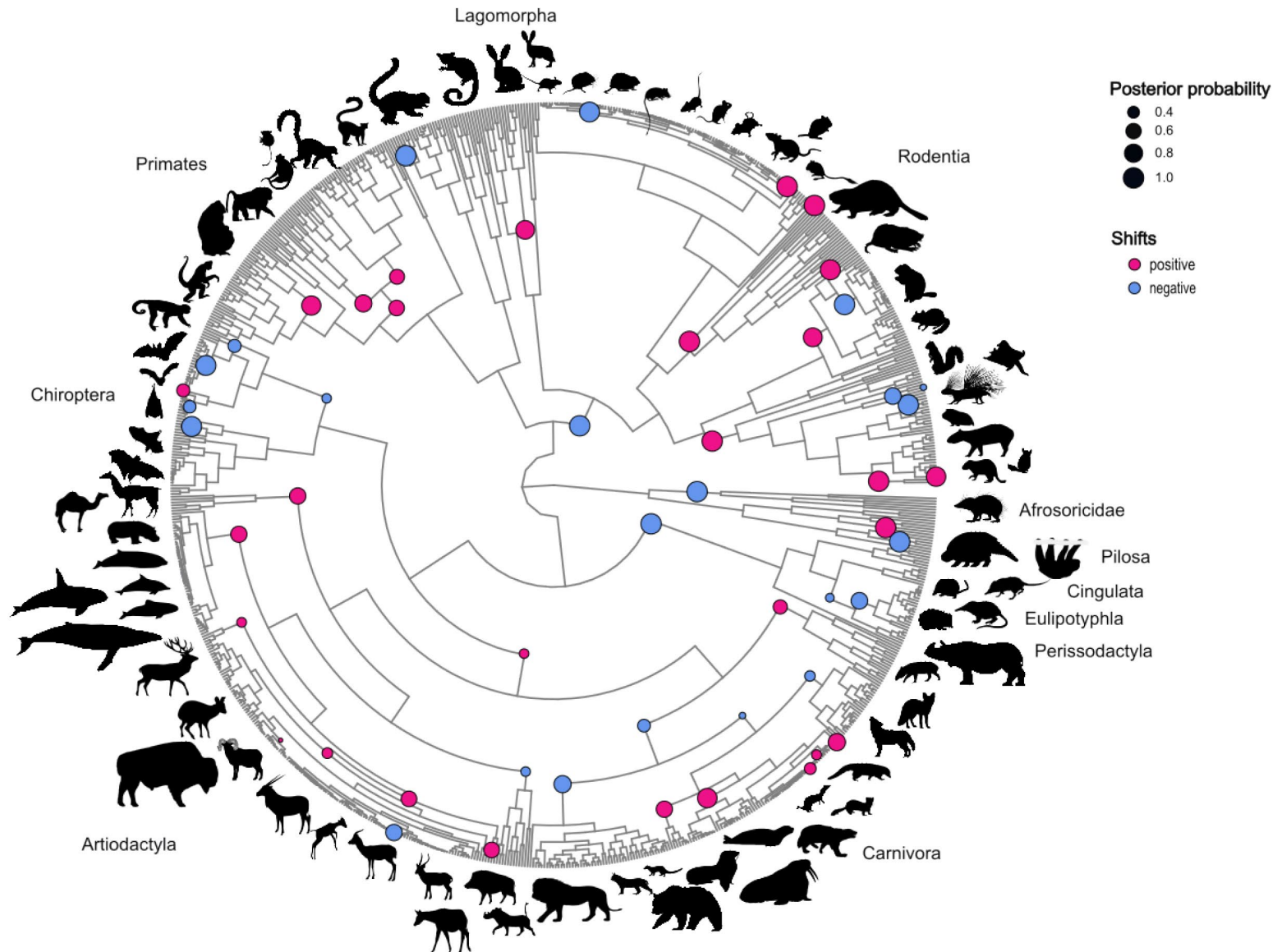
**Phylogenomics**
(NSF)

**The molecular foundations of the fungal lifestyle**
(NSF & NIH)

**The evolution of mammalian pregnancy**
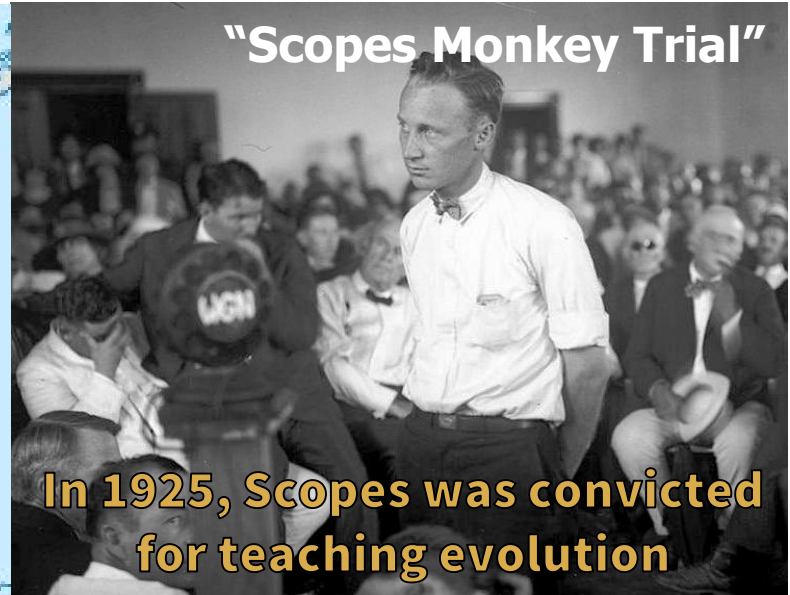(BWF & March of Dimes)

# The repeated evolution of fungal pathogenicity



*Rokas (2022) Nature Microbiol.*

# Evolution of variation in gestation length in mammals



*Danis & Rokas (2023) bioRxiv*

# ...But I'm also an evolutionist in Tennessee (USA)



"Scopes Monkey Trial"

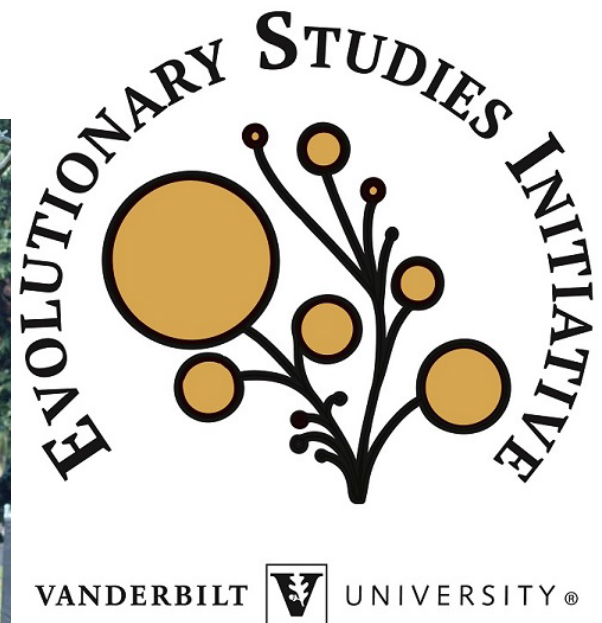In 1925, Scopes was convicted for teaching evolution

Foes of science faced ridicule at the Scopes trial. We're paying the price 95 years later.

The Washington Post

Opinion by **Max Boot**
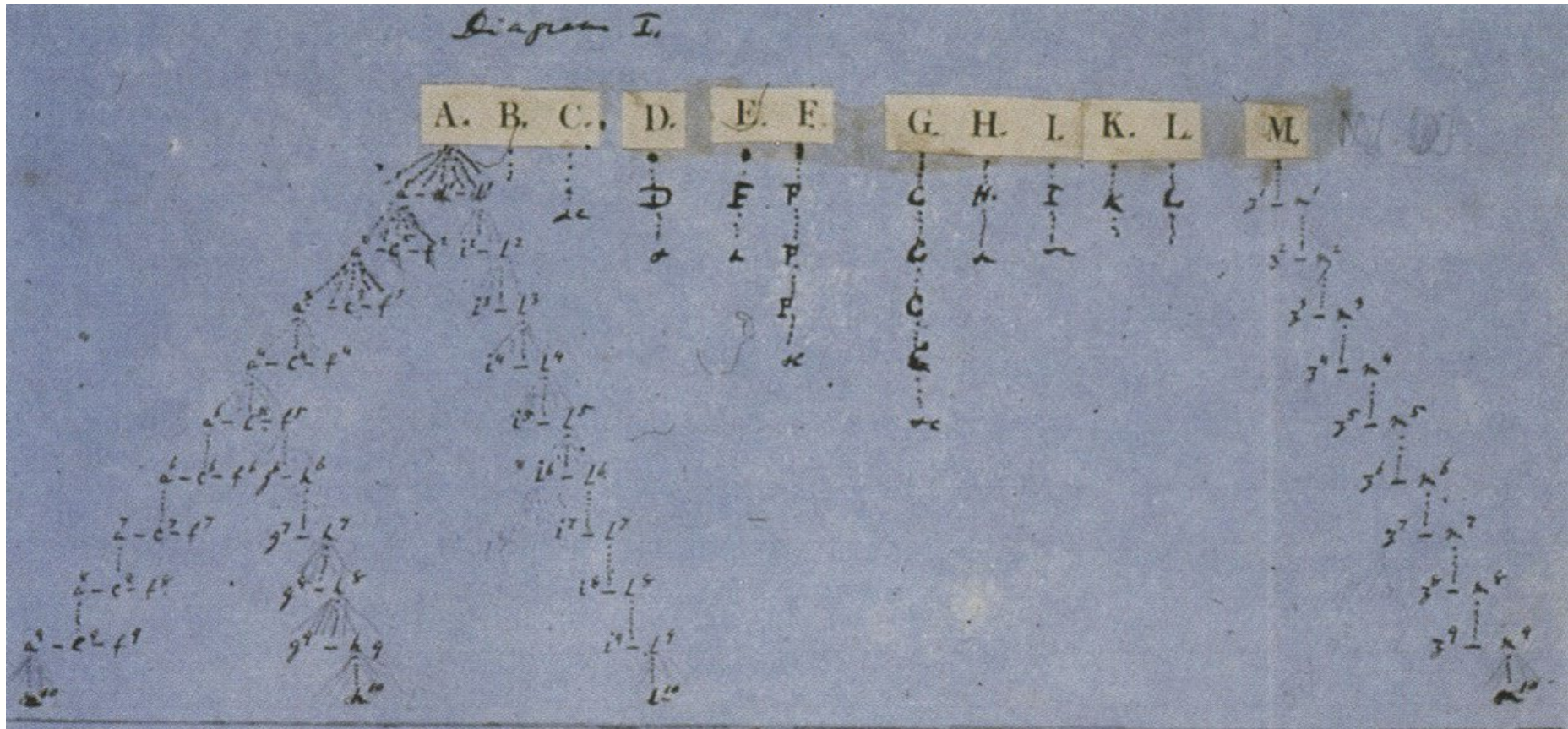Columnist

July 8, 2020 at 1:01 p.m. CDT

**www.vanderbilt.edu/evolution/**

❖ **Incongruence and its causes**

-------------------- **Coffee Break** --------------------------

❖ **Handling incongruence in phylogenomic data**
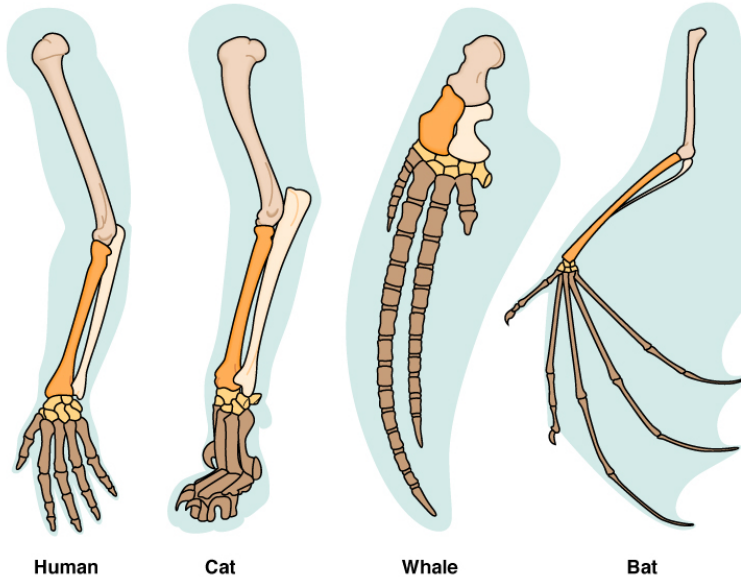
# Darwin's tree

and instinct as the summing up of many contrivances, each useful to the possessor, nearly in the same way as when we look at any great mechanical invention as the summing up of the labour, the experience, the reason, and even the blunders of numerous workmen; when we thus view each organic being, how far more interesting, I speak from experience, will the study of natural history become!
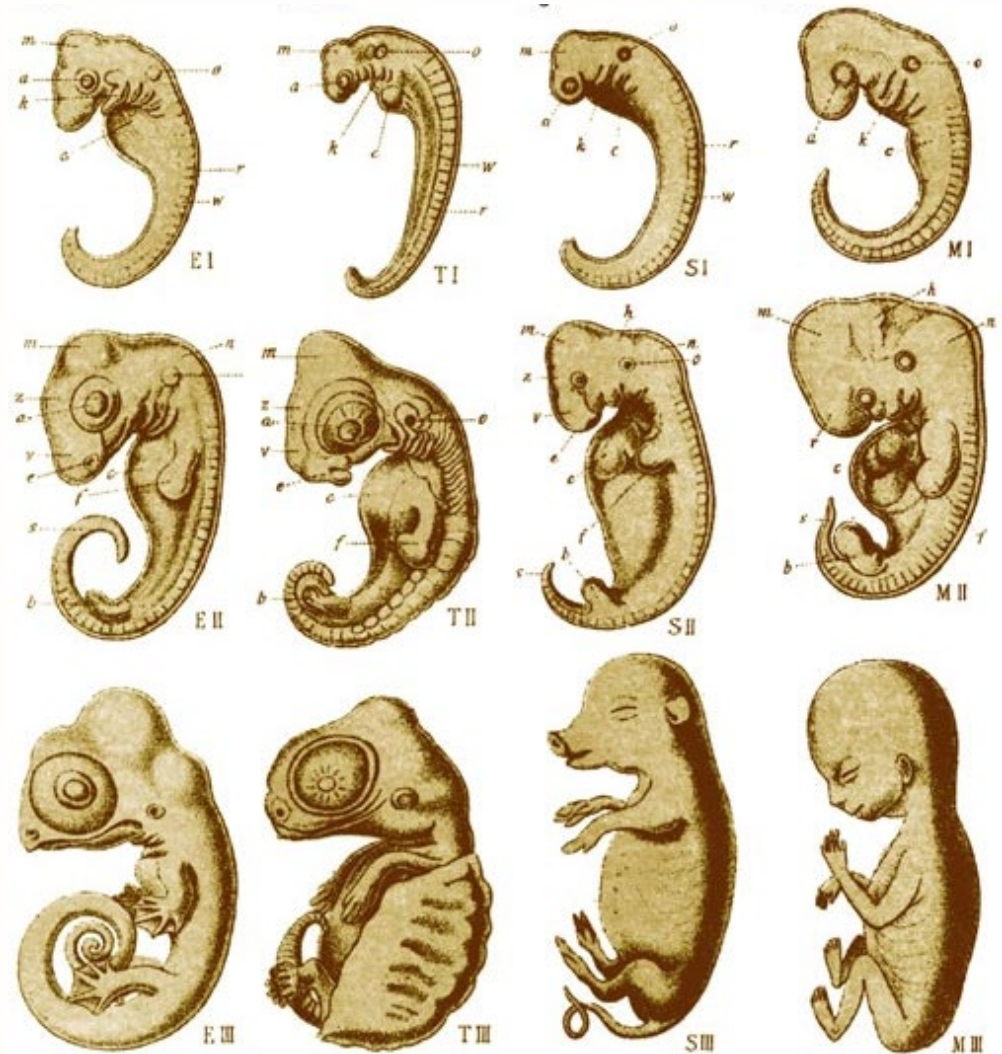
A grand and almost untrodden field of inquiry will be opened, on the causes and laws of variation, on correlation of growth, on the effects of use and disuse, on the direct action of external conditions, and so forth. The study of domestic productions will rise immensely in value. A new variety raised by man will be a far more important and interesting subject for study than one more species added to the infinitude of already recorded species. Our classifications will come to be, as far as they can be so made, genealogies; and will then truly give what may be called the plan of creation. The rules for classifying will no doubt become simpler when we have a definite object in view. We possess no pedigrees or armorial bearings; and we have to discover and trace the many diverging lines of descent in our natural genealogies, by characters of any kind which have long been inherited. Rudimentary organs will speak infallibly with respect to the nature of long-lost structures. Species and groups of species, which are called aberrant, and which may fancifully be called living fossils, will aid us in forming a picture of the ancient forms of life. Embryology will reveal to us the structure, in some degree obscured, of the prototypes of each great class.

When we can feel assured that all the individuals of the same species, and all the closely allied species of most genera, have within a not very remote period de-

Human    Cat    Whale    Bat

©1999 Addison Wesley Longman, Inc.

**St. George Jackson Mivart**

# Disagreement between phylogenies

**1865: SPINAL COLUMN**

- Homo
- Simia
- Troglodytes
- Hylobates
- Loris
- Nycticebus
- Perodicticus
- Arctocebus
- Other Primates

**1867: LIMBS**

- Homo
- Simia
- Hylobates
- Troglodytes
- Semnopithecinae
- Cynopithecinae
- Ateles
- Lagothrix
- Cebus
- Mycetes
- Other Primates

**St. George Jackson Mivart**

Darwin Correspondence Project letters 7718A & 7170

**"From the same facts, opposite conclusions are drawn; facts of the same kind will take us no further. […] Need we waste more effort in these vain and sophistical disputes. If facts of the old kind will not help, let us seek facts of a new kind."**
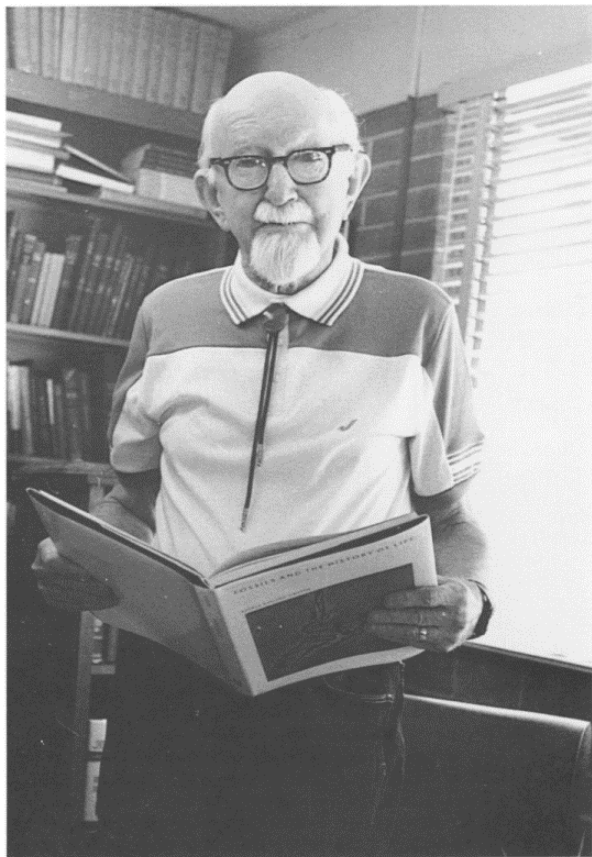
**William Bateson (1894)**
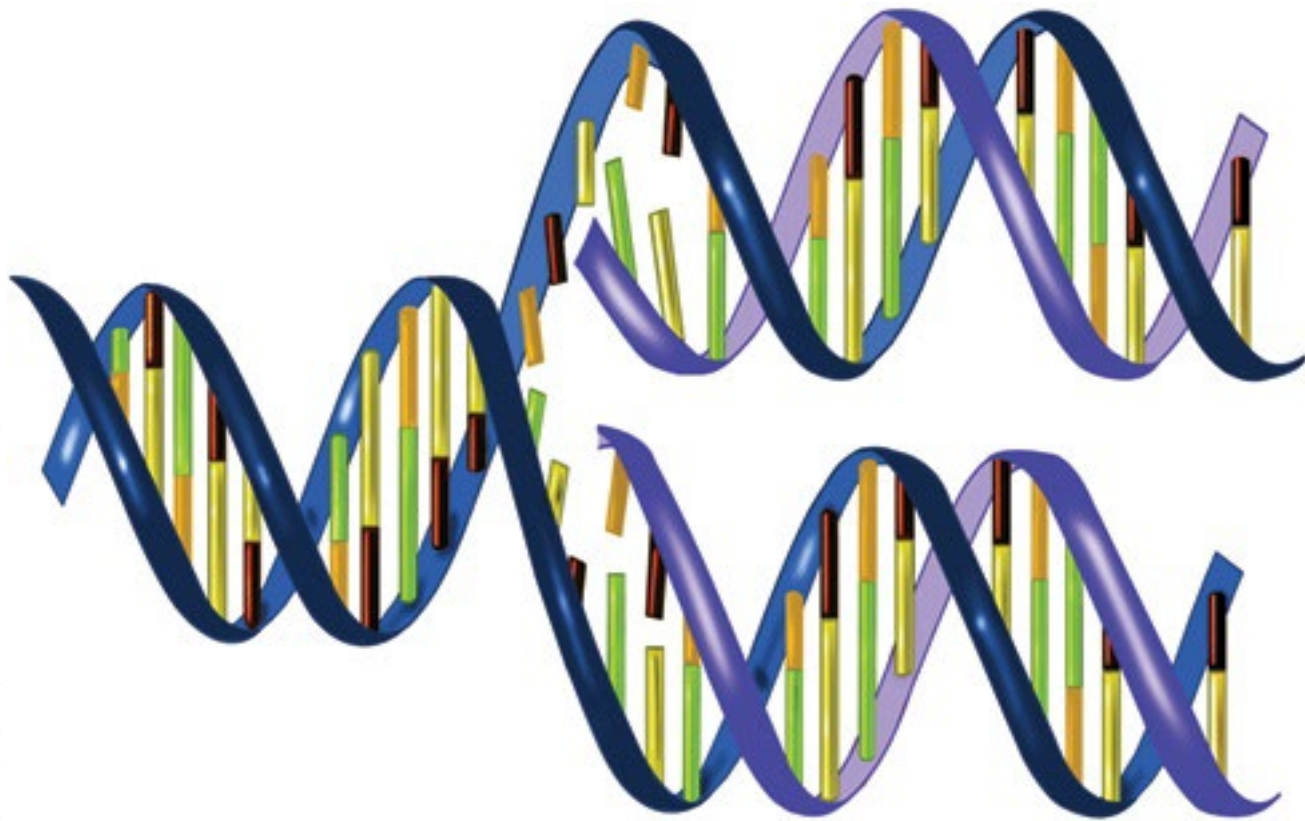
***Materials for the Study of Variation***



Courtesy of American Philosophical Society, Curt Stern Papers.
Noncommercial, educational use only.

George Gaylord Simpson

"The stream of heredity makes phylogeny; in a sense, it is phylogeny. Complete genetic analysis would provide the most priceless data for the mapping of this stream"

G. G. Simpson, 1945

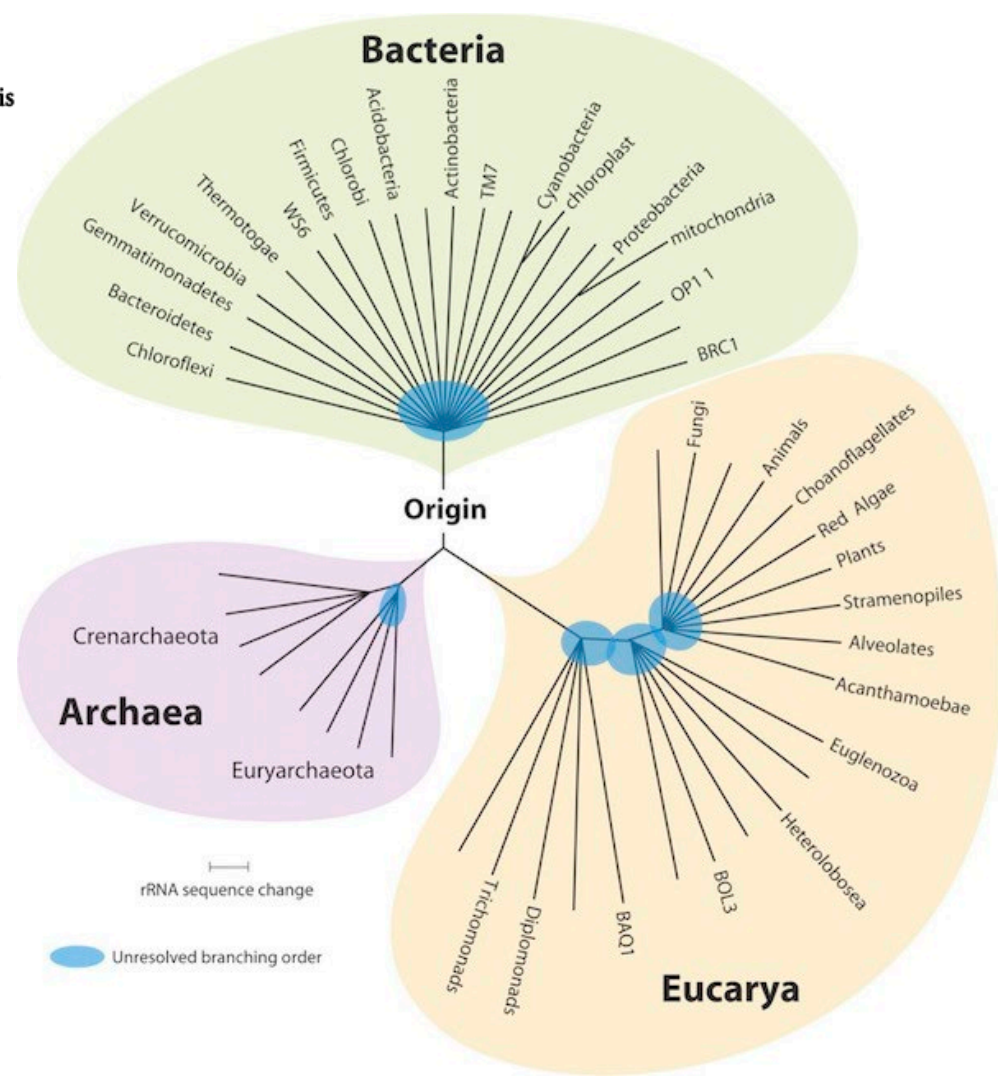# Phylogenetic structure of the prokaryotic domain: The primary kingdoms

(archaebacteria/eubacteria/urkaryote/16S ribosomal RNA/molecular phylogeny)
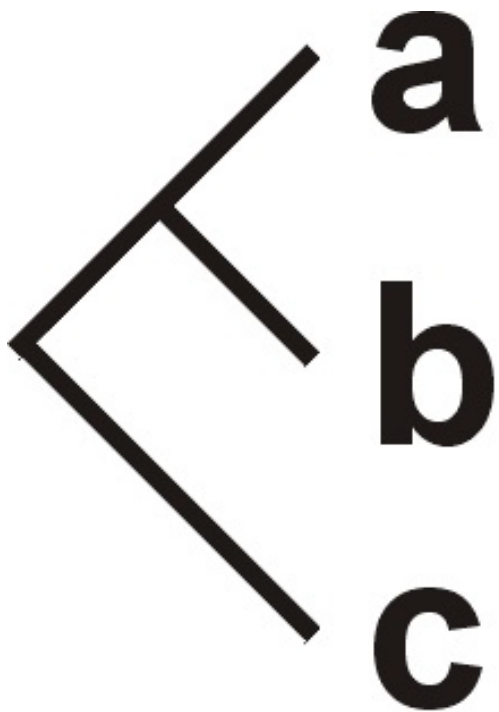
CARL R. WOESE AND GEORGE E. FOX*

Department of Genetics and Development, University of Illinois, Urbana, Illinois

ABSTRACT        A phylogenetic analysis based upon ribosomal RNA sequence characterization reveals that living systems represent one of three aboriginal lines of descent: (*i*) the eubacteria, comprising all typical bacteria; (*ii*) the archaebacteria, containing methanogenic bacteria; and (*iii*) the urkaryotes, now represented in the cytoplasmic component of eukaryotic cells.

*Incongruence / conflict / discordance*

a

b

c

**Gene X**

a

c

b

**Gene Y**

**Species phylogeny?**

# The genomics revolution



## Cost per Raw Megabase of DNA Sequence

Moore's Law

National Human Genome Research Institute

genome.gov/sequencingcosts

PhD    Faculty

# A systematic evaluation of single gene phylogenies

*S. cerevisiae*     *S. bayanus*

*S. paradoxus*     *S. castellii*

*S. mikatae*     *S. kluyveri*

*S. kudriavzevii*     *Candida glabrata*

S. cerevisiae
S. paradoxus
S. mikatae
S. kudriavzevii
S. bayanus
S. castellii
S. kluyveri
C. albicans

100/100
100

100/100
100

100/100
100

100/100
100

100/100
100

**ML / MP on nt**
**MP on aa**

# The use of many genes eliminates incongruence

ML / MP on nt
MP on aa

*Rokas et al. (2003) Nature*

# *The dawn of the phylogenomics era*

LETT

LETT

## Phylogenomic Analysis Resolves the Interordinal Relationships and Rapid Diversification of the Laurasiatherian Mammals

XUMING ZHOU, SHIXIA XU, JUNXIAO XU, BINGYAO CHEN, KAIYA ZHOU, AND GUANG YANG*

*Jiangsu Key Laboratory for Biodiversity and Biotechnology, College of Life Sciences, Nanjing Normal University, Nanjing 210046, China;
*Correspondence to be sent to: Jiangsu Key Laboratory for Biodiversity and Biotechnology, College of Life Sciences, Nanjing Normal University, Nanjing 210046, China; E-mail: gyang@njnu.edu.cn.

## Resolving the evolutionary relationships of molluscs with phylogenomic tools

nature

LETTERS

Stephen A. Smith[1,2], Nerida G. Wilson[3,4], Freya
Gonzalo Giribet[5] & Casey W. Dunn[1]

## Resolving Arthropod Phylogeny: Exploring Phylogenetic Signal within 41 kb of Protein-Coding Nuclear Gene Sequence

JEROME C. REGIER,[1] JEFFREY W. SHULTZ,[2] AUSTEN R. D. GANLEY,[3,6] APRIL HUSSEY,[1] DIANE SHI,[1]
BERNARD BALL,[3] ANDREAS ZWICK,[1] JASON E. STAJICH,[3,7] MICHAEL P. CUMMINGS,[4] JOEL W. MARTIN,[5]
AND CLIFFORD W. CUNNINGHAM[3]

## Toward Resolving t Tree: The Phylogen of Jakobids and Cercozoans

Yeast

An

## Toward Resolving Priors

## Prion-Like Proteins in the Fungal Kingdom

Edgar M. Medina · Gary W. Jones ·
David A. Fitzpatrick

## Towards

Renae C. Pratt,* Gillian C. Gibb,* Mary Morgan-Richards,* Matthew J. Phillips,† Michael
D. Hendy,* and David Penny*

Samuli Lehtonen

Department of Biology, U

*Allan Wilson Centre for Molecular Ecology and Evolution, Massey University, Palmerston North, New Zealand; and †Centre for Macroevolution and Macroecology, School of Botany and Zoology, Australian National University, Canberra ACT, Australia

# Have we eliminated incongruence?

## Biological factors

They lead to gene trees whose histories may differ from each other and from the species tree. Known factors include **stochastic lineage sorting**, **hidden paralogy**, **horizontal gene transfer**, **recombination, hybridization / introgression,** and **natural selection**

## Analytical factors

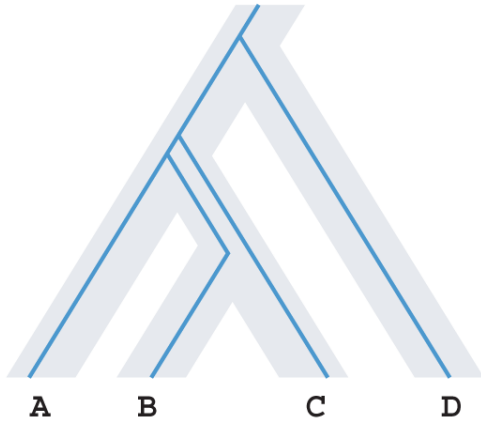They lead to failure in accurately inferring a gene tree; these can be either due to **stochastic error** (e.g., insufficient number of genes or taxa), **systematic error** (e.g., observed data deviate from model assumptions), or **treatment error** (e.g., excessive trimming)
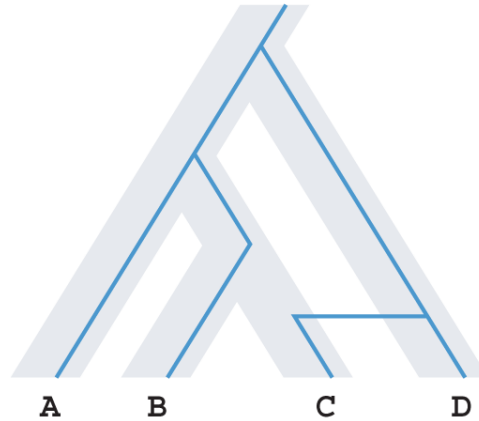
# Biological factors



Lineage sorting

Horizontal transfer

Hybridization / Introgression

Recombination

Duplication & Loss

Selection

**A**      **B**      **C**

**Genes' histories
can differ from
species ones**

**Speciation
of B and C
lineages**

**Splitting of A and
B alleles**

**Speciation
of A and BC
lineages**

**Splitting of AB
and C alleles**

# Lineage sorting in primates



**Informative Sites**

8,561 / 11,293
(~76%)

1,302 / 11,293
(~11.5%)

1,430 / 11,293
(~12.5%)

*Patterson et al. (2006) Nature*

**Exchange of genes between organisms other than through reproduction**

**A clade of yeasts acquired the enterobactin operon from Enterobacteria – organisms from both lineages co-occur in insect guts, where iron is a growth-limiting factor**

# Balancing selection can maintain "trans-species polymorph-isms", in which the alleles are more ancient than the species

**Best example: alleles at loci of the MHC – they have been retained by selection because they confer resistance to infection**

**Certain human MHC alleles appear to have diverged more than 65 million years ago (these alleles witnessed the extinction of dinosaurs!!!)**



Legend:
- MHC
- ABO
- ZC3HAV1
- TRIM5
- LAD1

*Pongo pygmaeus*   *Gorilla gorilla*   *Pan paniscus*   *Pan troglodytes*   *Homo sapiens*

**Phylogeny of *prestin*, a gene involved in echolocation**

# Gene duplication and loss

*S. eubayanus* was discovered in 2011 – until then, *S. bayanus* was thought to be a "pure" species

*S. cerevisiae* – *S. paradoxus* divergence ≈ human – mouse divergence
*S. cerevisiae* – *S. uvarum* divergence ≈ human – chicken divergence

**a** Taxon selection

Taxon 1   Taxon 2   Taxon 3   Taxon 4

Contributor of incongruence

- Insufficient taxon sampling
- Insufficient locus sampling
- Fast-evolving lineages
- Rogue taxa
- Outgroup choice

**b** Orthology inference

Taxon 1   Taxon 2   Taxon 3   Taxon 4

- Sequence length biases
- Erroneous orthologue inference (hidden paralogy and orthology)

**c** Alignment and site trimming

| Taxon 1 | MPSQP---VQ ... |
| Taxon 2 | MPSQP---VQ ... |
| Taxon 3 | MPSQPYVQVQ ... |
| Taxon 4 | M--QPYVQVQ ... |

| Taxon 1 | MGH--YEEN ... |
| Taxon 2 | M--LRY--- ... |
| Taxon 3 | MGHL-YEEN ... |
| Taxon 4 | M--LRYEEN ... |

| Taxon 1 | MSP-VKG-PR ... |
| Taxon 2 | MSPTVK--PR ... |
| Taxon 3 | MSPTVKGIPR ... |
| Taxon 4 | MS---KGI-R ... |

- Misalignment
- Excessive trimming
- Inappropriate recoding

**d** Selection of substitution model

Site-homogeneous model   Site-homogeneous with partitioning   Site-heterogeneous model

Taxon 1
Taxon 2
Taxon 3
Taxon 4

- Long-branch attraction
- Model misspecification
- Inadequate model complexity

**e** Method of tree inference

Concatenation   Coalescence

- Irreproducibility
- Single-locus accuracy

```
10 50

Cow      MAYPMQLGFQ DATSPIMEEL LHFHDHTLMI VFLISSLVLY IISLMLTTKL
Carp     MAHPTQLGFK DAAMPVMEEL LHFHDHALMI VLLISTLVLY IITAMVSTKL
Chicken  MANHSQLGFQ DASSPIMEEL VEFHDHALMV ALAICSLVLY LLTLMLMEKL
Human    MAHAAQVGLQ DATSPIMEEL ITFHDHALMI IFLICFLVLY ALFLTLTTKL
Loach    MAHPTQLGFQ DAASPVMEEL LHFHDHALMI VFLISALVLY VIITTVSTKL
Mouse    MAYPFQLGLQ DATSPIMEEL MNFHDHTLMI VFLISSLVLY IISLMLTTKL
Rat      MAYPFQLGLQ DATSPIMEEL TNFHDHTLMI VFLISSLVLY IISLMLTTKL
Seal     MAYPLQMGLQ DATSPIMEEL LHFHDHTLMI VFLISSLVLY IISLMLTTKL
Whale    MAYPFQLGFQ DAASPIMEEL LHFHDHTLMI VFLISSLVLY IITLMLTTKL
Frog     MAHPSQLGFQ DAASPIMEEL LHFHDHTLMA VFLISTLVLY IITIMMTTKL
```

**Long branch attraction**

male

Halictophagidae (Strepsiptera)

# The Strepsiptera Problem



p distance  HKY85  HKY85+GAMMA

**Multiple sequence alignment**

**Alignment trimming**

**Character recoding**

**Irreproducibility**

**…**

CLUSTAL W

MUSCLE

T-COFFEE

DIALIGN 2

MAFFT

DCA

PROBCONS

# Genes yielding irreproducible trees are less informative

We have not eliminated incongruence but now better understand the biological and analytical contributing factors and their impact; sources vary by lineage and depth

# Inference at shallow depths is easier:

analytical factors

## biological factors

# The evolution of hominids

# The phylogeny of primate genera

**Nomascus leucogenys**



NLE

**Hoolock leuconedys**



HLE

**Symphalangus syndactylus**



SSY

**Hylobates pileatus**



HPI

**Hylobates moloch**



HMO

# Inference in deep time can be more challenging:

## analytical factors

## biological factors

Gastropoda

Bivalvia

Scaphopoda

**Bilaterian animals**

**Ctenophores**

**Porifera**

❖ **Incongruence and its causes**

-------------------- **Coffee Break** --------------------------

❖ **Handling incongruence in phylogenomic data**

# An expanded yeast data matrix

**Yeast Gene Order Browser (YGOB)**



**Candida Gene Order Browser (CGOB)**



*Saccharomyces* lineage

*Candida* lineage

**1,070 genes
23 taxa
no missing data**

*Byrne & Wolfe (2005) Genome Res.*          *Fitzpatrick et al. (2010) BMC Genom.*

**Yeasts**    **Animals**    **Plants**

Amino acid substitutions / site: 0.0, 0.5, 1.0

*S. cerevisiae*, *S. uvarum*, *S. paradoxus*, *C. glabrata*, *K. lactis*, *Wy. anomalus*, *D. hansenii*, *C. albicans*, *B. bruxellensis*, *Y. lipolytica*, *L. starkeyi*

Human, Mouse, Python, Zebrafish, Lancelet, Starfish, Sea Urchin, Coral, Polyp, Sponge, Comb jelly, Roundworm

Thale cress (*A. thaliana*), Thale cress (*A. lyrata*), Mustard, Grape, Tomato, Apple, Rice, Orchid, Maize, Pine, Spikemoss, Green algae

***Saccharomyces*, *Candida*, *Kluyveromyces*, etc. are all polyphyletic genera**

# *The concatenation phylogeny is at least partly wrong*



- ❖ **5 genomic rearrangements that are uniquely shared by *S. cerevisiae* and *C. glabrata***

- ❖ **Much higher proportion of shared gene losses in *S. cerevisiae* and *C. glabrata***

- ❖ **Bias in the placement of *C. glabrata* as an outgroup of *S. cerevisiae* and *S. castellii***

# Gene trees are incongruent in most datasets



182 / 184

*Zhong et al. (2013) Trends Plant Sci.*

440 / 447

*Song et al. (2012) PNAS*

1,070 / 1,070

*Salichos & Rokas (2013) Nature*

14,536 / 14,536

*Jarvis et al. (2014) Science*

# The yeast phylogeny inferred by majority-rule consensus

**Gene Support Frequency (GSF): % of single gene trees supporting a given internode**

# New Methods to Calculate Concordance Factors for Phylogenomic Datasets

Bui Quang Minh (ID) ,[1,2] Matthew W. Hahn,[3,4] and Robert Lanfear*,[2]

[1]Research School of Computer Science, Australian National University, Canberra, ACT, Australia

[2]Department of Ecology and Evolution, Research School of Biology, Australian National University, Canberra, ACT, Australia

[3]Department of Biology, Indiana University, Bloomington, IN

[4]Department of Computer Science, Indiana University, Bloomington, IN

**Correspondence to:** *Corresponding author: E-mail: rob.lanfear@anu.edu.au.

**Associate editor:** Michael Rosenberg

## Abstract

We implement two measures for quantifying genealogical concordance in phylogenomic data sets: the gene concordance factor (gCF) and the novel site concordance factor (sCF). For every branch of a reference tree, gCF is defined as the percentage of "decisive" gene trees containing that branch. This measure is already in wide usage, but here we introduce a package that calculates it while accounting for variable taxon coverage among gene trees. sCF is a new measure defined as the percentage of decisive sites supporting a branch in the reference tree. gCF and sCF complement classical measures of branch support in phylogenetics by providing a full description of underlying disagreement among loci and sites. An easy to use implementation and tutorial is freely available in the IQ-TREE software package (http://www.iqtree.org/doc/Concordance-Factor, last accessed May 13, 2020).

*Key words:* phylogenetic inference, concordance factor, phylogenomics.

# Phylogenetic trees are sets of splits / bipartitions

**Division**

**Splits / Bipartitions**



Set of splits in reference tree: {A, B, C, D, E}    {F, G, H, I}



*Conflicting* set of splits: {I, B, C, D, E}    {F, G, H, A}

**Internode Certainty (IC)**: a measure of the support for a given internode (bipartition) by considering its frequency in a given set of trees jointly with that of the most prevalent conflicting bipartition in the same set of trees

**Tree Certainty (TC)**: the sum of IC values across all internodes

**Implemented in RAxML and QuartetScores (https://github.com/lutteropp/QuartetScores)**



Ratio of Supports for Two Conflicting Bipartitions

*Salichos et al. (2014) Mol. Biol. Evol.; Kobert et al. (2016) Mol. Biol. Evol.; Zhou et al. (2020) Syst. Biol.*

**Saccharomyces lineage**

*Kluyveromyces waltii* (Kwal)
*Kluyveromyces thermotolerans* (Kthe)
*Saccharomyces kluyveri* (Sklu)
*Kluyveromyces lactis* (Klac)
*Eremothecium gossypii* (Egos)
*Zygosaccharomyces rouxii* (Zrou)
*Kluyveromyces polysporus* (Kpol)
*Candida glabrata* (Cgla)
*Saccharomyces castellii* (Scas)
*Saccharomyces bayanus* (Sbay)
*Saccharomyces kudriavzevii* (Skud)
*Saccharomyces mikatae* (Smik)
*Saccharomyces paradoxus* (Spar)
*Saccharomyces cerevisiae* (Scer)

**62**

**% Support for most prevalent conflict:**
**#1: 6%**

**IC value: 0.59**

**52**

**% Support for most prevalent conflict:**
**#1: 29%**

**IC value: 0.06**

The yeast species phylogeny inferred using the STAR species tree method



**Coalescent units / IC**

# Coalescent branch lengths are correlated with GSF/IC

**Internode length: influences amount of phylogenetic signal (I)**
**Homoplasy: independent evolution of identical characters ( * , ● )**

# Certain recipes for handling incongruence didn't help

| Treatment | Tree Certainty | # of Internodes where IC increased \| decreased |
|---|---|---|
| Default analysis | 8.35 | n/a |
| *Removing sites containing gaps* | | |
| All sites with gaps excluded | 7.91 | 0 \| 7 |
| *Removing fast-evolving or unstable species* | | |
| *C. lusitaniae* | 8.15 | 1 \| 2 |
| *C. glabrata* | 8.30 | 2 \| 2 |
| *E. gossypii, C. glabrata, K. lactis* | 7.88 | 1 \| 3 |
| *Selecting genes that recover specific clades* | | |
| [*C. tropicalis, C. dubliniensis, C. albicans*] | 8.62 | 0 \| 0 |
| *Selecting the most slow-evolving genes* | | |
| *100 slowest-evolving genes* | 6.76 | 2 \| 9 |

# What do we do then?

| Treatment | Tree Certainty | # of Internodes where IC increased \| decreased |
|---|---|---|
| **Default analysis** | **8.35** | **n/a** |
| *Selecting genes whose bootstrap consensus trees have high average support* | | |
| **All genes with average BS ≥ 60%** | **8.59** | **4 \| 0** |
| **All genes with average BS ≥ 70%** | **9.18** | **14 \| 0** |
| **All genes with average BS ≥ 80%** | **9.92** | **15 \| 0** |

average BS ≥60%

100
- *S. cerevisiae*
- *S. castellii*
- *C. glabrata*

average BS ≥70%

73
- *S. cerevisiae*
- *S. castellii*
- *C. glabrata*

average BS ≥80%

95
- *S. cerevisiae*
- *C. glabrata*
- *S. castellii*

# Selecting specific bipartitions dramatically improves phylogeny



| Treatment | Tree Certainty | # of Internodes where IC increased \| decreased |
|:---:|:---:|:---:|
| **Default analysis** | **8.35** | **n/a** |
| *Selecting genes whose bootstrap consensus trees have high average support* | | |
| **All bipartitions with BS ≥ 60%** | **10.11** | **14 \| 0** |
| **All bipartitions with BS ≥ 70%** | **10.70** | **16 \| 0** |
| **All bipartitions with BS ≥ 80%** | **11.32** | **15 \| 0** |

# The status of the yeast phylogeny



Supported by Rare Genomic Characters

# Vertebrates
**(1,086 genes, 18 taxa)**

# Animals
**(225 genes, 21 taxa)**

# Mosquitoes
**(2,007 genes, 20 taxa)**

**Hypothesis:** these debates concern internodes that are poorly supported by individual gene trees

**Test:** measure the phylogenetic signal in contentious branches of the tree of life

**A measure of the statistical dependence among species' trait values due to their phylogenetic relationships / the tendency of related species to resemble each other more than species drawn at random from the same tree**

**Revell et al. (2008) Syst. Biol.**
**Münkemüller et al. (2012) Methods Ecol. Evol.**

**The amount of support for a particular topology, e.g., the relative number of resolved internodes in a consensus tree**

**Sanderson (2008) Science**

**A measure of the substitutions occurring along a given branch of the evolutionary tree. In parsimony methods, the signal is encoded in shared derived characters. In probabilistic methods, the amount of phylogenetic signal actually extracted from a given dataset depends on the model and is expected to increase with the fit of the model to the data**

**Philippe et al. (2011) PLoS Biol.**
**Townsend et al. (2012) Syst. Biol.**

**Maximum Likelihood tree**

**(T1)**

**Conflicting tree**

**(T2)**



$$\ln(T_1|X_i) = -100$$

$$\ln(T_2|X_i) = -150$$

$$Phylogenetic\ Signal = -(\ln(T_1|X_i) - \ln(T_2|X_i))$$

**1,080 genes from 36 animal taxa**

# Signal of genes in a phylogenomic data matrix

# Signal of genes in multiple phylogenomic data matrices



**Phylogenetic Signal** (y-axis)

**Genes** (x-axis)

Panels: Borowiec_total1080 (1,080); Borowiec_best108 (108); Ryanl_EST_Opisthokonta (406); Ryan_EST_Choanoflagellata (406); Whelan_Dataset1_Opisthokonta (251); Whelan_Dataset1_Choanoflagellata (251); Whelan_Dataset16_Opisthokonta (89); Whelan_Dataset16_Choanoflagellata (89)

# Summarizing phylogenetic signal across genes and sites

# Testing several contentious branches of the tree of life

| Clade | ML Tree (T1) | Conflicting Tree (T2) |
|---|---|---|
| **Plants** | *Amborella* as sister to all other flowering plants | *Amborella* + *Nuphar* as sister to all other flowering plants |
| | Magnoliids as sister to Eudicots + Chloranthales | Eudicots as sister to Magnoliids + Chloranthales |
| | Hornworts as sister to all other land plants, followed by a mosses + liverworts clade | Hornworts as sister to a mosses + liverworts clade |
| | Gnetales as sister to the Pinaceae, nested within the Coniferales | Gnetales as sister to the Coniferales |
| | Zygnematophyceae as sister to all land plants | Charales as sister to all land plants |
| **Vertebrates** | Gymnophiona as sister to all other amphibians | Anura as sister to all other amphibians |
| | Atlantogenata (Afrotheria + Xenarthra) as sister to all other placental mammals | Afrotheria as sister to all other placental mammals |
| | Lungfishes as sister to all tetrapods | Lungfishes + coelacanths as sister to all tetrapods |
| | Pigeons as sister to all other Neoaves | Falcons as sister to all other Neoaves |
| | Elopomorpha + Osteoglossomorpha as sister to all other teleosts | Osteoglossomorpha alone as sister to all other teleosts |
| | Turtles as sister to archosaurs (birds + crocodiles) | Turtles as sister to crocodiles |
| **Yeasts** | Ascoideaceae as sister to Phaffomycetaceae + Saccharomycetaceae | Ascoideaceae as sister to a clade comprising Pichiaceae, Debaryomycetaceae, Phaffomycetaceae, and Saccharomycetaceae |
| | *Candida glabrata* rather than *Naumovozyma castellii* as sister to Saccharomyces sensu stricto yeasts | *Naumovozyma castellii* rather than *Candida glabrata* sister to Saccharomyces sensu stricto yeasts |
| | *Hyphopichia burtonii* as sister to *Candida auris* + *Metschnikowia bicuspidata* | *Hyphopichia burtonii* as sister to *Debaryomyces hansenii* |
| | *Zygosaccharomyces rouxii* as sister to all other yeasts with occurring whole-genome duplication event | *Vanderwaltozyma polyspora* as sister to all other yeast with occurring whole-genome duplication event |
| | *Meyerozyma guilliermondii* as sister to *Debaryomyces hansenii* | *Meyerozyma guilliermondii* as sister to *Hyphopichia burtonii* + *Candida auris* |
| | *Candida tanzawaensis* as sister to *Pichia stipiti* + *Candida maltosa* | *Pichia stipiti* as sister to *Candida tanzawaensis* + *Candida maltosa* |

**1233 genes, 86 yeast taxa**



**Difference in phy-logenetic signal** (y-axis): 300, 200, 100, 0

**Removal of this gene switches support from T1 to T2**

T1

T2

**Genes** (x-axis)

# What happens if we remove that one gene?

Quantifying the impact of removing opinionated genes

*Shen et al. (2017) Nature Ecol. Evol.*

Which branches are resolved and which are unresolved?

*Shen et al. (2017) Nature Ecol. Evol.*

**Explanation #1: Biological factors (parts of the tree of life are bush-like / network-like rather than tree-like)**

**Explanation #2: Analytical factors (systematic error due to the bad fit of our models to our data)**

# Genome-scale phylogeny of 332 yeasts

*Saccharomyces*

WGD
clade

*Candida
albicans*

332 taxa
2,408 genes

*Shen, Opulente,
Kominek, Zhou et
al. (2018) Cell*

# The 32 conflicting branches in the yeast phylogeny



Concatenation Coalescence

Internal branches

Data matrices

~10% (32 / 331) of internal branches show conflict between analyses

Legend:
- Strong T1 (dark green)
- Weak T1 (light green)
- Strong T2 (red)
- Weak T2 (orange)
- Strong T3 (dark blue)
- Weak T3 (light blue)

*Shen, Opulente, Kominek, Zhou et al. (2018) Cell*

# Distribution of conflict on the yeast phylogeny



Lipomycetaceae

Trigonopsidaceae

Dipodascaceae / Trichomonascaceae

Alloascoideaceae

*Sporopachydermia* clade

CUG-Ala clade

Pichiaceae

CUG-Ser1 clade

CUG-Ser2 clade

Phaffomycetaceae

Saccharomycodaceae

Saccharomycetaceae

*Shen, Opulente, Kominek, Zhou et al. (2018) Cell*

*Lipomyces starkeyi* — **Lipomycetaceae**
*Tortispora caseinolytica* — **Trigonopsidaceae**
*Nadsonia fulvescens*
*Yarrowia lipolytica*
*Geotrichum candidum* CLIB 918 — **Yarrowia clade**
*Saprochaete clavata*
*Blastobotrys adeninivorans*
*Candida apicola*
*Starmerella bombicola*
*Pachysolen tannophilus* — **Pichiaceae**
*Komagataella pastoris* — **Komagataella clade**
*Kuraishia capsulata*
*Candida boidinii*
*Candida arabinofermentans*
*Ogataea polymorpha* — **Pichiaceae**
*Ogataea parapolymorpha*
*Brettanomyces bruxellensis*
*Brettanomyces anomalus*
*Pichia membranifaciens*
*Pichia kudriavzevii*
*Babjeviella inositovora*
*Candida tenuis*
*Hyphopichia burtonii*
*Candida auris*
*Clavispora lusitaniae*
*Metschnikowia bicuspidata*
*Metschnikowia fructicola*
*Debaryomyces hansenii* — **Debaryomycetaceae Metschnikowiaceae**
*Meyerozyma guilliermondii*
*Meyerozyma caribbica*
*Suhomyces tanzawaensis*
*Scheffersomyces stipitis*
*Spathaspora arborariae*
*Spathaspora passalidarum*
*Lodderomyces elongisporus*
*Candida parapsilosis*
*Candida orthopsilosis*
*Candida dubliniensis*
*Candida albicans*
*Candida maltosa*
*Candida sojae*
*Candida tropicalis*
*Ascoidea rubescens* — **Ascoideaceae**
*Wickerhamomyces ciferrii*
*Wickerhamomyces anomalus*
*Cyberlindnera fabianii* — **Phaffomycetaceae**
*Cyberlindnera jadinii*
*Hanseniaspora vineae*
*Hanseniaspora uvarum* — **Saccharomycodaceae**
*Hanseniaspora valbyensis*
*Lachancea kluyveri*
*Lachancea lanzarotensis*
*Lachancea waltii*
*Lachancea thermotolerans*
*Eremothecium coryli*
*Eremothecium cymbalariae*
*Eremothecium gossypii*
*'Ashbya' aceri*
*Kluyveromyces aestuarii*
*Kluyveromyces wickerhamii*
*Kluyveromyces marxianus*
*Kluyveromyces dobzhanskii*
*Kluyveromyces lactis*
*Torulaspora delbrueckii*
*Zygosaccharomyces bailii* — **Saccharomycetaceae**
*Zygosaccharomyces rouxii*
*Tetrapisispora blattae*
*Vanderwaltozyma polyspora*
*Tetrapisispora phaffii*
*Naumovozyma castellii*
*Naumovozyma dairenensis*
*Kazachstania africana*
*Kazachstania naganishii*
*Candida castellii*
*Nakaseomyces bacillisporus*
*Candida glabrata*
*Candida bracarensis*
*Nakaseomyces delphensis*
*Candida nivariensis*
*Saccharomyces uvarum*
*Saccharomyces eubayanus*
*Saccharomyces arboricola*
*Saccharomyces kudriavzevii*
*Saccharomyces mikatae*
*Saccharomyces paradoxus*
*Saccharomyces cerevisiae*

CUG

WGD

**Unresolved**

**1,233-gene, 86-taxon data matrix**

**~13% (11 / 85) of internal branches conflict between analyses**

**Despite increasing # internal branches ~4X, (85 -> 331), conflict decreased**

*Shen et al. (2016) G3*

A single gene governs the placement of Ascoideaceae

1,233 genes, 86 yeast taxa

Removal of this gene (*DPM1*) switches support from T1 to T2

Ascoidea asiatica

Ascoidea rubescens

Saccharomycopsis malanga

Saccharomycopsis capsularis

**Phaffomycetaceae**

**Saccharomycetaceae / Saccharomycodaceae**

# Genomfart? The way forward?

**Multiple sequence alignment / data matrix reconstruction (more taxa, more genes)**

→

**Apply different phylogenetic analyses (diff. optimality criteria / diff. approaches / different treatments)**

**Assess conflict (e.g., use internode certainty / concordance factors)**

**Investigate alternative hypotheses for branches showing conflict / assess sensitivity of results (keep biology of lineage in mind!)**

**Only report resolution of branches that you have support for**

**"One can use the most sophisticated audio equipment to listen, for an eternity, to a recording of white noise and still not glean a useful scrap of information"**

**Rodrigo et al. (1994)**
**Chapter in: Sponge in Time and Space; Biology, Chemistry, Paleontology**

# Acknowledgements



Xing-Xing Shen

Xiaofan Zhou

Leonidas Salichos — NEW YORK INSTITUTE OF TECHNOLOGY

Jacob Steenwyk — Berkeley UNIVERSITY OF CALIFORNIA

HITS — Alexis Stamatakis

National Science Foundation — WHERE DISCOVERIES BEGIN

Chris Todd Hittinger

Marizeth Groenewald

NIH — National Institute of Allergy and Infectious Diseases

BURROUGHS WELLCOME FUND

WISCONSIN UNIVERSITY OF WISCONSIN-MADISON

WESTERDIJK FUNGAL BIODIVERSITY INSTITUTE

http://www.rokaslab.org/        @RokasLab