# Era of Phylogenomics: Influx of genomic data



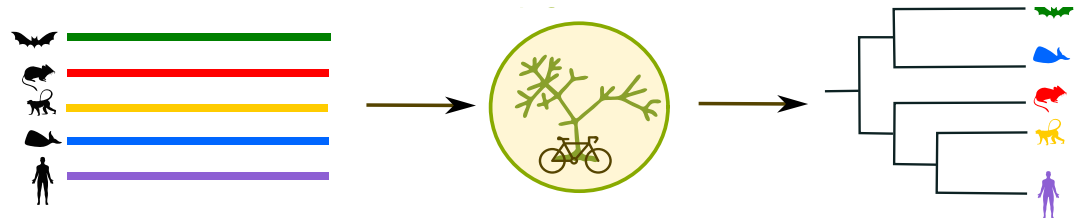Image source: G A Bravo et al., AR Ecology, Evolution, and Systematics, 2021

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# Large-scale Genomic Sequencing

- Multiple consortiums aiming to sequence thousands to millions of species

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# Large-scale Genomic Sequencing

- Multiple consortiums aiming to sequence thousands to millions of species

- Phylogenomic analyses of huge datasets solves various questions related to Tree of Life

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego

JACOBS SCHOOL OF ENGINEERING

A — Stiller at al. (2024)
B — Prum et al. (2015)
C — Jarvis et al. (2014)
D — Kuhl et al. (2021)

# Evolutionary Trees are still debated!

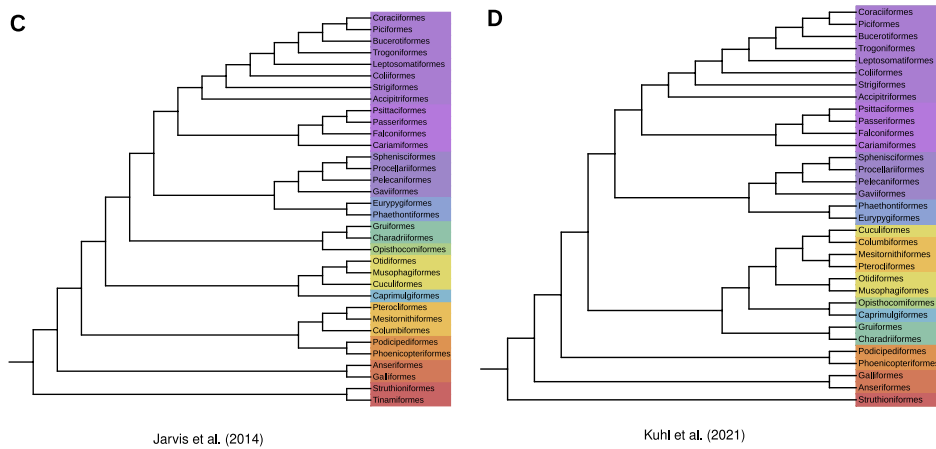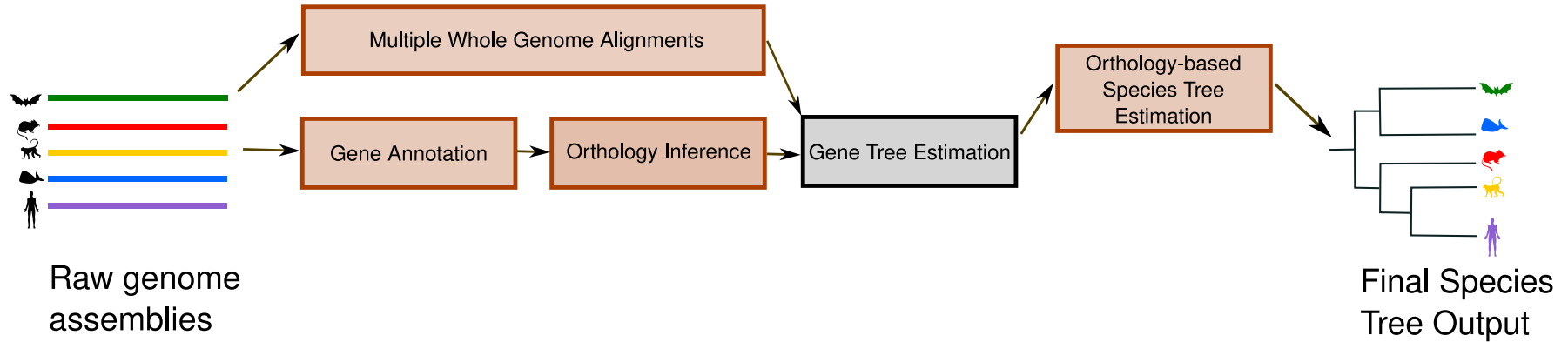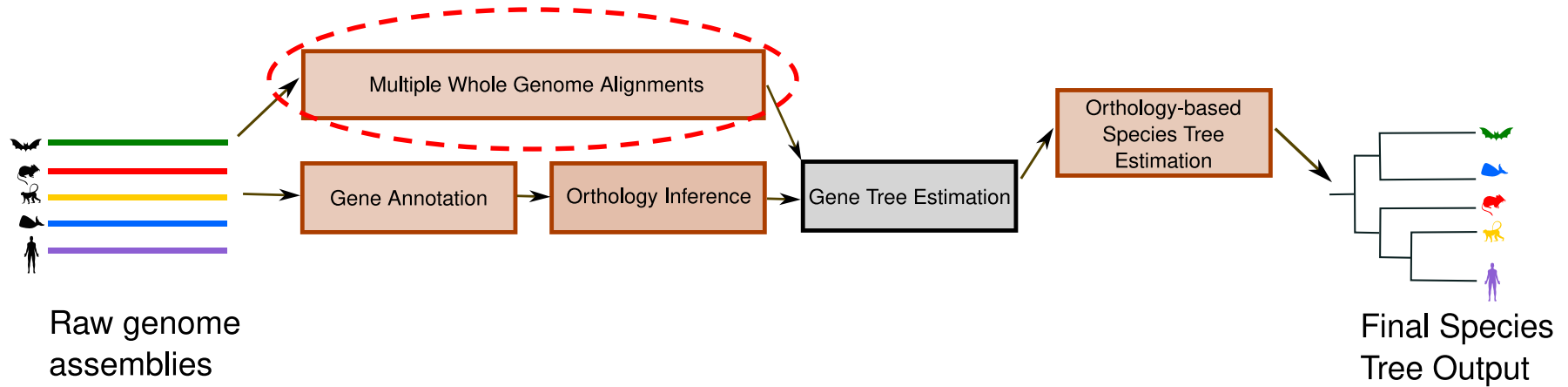Figure shows different avian phylogenies proposed by various groups

Image source:

A. Stiller, J. et al. Complexity of avian evolution revealed by family-level genomes. Nature (2024).
B. Prum, R., Berv, J., Dornburg, A. et al. A comprehensive phylogeny of birds (Aves) using targeted next-generation DNA sequencing. Nature (2015).
C. Jarvis, E. D. et al. Whole-genome analyses resolve early branches in the tree of life of modern birds. Science (2014).
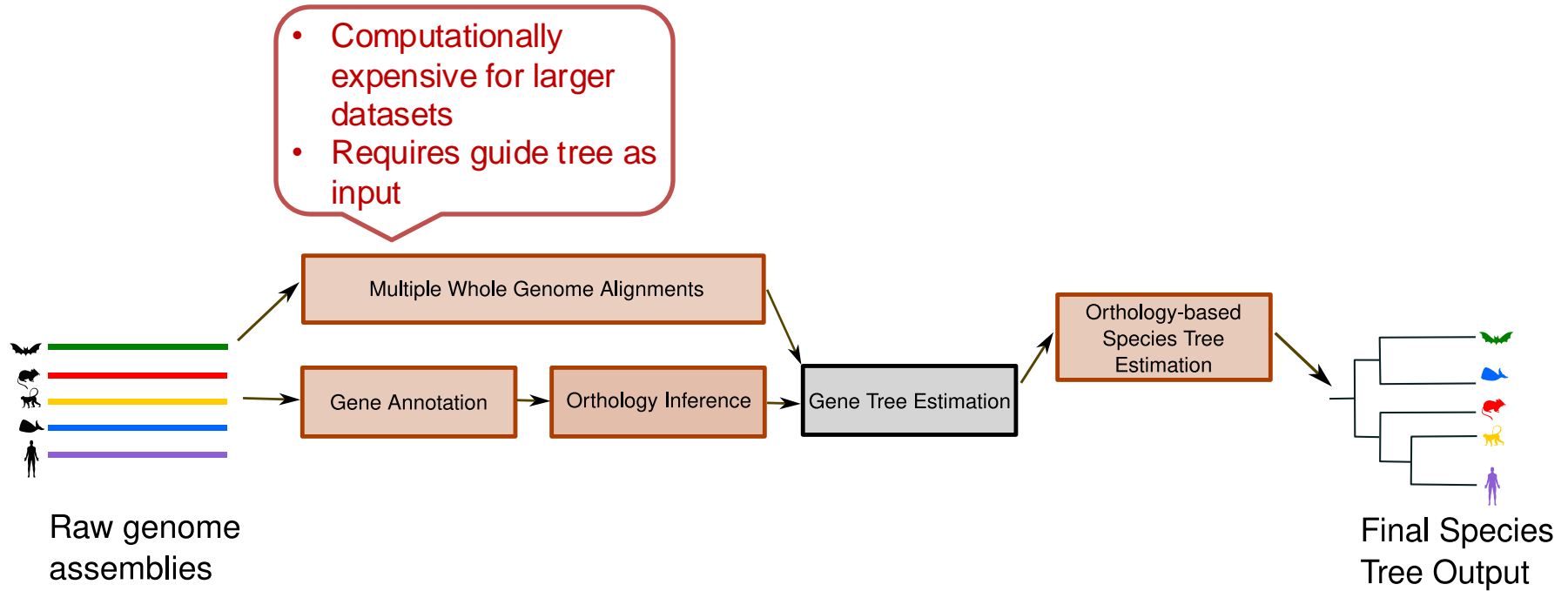D. Kuhl et al., An Unbiased Molecular Approach Using 3'-UTRs Resolves the Avian Family-Level Tree of Life, MBE (2021).

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# Existing approaches to estimate species tree



Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# Existing approaches to estimate species tree



Raw genome assemblies

Multiple Whole Genome Alignments

Gene Annotation

Orthology Inference

Gene Tree Estimation

Orthology-based Species Tree Estimation

Final Species Tree Output

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego

**JACOBS SCHOOL OF ENGINEERING**

# Existing approaches to estimate species tree



Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# Existing approaches to estimate species tree



Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

# Existing approaches to estimate species tree



Raw genome assemblies

Multiple Whole Genome Alignments

Gene Annotation

Orthology Inference

Gene Tree Estimation

Orthology-based Species Tree Estimation

Final Species Tree Output

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# Existing approaches to estimate species tree



Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

# Existing approaches to estimate species tree



Raw genome assemblies

Final Species Tree Output

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego

**JACOBS SCHOOL OF ENGINEERING**

# Gene Tree Discordances

- Different parts of the genome can infer different phylogenies



Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

# Gene tree discordance aware tool - ASTRAL[1]



(probabilities are made-up and wrong in this table)

- Estimates an unrooted species tree given a set of unrooted gene trees.

- It is statistically consistent under the multi-species coalescent model

[1]Mirarab et al., ASTRAL: genome-scale coalescent-based species tree estimation, Bioinformatics (2014).

Image source: UCSD ECE 208 Lecture Slides

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

# Existing approaches to estimate species tree



Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

# Existing approaches to estimate species tree



Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# No automated yet accurate tool exists to infer phylogeny directly from raw genomes

- Separate tool exists for individual steps

- Accuracy depends on input tree/alignment -> error-prone

- Relies on domain expertise

- Takes months to complete

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

# Objectives

To develop a

- Reference-free

- Orthology-free

- Alignment-free

- Discordance-aware

Approach for Estimation of Species tree

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

**UC San Diego**

**JACOBS SCHOOL OF ENGINEERING**

# Objectives

To develop a

- **R**eference-free

- **O**rthology-free

- **A**lignment-free

- **DI**scordance-aware

Approach for **E**stimation of **S**pecies tree

**ROADIES** – **R**eference-free **O**rthology-free **A**lignment-free **DI**scordance-aware **E**stimation of **S**pecies Tree

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# What is ROADIES?

Automated tool which takes raw genomic assemblies as input and outputs species tree



Input Genomic Sequences

ROADIES

Final Species Tree Output

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

**UC San Diego**
**JACOBS SCHOOL OF ENGINEERING**

# What ROADIES does?



Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# What ROADIES does?



Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

**UC San Diego**
**JACOBS SCHOOL OF ENGINEERING**

# Species Tree Estimation by ASTRAL-Pro[1]



Gorilla  Human  Chimp  Orang.  Dog
(probabilities are made-up and wrong in this table)

- ASTRAL Pro → ASTRAL for PaRalogs and Orthologs

- Statistically consistent discordance aware tool

- Finds the best tree with maximum dominant quartets

- Does not require separation of orthologs and paralogs

[1]Zhang et al. "ASTRAL-Pro: Quartet-Based Species-Tree Inference despite Paralogy", MBE 2020

Image source: UCSD ECE 208 Lecture Slides

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# What ROADIES does?



Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

**UC San Diego**

**JACOBS SCHOOL OF ENGINEERING**

# What ROADIES does?



Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego

JACOBS SCHOOL OF ENGINEERING

# What ROADIES does?



Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# What ROADIES does?



Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

**UC San Diego**
**JACOBS SCHOOL OF ENGINEERING**

# ROADIES pipeline



Random Sampling of genes from input genomic assemblies

Genes

Whole genome assemblies

Pairwise alignment of sampled genes to other genomes

Sampled Genes

Pairwise whole genome alignment

Filter low quality pairwise alignments and repeats

Discarded alignments

Retained alignments

Perform Multiple sequence alignment of filtered genes across all species

Filtered genes

MSA

Estimate genes trees from MSAs

Gene trees

Estimate species tree from lists of gene trees

Estimated species tree

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

**UC San Diego**

**JACOBS SCHOOL OF ENGINEERING**

# ROADIES pipeline

Random Sampling of
genes from input
genomic assemblies

Genes

Whole genome
assemblies

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

**UC San Diego**

**JACOBS SCHOOL OF ENGINEERING**

# ROADIES pipeline

Random Sampling of genes from input genomic assemblies
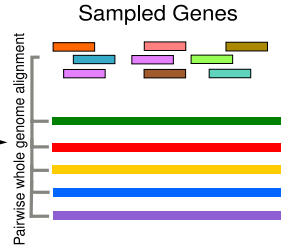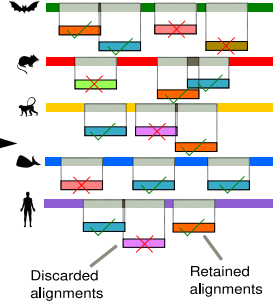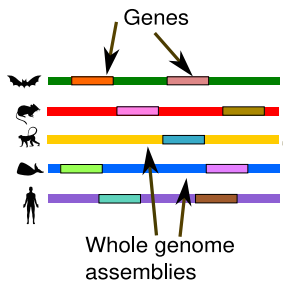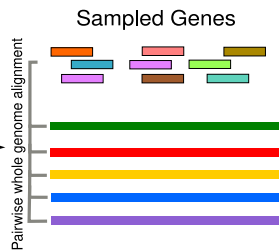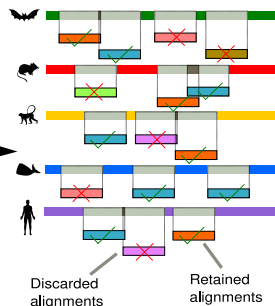
Pairwise alignment of sampled genes to other genomes

Genes

Sampled Genes

Whole genome assemblies

Pairwise whole genome alignment

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# ROADIES pipeline



Random Sampling of genes from input genomic assemblies

Pairwise alignment of sampled genes to other genomes

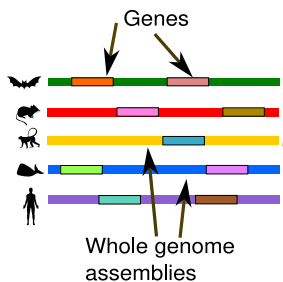Filter low quality pairwise alignments and repeats

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES
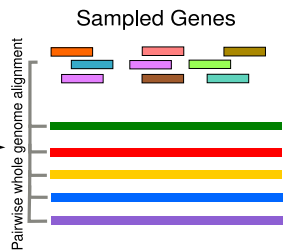
UC San Diego
JACOBS SCHOOL OF ENGINEERING
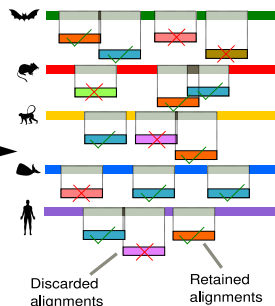
# ROADIES pipeline


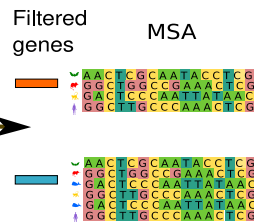
Random Sampling of genes from input genomic assemblies
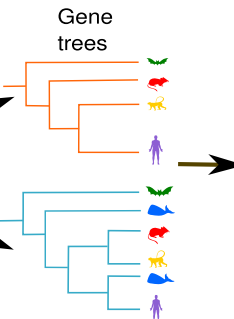
Pairwise alignment of sampled genes to other genomes

Filter low quality pairwise alignments and repeats

Perform Multiple sequence alignment of filtered genes across all species

Genes

Whole genome assemblies

Pairwise whole genome alignment

Sampled Genes

Discarded alignments

Retained alignments

Filtered genes

MSA

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

**UC San Diego**
**JACOBS SCHOOL OF ENGINEERING**

# ROADIES pipeline



Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

**UC San Diego**
**JACOBS SCHOOL OF ENGINEERING**

# ROADIES pipeline



Random Sampling of genes from input genomic assemblies

Pairwise alignment of sampled genes to other genomes

Filter low quality pairwise alignments and repeats

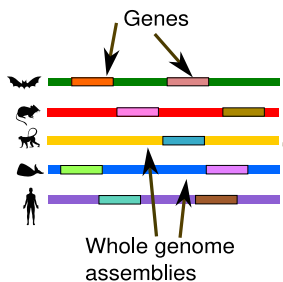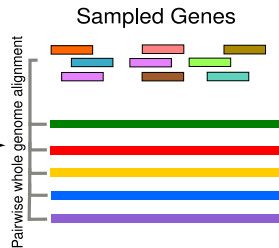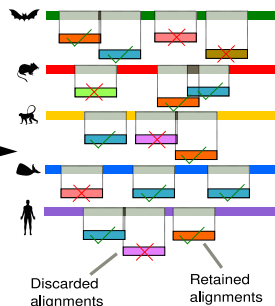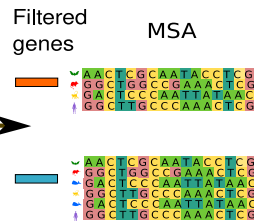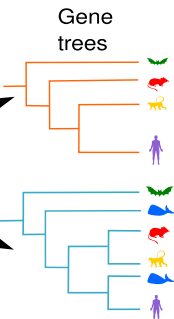Perform Multiple sequence alignment of filtered genes across all species

Estimate genes trees from MSAs

Estimate species tree from lists of gene trees

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

**UC San Diego**
**JACOBS SCHOOL OF ENGINEERING**

# Problem: How many genes to start with?



Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# ROADIES converges into accurate tree with more gene trees



Random Sampling of genes from input genomic assemblies

Pairwise alignment of sampled genes to other genomes

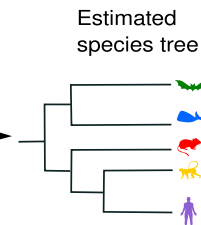Filter low quality pairwise alignments and repeats

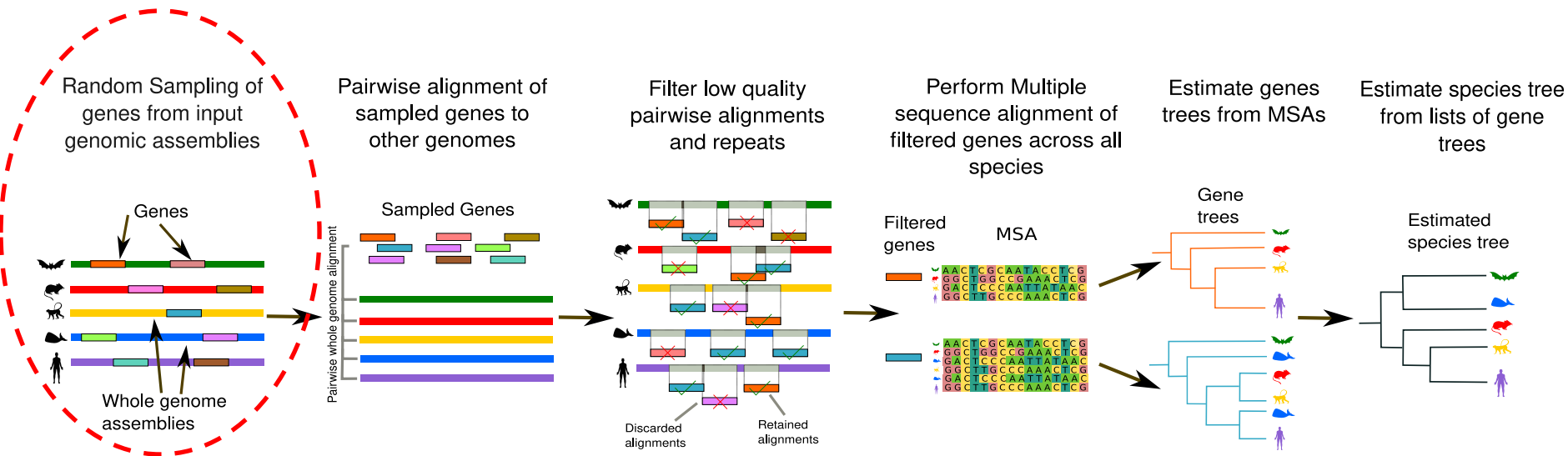Perform Multiple sequence alignment of filtered genes across all species
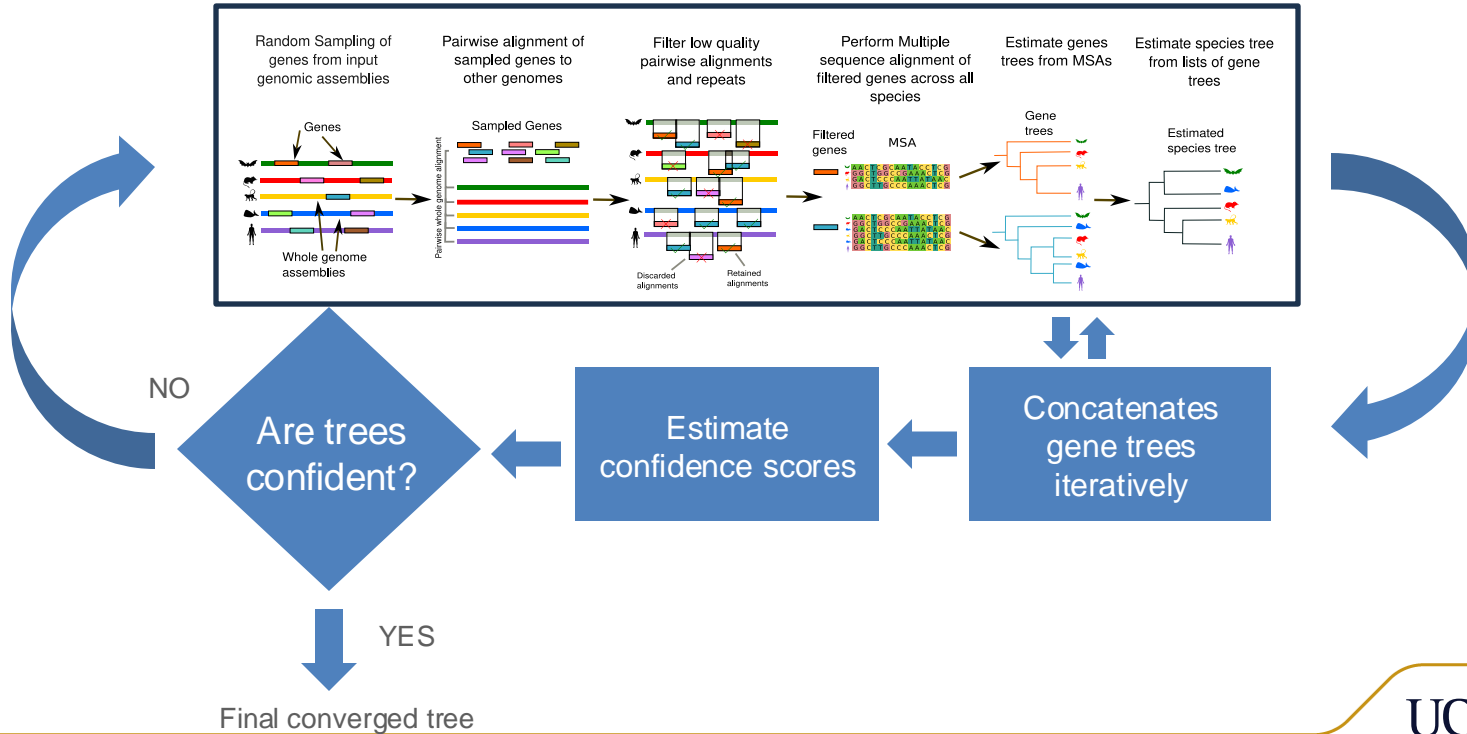
Estimate genes trees from MSAs
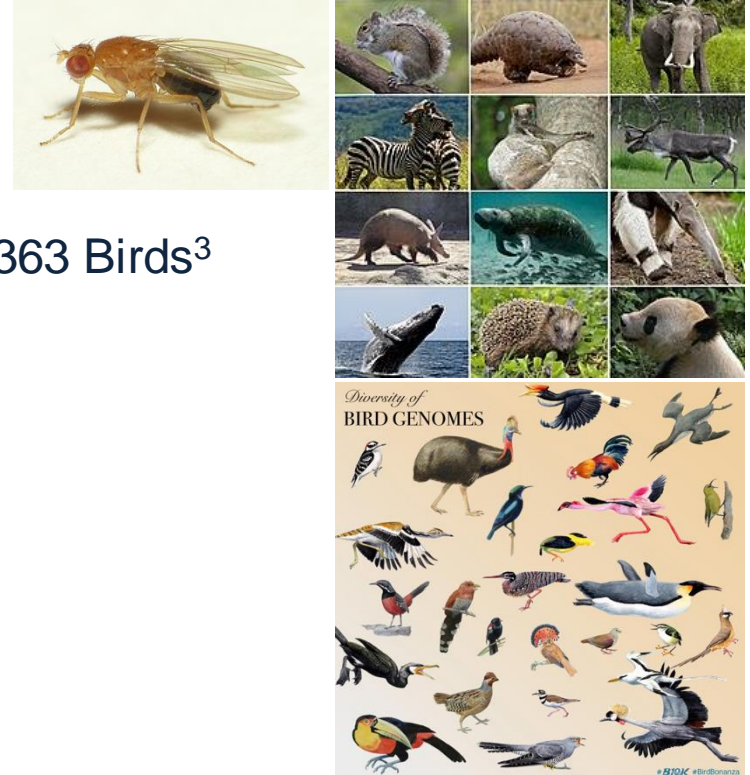
Estimate species tree from lists of gene trees

Are trees confident?

Estimate confidence scores

Concatenates gene trees iteratively

NO

YES

Final converged tree

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

**UC San Diego**

**JACOBS SCHOOL OF ENGINEERING**
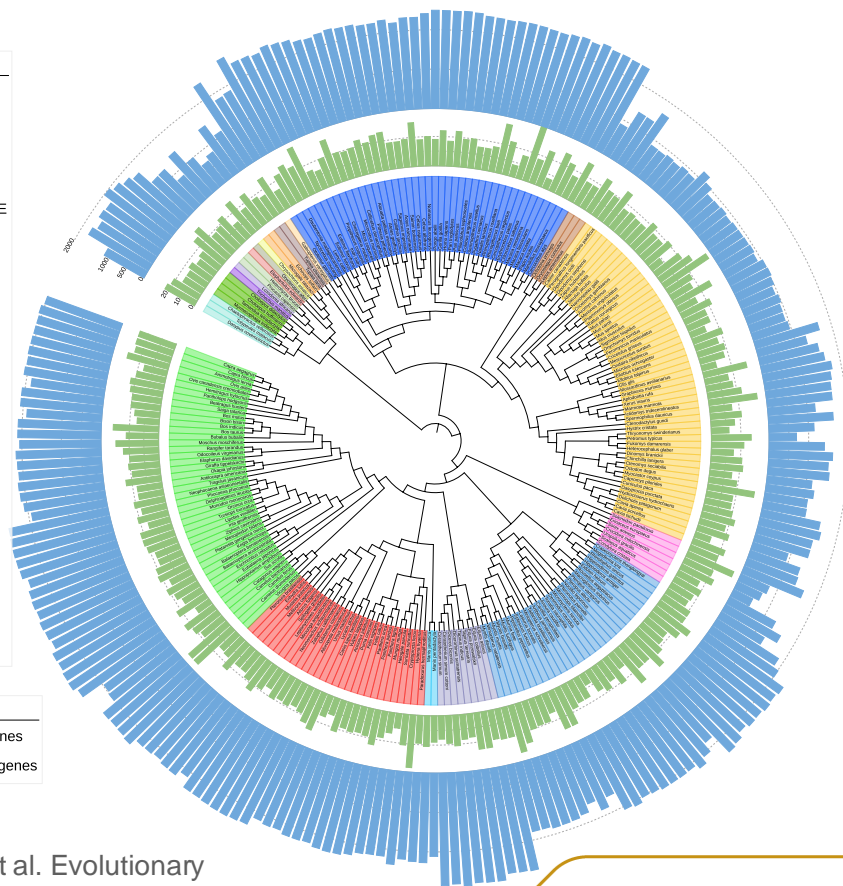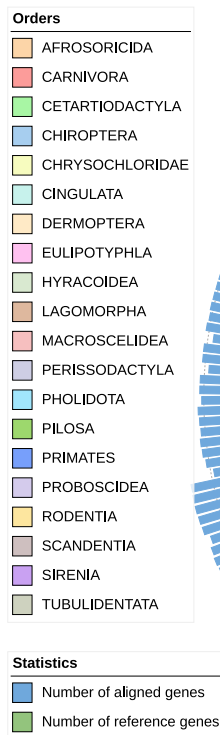
# Methodology



- **Datasets** – 240 Mammals[1], 100 Flies[2] and 363 Birds[3]

- **Normalized Robinson-Foulds distance**
  - Species-level
  - Orders/Group-level

- **Tree confidence metric**
  - Local posterior probability
  - Quartet scores

1 – Zoonomia 2020, 2- Kim et al. 2021, 3 – Birds 10k Genome Project (Feng et al. 2020), Image source: Wikipedia, Researchers Sequence Genomes of 363 Bird Species - SciNews
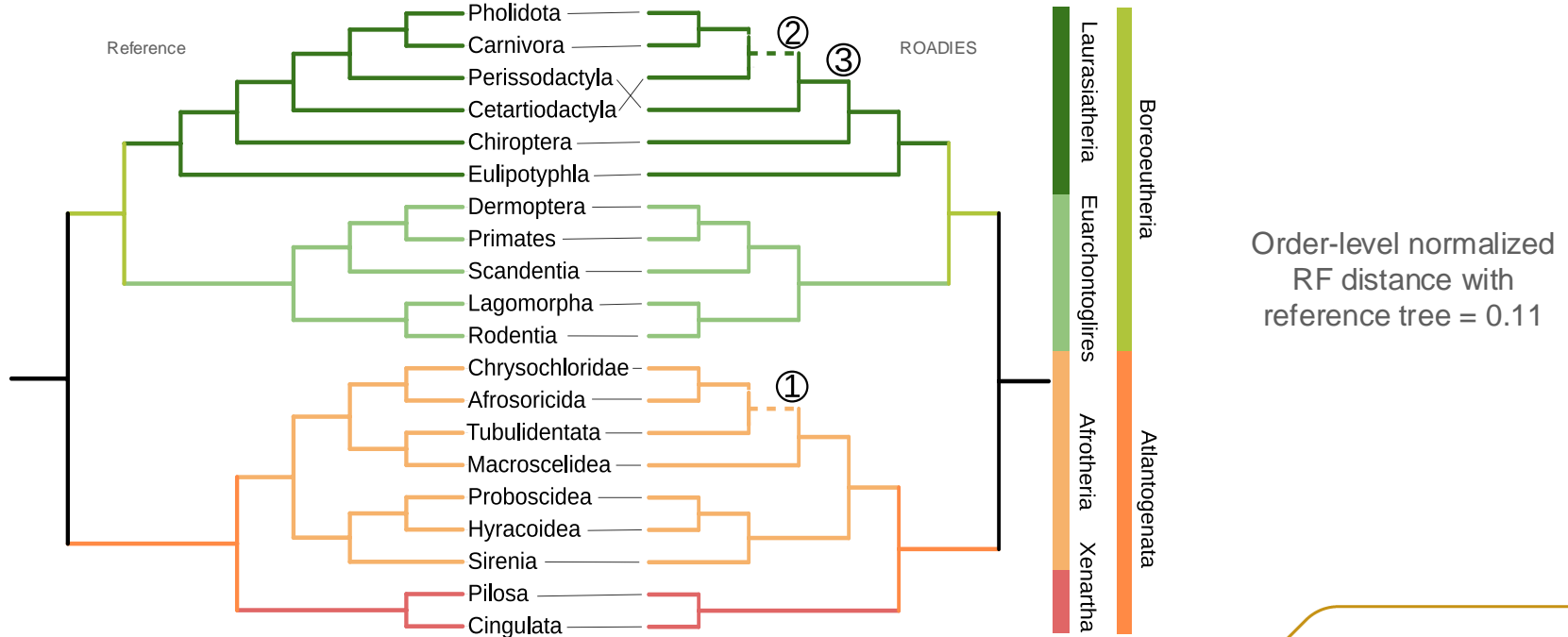
Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

# ROADIES estimates accurate phylogeny of 240 placental mammals

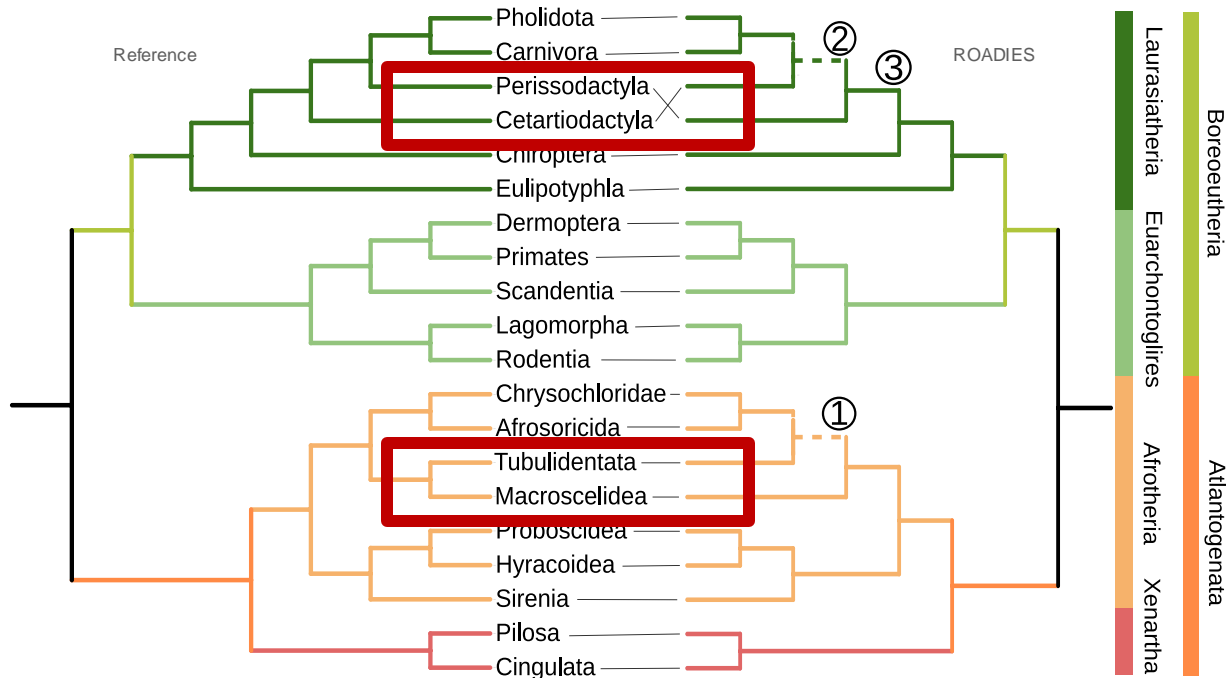Species-level normalized RF distance
with reference tree = 0.037

Reference tree and datasets taken from Zoonomia project (Christmas et al. Evolutionary constraint and innovation across hundreds of placental mammals, Science 2023).



**Orders**
- AFROSORICIDA
- CARNIVORA
- CETARTIODACTYLA
- CHIROPTERA
- CHRYSOCHLORIDAE
- CINGULATA
- DERMOPTERA
- EULIPOTYPHLA
- HYRACOIDEA
- LAGOMORPHA
- MACROSCELIDEA
- PERISSODACTYLA
- PHOLIDOTA
- PILOSA
- PRIMATES
- PROBOSCIDEA
- RODENTIA
- SCANDENTIA
- SIRENIA
- TUBULIDENTATA

**Statistics**
- Number of aligned genes
- Number of reference genes

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# ROADIES estimates accurate phylogeny of 240 placental mammals at order-level



Order-level normalized RF distance with reference tree = 0.11

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# ROADIES estimates accurate phylogeny of 240 placental mammals at order-level



Order-level normalized RF distance with reference tree = 0.11

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# ROADIES estimates accurate phylogeny of 240 placental mammals at order-level

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

# ROADIES converges with more gene trees



Experiments run on AWS R6a 16-core instances
Runtime is calculated as wall clock time
High support nodes: Nodes with localPP >= 0.95

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# ROADIES estimates accurate phylogeny of 100 Drosophilid genomes (fruit flies)



**Groups**
- Cardini
- Funebris
- hawaiian_drosophila
- Immigrans
- Repleta
- Scaptomyza
- Tumiditarsus
- Virilis
- Zaprionus
- Lordiphosa
- Melanogaster
- Obscura
- Saltans
- Willstoni
- Chymomyza
- Leucophenga

**Statistics**
- Number of reference genes
- Number of aligned genes

Species-level normalized RF distance with reference tree = 0.062

Reference tree and datasets taken from the paper: Kim, B. Y. et al. Highly contiguous assemblies of 101 drosophilid genomes. ELife, 2021.

**UC San Diego**
**JACOBS SCHOOL OF ENGINEERING**

# ROADIES estimates accurate phylogeny of 363 avian species

Species-level normalized RF distance with reference tree = 0.037

Reference tree and datasets taken from Stiller, J. et al. Complexity of avian evolution revealed by family-level genomes. Nature (2024) doi:10.1038/s41586-024-07323-1



Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# ROADIES is 176x faster than conventional approaches



Experiments run on AWS R6a 16-core instances
Runtime is calculated as wall clock time
Experiments are performed with mammals datasets

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

**UC San Diego**
**JACOBS SCHOOL OF ENGINEERING**

# ROADIES scales well with increasing system cores and species count



ROADIES scales between linear and quadratic with more species count

Experiments run on AWS R6a instances
Runtime is calculated as wall clock time

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# ROADIES aims to support diverse evolutionary timescales

- Mammals ~ 100 million years

- Flies ~ 40 million years

- Birds ~ 150 million years

- Fish ~ 500 million years

- Bacteria ~ 3 billion years

- SARS Cov2 ~ 3-4 years



Image Source: geologyscience.com/geology-branches/paleontology/geologic-time-scale/

UC San Diego

JACOBS SCHOOL OF ENGINEERING

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

# Summary

- ROADIES is first-of-its-kind tool which automates species tree inference directly from raw genome assemblies

- Highly configurable and scalable

- 176x faster than conventional methods (for mammals dataset)

- Accurate results for mammals, flies and, birds dataset

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

**UC San Diego**

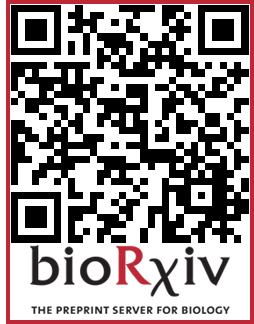**JACOBS SCHOOL OF ENGINEERING**

# Acknowledgments

**Thank you for your attention.**

**Questions?**

**Advisors:**

- Yatish Turakhia (UCSD)

- Siavash Mirarab (UCSD)

**Collaborators:**

- Tian (Kevin) Liu (UCSD)

- Hiram Clawson (UCSC)

- Guojie Zhang (Zhejiang University)

- Yulong Xie (Zhejiang University)

- Benedict Paten (UCSC)

bioRχiv
THE PREPRINT SERVER FOR BIOLOGY

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

**UC San Diego**
**JACOBS SCHOOL OF ENGINEERING**

# Additional Slides

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

**UC San Diego**
**JACOBS SCHOOL OF ENGINEERING**

# Modes of operation

**Accurate mode:**

Sampling of genes from species → LASTZ → Filtering → PASTA → RAxML-NG → ASTRAL-Pro

**Balanced mode:**

Sampling of genes from species → LASTZ → Filtering → PASTA → FastTree → ASTRAL-Pro

**Fast mode:**

Sampling of genes from species → LASTZ → Filtering → Mashtree → ASTRAL-Pro

Higher accuracy and runtime

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

**UC San Diego**
**JACOBS SCHOOL OF ENGINEERING**

# ROADIES is more accurate than MashTree



Accuracy of Final Tree of MashTree with all datasets

| Dataset | Normalized RF distance at Species-level | Normalized RF distance at Order-level |
|---------|------------------------------------------|----------------------------------------|
| Flies | 0.15 | 0.14 |
| Mammals | 0.2 | 0.61 |
| Birds | 0.25 | 0.71 |

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING
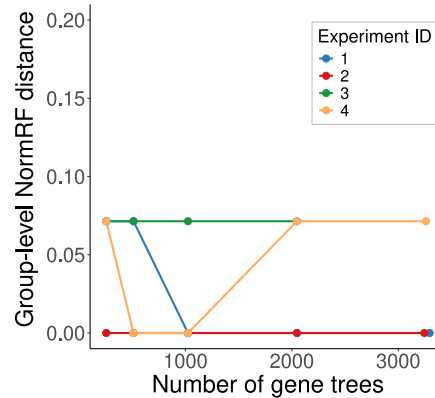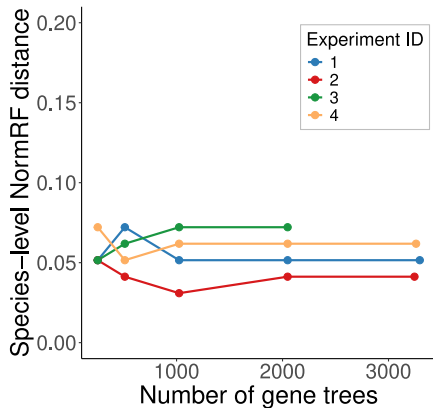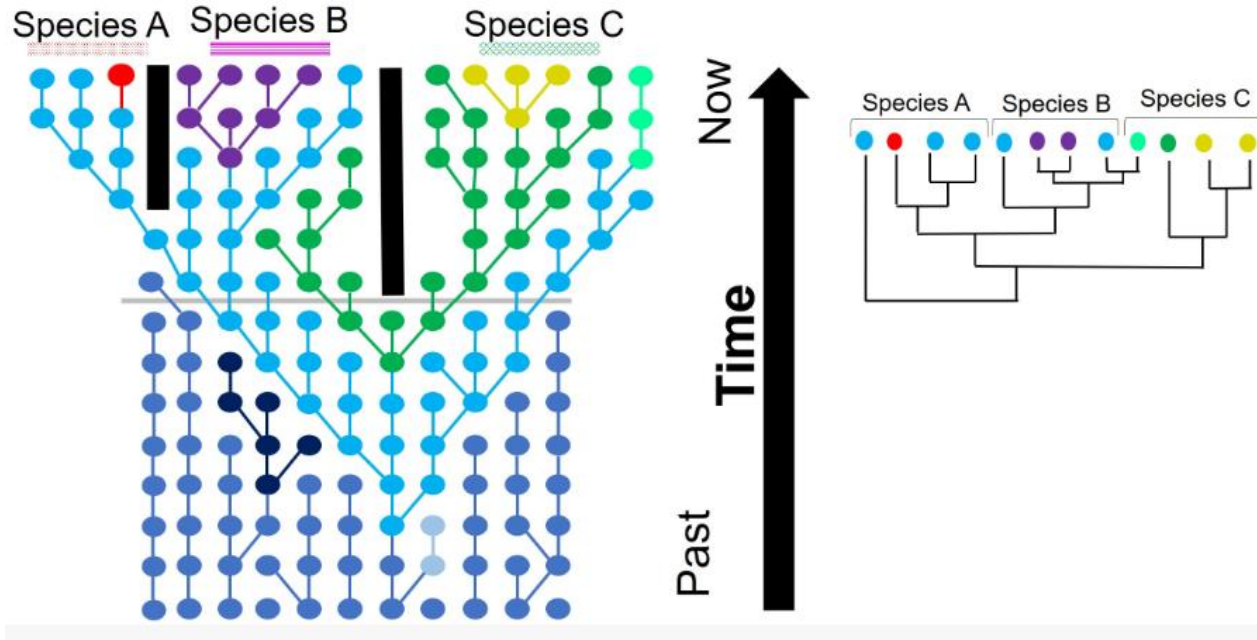
# ROADIES gives stable results

Even if ROADIES randomly samples genes, results are consistent



Experiments run on AWS R6a instances
Runtime is calculated as wall clock time
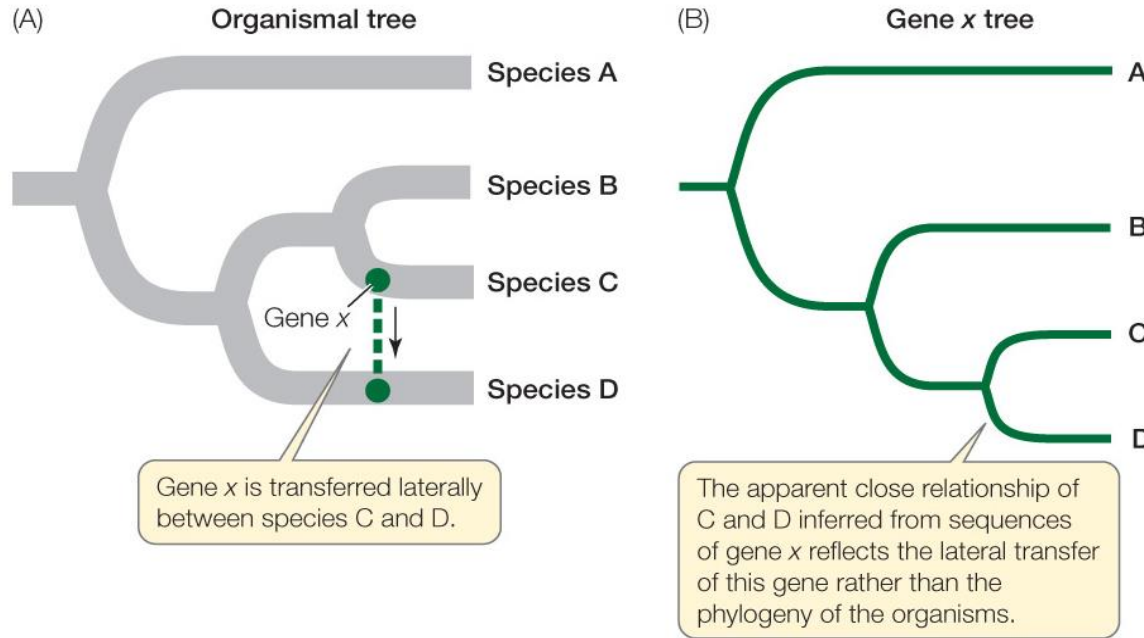Variance experiments are tested with Drosophila datasets

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# Causes of Gene Tree Discordance

Incomplete Lineage Sorting



Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

**UC San Diego**
**JACOBS SCHOOL OF ENGINEERING**

# Causes of Gene Tree Discordance

Horizontal Gene Transfer



(A) Organismal tree

Species A
Species B
Species C
Species D

Gene x

Gene x is transferred laterally between species C and D.

(B) Gene x tree

A
B
C
D

The apparent close relationship of C and D inferred from sequences of gene x reflects the lateral transfer of this gene rather than the phylogeny of the organisms.

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES
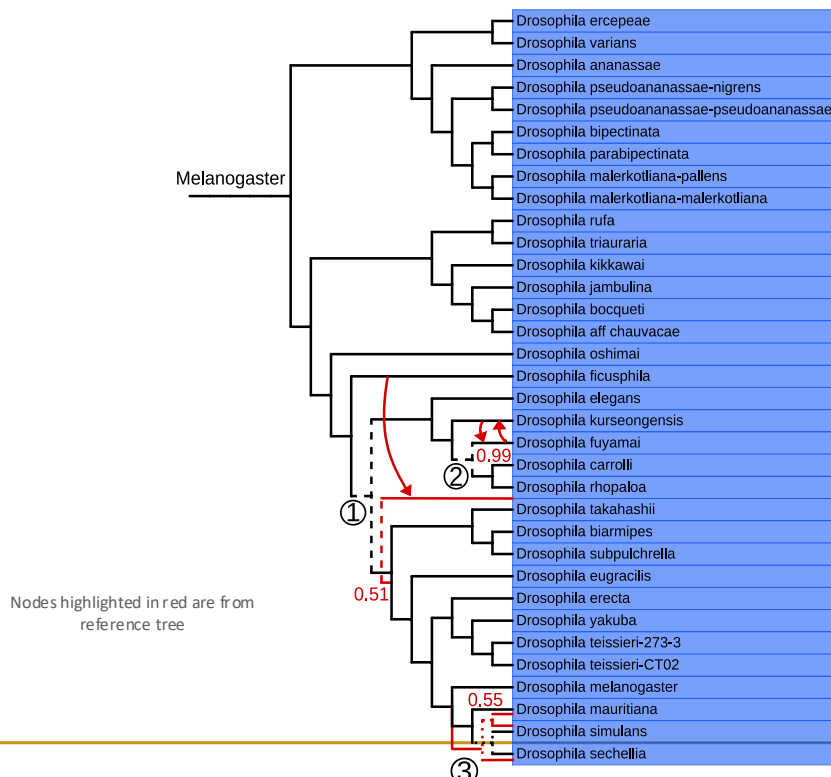
UC San Diego
JACOBS SCHOOL OF ENGINEERING

# ROADIES estimates exact same phylogeny of 100 Drosophilid genomes at group-level



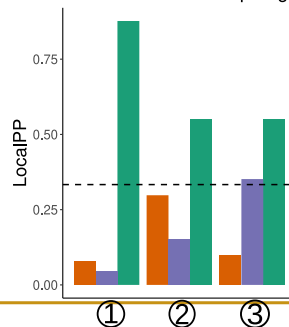Group-level normalized RF distance with reference tree = 0

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# ROADIES' differences with reference is limited to low-confident branches



Nodes highlighted in red are from reference tree

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego

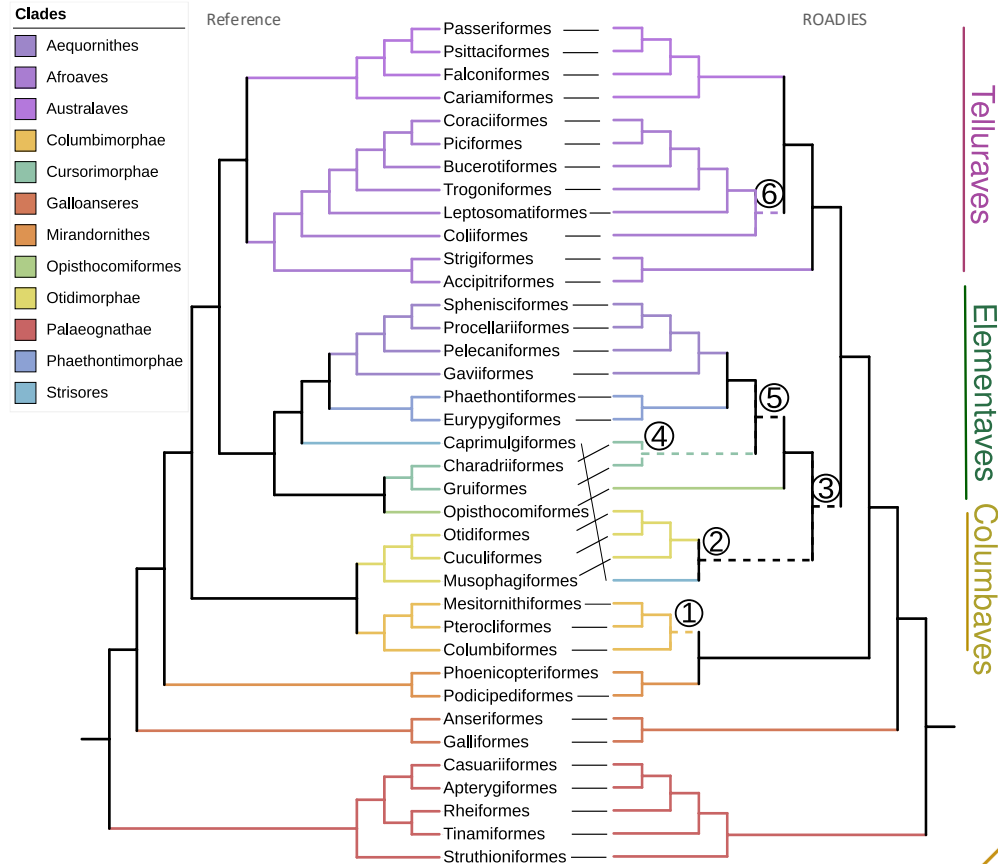JACOBS SCHOOL OF ENGINEERING

# ROADIES convergence results



Experiments run on AWS R6a 16-core instances
Runtime is calculated as wall clock time

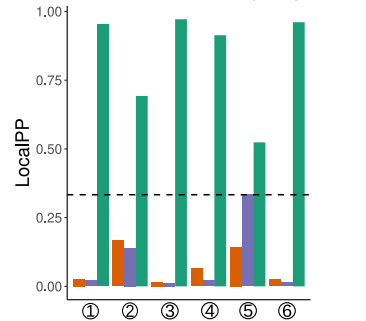Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING

# Order-level phylogeny



Order-level normalized RF distance with reference tree = 0.28

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES
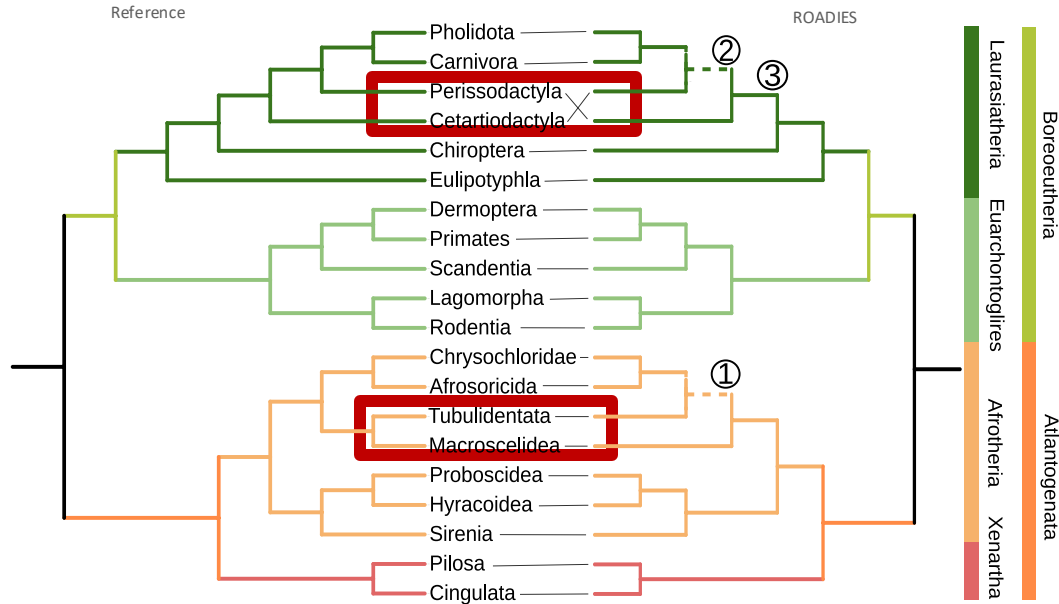
**UC San Diego**
**JACOBS SCHOOL OF ENGINEERING**

# ROADIES convergence results



Experiments run on AWS R6a 16-core instances
Runtime is calculated as wall clock time

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

**UC San Diego**
**JACOBS SCHOOL OF ENGINEERING**

# Low-confident branches are debatable



Quartet scores of different topologies

LocalPP of different topologies

Accurate, scalable, and fully automated inference of species trees from raw genome assemblies using ROADIES

UC San Diego
JACOBS SCHOOL OF ENGINEERING