

# Multi-locus multi-species coalescent phylogenetic inference

Huw A. Ogilvie

Research School of Biology, Australian National University

Centre for Computational Evolution, University of Auckland

February 2017



# Introduction

## (Re-)introducing the multispecies coalescent

The multi-species  
coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

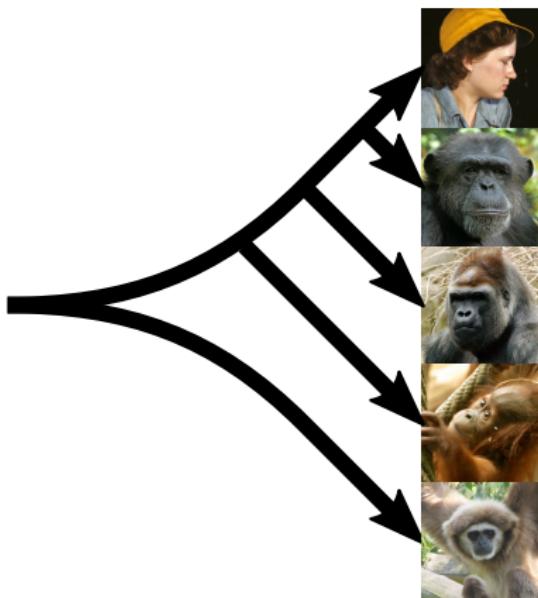
MSC accuracy

MSC implementations

Multi-species-network  
coalescent

Acknowledgements

# Speciation processes



- ▶ Evolution of species
- ▶ Birth-death process  
(speciation & extinction)
- ▶ Diversification rate  
 $\text{Birth} - \text{Death} = \lambda$
- ▶ Extinction ratio  
 $\text{Death} \div \text{Birth} = \nu$
- ▶ Forward in time

The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

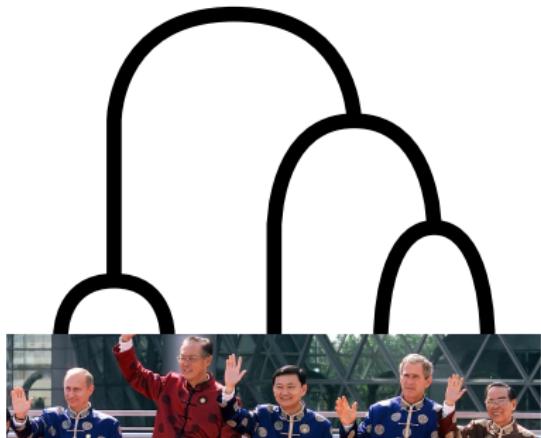
MSC accuracy

MSC implementations

Multi-species-network coalescent

Acknowledgements

# Coalescent processes



- ▶ Evolution of genes
- ▶ Dependent on effective population size  $N_e$
- ▶ Constant, linear, exponential or stepwise
- ▶ Backwards in time

The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

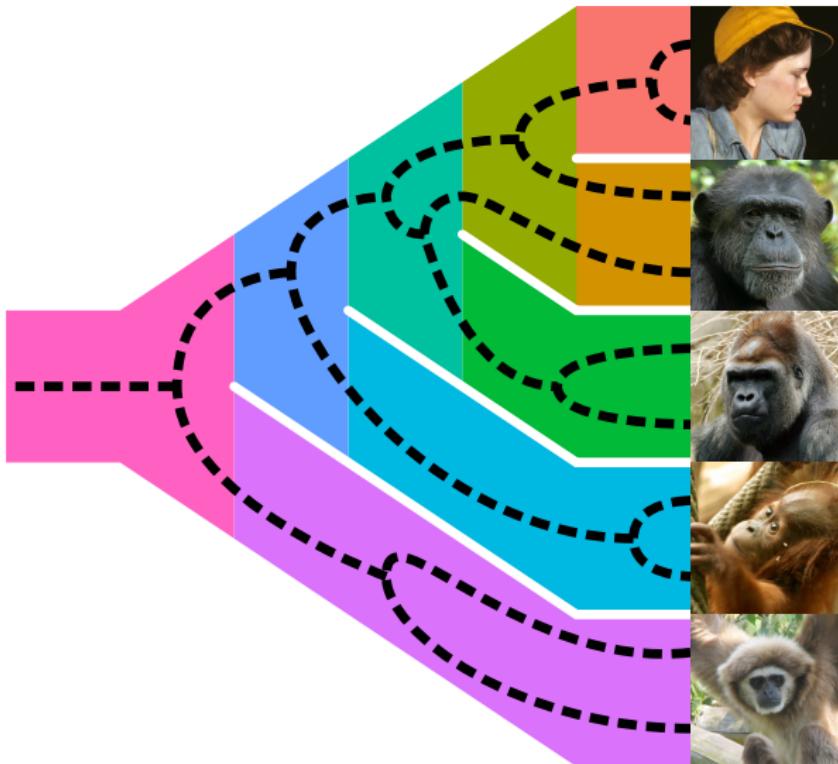
MSC accuracy

MSC implementations

Multi-species-network coalescent

Acknowledgements

# Putting them together



The multi-species  
coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

MSC accuracy

MSC implementations

Multi-species-network  
coalescent

Acknowledgements

# Why use the MSC

Why the MSC is important  
*or*  
reasons to use MSC methods

The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

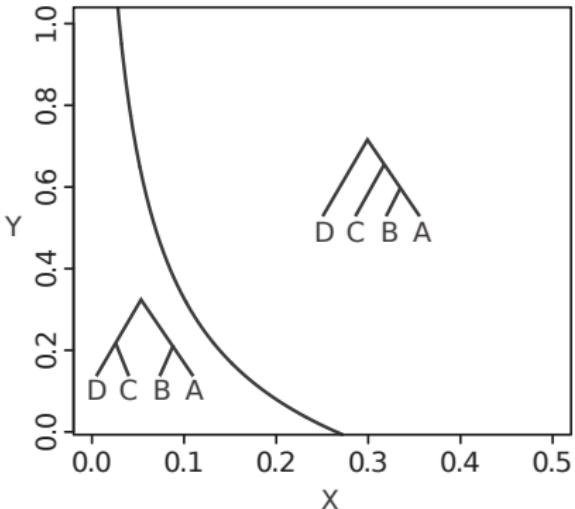
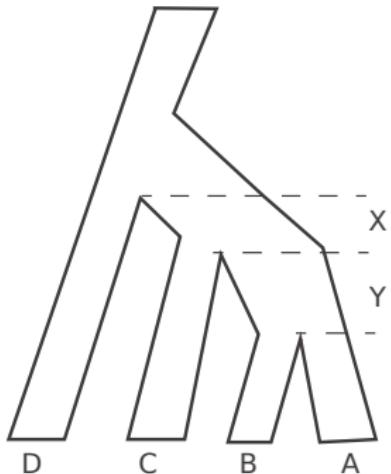
MSC accuracy

MSC implementations

Multi-species-network coalescent

Acknowledgements

# Welcome to the anomaly zone



Linkem, Minin and Leaché (2016) Systematic Biology

The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

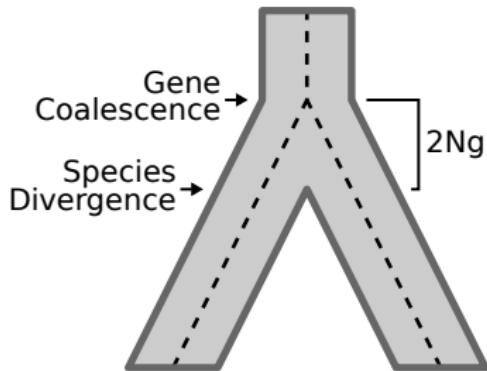
MSC accuracy

MSC implementations

Multi-species-network coalescent

Acknowledgements

# Impact on estimated divergence times



- ▶ Human-chip divergence time  $\approx 4\text{my}$
- ▶ Ancestral effective population size  $N \approx 50000$
- ▶ Ancestral generation time  $g \approx 20\text{y}$
- ▶  $2Ng = 2 \cdot 50000 \cdot 20\text{y} = 2\text{my}$
- ▶ Error in estimate will be  $4\text{my} + 2\text{my} = 6\text{my}$

Hobolth *et al.* (2011) Genome Research

The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

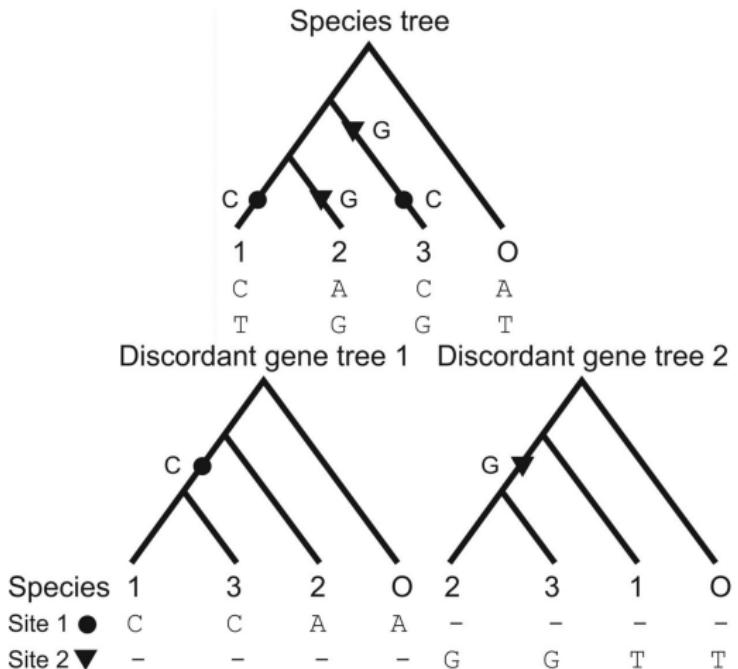
MSC accuracy

MSC implementations

Multi-species-network coalescent

Acknowledgements

# It gets worse



The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

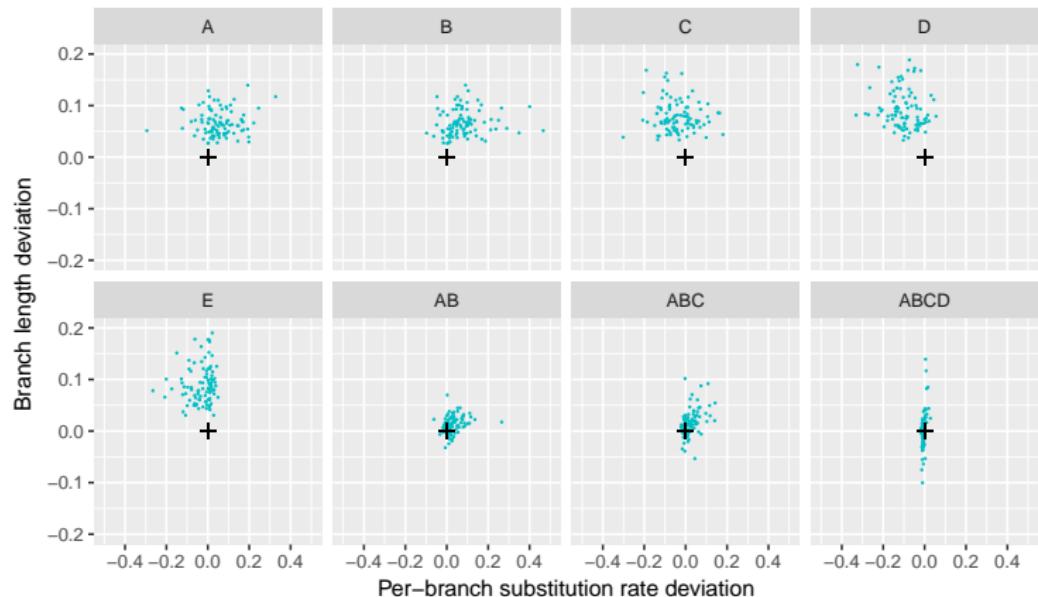
MSC accuracy

MSC implementations

Multi-species-network coalescent

Acknowledgements

# The impact of SPILS



The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

MSC accuracy

MSC implementations

Multi-species-network coalescent

Acknowledgements

Ogilvie, Bouckaert & Drummond (2017, if the journal gods are kind)

## Estimating genes trees and species trees

The multi-species  
coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

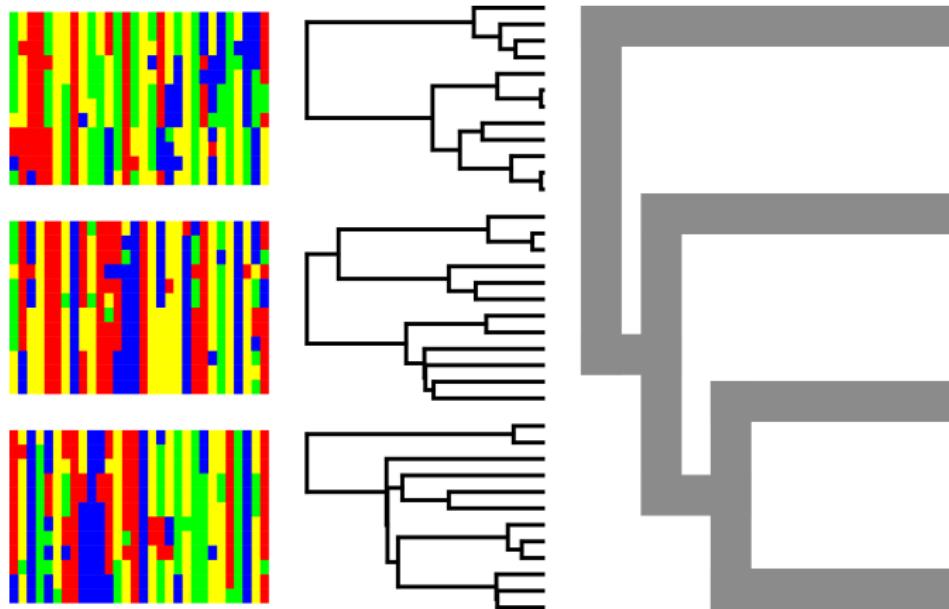
MSC accuracy

MSC implementations

Multi-species-network  
coalescent

Acknowledgements

# Sequence alignments to species trees



The multi-species  
coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

MSC accuracy

MSC implementations

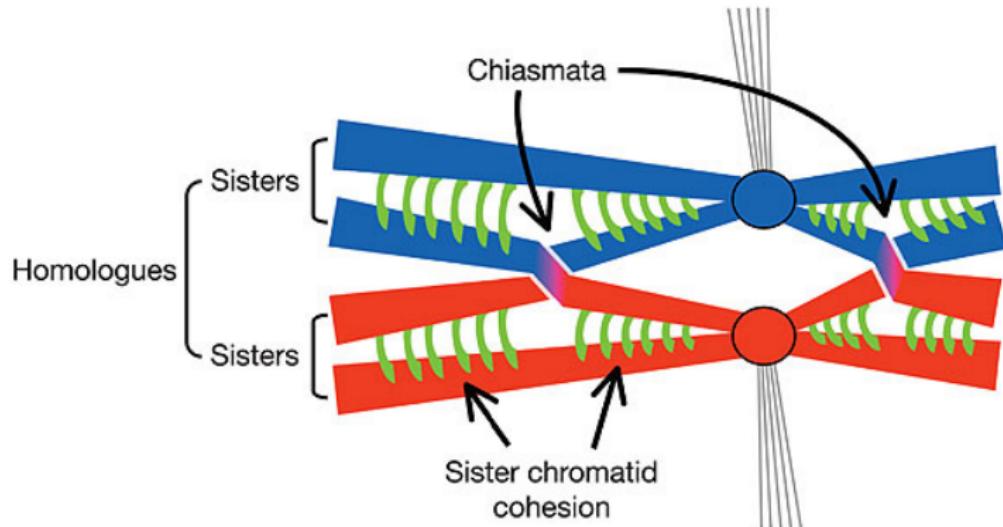
Multi-species-network  
coalescent

Acknowledgements

$$\Pr(S|D) = \frac{\Pr(D|G) \cdot \Pr(G|S) \cdot \Pr(S)}{\Pr(D)} \quad (1)$$

# Recombination I

Multi-locus MSC



Neale & Keeney (2006) Nature

The multi-species  
coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

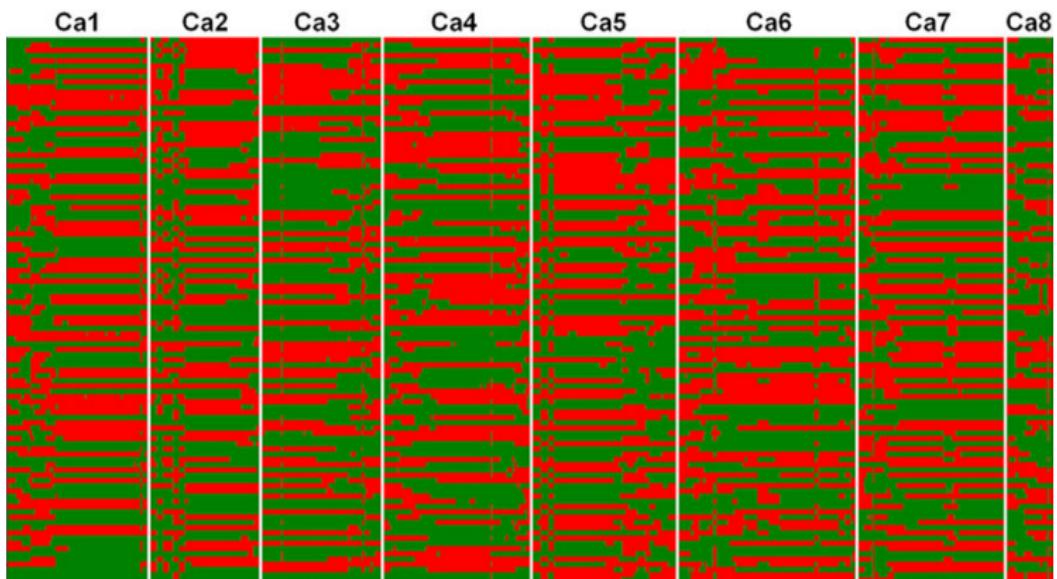
MSC accuracy

MSC implementations

Multi-species-network  
coalescent

Acknowledgements

# Recombination II



The multi-species  
coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

MSC accuracy

MSC implementations

Multi-species-network  
coalescent

Acknowledgements

Kale *et al.* (2015) Scientific Reports

# Multi-locus MSC assumptions

- ▶ Sites within a locus are linked
  - ▶ Sequences should be short enough that recombination is unlikely within a locus
  - ▶ Deeper trees → more recombination → shorter sequences
- ▶ No linkage between loci
  - ▶ Distances between loci should be large enough that recombination is common
  - ▶ Deeper trees → more recombination → shorter distances

The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

MSC accuracy

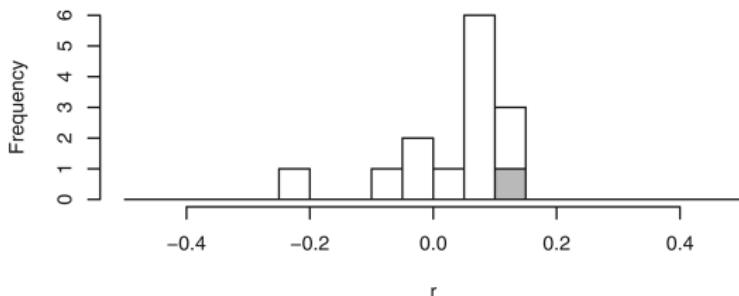
MSC implementations

Multi-species-network coalescent

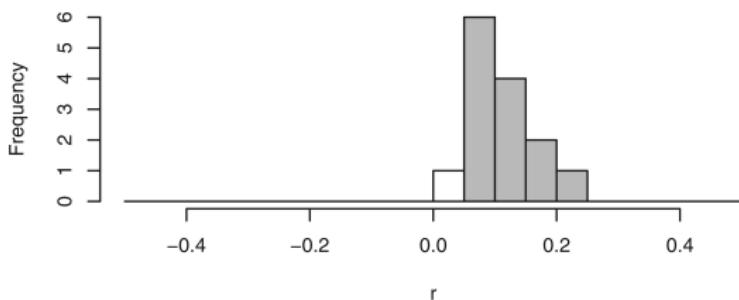
Acknowledgements

# Correlation of phylogenetic signal

Two exons from a single gene



Two halves of a single exon



The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

MSC accuracy

MSC implementations

Multi-species-network coalescent

Acknowledgements

# Species tree clocks

## Species tree clocks in StarBEAST2

The multi-species  
coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

MSC accuracy

MSC implementations

Multi-species-network  
coalescent

Acknowledgements

# Per-species substitution rates

- ▶ Differences exist in substitution (clock) rates between loci
- ▶ Differences also exist in substitution rates between **species**
- ▶ But phylogenetic likelihood requires clock rates be applied to **gene** branches
- ▶ Species tree relaxed clocks extend the MSC to derive gene tree branch rates from species tree branch rates

The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

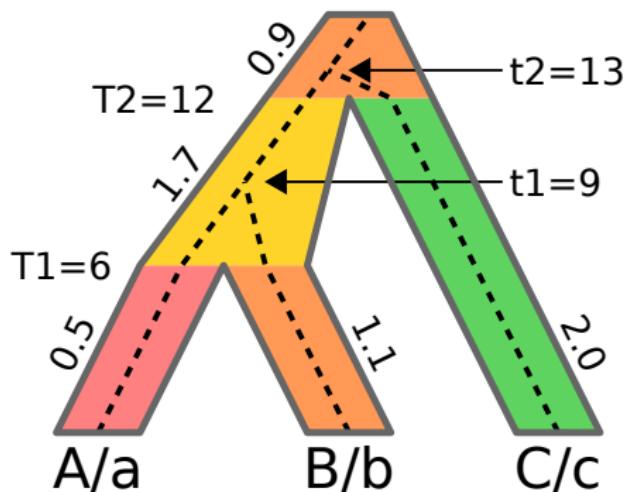
MSC accuracy

MSC implementations

Multi-species-network coalescent

Acknowledgements

# Computing clock rates I



$$\text{rate}_{\text{branch}} = \text{rate}_{\text{locus}} \sum^{\text{species}} \left( \text{rate}_{\text{species}} \times \frac{\text{overlap}_{\text{species,branch}}}{\text{length}_{\text{branch}}} \right) \quad (2)$$

The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

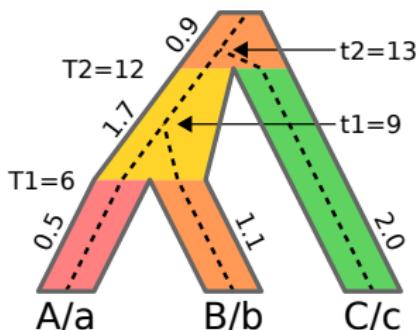
MSC accuracy

MSC implementations

Multi-species-network coalescent

Acknowledgements

# Computing clock rates II



The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

MSC accuracy

MSC implementations

Multi-species-network coalescent

Acknowledgements

$$\text{rate}_a = 10^{-9} \left( 0.5 \times \frac{6.0}{9.0} + 1.7 \times \frac{3.0}{9.0} \right) = 0.9 \times 10^{-9} \quad (3)$$

$$\text{rate}_b = 10^{-9} \left( 1.1 \times \frac{6.0}{9.0} + 1.7 \times \frac{3.0}{9.0} \right) = 1.3 \times 10^{-9} \quad (4)$$

$$\text{rate}_{ab} = 10^{-9} \left( 1.7 \times \frac{3.0}{4.0} + 0.9 \times \frac{1.0}{4.0} \right) = 1.5 \times 10^{-9} \quad (5)$$

# MSC accuracy

## The power of the multi-locus MSC

The multi-species  
coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

MSC accuracy

MSC implementations

Multi-species-network  
coalescent

Acknowledgements

# Simulation details

- ▶ 21 species, two phased sequences per species
- ▶ 600nt long sequence alignments
- ▶ Substitution rate  $\approx 10^{-3}$  per million years
- ▶  $N_e \approx 2$  million assuming annual generation time
- ▶ Species tree height  $\approx 34$  million years
- ▶ “StarBEAST2 brings faster species tree inference and accurate estimates of substitution rates” on bioRxiv

The multi-species  
coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

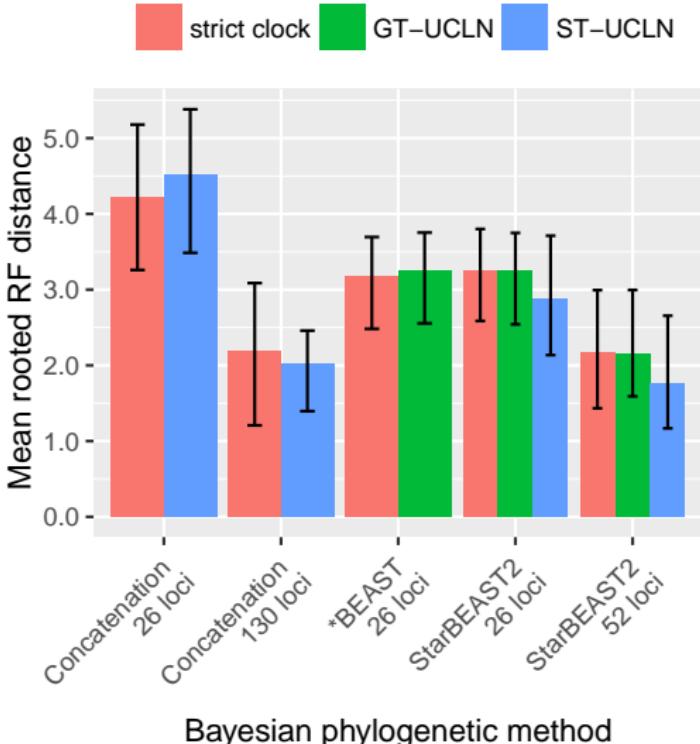
MSC accuracy

MSC implementations

Multi-species-network  
coalescent

Acknowledgements

# Topological accuracy



The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

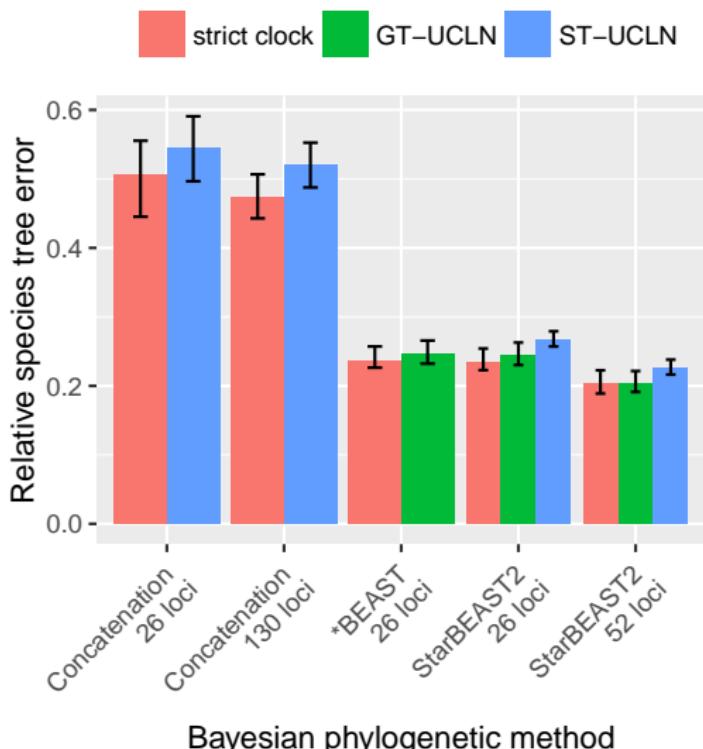
MSC accuracy

MSC implementations

Multi-species-network coalescent

Acknowledgements

# Branch length accuracy



The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

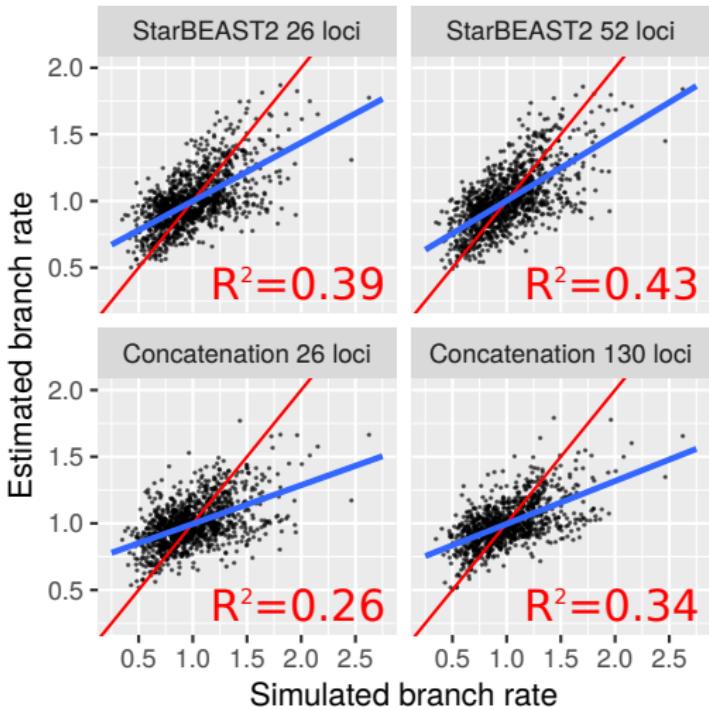
MSC accuracy

MSC implementations

Multi-species-network coalescent

Acknowledgements

# Substitution rate accuracy



The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

MSC accuracy

MSC implementations

Multi-species-network coalescent

Acknowledgements

# MSC implementations

## Implementations of the multi-locus MSC

The multi-species  
coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

MSC accuracy

MSC implementations

Multi-species-network  
coalescent

Acknowledgements

# StarBEAST2

- ▶ Successor to \*BEAST, a method originally developed by Joseph Heled and Alexei Drummond
- ▶ 10 to 30 × faster than \*BEAST — depending on the data and model
- ▶ Can integrate out population sizes analytically — slightly faster if population sizes are nuisance parameters
- ▶ The first method to implement species tree relaxed clocks

The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

MSC accuracy

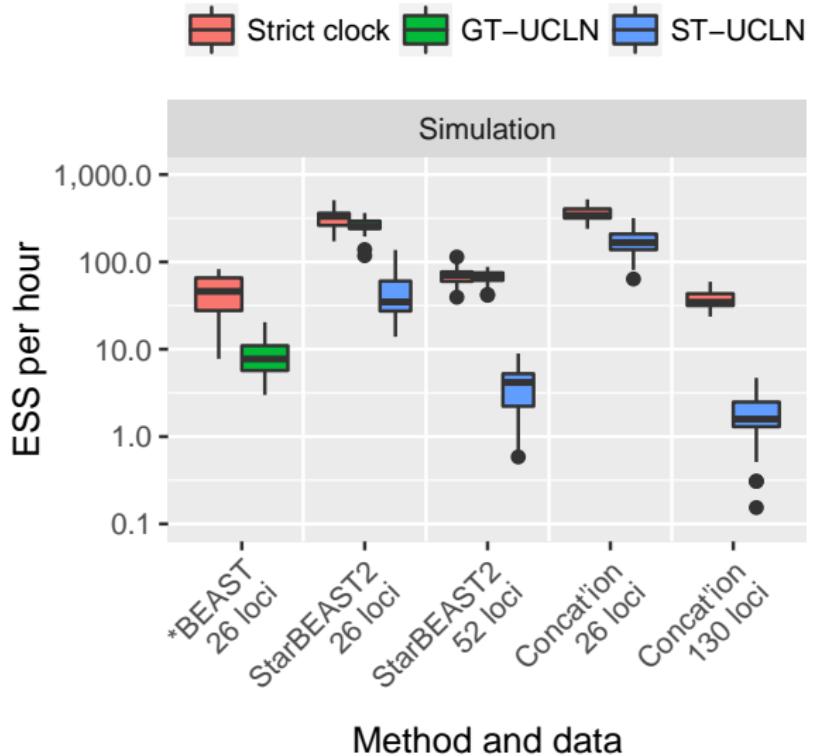
MSC implementations

Multi-species-network coalescent

Acknowledgements

# StarBEAST2 performance

Multi-locus MSC



The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

MSC accuracy

MSC implementations

Multi-species-network coalescent

Acknowledgements

# Other methods

- BPP Simultaneous species delimitation and species tree inference. Limited to Jukes-Cantor, strict clocks and gamma priors.
- RevBayes A challenger approaches!
- ASTRAL Can estimate topology but not branch lengths  
(except in coalescent units which are useless)
- MP-EST Like ASTRAL but does a terrible job at estimating topology (what's the point???)

The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

MSC accuracy

MSC implementations

Multi-species-network coalescent

Acknowledgements

# Multi-species-network coalescent

The multi-species  
coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

MSC accuracy

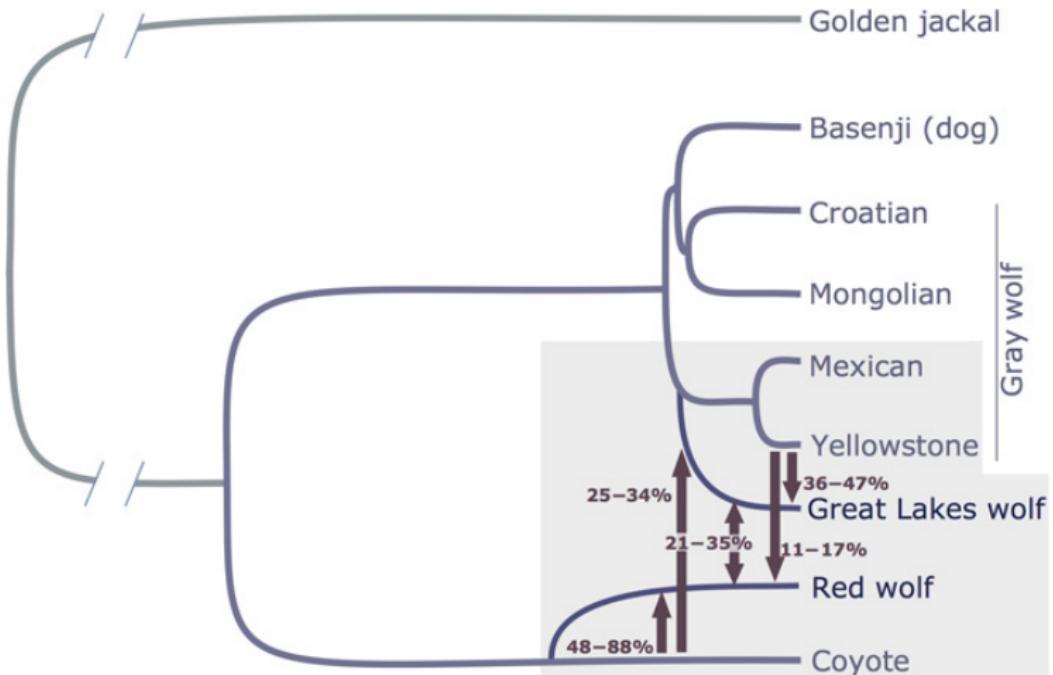
MSC implementations

Multi-species-network  
coalescent

Acknowledgements

Because biology likes to play  
tricks on biologists

# North American canids



The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

MSC accuracy

MSC implementations

Multi-species-network coalescent

Acknowledgements

vonHoldt *et al.* (2016) Science Advances

# MSNC implementations

Multi-locus MSNC: joint estimation of gene trees, species networks and inheritance probabilities

The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

MSC accuracy

MSC implementations

Multi-species-network coalescent

Acknowledgements

PhyloNet Limited range of substitution models and prior distributions implemented

BEAST2 Species Networks Chi Zhang, Alexei Drummond, Tanja Stadler and me

# Acknowledgements

## Auckland

Alexei Drummond  
Tim Vaughan  
Remco Bouckaert  
Joseph Heled

## Canberra

Craig Moritz  
Jason Bragg

## Basel

Tanja Stadler  
Chi Zhang

The multi-species coalescent

Why use the MSC

Multi-locus MSC

Species tree clocks

MSC accuracy

MSC implementations

Multi-species-network coalescent

Acknowledgements

## Photo credits:

Human, Howard R. Hollem; public domain U.S. Government work  
Chimpanzee, Clément Bardot; CC-BY-SA  
Gorilla, Kabir Bakie; CC-BY-SA  
Orangutan, "Winkelbohrer"; CC-BY-SA  
Gibbon, Kabir Bakie; CC-BY-SA  
APEC group photo, Stephen Jaffe/Agence France-Presse/Getty Images