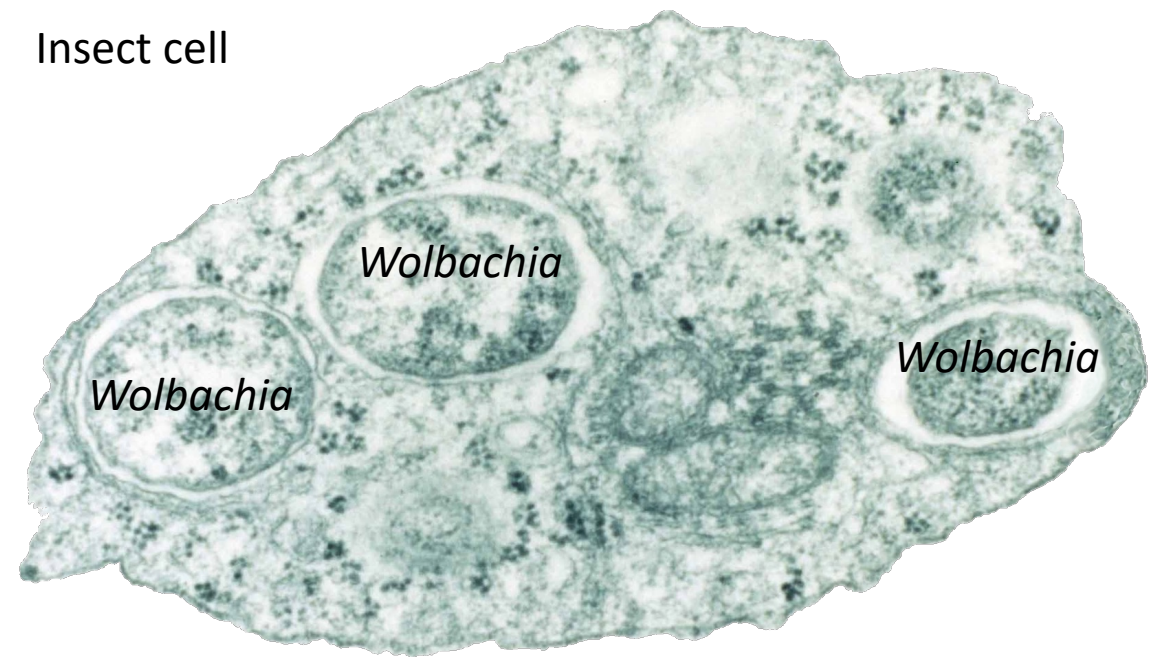# Phylogenetic Network Estimation for *Wolbachia*

*Paul Zaharias, Danny Rice, Tandy Warnow, and Irene Newton*

# The parasitic bacteria *Wolbachia*

- Infects the reproductive tracts of arthropods and nematods

- 25 to 70% of all insect species are potential hosts

- Impacts on human health (e.g., Dengue fever)

Insect cell

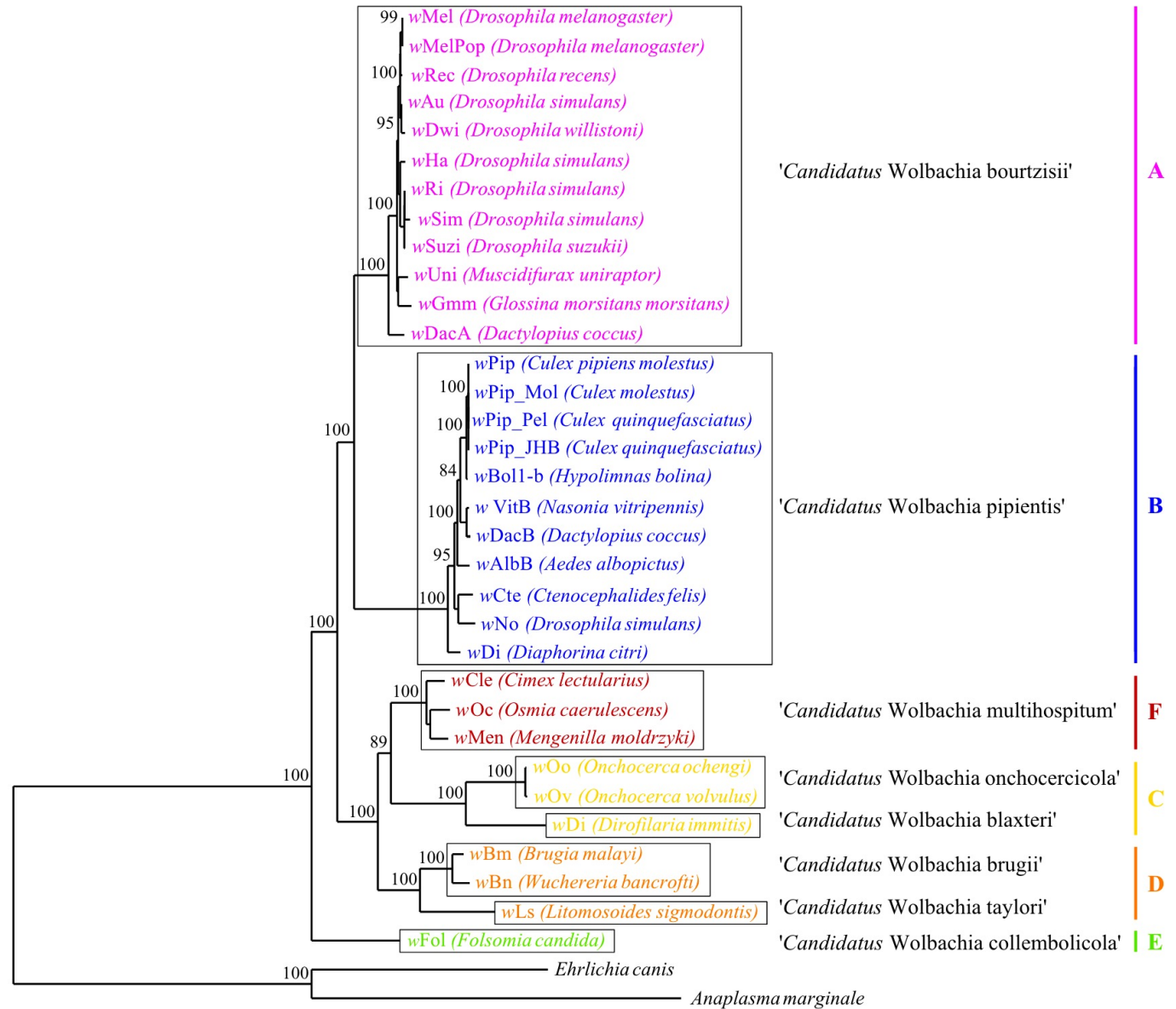*Wolbachia*

*Wolbachia*

*Wolbachia*

# Questions of interest

- What genes are under selection?

- Can we improve species delimitation?

- What is the percentage of horizontally-transferred genes in *Wolbachia?*

# One of the trees from the literature

Ramírez-Puebla et al.,
Syst. and Appl. Microbiology, 2016



Ramírez-Puebla, Shamayim T., et al. "A response to Lindsey et al." Wolbachia pipientis should not be split into multiple species: A response to Ramírez-Puebla et al."." *Systematic and applied microbiology* 39.3 (2016): 223-225.
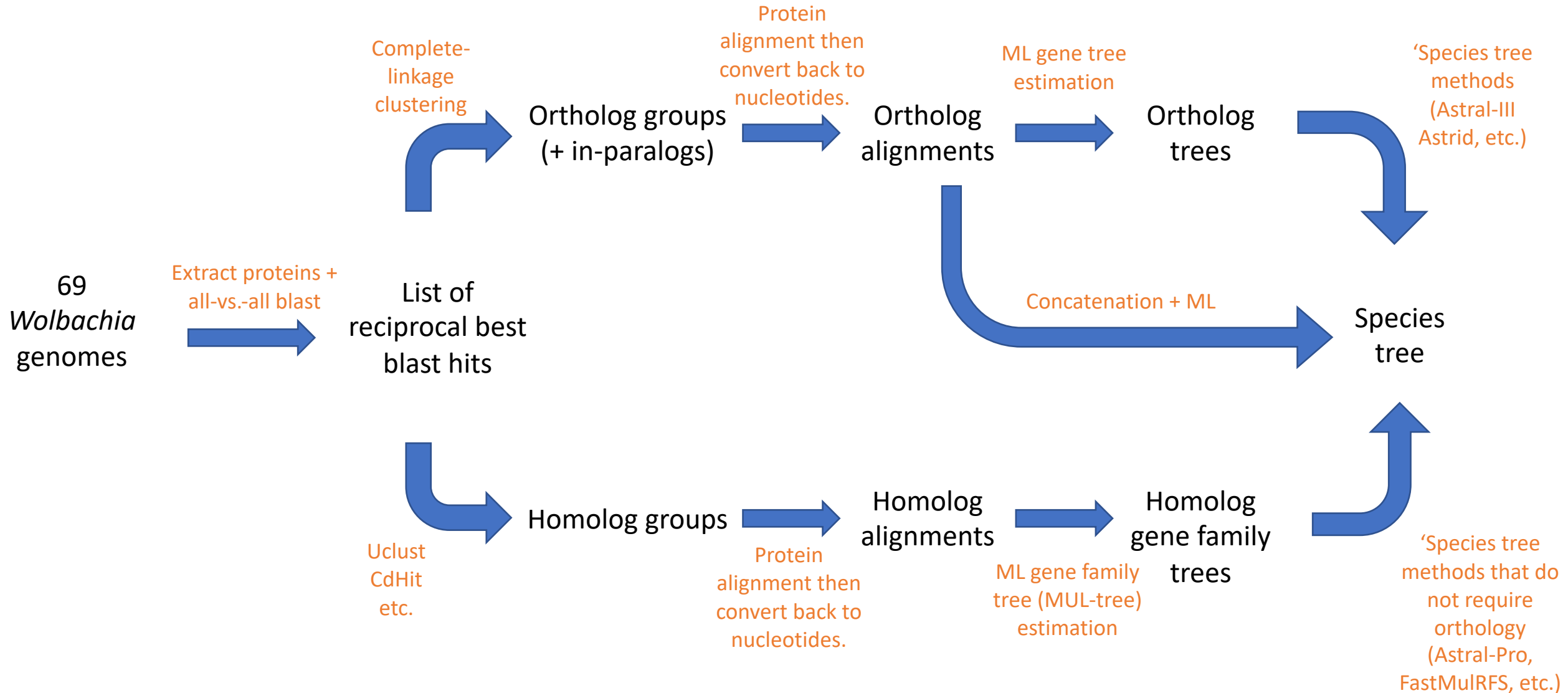
# Species tree estimation options

- From orthologs:
  - Concatenation using maximum likelihood: statistically inconsistent, computationally intensive, can be accurate
  - Summary methods (that combine gene trees), such as ASTRAL, ASTRID: faster than concatenation, statistically consistent if given true gene trees, can be accurate
  - Site-based methods, such as LilyQ and SVDquartets/SVDquest: statistically consistent, limited to perhaps 100 leaves due to computational challenges
- The choice between methods depends on the data

# Species tree estimation options

- From homologs:
  - Construct gene family trees and combine
    - ASTRAL-Pro (variant of ASTRAL to handle gene family trees)
    - DISCO+ASTRID (decomposes gene family trees into single copy trees, then runs ASTRID)
    - FastMulRFS (supertree method applied to MUL-trees)
  - Variant of concatenation
    - Construct gene family trees, decompose using DISCO, then use concatenation
- Statistical consistency without knowing orthology is achieved in some cases
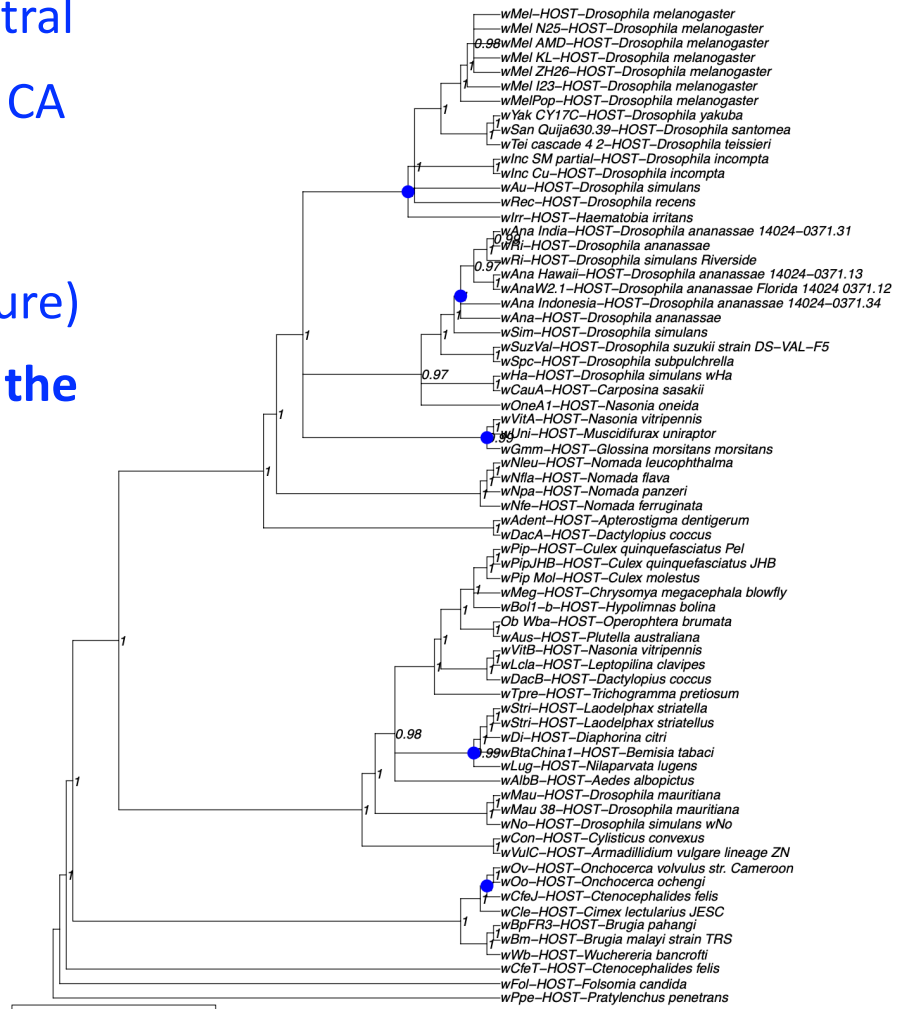- The choice between methods depends on the data
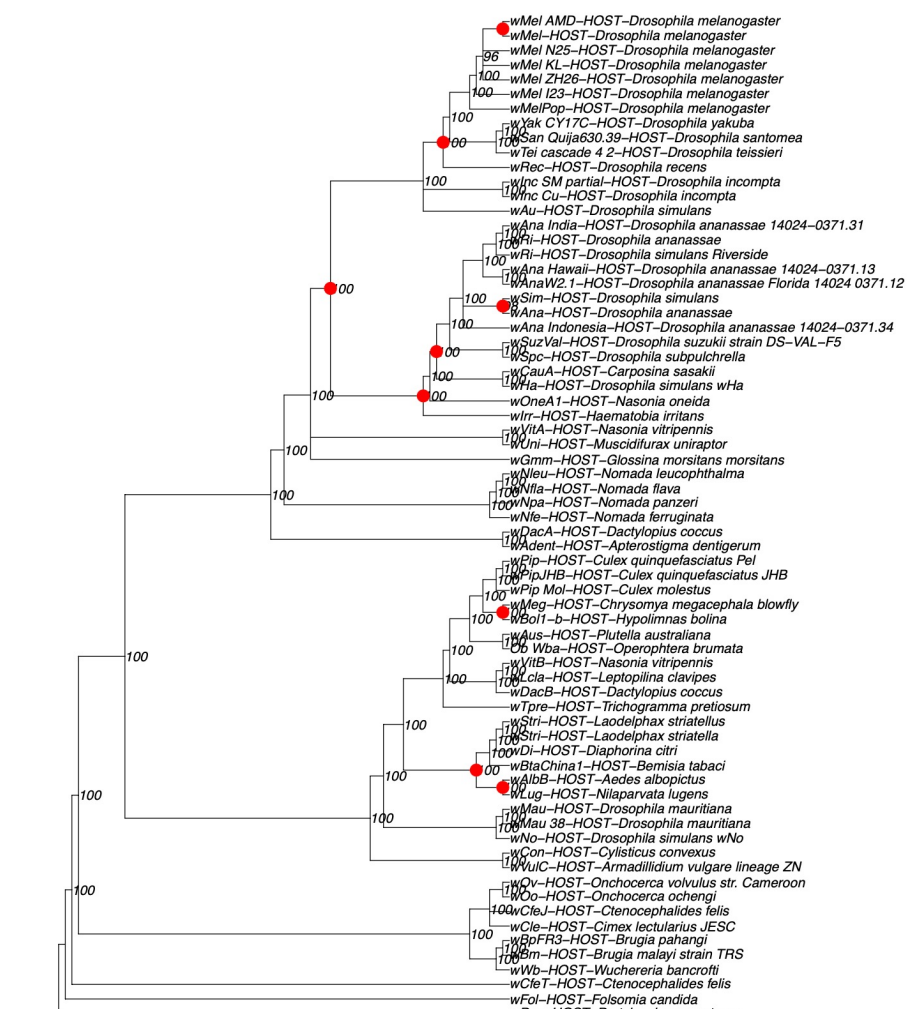
# What we did: genomes to species trees

# ASTRAL and Concatenation-IQTree: very similar

## Astral

## CA (IQ-TREE)

- Collapsed pp < 0.95 for Astral

- Collapsed UfBoot < 95 for CA

- Both trees are very well supported, but with some differences (marked in figure)

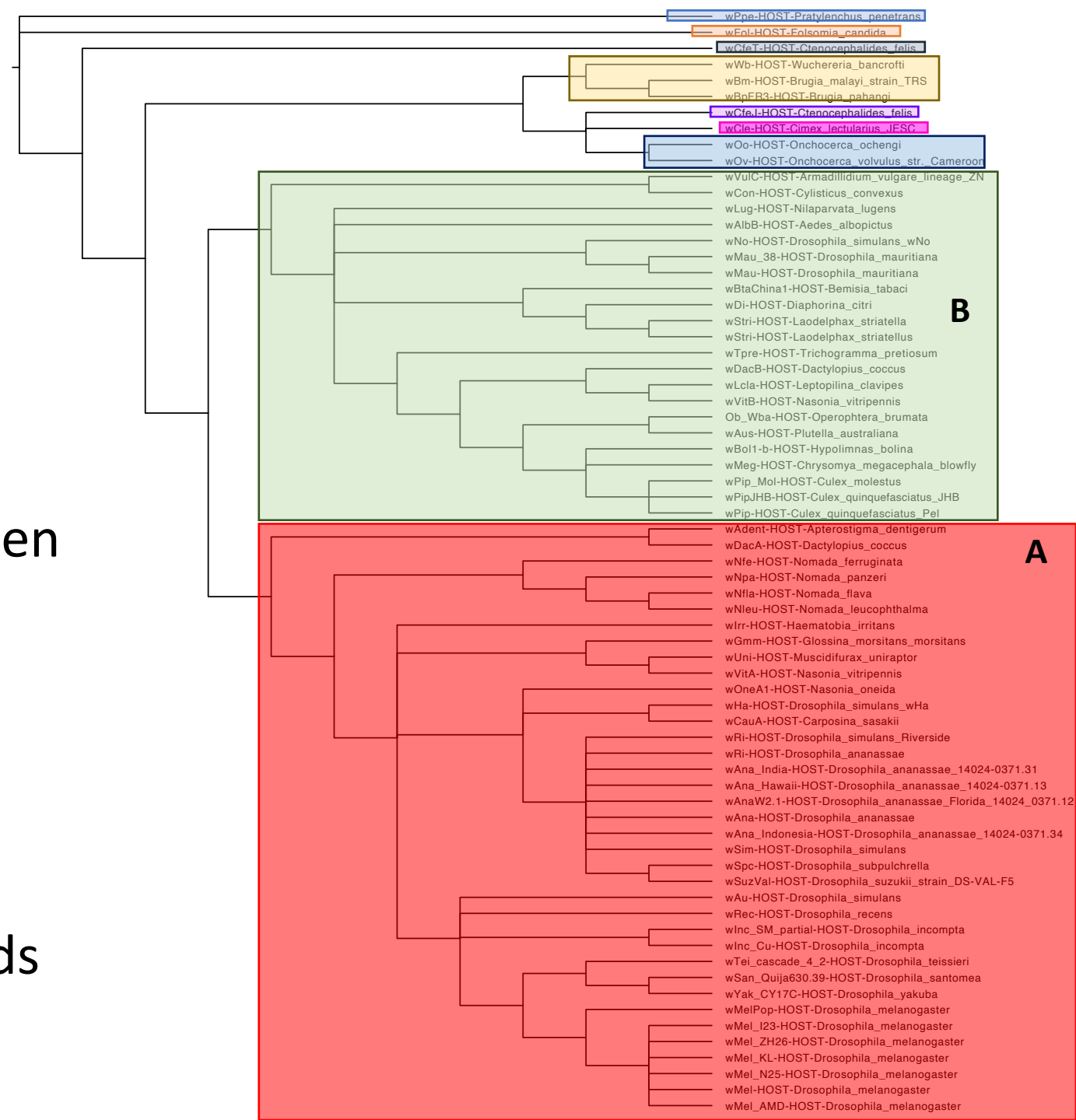- **Are this incompatibilities the result of HGT events or something else?**
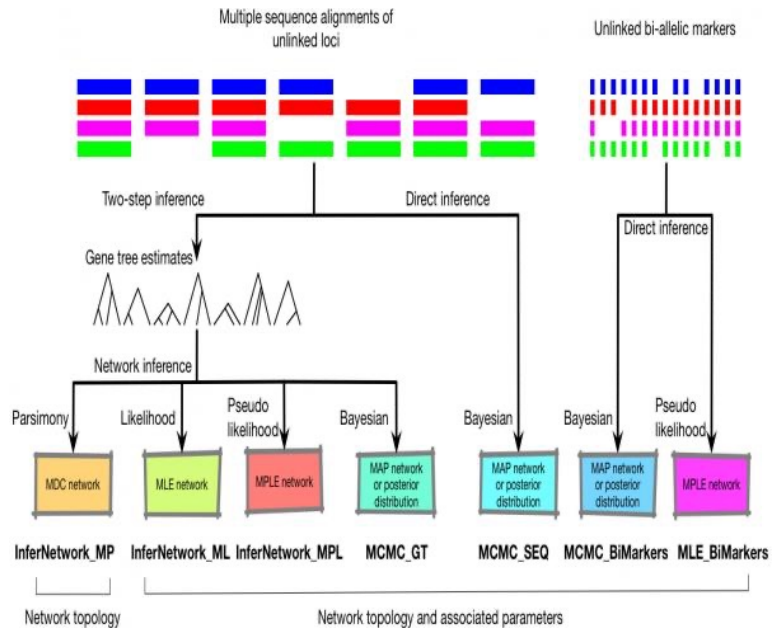
# Reticulation events

- Consensus tree

- Between supergroups ( = 'species') ? In particular, between supergroup A and B (most studied)

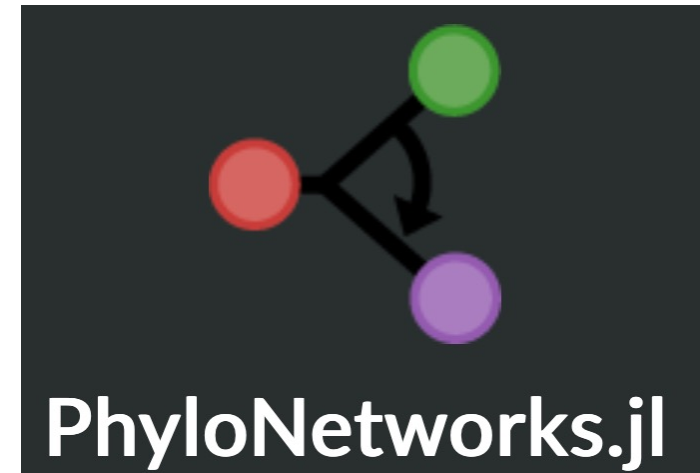- Within 'species' ?

- → Phylogenetic network methods

# Phylonet (Nakhleh) and PhyloNetworks (Solis-Lemus)

## Phylonet (Luay Nakhleh)

## PhyloNetworks (Claudia Solís-Lemus)

# Phylogenetic Network Pipeline

- In all analyses, assumption is you have gene trees and you are seeking a phylogenetic network

- Three variants of the problem:
  - Given an underlying species tree, add reticulations (extra edges) to create network (easier)
  - Given an underlying species tree, add reticulations and allow tree to change (a bit harder)
  - Construct the phylogenetic network directly (much harder)

# Phylogenetic Network Pipeline

- In all analyses, assumption is you have gene trees and you are seeking a phylogenetic network

- Three variants of the problem:
  - Given an underlying species tree, add reticulations (extra edges) to create network (easier)
  - Given an underlying species tree, add reticulations and allow tree to change (a bit harder)
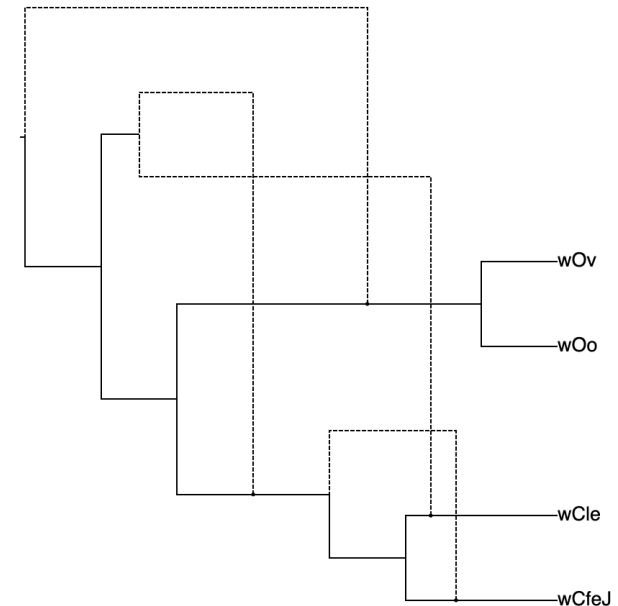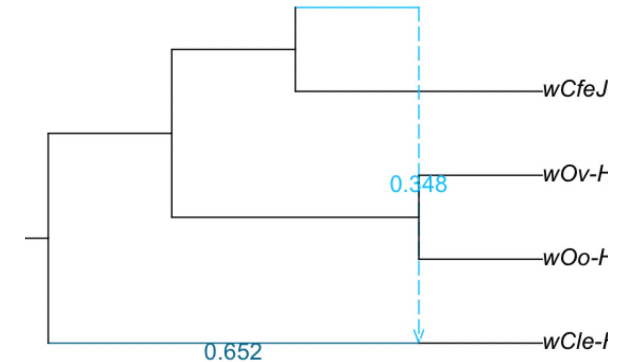  - Construct the phylogenetic network directly (much harder)

# How do you proceed?

- Maximum Parsimony, Maximum Likelihood, Maximum Pseudo-Likelihood:
  - Fix a maximum number of 'hybridization' events and optimize a criterion
  - Choosing between different networks with different numbers of hybridization events is not well solved
  - Knowing the number of hybridization events is also not well solved

- Bayesian approaches: deal with the model complexity but limited to very small datasets (maybe 4 or 5 leaves)
  - **MCMC_GT** : input rooted gene trees
  - **MCMC_SEQ** : input alignments (direct inference)

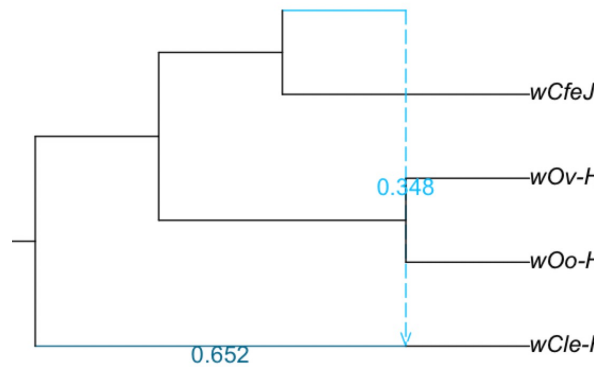# Problem n°1 : different methods return different number of reticulations

Only considering one 4-taxon 'species':

- SNaQ → 1 reticulation
  (too conservative?)

- MCMC_SEQ (direct Bayesian inference) →
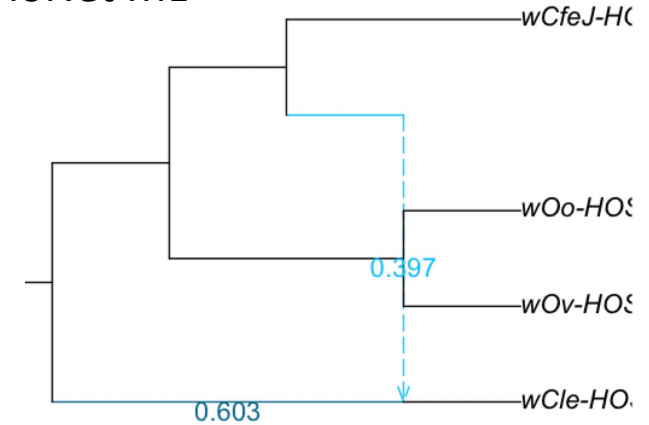  converges to 4 reticulations
  (too aggressive?)

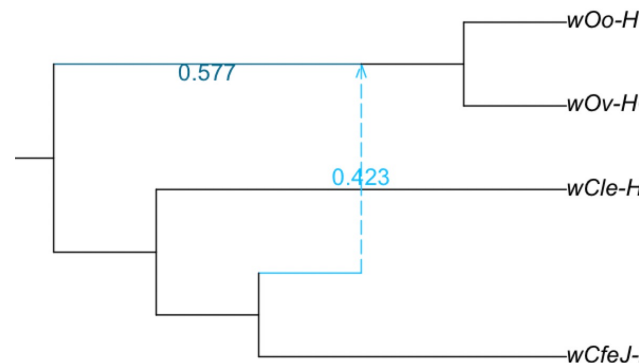# Problem n°2: For fixed number of reticulations, different methods return different networks
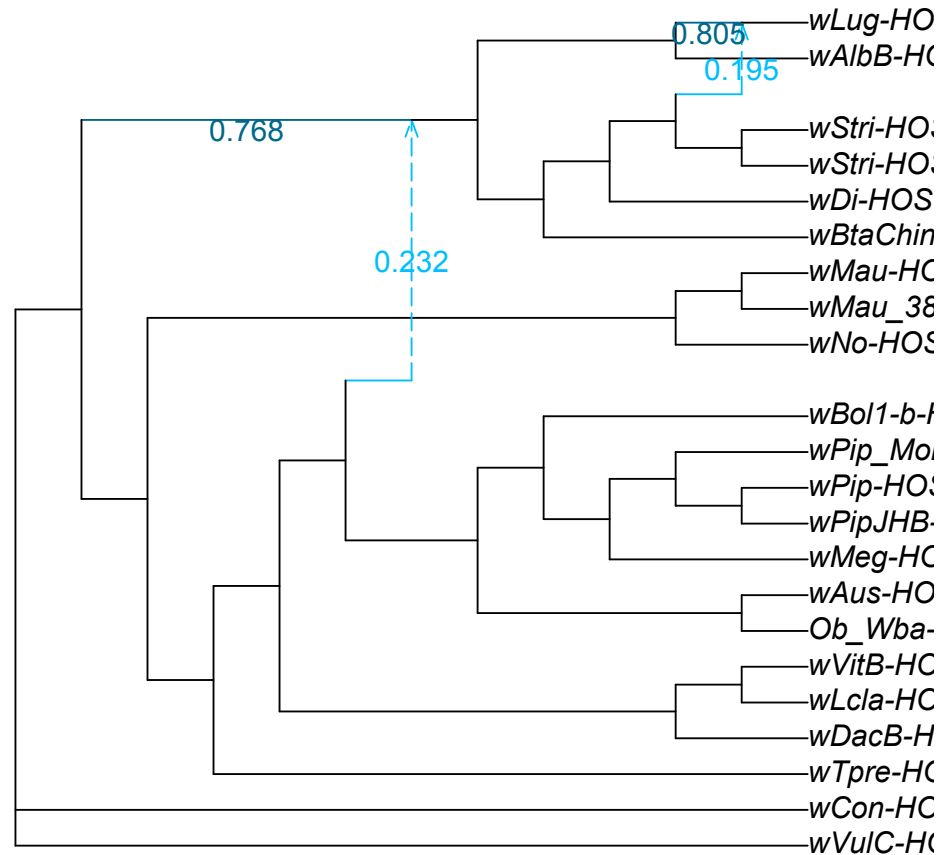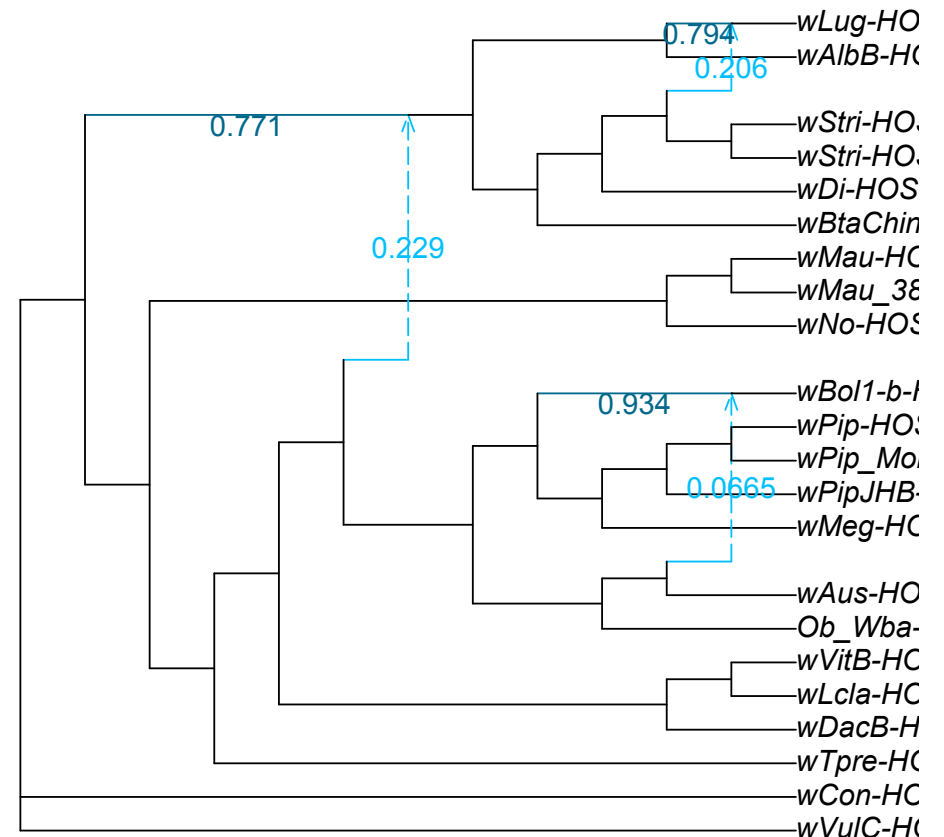


SNaQ

PhyloNet ML

PhyloNet MPL

→ How do you score networks?
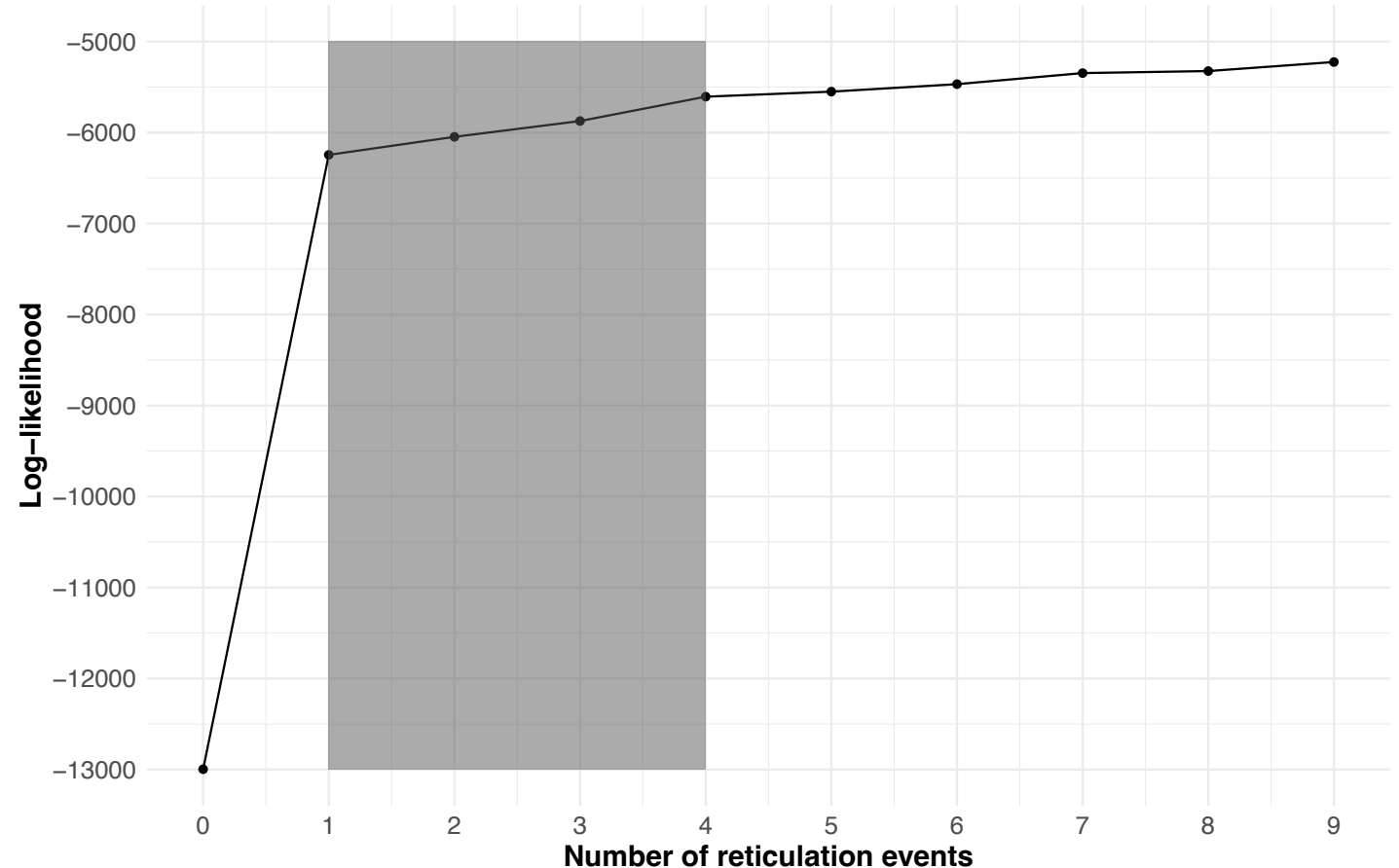
# Problem n°3: how to choose number of reticulations?



SNaQ analysis with 2 vs. 3 hybridization events

# (Pseudo) Log-likelihood improves as you increase the number of reticulation events

- Risk of falsely detecting extra reticulations

- Choice of "best" number of hybridization events usually done in a hand-wavy way (i.e. within the greyzone hereafter)

- Few statistical tests available (but see Cai & Ané, 2020)

# Conclusion

- Species tree inference is hard, as it requires careful processing of the genomes (assembly pipeline), alignment procedure and tree estimation

- Phylogenetic network inference is much harder, as we are limited by the number of methods and most of them are highly time-consuming

- Wolbachia seems to have a limited amount of <u>major</u> hybridization events but we cannot quantify it with certainty