# Inferring Bacterial Recombination Graphs using BEAST2

Lecturer: Tim Vaughan

Department of Biosystems Science and Enginerring
ETH Zürich, Switzerland

Taming the BEAST, London, July 2017
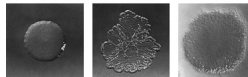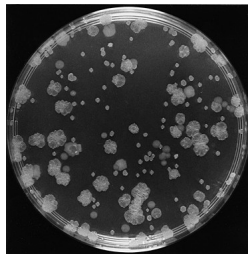
# Why study bacterial phylogenetics?

# Why study bacterial phylogenetics?

- ▶ Bacteria play important roles (both positive and negative) in the health of humans, animals and plants.
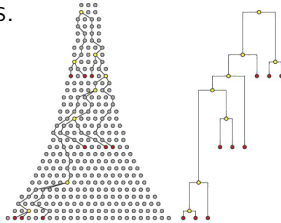
# Why study bacterial phylogenetics?

- ► Bacteria play important roles (both positive and negative) in the health of humans, animals and plants.

- ► Many bacteria possess interesting and *experimentally accessible* evolutionary dynamics.
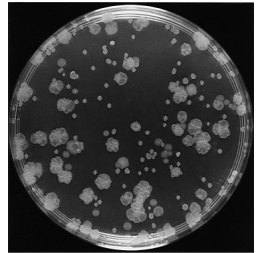






Rainey & Travisano, Nature (1998)

# Why study bacterial phylogenetics?

- ► Bacteria play important roles (both positive and negative) in the health of humans, animals and plants.

- ► Many bacteria possess interesting and *experimentally accessible* evolutionary dynamics.

- ► Bacterial genomes are measurably evolving over relatively short study periods.



Drummond & Rambaut, TIEE (2003)

Rainey & Travisano, Nature (1998)
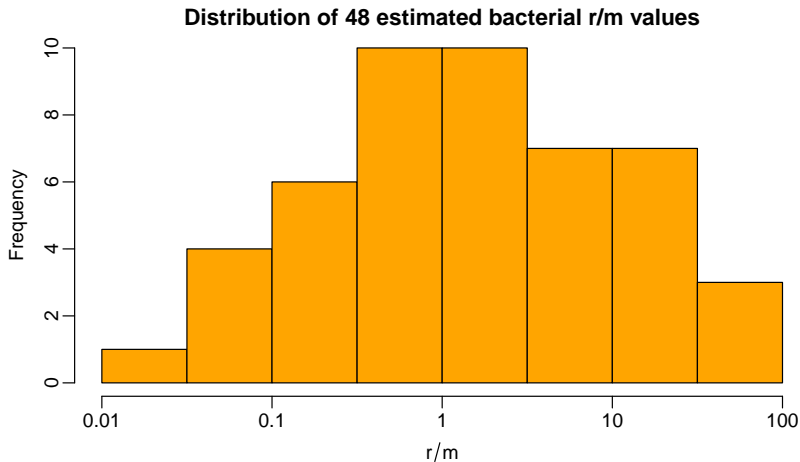
# Bacterial Recombination

- ▶ Bacteria reproduce clonally via binary fission.
- ▶ Multiple mechanisms allow for non-vertical transfer of genetic information:
  - ▶ Conjugation
  - ▶ Natural transformation
  - ▶ Phage-mediated transduction
- ▶ The frequency at which these events occur depends on the bacterial species (i.e. depends on the genome: a strange loop!)
- ▶ The effect of these events can be:
  - ▶ Plasmid transfer
  - ▶ Insertion
  - ▶ Homologous recombination
- ▶ Focus solely on homologous recombination: only event which doesn't alter the length of the sequence.

# The Problem for Phylogenetic Inference

For many bacteria, the ratio between the recombination rate and the mutation rate is very high.

# The Problem for Phylogenetic Inference

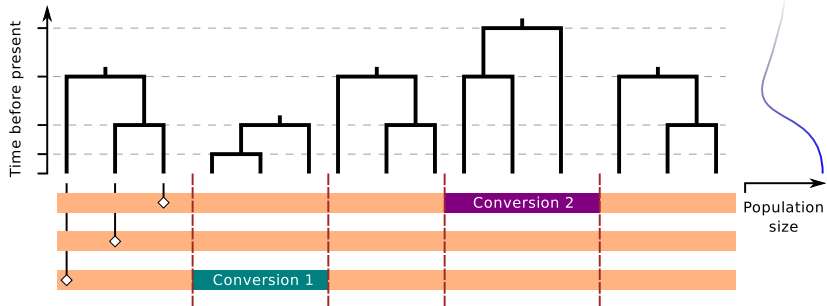For many bacteria, the ratio between the recombination rate and the mutation rate is very high.



**Distribution of 48 estimated bacterial r/m values**

Vos and Didelot (2009)

# The Benefit for Demographic Inference

# The Benefit for Demographic Inference



▶ Relationship used by Li and Durbin (2011) to infer human demographic history from pairs of autosomes.

# Existing solutions

# Existing solutions

- ▶ Pre-processing of data to identify and remove non-vertically inherited material. (eg. START: Jolley et al. (2001))

# Existing solutions

- ▶ Pre-processing of data to identify and remove non-vertically inherited material. (eg. START: Jolley et al. (2001))

| Pros | Cons |
|------|------|
| • Can use standard tools for phylogenetic inference. | • Data is being thrown away.<br>• Ad hoc, may bias results. |

# Existing solutions

- ▶ Pre-processing of data to identify and remove non-vertically inherited material. (eg. START: Jolley et al. (2001))

| Pros | Cons |
|---|---|
| • Can use standard tools for phylogenetic inference. | • Data is being thrown away. <br> • Ad hoc, may bias results. |

- ▶ Explicit modelling of bacterial recombination. (eg. ClonalFrame and ClonalOrigin: Didelot and Falush (2007); Didelot et al. (2010))

# Existing solutions

- Pre-processing of data to identify and remove non-vertically inherited material. (eg. START: Jolley et al. (2001))

| Pros | Cons |
|---|---|
| • Can use standard tools for phylogenetic inference. | • Data is being thrown away. <br> • Ad hoc, may bias results. |

- Explicit modelling of bacterial recombination. (eg. ClonalFrame and ClonalOrigin: Didelot and Falush (2007); Didelot et al. (2010))

| Pros | Cons |
|---|---|
| • Can make use of all data. <br> • Can infer additional parameters such as recombination rates. <br> • May yield increased confidence in estimates | • Models can be complex, with many parameters. <br> • Both computationally and statistically challenging. <br> • Existing implementations are too restrictive. |

# The coalescent with gene conversion



Generations

Wiuf (2000); Wiuf and Hein (2000)

# The coalescent with gene conversion



Generations

Wiuf (2000); Wiuf and Hein (2000)

# The coalescent with gene conversion



Generations

Wiuf (2000); Wiuf and Hein (2000)

# The coalescent with gene conversion



Wiuf (2000); Wiuf and Hein (2000)

# The coalescent with gene conversion



Generations

Wiuf (2000); Wiuf and Hein (2000)

# The coalescent with gene conversion



$(N_e(t)g)^{-1}$    Coalescence rate
$\rho_s$    Conversion rate
$\delta$    Expected tract length

clonal frame

recombinant edges

Generations

Wiuf (2000); Wiuf and Hein (2000)

The space of possible ancestral recombination graphs is extremely large. Even two samples have infinitely many distinct ancestries!

Full ARG

# Approximation 1: ClonalFrame

# Approximation 2: ClonalOrigin

# Inference under the ClonalOrigin model

Inference follows the standard Bayesian phylogenetic tradition:

$$f(G, N, \mu, \rho, \delta | A) \propto P_F(A|G, \mu) f_{CGC}(G|N, \rho, \delta) f_{prior}(N, \mu, \rho, \delta)$$

where

      $A$ is the sequence alignment,

      $\mu$ are the substitution model parameters, and

      $G$ is the full sample genealogy including clonal frame $T$ and $M$ conversions $\{C_i\}_{i \in [1 \dots M]}$.

The genealogy density under ClonalOrigin model can be expanded

$$f_{CGC}(G|\rho', \delta, N) = \left( \prod_{i=1}^{M} f(C_i|T, N, \delta) \right) P(M|T, \rho) f_C(T|N)$$

# Identifiability

Despite using a simplified model, an infinite number of ARGs still possess the same likelihood given a sequence alignment.

# Identifiability

Despite using a simplified model, an infinite number of ARGs still possess the same likelihood given a sequence alignment.

# Identifiability

Despite using a simplified model, an infinite number of ARGs still possess the same likelihood given a sequence alignment.



*Very important for an MCMC algorithm to propose state changes which minimize effect on likelihood.*

# ClonalOrigin in BEAST



- BEAST package that performs inference under the ClonalOrigin model.
- Joint inference of clonal frame and recombinant edges.
- Can be combined with usual variety of substitution models and parametric population models.
- Straight-forward usage via BEAUti.

`http://tgvaughan.github.io/bacter`

# Setting up a Bacter analysis

# Setting up a Bacter analysis

# Setting up a Bacter analysis

# Setting up a Bacter analysis

# Setting up a Bacter analysis

True ARG:

Randomly-selected ARG from MCMC:

Locus
■ l1
■ l2

# Summary networks

# Summary networks

# Summary networks

# Summary networks

# Summary networks



Converted fraction

1

0

Threshold

Site

(A,B,C)

(A,B)

A    B    C

Conversions
from distinct
states sampled
from posterior

# Summary networks

# ACGAnnotator

# ACGAnnotator

True ARG:

Summary ARG from MCMC:

# Bacter limitations

- Computational complexity scales with the number of proposed conversions.
  - This can be huge, even for small sample sets!
  - Can't use BEAGLE to speed things up: doesn't perform well due to peculiarities of ARG likelihood computation.
- MCMC algorithm used does not intelligently locate conversions.
  - Looking at fixing this in the near future.
- Summary algorithm can produce peculiar results.
  - More research needs to be done to find a better algorithm.

# baCter

## Bacterial phylogenetics in BEAST 2

[ Installation & Usage Tutorial ]  [ Source Code ]

Bacter is a BEAST 2 package which facilitates joint Bayesian inference of ancestral recombination graphs (ARG) and related parameters from bacterial sequence alignments under the ClonalOrigin model.

## ARGs

Bacter allows you to infer the ARG that gave rise to your sampled data, including both clonal frame and the recombinant edges. Loci and sites affected by each conversion are also inferred.

## Parameters

For informative data sets, Bacter can sample posterior distributions for parameters such as the recombination rate and expected converted tract length jointly with the ARG. For serially sampled data, joint estimation of the molecular clock rate is also possible.

## Population dynamics

Bacter includes support for producing Bayesian Skyline Plot (BSP) inferences of bacterial population dynamics. This inference may actually be more powerful when recombination is present in your data set.

tgvaughan.github.io/bacter

# Bacter Tutorial

1. Open Bacter tutorial at
   taming-the-beast.github.io/tutorials/Bacter-Tutorial
2. Begin the tutorial!

# References

Didelot, X. and Falush, D. (2007). Inference of bacterial microevolution using multilocus sequence data. *Genetics*, 175:1251.

Didelot, X., Lawson, D., Daarling, A., and Falush, D. (2010). Inference of homologous recombination in bacteria using whole-genome sequences. *Genetics*, 186:1435.

Jolley, K. A., Feil, E. J., Chan, M. S., and Maiden, M. C. (2001). Sequence type analysis and recombinational tests (start). *Bioinformatics (Oxford, England)*, 17:1230–1231.

Li, H. and Durbin, R. (2011). Inference of human population history from individual whole-genome sequences. *Nature*, 475(7357):493–496.

Vos, M. and Didelot, X. (2009). A comparison of homologous recombination rates in bacteria and archaea. *ISME J*, 3(2):199–208.

Wiuf, C. (2000). A coalescence approach to gene conversion. *Theor Popul Biol*, 57(4):357–367.

Wiuf, C. and Hein, J. (2000). The coalescent with gene conversion. *Genetics*, 155(1):451–462.