

# **Non-tree like evolution: Dealing with recombination in phylogenetics**

**David Rasmussen**

Department of Entomology and Plant Pathology  
Bioinformatics Research Center  
North Carolina State University

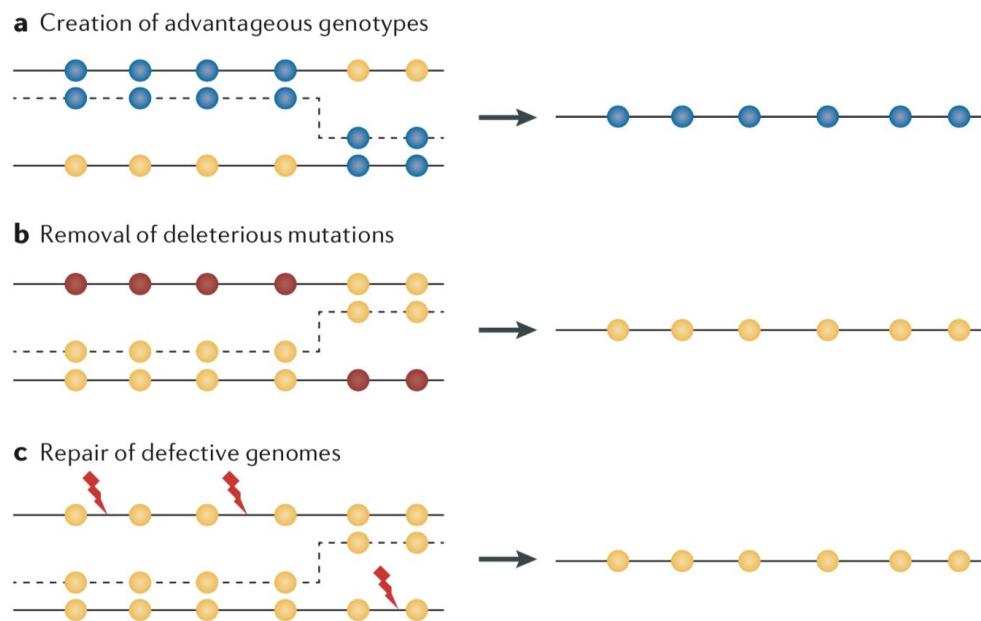
Taming the BEAST, Squamish  
August 17<sup>th</sup>, 2023

**Recombination is a  
major force shaping  
the evolution of  
nearly all populations**

# The advantages of recombination

Recombination can shuffle parental genetic material to:

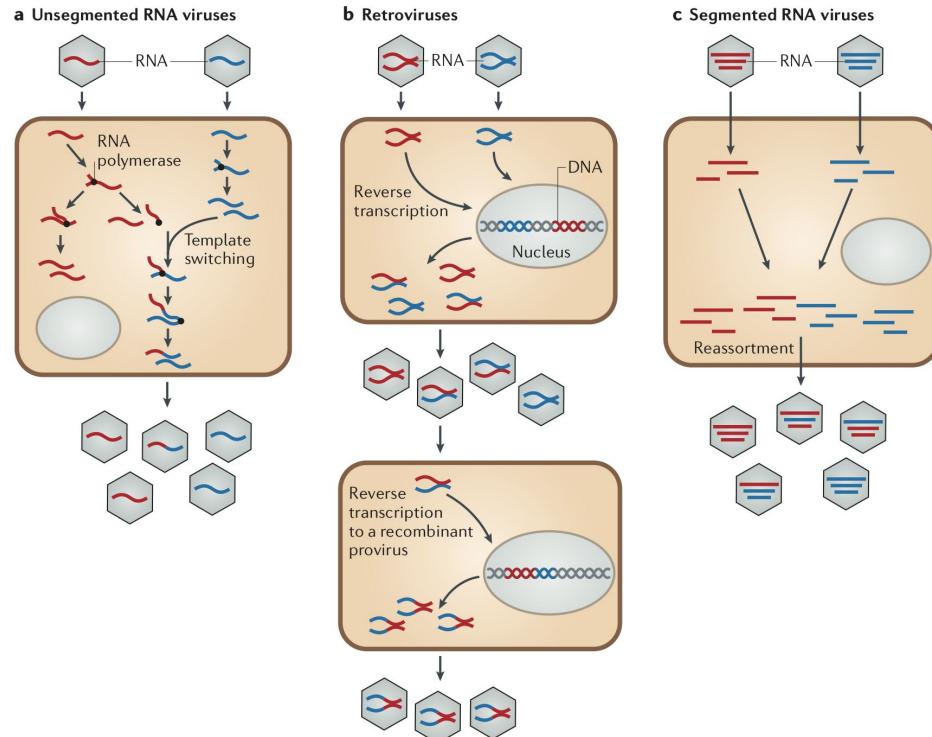
- Combine beneficial mutations
- Purge deleterious mutations
- Repair defective genomes



# Recombination creates novel pathogens

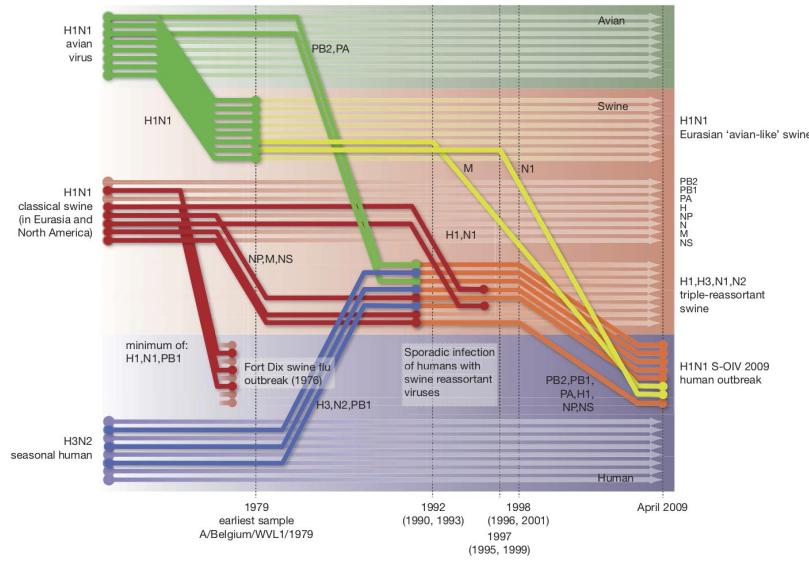
Co-infection of a cell by genetically distinct viral strains can lead to the generation of recombinant viruses.

End result: progeny inherit genetic material from both parents.



# Recombination creates novel pathogens

The 2009 pandemic H1N1 resulted from a triple reassortment event that exchanged genome segments between human, swine and avian viruses.



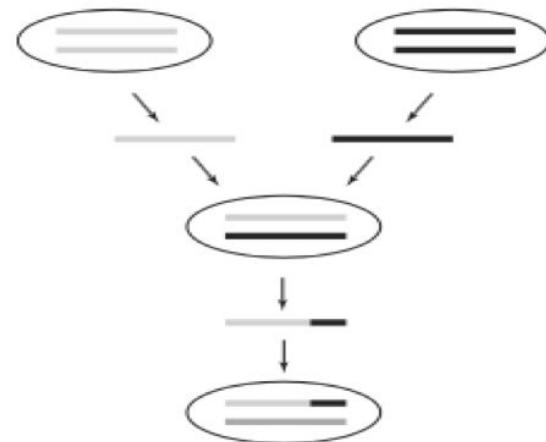
**If recombination is so  
good for pathogens,  
why is it so bad for  
phylogenetics?**

# Recombination creates mosaic ancestry

Without any recombination, the entire genome of an individual will share the same ancestry (i.e. phylogenetic history).

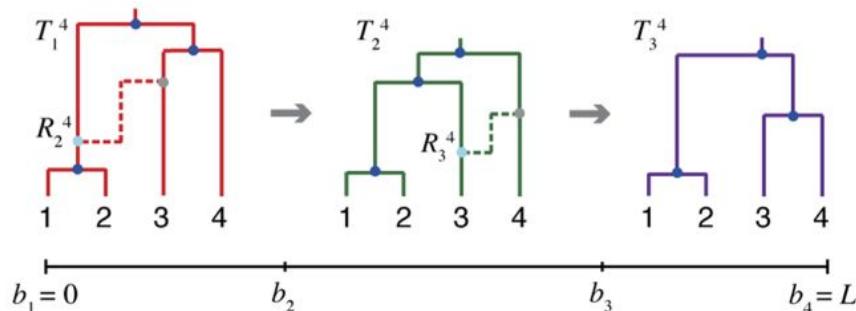
With recombination, genomes become mosaics where different segments descend from different ancestors.

No single phylogenetic tree can therefore describe the genetic ancestry of a sample of recombining sequences.



# Recombination creates mosaic ancestry

Different regions of the genome will have different phylogenetic histories:



**C**

**$D^4$**

1	C	G	A
2	A C	A C	T A
3	T G	C G	A A
4	T A	A G	T T

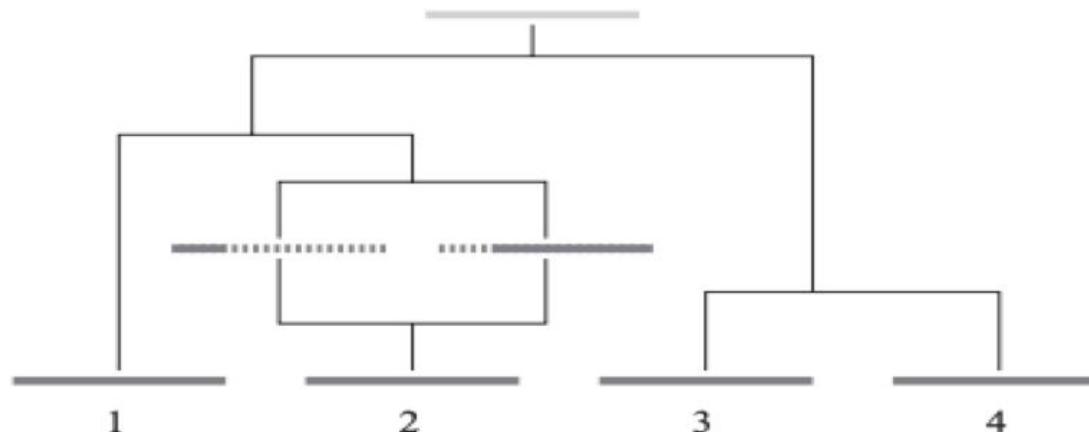
# Effect of a single recombination event

A single recombination event between two sampled lineages will have one of three possible effects on the phylogeny:

- No effect
- Effect only the branch lengths
- Effect the tree topology

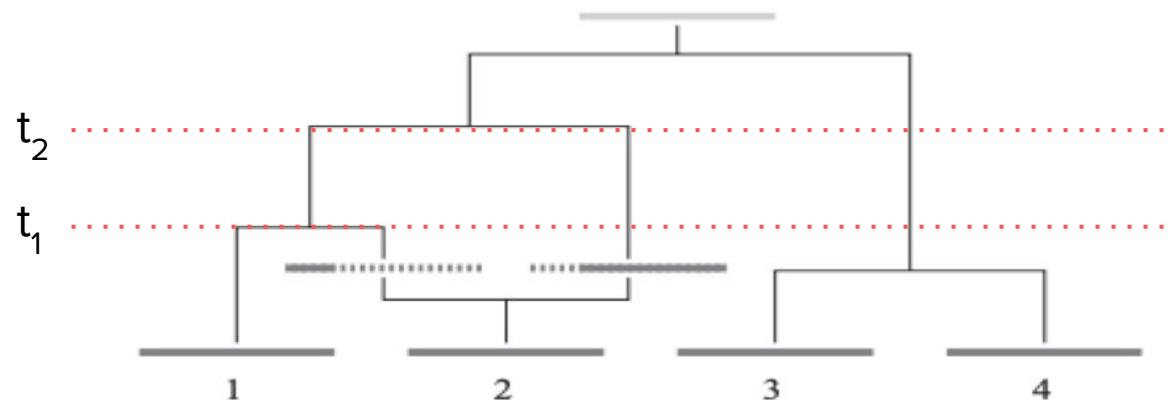
# Effect of a single recombination event

If the recombining parent lineages coalesce before they coalesce with any other lineage, the recombination event will have **no effect** on the phylogeny.



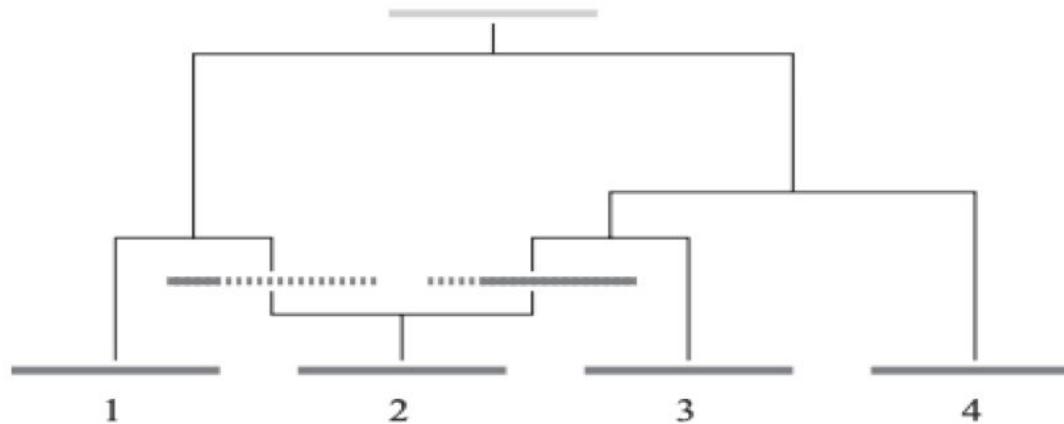
# Effect of a single recombination event

Only **branch lengths** will change if one of two recombining parent lineages merges with another sequence before coalescing with the other parent..



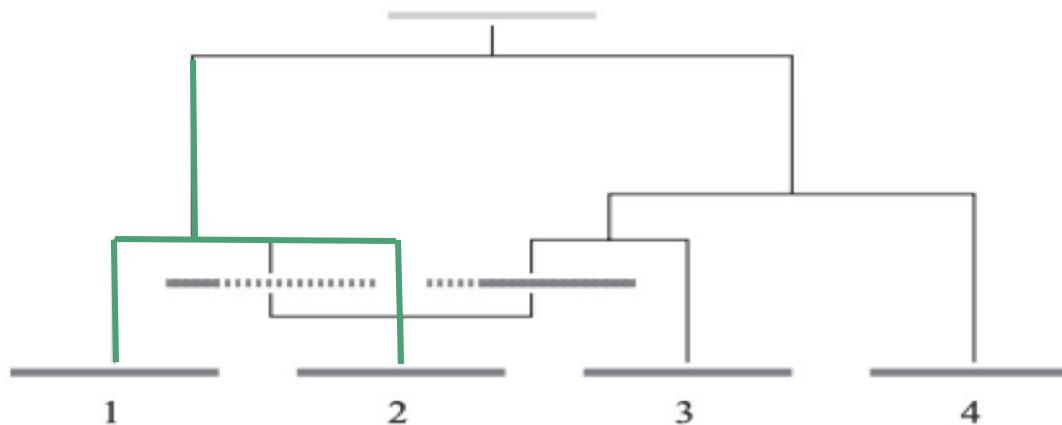
# Effect of a single recombination event

The **tree topology will change** if the recombining parent lineages coalesce with other sequences before the two parent lineages coalesce.



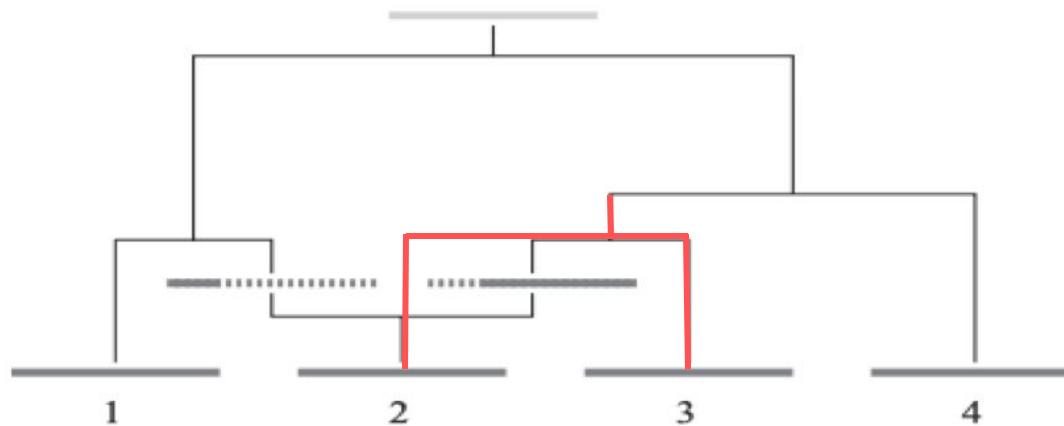
# Effect of a single recombination event

The **tree topology will change** if the recombining parent lineages coalesce with other sequences before the two parent lineages coalesce.



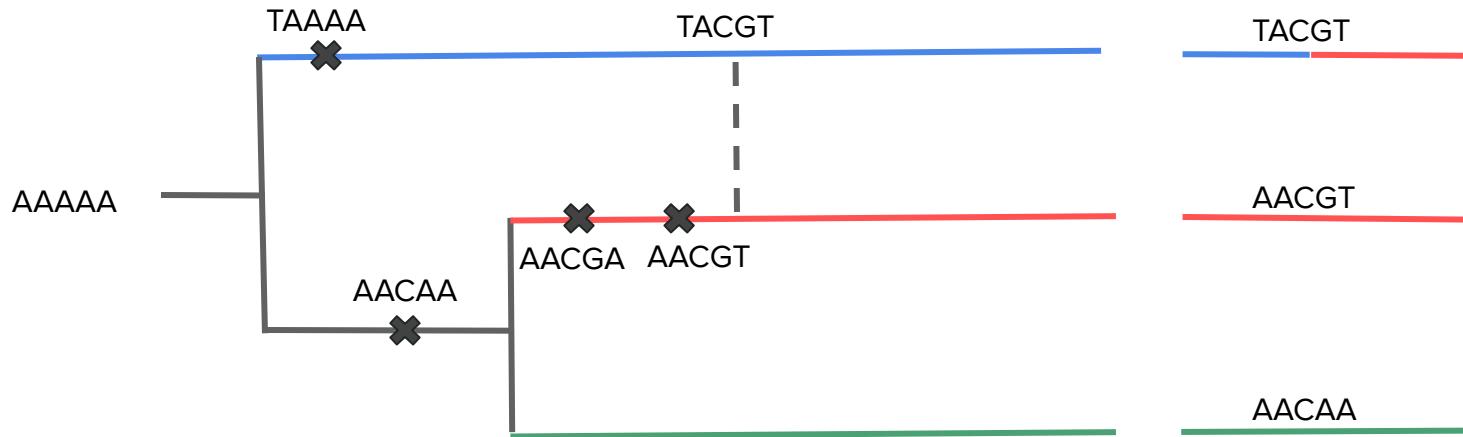
# Effect of a single recombination event

The **tree topology will change** if the recombining parent lineages coalesce with other sequences before the two parent lineages coalesce.



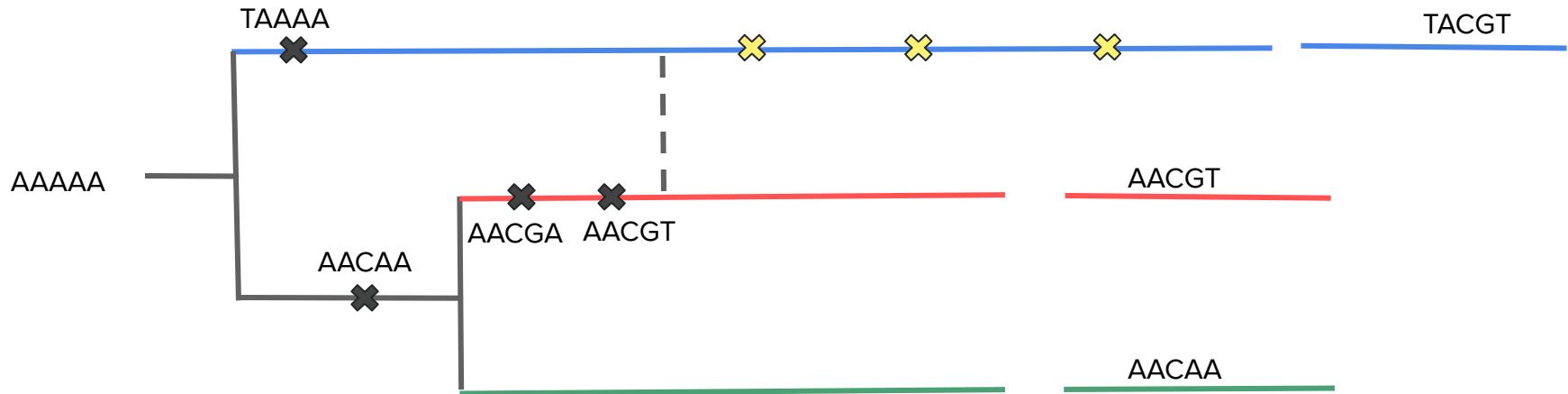
# Effect of a single recombination event

A recombination event between two sequences can generate recombinant sequences that are quite genetically divergent from the parent sequences.



# Effect of a single recombination event

This will result in abnormally long branches leading to recombinant sequences if recombination is ignored when reconstructing the phylogeny.



# Effect of many recombination events

In the presence of multiple recombination events, phylogenies:

- Have longer terminal branches
- Become more star-like
- Behave less clock-like\*\*\*

\*\*\* Wreaks havoc on estimating the molecular clock rate

**But how do we actually  
deal with recombination  
in phylogenetic  
analyses?**

# Some potential options

Remove recombinant sequences from alignments.

Remove recombinant genomic regions and reconstruct local trees from recombination-free blocks.

Assume evolution is mostly tree-like and reconstruct a clonal frame

Reconstruct a full ancestral recombination graph

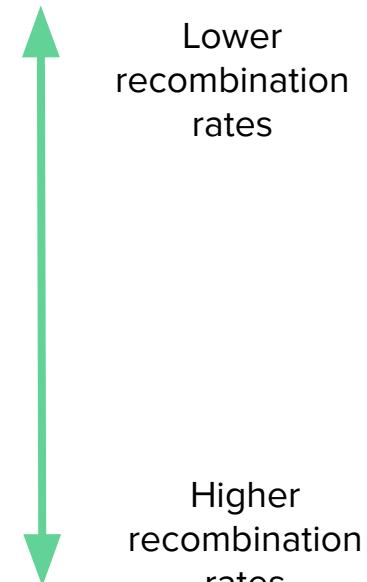
# Some potential options

Remove recombinant sequences from alignments.

Remove recombinant genomic regions and reconstruct local trees from recombination-free blocks.

Assume evolution is mostly tree-like and reconstruct a clonal frame

Reconstruct a full ancestral recombination graph



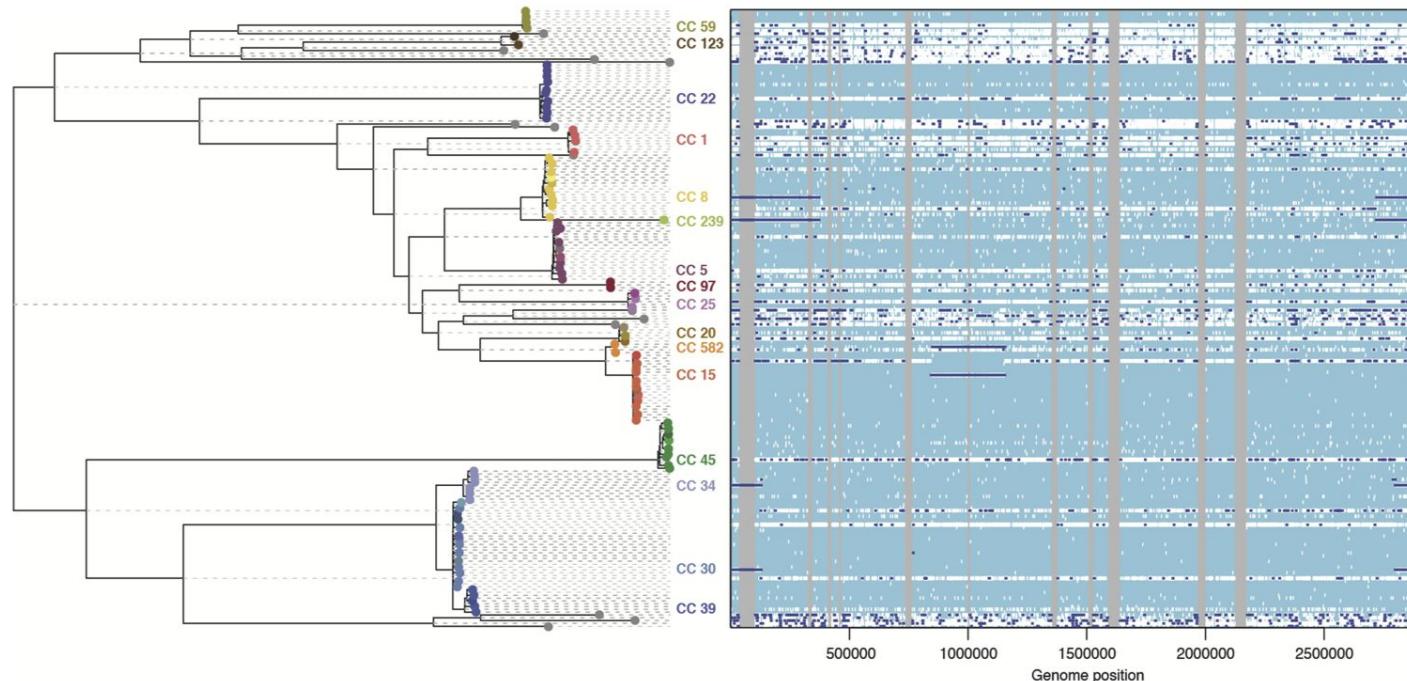
# Clonal frames

A **clonal frame** attempts to describe the true ancestral relationships among sampled individuals as a single tree.

Assumes the majority of the genome is inherited clonally while accounting for recombination within certain regions of the genome

Clonal frames are a popular choice for bacteria where the majority of the genome is assumed to be inherited clonally (i.e. the core genome) but gene conversion and other horizontal transfers overwrites small portions of the genome.

# ClonalFrame of *Staphylococcus aureus*



Dark blue = recombinant regions to be masked

Didelot et al. (PLoS Comp Bio, 2015)

# Bacter: Clonal frames in BEAST 2

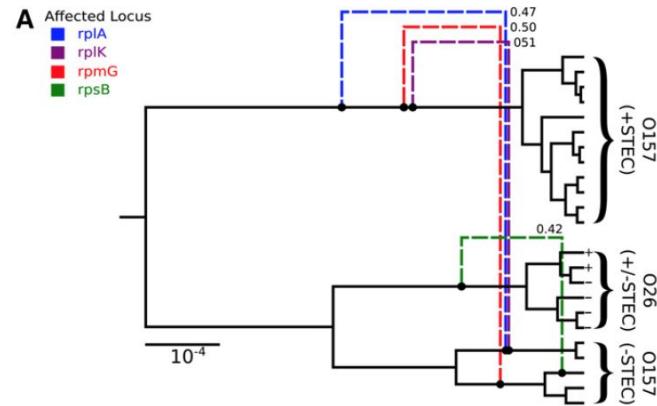
GENETICS | INVESTIGATION ■■■

## Inferring Ancestral Recombination Graphs from Bacterial Genomic Data

Timothy G. Vaughan,<sup>\*,†,‡</sup> David Welch,<sup>\*,†</sup> Alexei J. Drummond,<sup>\*,†</sup> Patrick J. Biggs,<sup>‡</sup> Tessy George,<sup>‡</sup> and Nigel P. French<sup>†</sup>

\*Centre for Computational Evolution, and <sup>†</sup>Department of Computer Science, The University of Auckland, 1010, New Zealand,

and <sup>‡</sup>Molecular Epidemiology and Public Health Laboratory, Infectious Disease Research Centre, Hopkirk Research Institute, Massey University, Palmerston North 4442, New Zealand



<https://taming-the-beast.org/tutorials/Bacter-Tutorial/>

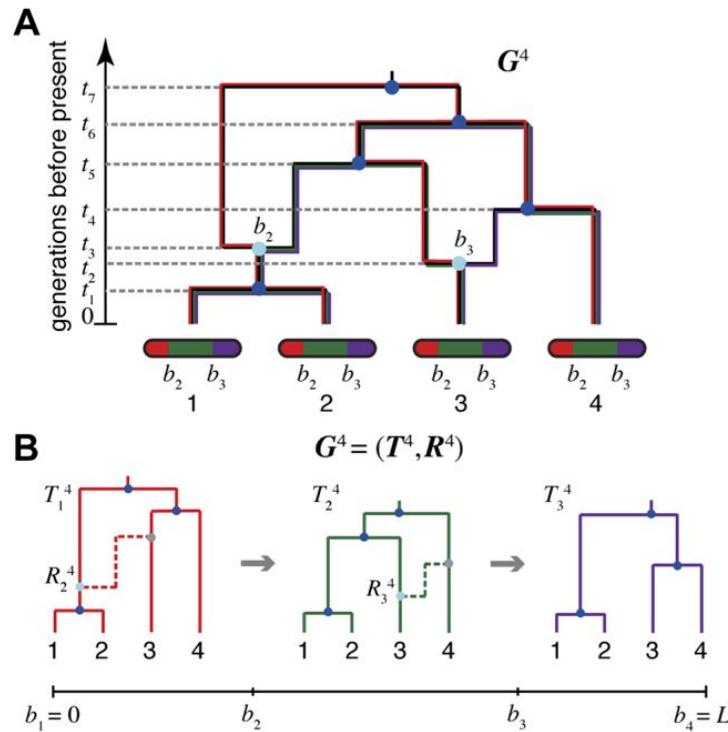
# Ancestral recombination graphs

ARGs provide a complete record of the ancestry of all sequences as a graph/network.

This graph includes all recombination and coalescent events in the history of the sample as well as information about the location of recombination breakpoints.

The local phylogeny at each genomic position is embedded in the full ARG

# A hypothetical ARG



# Ancestral recombination graphs

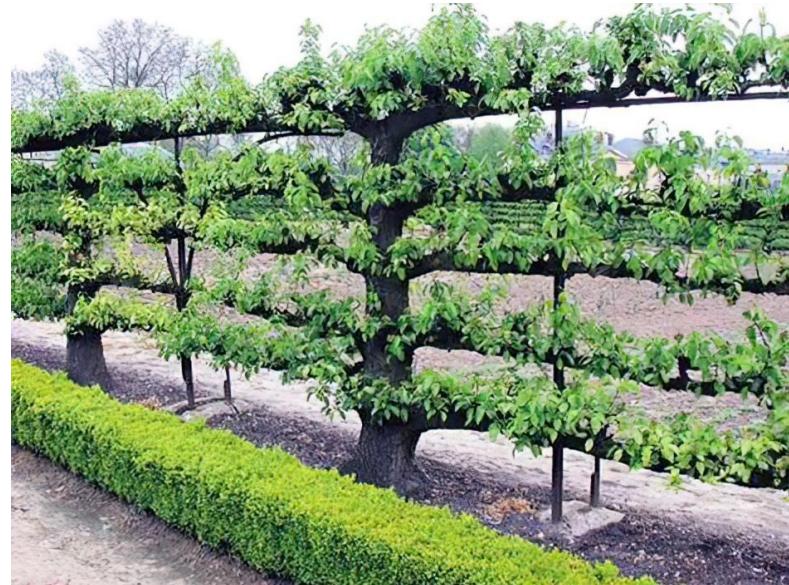
ARGs are in theory the ideal way to represent the history of sequences with recombination.

However, even state-of-the-art methods like *ARGweaver* (Rasmussen et al., 2014) that employ very efficient HMM methods work with at most dozens of sequences.

Notoriously difficult to infer full ARGs...

# ARG reconstruction using Espalier

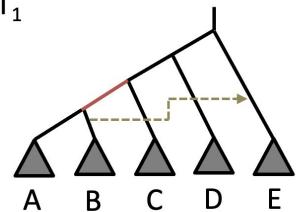
***Espalier (noun):*** the ancient agricultural practice of controlling woody plant growth by pruning and training branches to a frame.



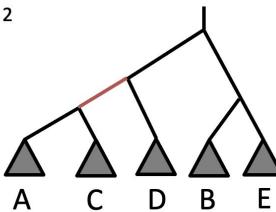
# Maximum agreement forests

Given a pair of discordant trees, a MAF provides the smallest possible set of subtrees that are all topologically consistent between the pair.

$T_1$



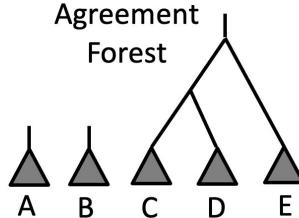
$T_2$



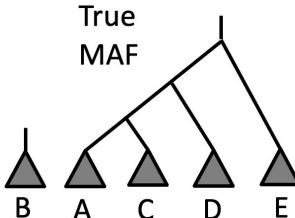
Edge bipartition = AB|CDE

Edge bipartition = AC|DBE

Agreement Forest



True  
MAF



# Maximum agreement forests

MAFs have incredibly useful properties for studying recombination!

1. The number of subtrees in a MAF gives the subtree-prune-regraft (SPR) distance between two trees, which tells us the number of recombination events required to reconcile the trees.

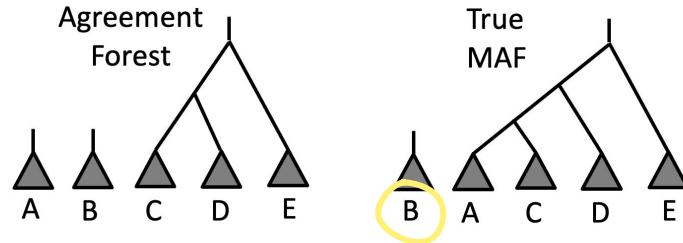
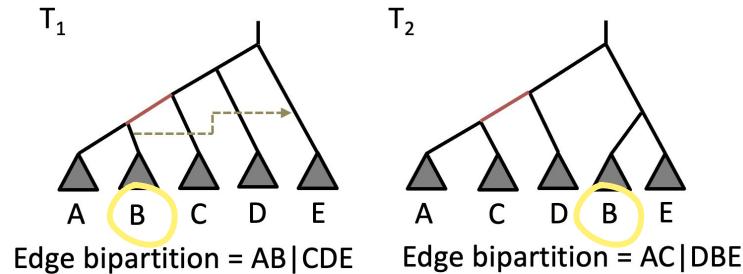
# Maximum agreement forests

MAFs have incredibly useful properties for studying recombination!

1. The number of subtrees in a MAF gives the subtree-prune-regraft (SPR) distance between two trees, which tells us the number of recombination events required to reconcile the trees.
2. The subtrees in a MAF indicate which lineages likely recombined.

# Maximum agreement forests

Given a pair of discordant trees, a MAF provides the smallest possible set of subtrees that are all topologically consistent between the pair.



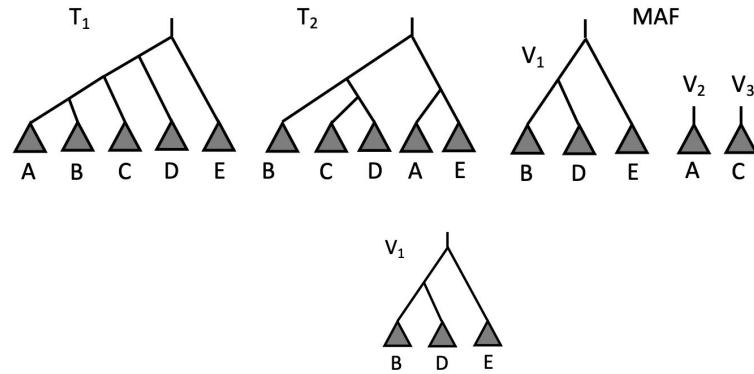
# Maximum agreement forests

MAFs have incredibly useful properties for studying recombination!

1. The number of subtrees in a MAF gives the subtree-prune-regraft (SPR) distance between two trees, which tells us the number of recombination events required to reconcile the trees.
2. The subtrees in a MAF indicate which lineages likely recombined.
3. The MAF can be used to reconcile discordant trees by re-grafting subtrees back on to the starting trees.

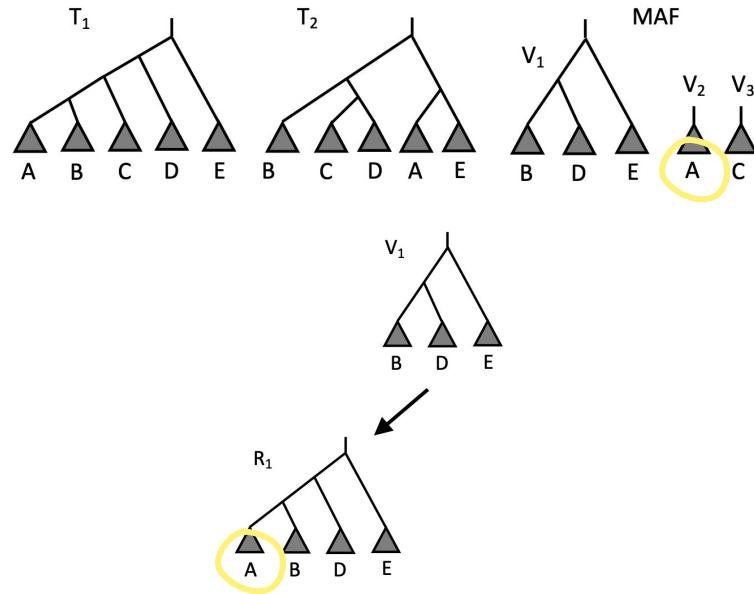
# Reconciliation through iterative regrafting

Starting with a MAF, pruned subtrees are regrafted back to their positions in the original trees.



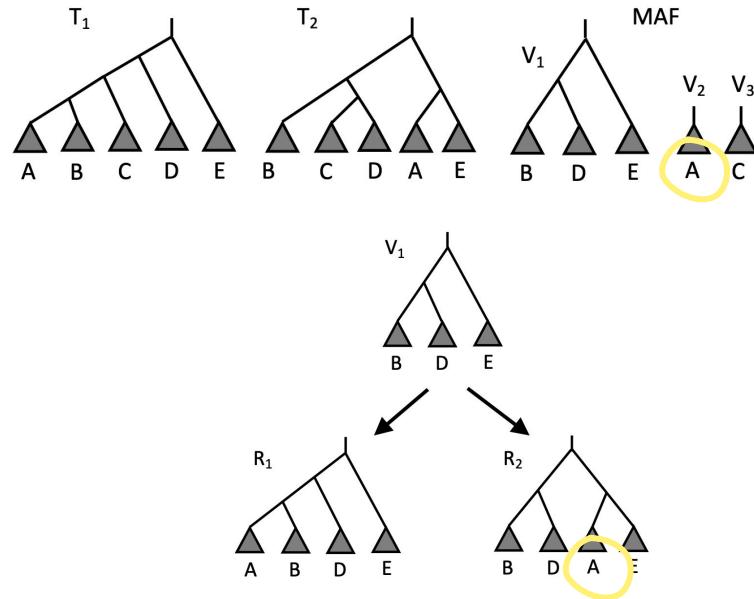
# Reconciliation through iterative regrafting

Starting with a MAF, pruned subtrees are regrafted back to their positions in the original trees.



# Reconciliation through iterative regrafting

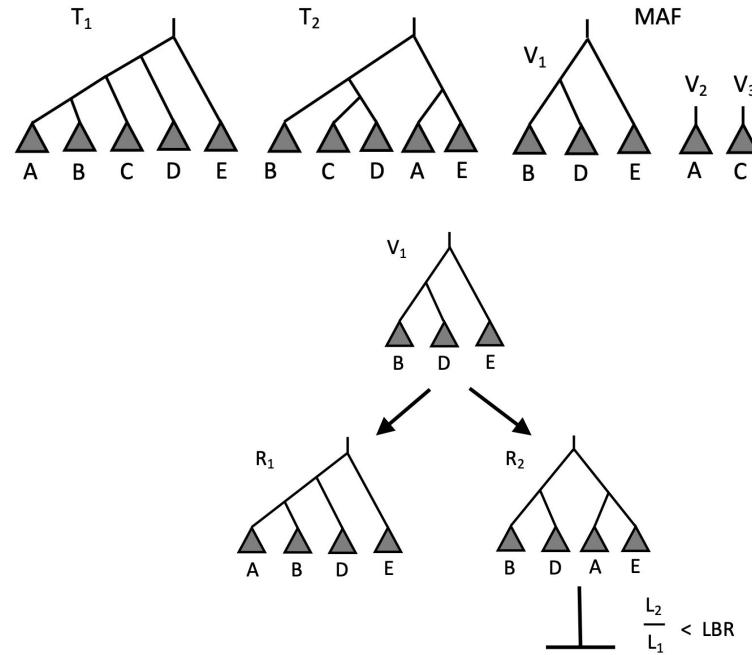
Starting with a MAF, pruned subtrees are regrafted back to their positions in the original trees.



# Reconciliation through iterative regrafting

Starting with a MAF, pruned subtrees are regrafted back to their positions in the original trees.

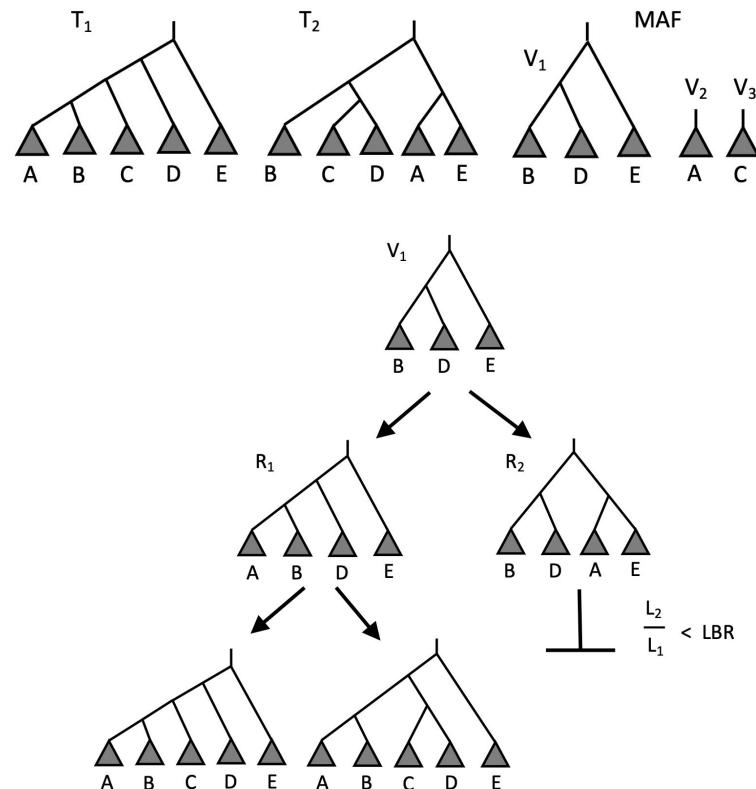
We then accept one or both regrafted trees based on the likelihood ratio of the sequence data given the trees.



# Reconciliation through iterative regrafting

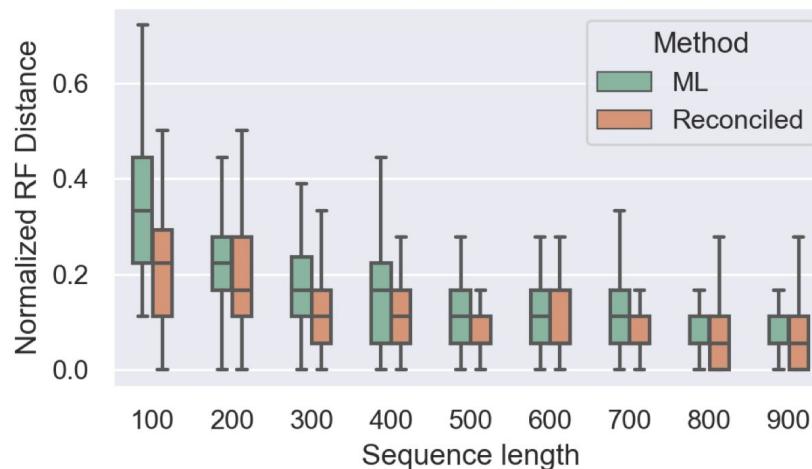
Starting with a MAF, pruned subtrees are regrafted back to their positions in the original trees.

We then accept one or both regrafted trees based on the likelihood ratio of the sequence data given the trees.



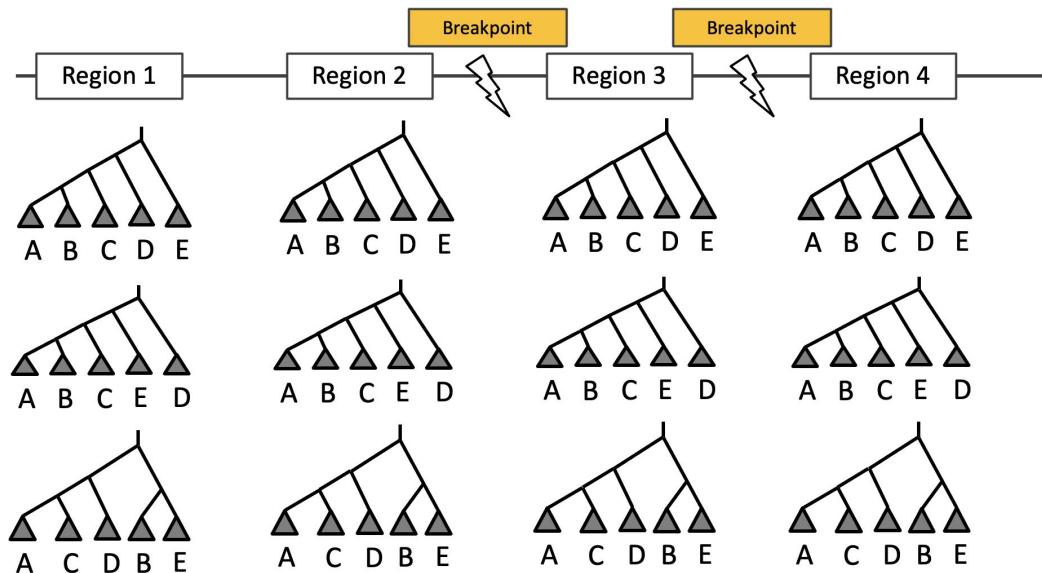
# Reconciliation through iterative regrafting

Reconciliation removes discordances between trees not strongly supported by the sequence data while retaining those that are.



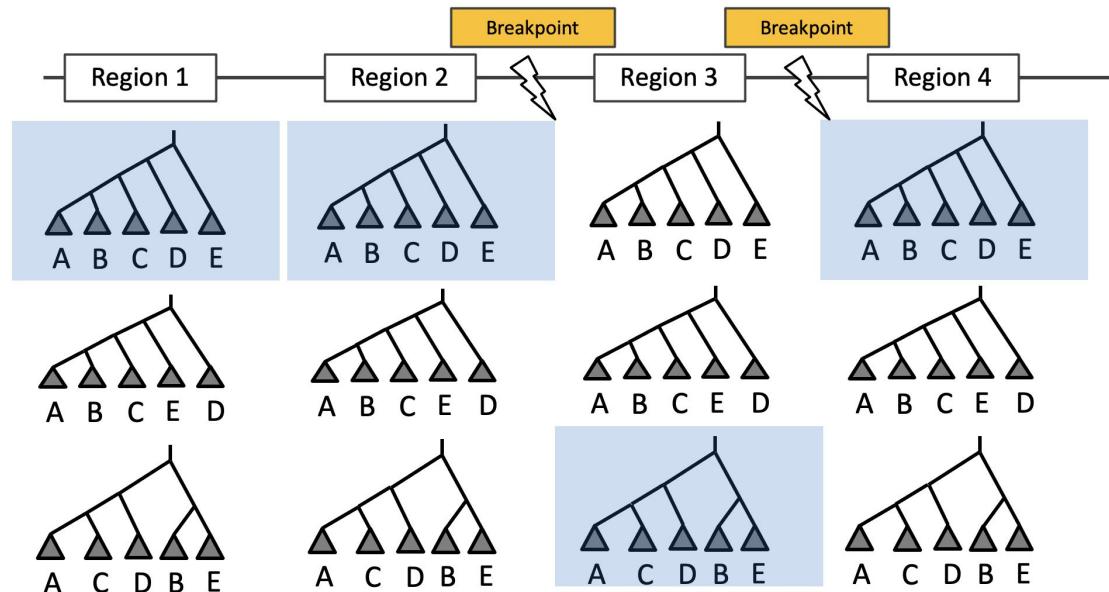
# Selecting a local tree path

A local tree path is sampled that maximizes the likelihood of the genomic sequence data while minimizing the discordance between trees.



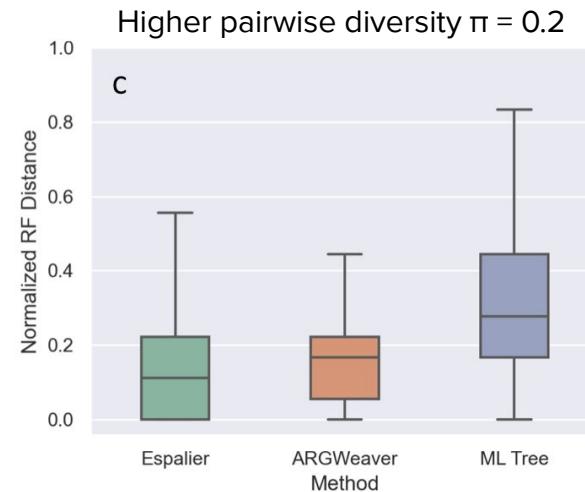
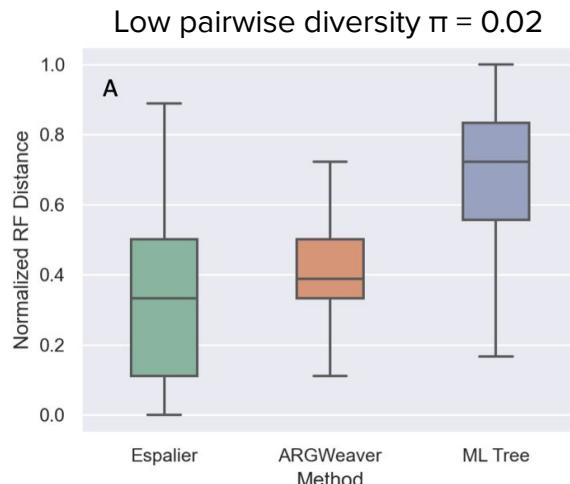
# Selecting a local tree path

A local tree path is sampled that maximizes the likelihood of the genomic sequence data while minimizing the discordance between trees.



# ARG reconstruction performance

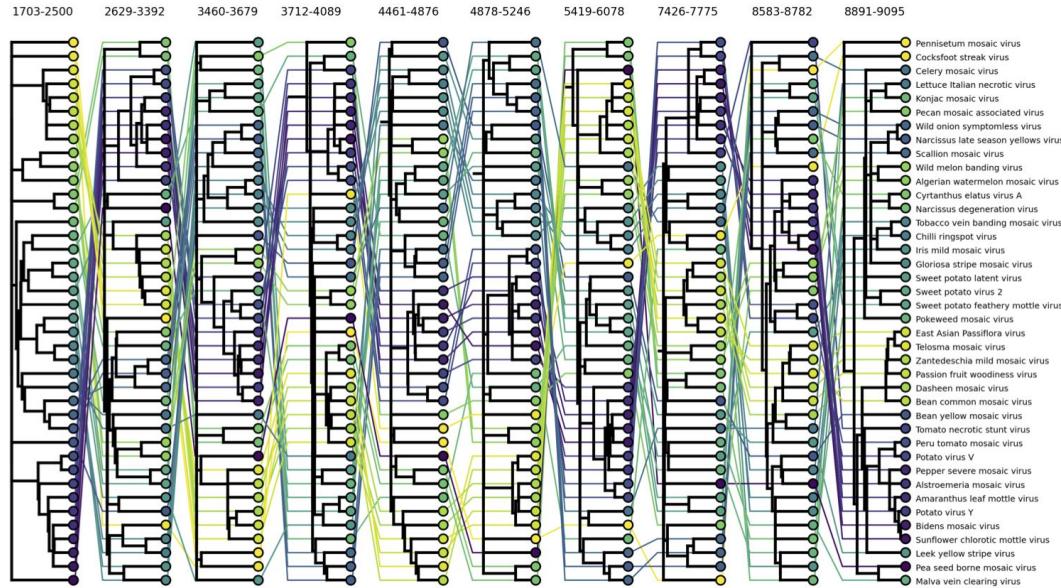
To test reconstruction accuracy, we compared the distance between local trees in reconstructed ARGs to the local trees in the true (simulated) ARGs.



# Disentangling recombination with Espalier

Potyvirus phylogenies show widespread phylogenetic conflict across different regions of the viral genome.

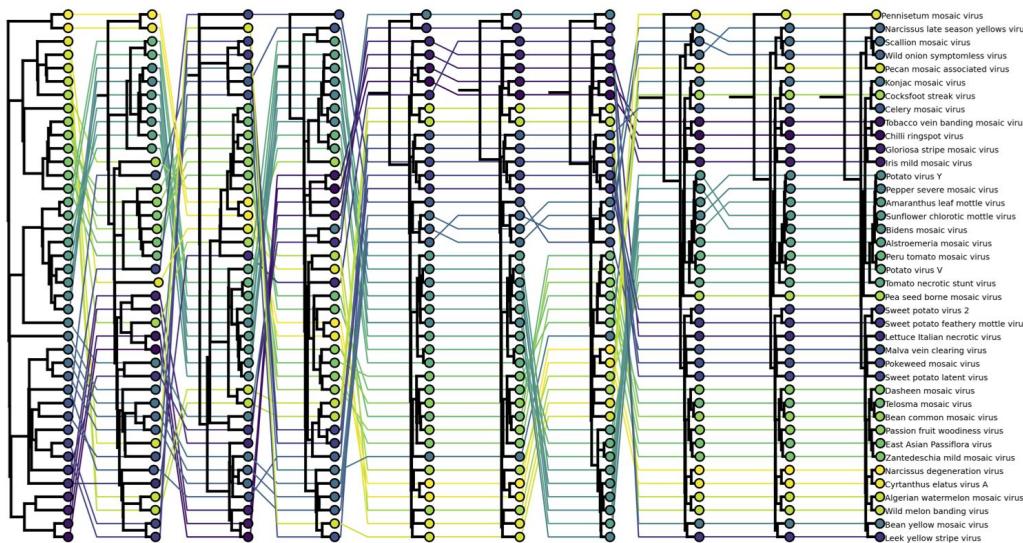
Tanglegram of reconstructed local ML trees



# Disentangling recombination with Espalier

Espalier “smooths” discordances between trees by removing conflicts likely due to phylogenetic error rather than true recombination events.

Tanglegram of local  
trees in Espalier  
ARG

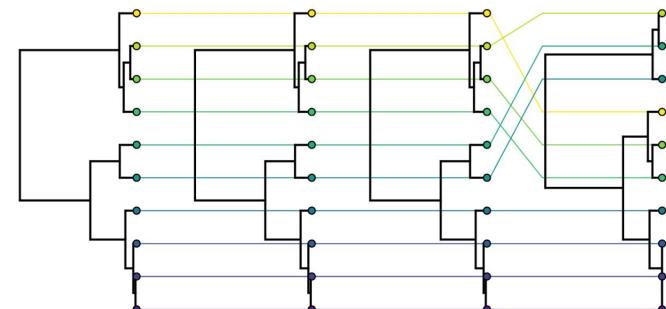


# Espalier: the package

Espalier is a python package and command line tool for computing MAFs, reconciling trees and reconstructing approximate ARGs.

The screenshot shows the Espalier documentation homepage. The top navigation bar includes a logo, the text "Espalier latest", and a "Search docs" input field. Below this is a "CONTENTS:" section with a "Introduction" link. Under "Espalier Primer", there are several links: "Installation", "Computing Maximum Agreement Forests", "Computing SPR distances", "Tree reconciliation through MAFs", "Reconstructing Ancestral Recombination Graphs" (which is highlighted in grey), "Automated breakpoint detection", and "Command Line Interface". At the bottom of the page are "Read the Docs" and "v: latest" buttons.

```
# Write local trees in reconstructed ARG to files
ARG_tree_files = ["ARG_example_ARGLocalTree" + str(i) + ".tre" for i in range(segments)]
for idx,tr in enumerate(tree_path):
    tr.write(path=ARG_tree_files[idx],schema='newick',suppress_annotations=True, suppress_rooting=True)
tanglegram_fig_name = 'ARGLocalTree-tanglegram.png'
PlotTanglegrams.plot(ARG_tree_files, tanglegram_fig_name, numerical_taxa_names=True)
```



Comparing the tanglegram for the reconstructed ARG to the tanglegram for the true ARG, we see that Espalier does a pretty good job of reconstructing the topology of the local trees. Some of the

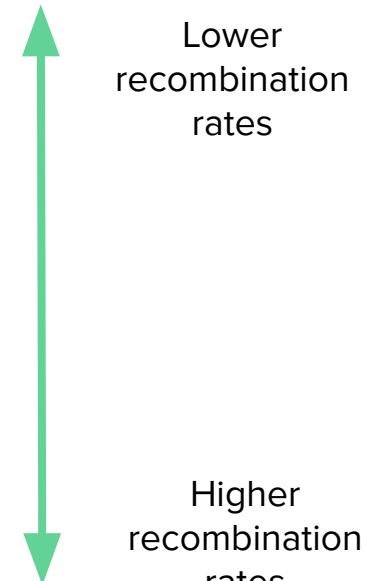
# Some potential options

Remove recombinant sequences from alignments.

Remove recombinant genomic regions and reconstruct local trees from recombination-free blocks.

Assume evolution is mostly tree-like and reconstruct a clonal frame

Reconstruct a full ancestral recombination graph



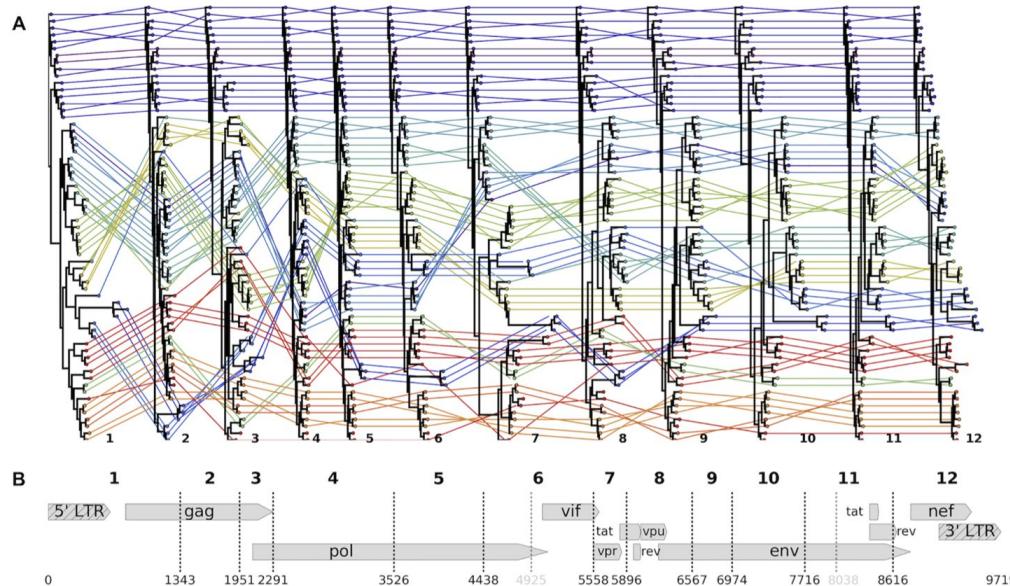
**Bonus slides: how do  
we detect  
recombination in the  
first place?**

# How do we detect recombination?

- Phylogenetic discordance between loci
- Linkage disequilibrium maps
- Substitution distribution/mosaic tests

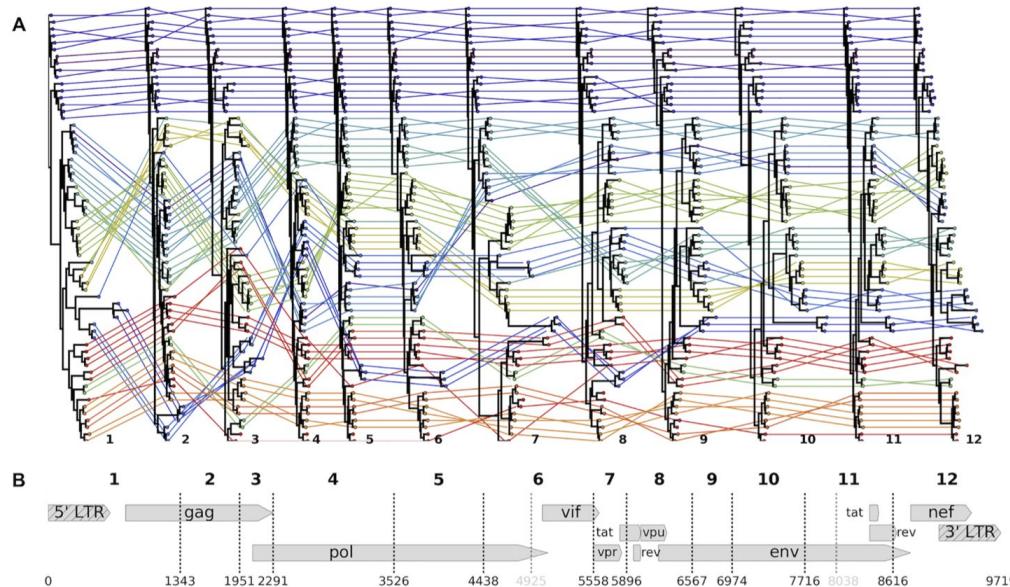
# Phylogenetic discordance

Phylogenetic discordance between ‘local’ trees can be used to detect recombination but may also arise due to errors in reconstruction.



# Phylogenetic discordance

Phylogenetic discordance between ‘local’ trees can be used to detect recombination but may also arise due to errors in reconstruction.



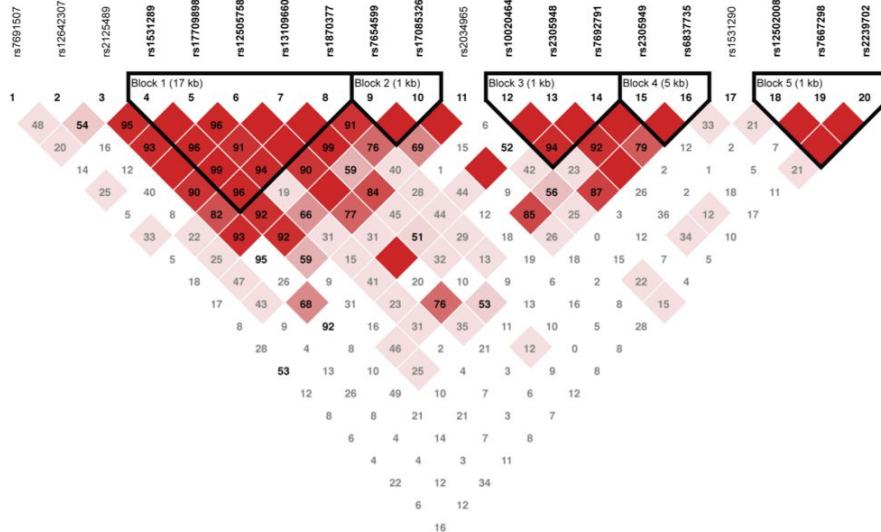
Phylogenetic recombination detection methods like **GARD** (Pond *et al.*, 2006) allow for statistical tests of discordance.

# How do we detect recombination?

- Phylogenetic discordance between loci
- Linkage disequilibrium maps
- Substitution distribution/mosaic tests

# Linkage disequilibrium maps

Sharp changes in linkage disequilibrium -- correlations in the presence/absence of alleles -- can indicate recombination in the history of the sample



**Linkage disequilibrium:**  
correlations between  
sites in the presence or  
absence of alleles.

# How do we detect recombination?

- Phylogenetic discordance between loci
- Linkage disequilibrium maps
- Substitution distribution/mosaic tests

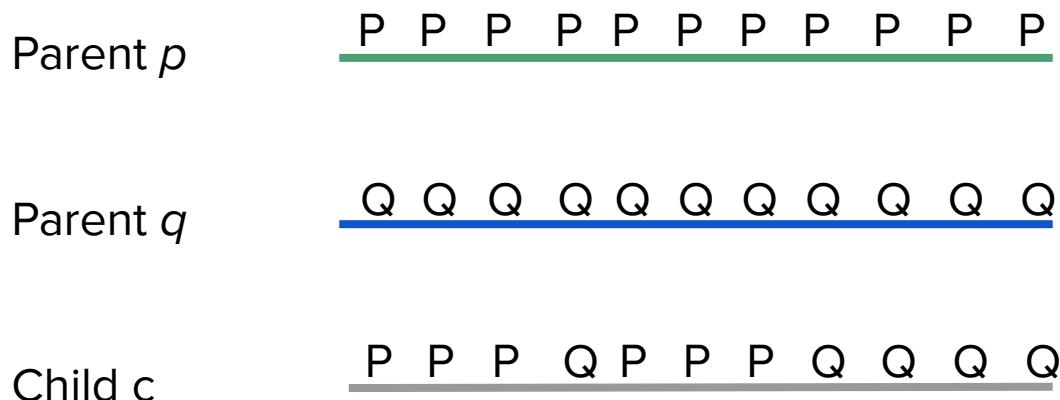
# How do we detect recombination?

Many methods try to detect a mosaic structure in the distribution of substitutions (i.e. polymorphic sites) across the genome.

For example, MaxChi uses a  $\chi^2$  to test whether there are significantly more substitutions between a pair sequences to the left or right of a potential breakpoint than would be expected by chance (Smith, 1992)

We'll consider 3SEQ (Boni *et al.*, 2007) which compares three sequences, one is assumed to be a potential child sequence that could have arisen by the two other "parent" sequences recombining.

# The 3SEQ triplet test



Let the  $P$ 's be mutations that the child shares in common with parent  $p$  and the  $Q$ 's be mutations the child shares with parent  $q$

# The 3SEQ triplet test

We can think of the mutations as up and down steps in a discrete random walk.

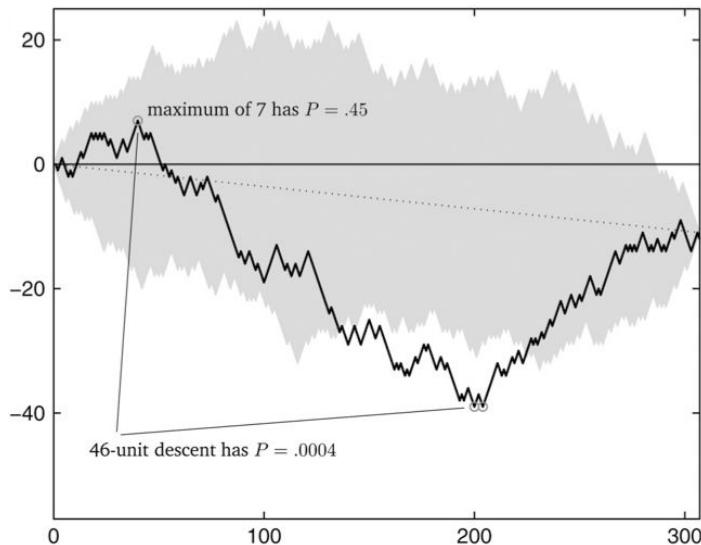
Let the  $P$ 's be thought of as up steps in the random walk.

And the  $Q$ 's as down steps.

A hypergeometric random walk model can be used to test whether the distribution (order) of  $P$ 's and  $Q$ 's is nonrandom based on the height of the random walk.

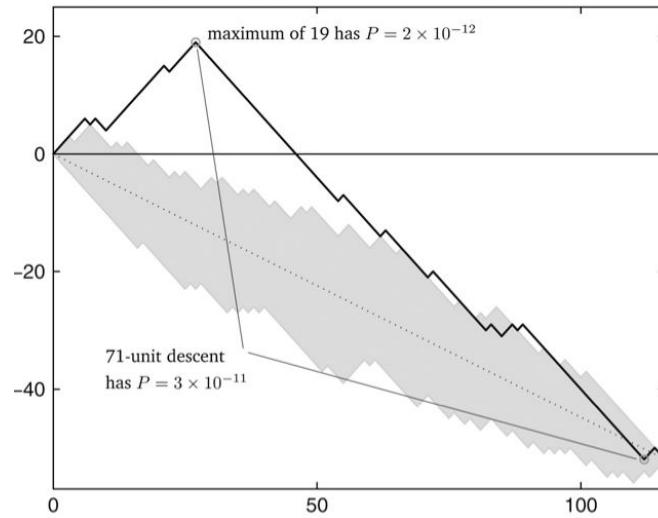
# The 3SEQ test for 1918 Spanish influenza

Small deviations from plausible random walks provide weak evidence for recombination



# The 3SEQ test for *Neisseria*

A recombinant will have a statistically improbable heights with its up steps clustered towards one end and down steps clustered towards the other end.

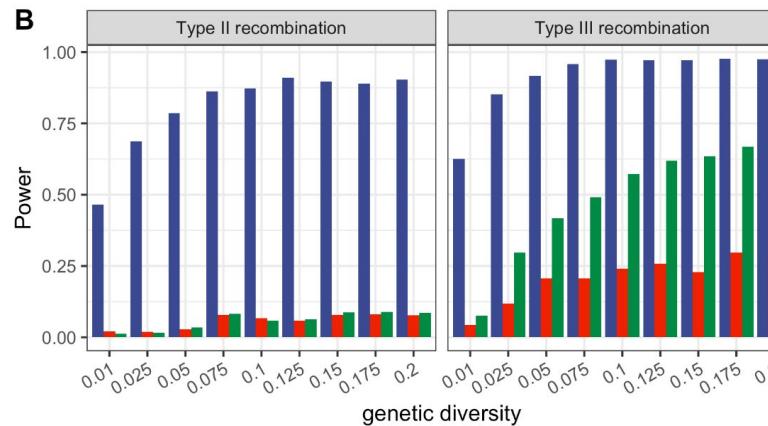
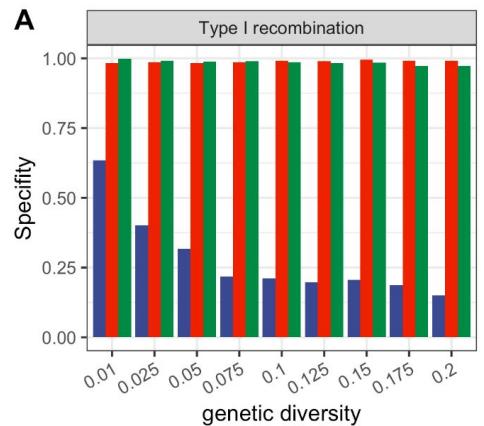


Boni et al. (Genetics, 2007)

**But which methods  
work best for detecting  
and localizing  
recombination  
breakpoints?**

# Sensitivity versus specificity

Detection power increases with genetic diversity but there is a tradeoff between power (sensitivity) and specificity.



**Specificity** = True negative rate

**Power** = True positive rate

method

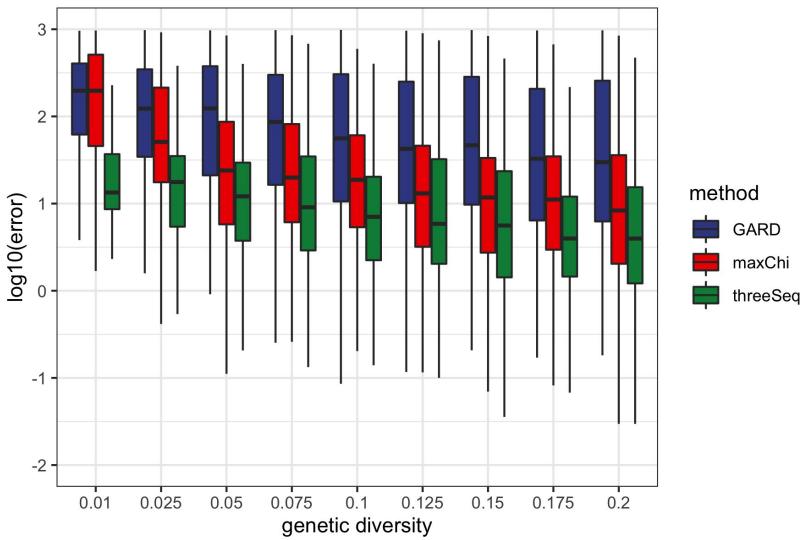
- GARD
- maxChi
- threeSeq



Shi Cen

# Breakpoint location accuracy

3SEQ performs best in accurately locating breakpoints but is still highly dependent on patterns of genetic polymorphisms in the sequences.



	Type I	Type II & Type III
GARD	$364.66 \pm 4.47$	$195.72 \pm 4.01$
maxChi	$395.27 \pm 22.40$	$72.41 \pm 4.12$
3SEQ	$203.40 \pm 19.19$	$33.69 \pm 2.07$

Localization accuracy (mean $\pm$ SEM) of three detection methods



Shi Cen

# Recombination vs. mutation rates

Whether or not it is possible to infer phylogenies ultimately depends of the ratio of the recombination rate  $r$  to the mutation rate  $m$ .

If  $r/m \ll 1$ , most changes in the genome occur due to mutation and it will generally be possible to infer local phylogenies within non-recombining regions.

If  $r/m > 1$ , most changes occur by recombination and there will not be enough mutations between recombination breakpoints to reliably reconstruct phylogenies.

# Recombination vs. mutation rates

The ratio r/m varies widely among different microbial pathogens

**Table 1** The ratio of nucleotide changes as the result of recombination relative to point mutation ( $r/m$ ) for different bacteria and archaea estimated from MLST data using ClonalFrame

Species	Phylum/division	Ecology	n STs	n loci	r/m	95% CI	Reference
<i>Flavobacterium psychrophilum</i>	Bacteroidetes	Obligate pathogen	33	7	63.6	32.8–82.8	Nicolas <i>et al.</i> (2008)
<i>Plagibacter ubique</i> (SAR 11)	z-proteobacteria	Free-living, marine	9	8	63.1	47.6–81.8	Vergin <i>et al.</i> (2007)
<i>Vibrio parahaemolyticus</i>	y-proteobacteria	Free-living, marine (OP)	20	7	39.8	27.4–48.2	Gonzalez-Escalona <i>et al.</i> (2008) web.mpiib-berlin.mpg.de/mlst
<i>Salmonella enterica</i>	y-proteobacteria	Commensal	50	7	30.2	21.0–36.5	Bisharat <i>et al.</i> (2007)
<i>Vibrio vulnificus</i>	y-proteobacteria	Free-living, marine (OP)	41	5	26.7	19.4–33.3	Hanage <i>et al.</i> (2005)
<i>Streptococcus pneumoniae</i>	Firmicutes	Commensal (OP)	52	6	23.1	16.7–29.0	Tanabe <i>et al.</i> (2007)
<i>Microcystis aeruginosa</i>	Cyanobacteria	Free-living, aquatic	79	7	18.3	13.7–21.2	Enright <i>et al.</i> (2001)
<i>Streptococcus pyogenes</i>	Firmicutes	Commensal (OP)	50	7	17.2	6.8–24.4	pubmlst.org
<i>Helicobacter pylori</i>	z-proteobacteria	Commensal (OP)	117	8	13.6	12.2–15.5	Jolley <i>et al.</i> (2005)
<i>Moraxella catarrhalis</i>	y-proteobacteria	Commensal (OP)	50	8	10.1	4.5–18.6	Salerno <i>et al.</i> (2007)
<i>Neisseria meningitidis</i>	β-proteobacteria	Commensal (OP)	83	7	7.1	5.1–9.5	Olvera <i>et al.</i> (2006)
<i>Plesiomonas shigelloides</i>	y-proteobacteria	Free-living, aquatic	58	5	7.1	3.8–13.0	pubmlst.net
<i>Neisseria lactamica</i>	β-proteobacteria	Commensal	180	7	6.2	4.9–7.4	Vos & Velter (2008)
<i>Myxococcus xanthus</i>	δ-proteobacteria	Free-living, terrestrial	57	5	5.5	1.9–11.3	Meats <i>et al.</i> (2003)
<i>Haemophilus influenzae</i>	y-proteobacteria	Commensal (OP)	50	7	3.7	2.6–5.4	Baldo <i>et al.</i> (2006)
<i>Wolbachia</i> b complex	z-proteobacteria	Endosymbiont	16	5	3.5	1.8–6.3	Stoddard <i>et al.</i> (2007)
<i>Campylobacter insulaeigrae</i>	z-proteobacteria	Commensal (OP)	59	7	3.2	1.9–5.0	Mayor <i>et al.</i> (2007)
<i>Mycoplasma pneumoniae</i>	Firmicutes	Commensal (OP)	33	7	3.0	1.1–5.8	Olvera <i>et al.</i> (2006)
<i>Haemophilus parasuis</i>	y-proteobacteria	Commensal (OP)	79	7	2.7	2.1–3.6	Papke <i>et al.</i> (2004)
<i>Campylobacter jejuni</i>	z-proteobacteria	Commensal (OP)	110	7	2.2	1.7–2.8	pubmlst.org
<i>Halorubrum</i> sp.	Halobacteria (Archaea)	Halophile	28	4	2.1	1.2–3.3	Sorokin <i>et al.</i> (2006)
<i>Pseudomonas viridisflava</i>	y-proteobacteria	Free-living, plant pathogen	92	3	2.0	1.2–2.9	Goss <i>et al.</i> (2005)
<i>Bacillus weihenstephanensis</i>	Firmicutes	Free-living, terrestrial	36	6	2.0	1.3–2.8	de Las Rivas <i>et al.</i> (2004)
<i>Pseudomonas syringae</i>	y-proteobacteria	Free-living, plant pathogen	95	4	1.5	1.1–2.0	Whittaker <i>et al.</i> (2005)
<i>Sulfolobus islandicus</i>	Thermoprotei (Archaea)	Thermacidiphile	17	5	1.2	0.1–4.5	Castillo and Greenberg (2007)
<i>Ralstonia solanacearum</i>	β-proteobacteria	Plant pathogen	58	7	1.1	0.7–1.6	Homan <i>et al.</i> (2002)
<i>Enterococcus faecium</i>	Firmicutes	Commensal (OP)	15	7	1.1	0.3–2.5	Miller <i>et al.</i> (2007)
<i>Mastigocladus laminosus</i>	Cyanobacteria	Thermophile	34	4	0.9	0.5–1.5	Coscolla and Gonzalez-Candelas (2007)
<i>Legionella pneumophila</i>	y-proteobacteria	Protozoa pathogen	30	2	0.9	0.2–1.9	Ludders <i>et al.</i> (2005)
<i>Microcoleus chthonoplastes</i>	Cyanobacteria	Free-living, marine	22	2	0.8	0.2–1.9	Sorokin <i>et al.</i> (2006)
<i>Bacillus thuringiensis</i>	Firmicutes	Insect pathogen	22	6	0.8	0.4–1.3	de Las Rivas <i>et al.</i> (2004)
<i>Bacillus cereus</i>	Firmicutes	Free-living, terrestrial (OP)	13	6	0.7	0.2–1.6	Salcedo <i>et al.</i> (2003)
<i>Oenococcus oeni</i>	Firmicutes	Free-living, terrestrial	17	5	0.7	0.2–1.7	Ruiz-Garbajosa <i>et al.</i> (2006)
<i>Escherichia coli</i> ET-1 group	y-proteobacteria	Commensal (free-living?)	44	7	0.7	0.03–2.0	Enersen <i>et al.</i> (2006)
<i>Listeria monocytogenes</i>	Firmicutes	Free-living, terrestrial (OP)	34	7	0.7	0.4–1.1	Pannekoek <i>et al.</i> (2008)
<i>Enterococcus faecalis</i>	Firmicutes	Commensal (OP)	37	7	0.6	0.0–3.2	Diancourt <i>et al.</i> (2007)
<i>Porphyromonas gingivalis</i>	Bacteroidetes	Obligate pathogen	99	7	0.4	0.0–3.4	Thiaguppanai <i>et al.</i> (2007)
<i>Yersinia pseudotuberculosis</i>	y-proteobacteria	Obligate pathogen	43	7	0.3	0.0–1.1	Vos & Didelot (ISME, 2008)
<i>Chlamydia trachomatis</i>	Chlamydiaceae	Obligate pathogen	14	7	0.3	0.0–1.8	Diancourt <i>et al.</i> (2005)
<i>Klebsiella pneumoniae</i>	y-proteobacteria	Free-living, terrestrial (OP)	45	7	0.3	0.0–2.1	Diavatopoulos <i>et al.</i> (2005)
<i>Bordetella pertussis</i>	β-proteobacteria	Obligate pathogen	32	7	0.2	0.0–0.7	Rasback <i>et al.</i> (2007)
<i>Brachyspira</i> sp.	Spirochaetes	Commensal (OP)	36	7	0.2	0.1–0.4	Lemee <i>et al.</i> (2004)
<i>Clostridium difficile</i>	Firmicutes	Commensal (OP)	34	6	0.2	0.0–0.5	Arvand <i>et al.</i> (2007)
<i>Bartonella henselae</i>	z-proteobacteria	Obligate pathogen	14	7	0.1	0.0–0.7	Thiaguppanai <i>et al.</i> (2007)
<i>Lactobacillus casei</i>	Firmicutes	Commensal	32	7	0.1	0.0–0.5	Enright <i>et al.</i> (2000)
<i>Staphylococcus aureus</i>	Firmicutes	Commensal (OP)	53	7	0.1	0.0–0.6	Thiaguppanai <i>et al.</i> (2007)
<i>Rhizobium gallicum</i>	z-proteobacteria	Free-living, terrestrial	33	3	0.1	0.0–0.3	Silva <i>et al.</i> (2005)
<i>Leptospira interrogans</i>	Spirochaetes	Commensal (OP)	61	7	0.02	0.0–0.1	Vos & Didelot (ISME, 2008)