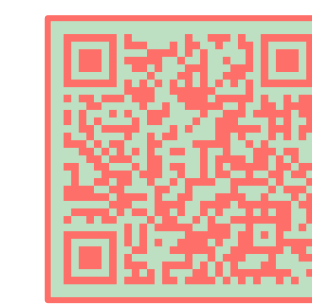


# Tools for analysis of clonal populations in R



More about poppr

[github.com/grunwaldlab/poppr#readme](https://github.com/grunwaldlab/poppr#readme)

Code for this poster

[github.com/grunwaldlab/poppr-poster-aps-2016](https://github.com/grunwaldlab/poppr-poster-aps-2016)



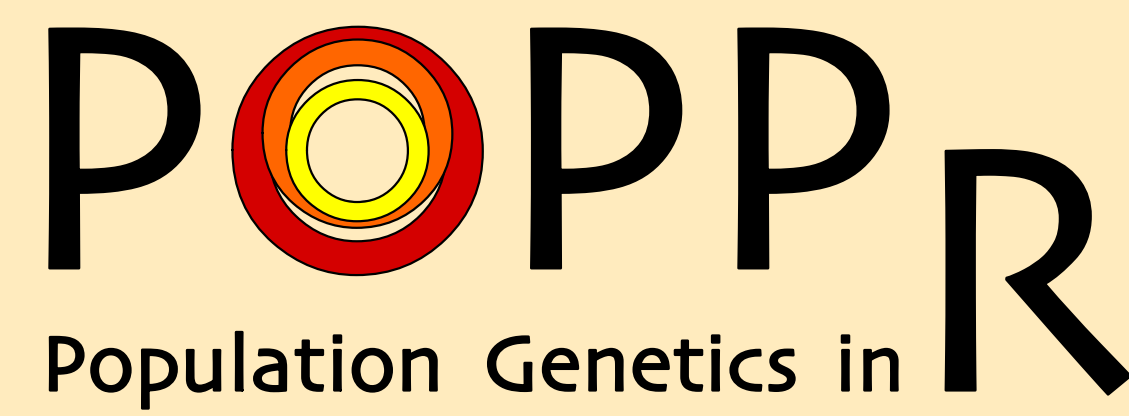
Zhian N. Kamvar<sup>1</sup> Jonah C. Brooks<sup>2</sup> Niklaus J. Grünwald<sup>1,3</sup>

<sup>1</sup>: Department of Botany and Plant Pathology, Oregon State University  
<sup>2</sup>: College of Electrical Engineering and Computer Science, Oregon State University  
<sup>3</sup>: Horticultural Crops Research Laboratory, USDA Agricultural Research Service, Corvallis, OR, USA

## Motivation

Knowledge of the population dynamics and evolution of plant pathogens allows inferences on evolutionary processes involved in their adaptation to hosts, pesticides, and other environmental pressures. With the advent of high-throughput sequencing technologies, obtaining genome-wide data has become easier and cheaper than ever before with techniques such as genotyping-by-sequencing (GBS).

We present here an overview of the **novel tools in the R package poppr** as it pertains to traditional and high-throughput population genetic data with select applications (Kamvar et al., 2015).

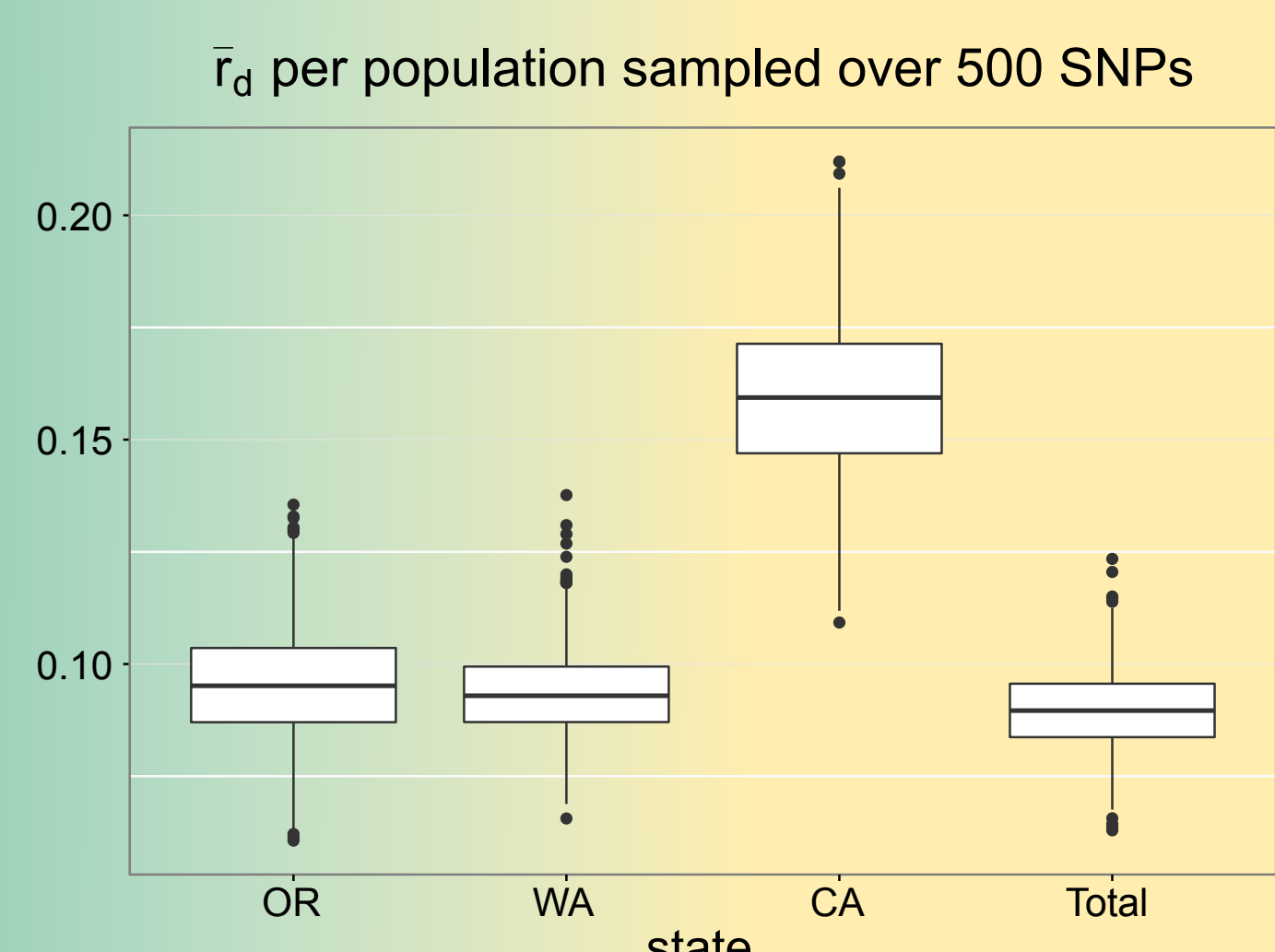
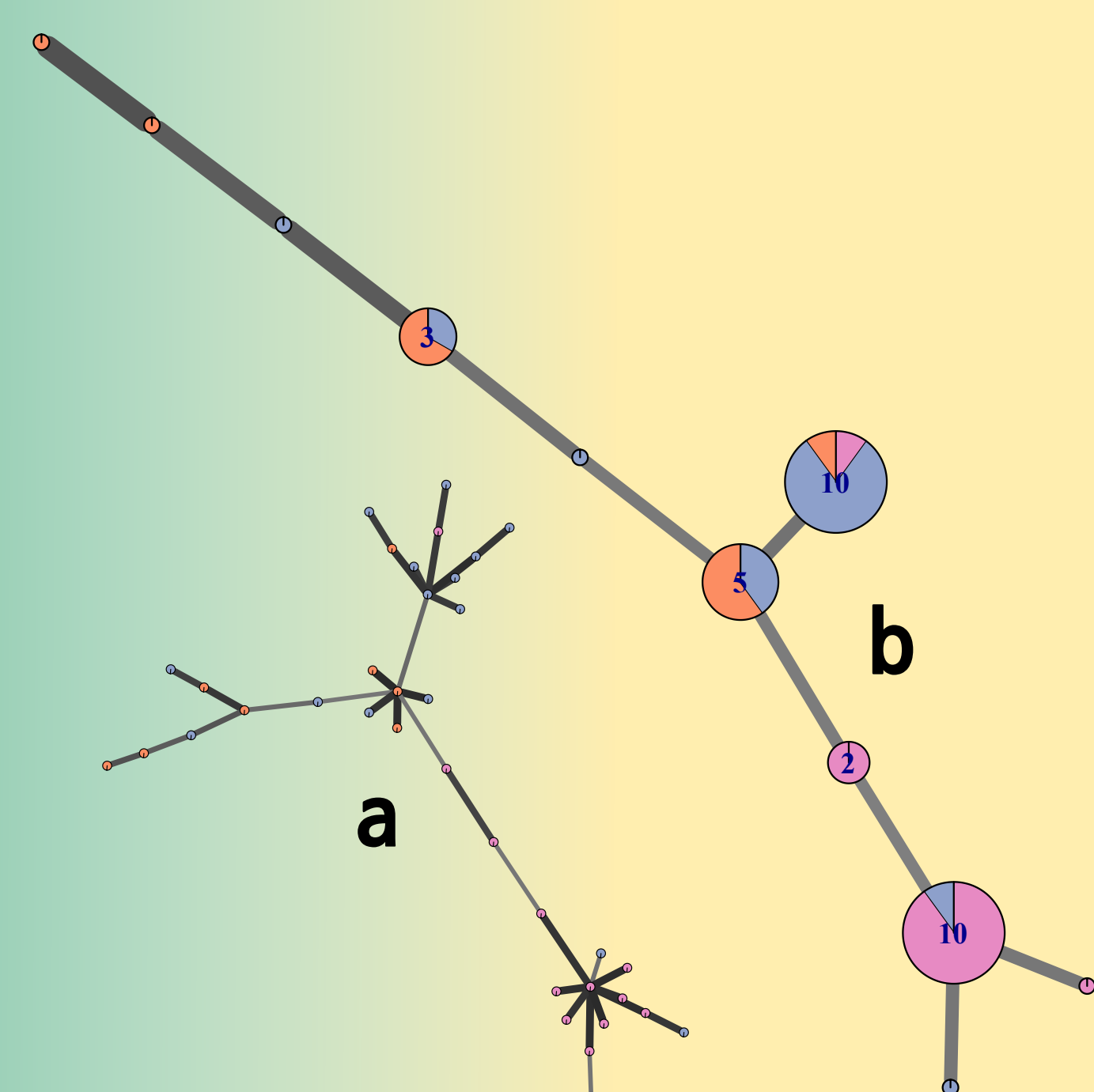
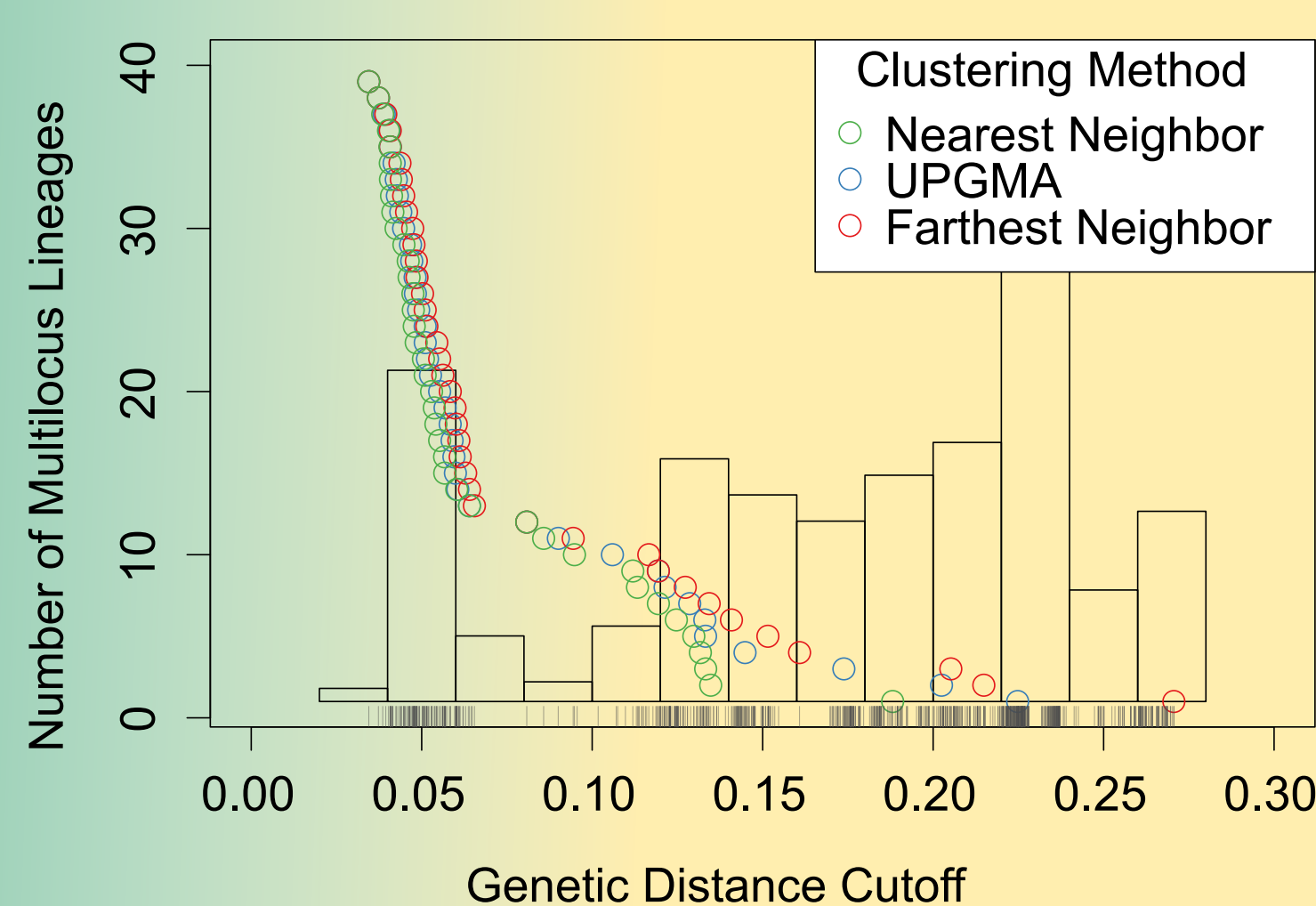
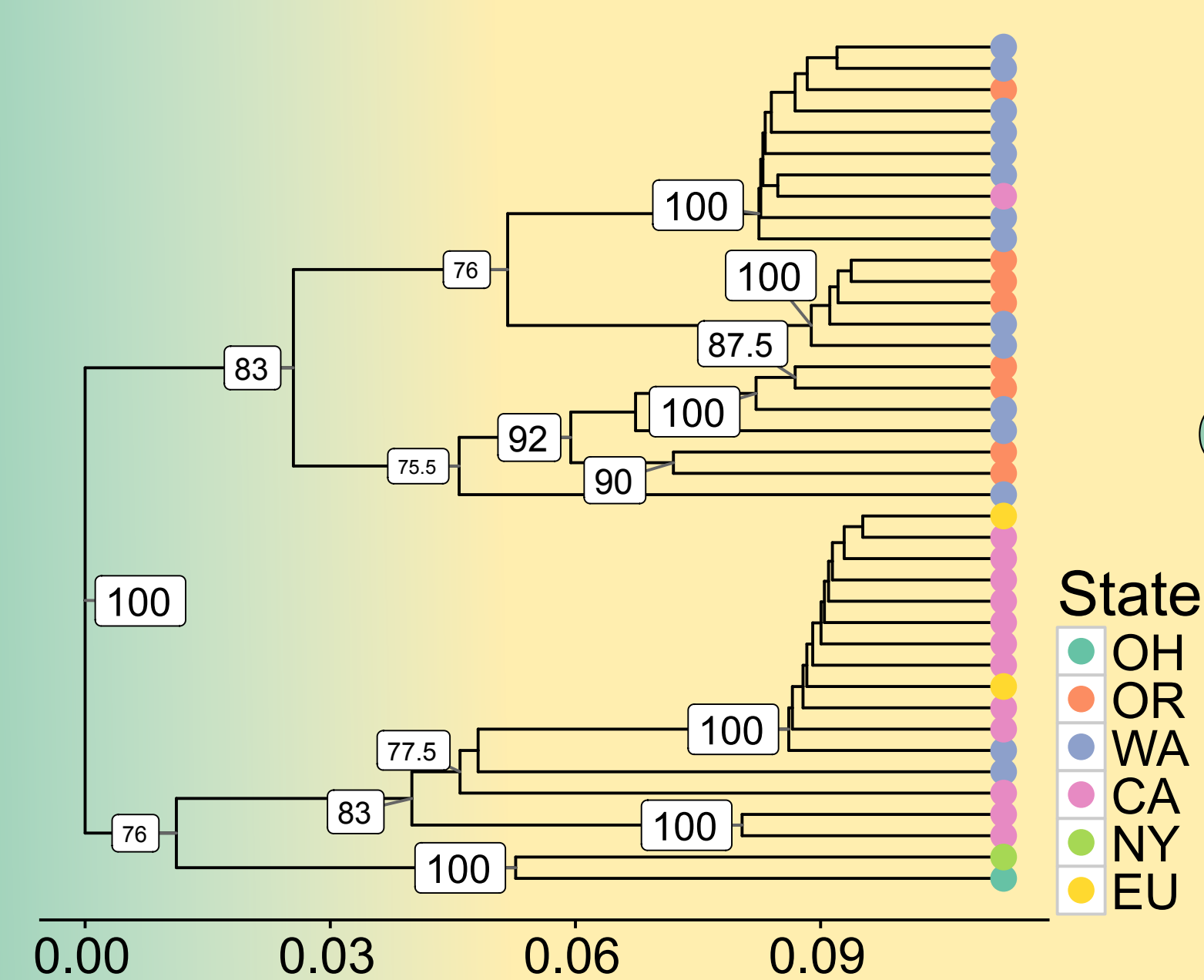


## Genotype Accumulation Curve

One recent addition to poppr is the genotype accumulation curve (GAC). When starting analysis with a handful of markers, it's important to check that the number of markers you have is sufficient to encompass the diversity of multilocus genotypes (MLGs). This is visualised by the GAC (S1), which randomly samples  $n$  loci from your data (x axis), reporting the number of MLGs observed (y axis).

## GBS (SNP) Data

*Phytophthora rubi*: 40 genotypes, 43,414 SNPs  
(JF Tabima & NJ Grünwald, unpublished)



## Flexible bootstrap analysis

Bootstrapping dendrograms is a way to assess the internal consistency of your data. Poppr provides the novel function `about` ("any boot") for quick and versatile bootstrap analysis. This function will create and bootstrap dendrograms for any data type including SNPs (G1) and SSRs (S2), with any genetic distance, from individuals to populations.

## Define Multilocus Lineages by Genetic Distance

Multilocus genotypes (MLGs) are initially defined as the unique combination of alleles across all loci. New MLGs are created when this combination differs, which could be caused by genotyping error. New in poppr v2 is the ability to use genetic distance to identify and collapse similar MLGs into multilocus lineages (MLLs) for any data type using three different clustering methods, Nearest Neighbor, Average Neighbor (UPGMA), and Farthest Neighbor (G2, S3).

## Minimum Spanning Networks with filtered MLGs

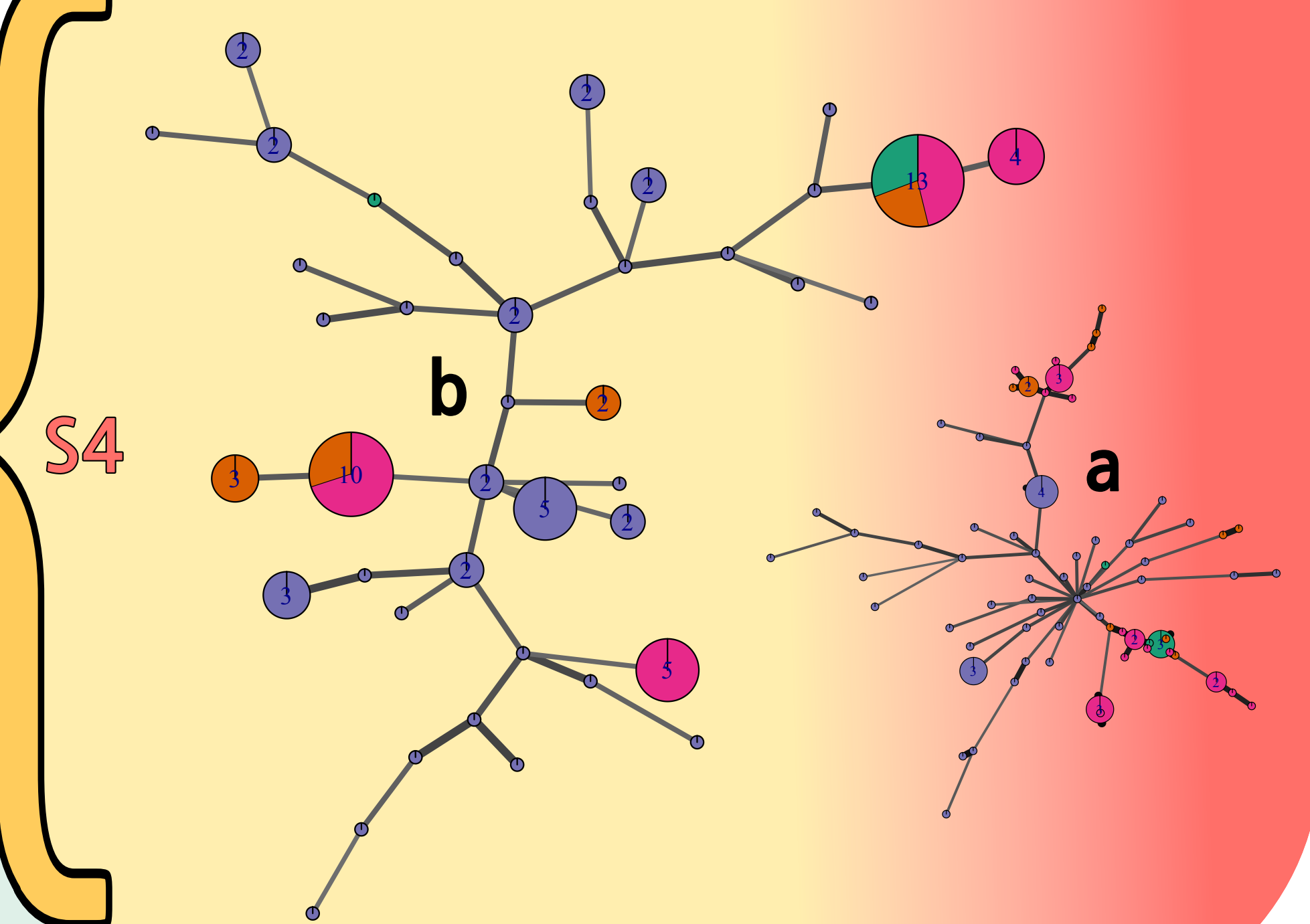
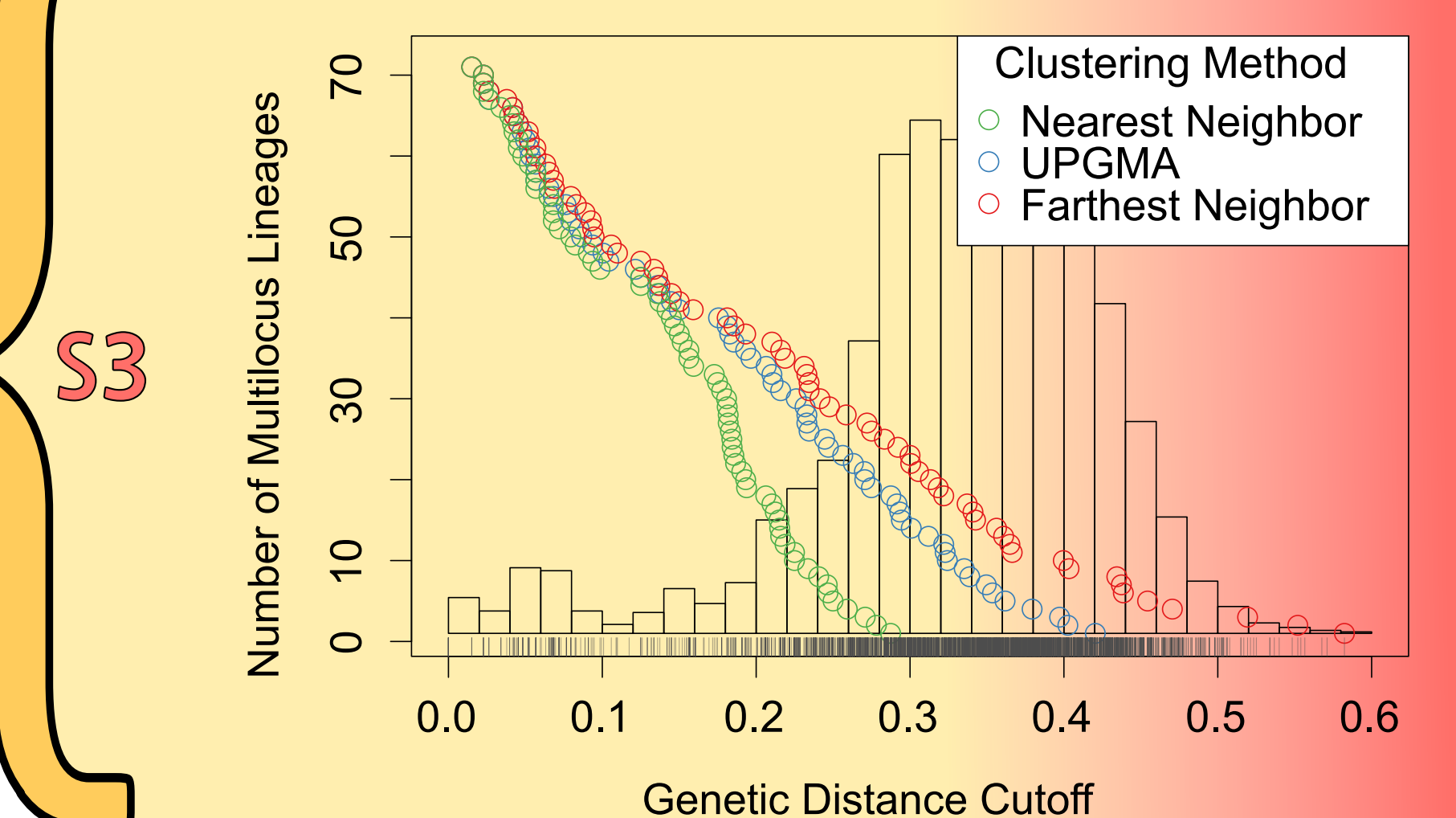
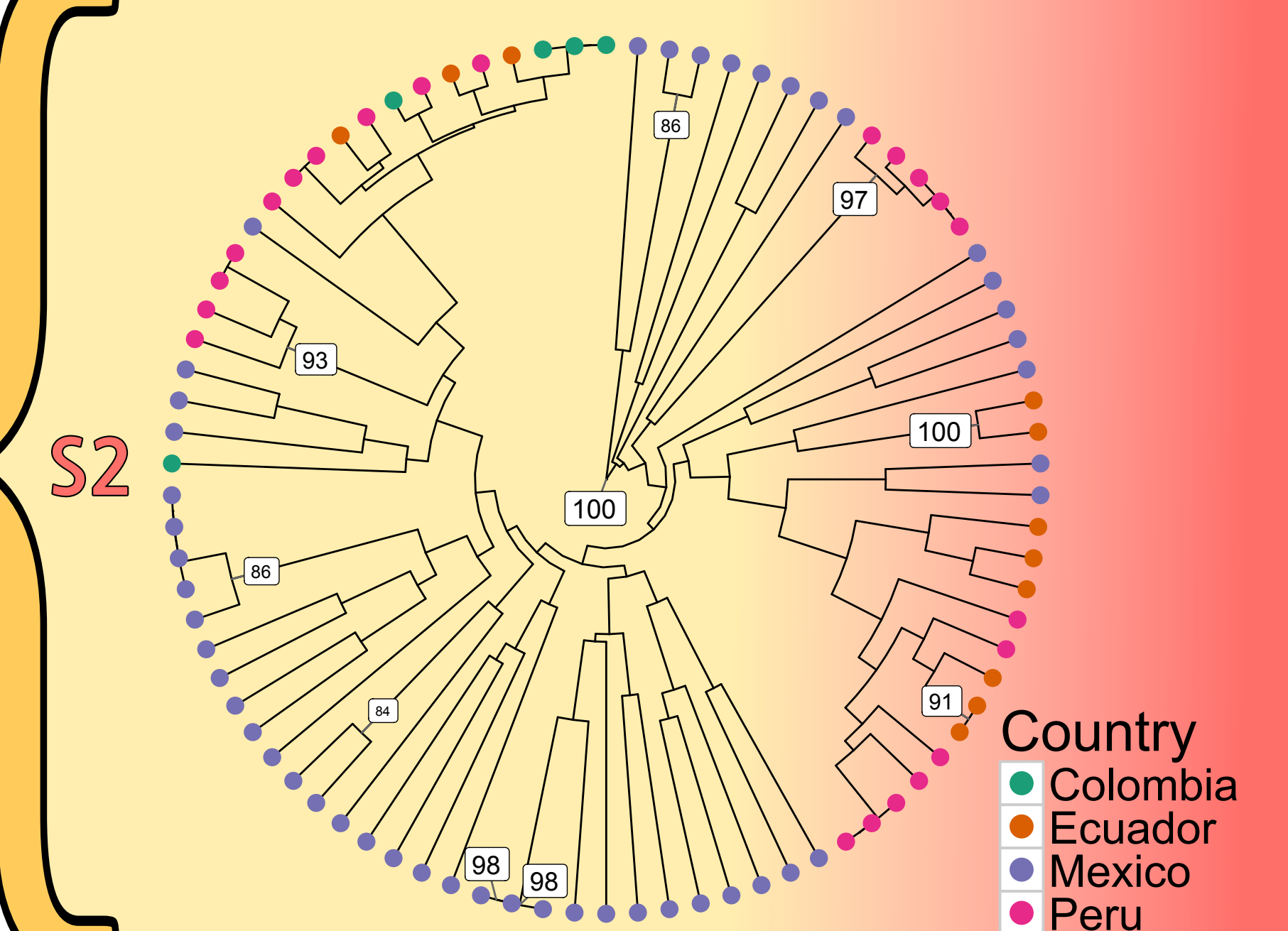
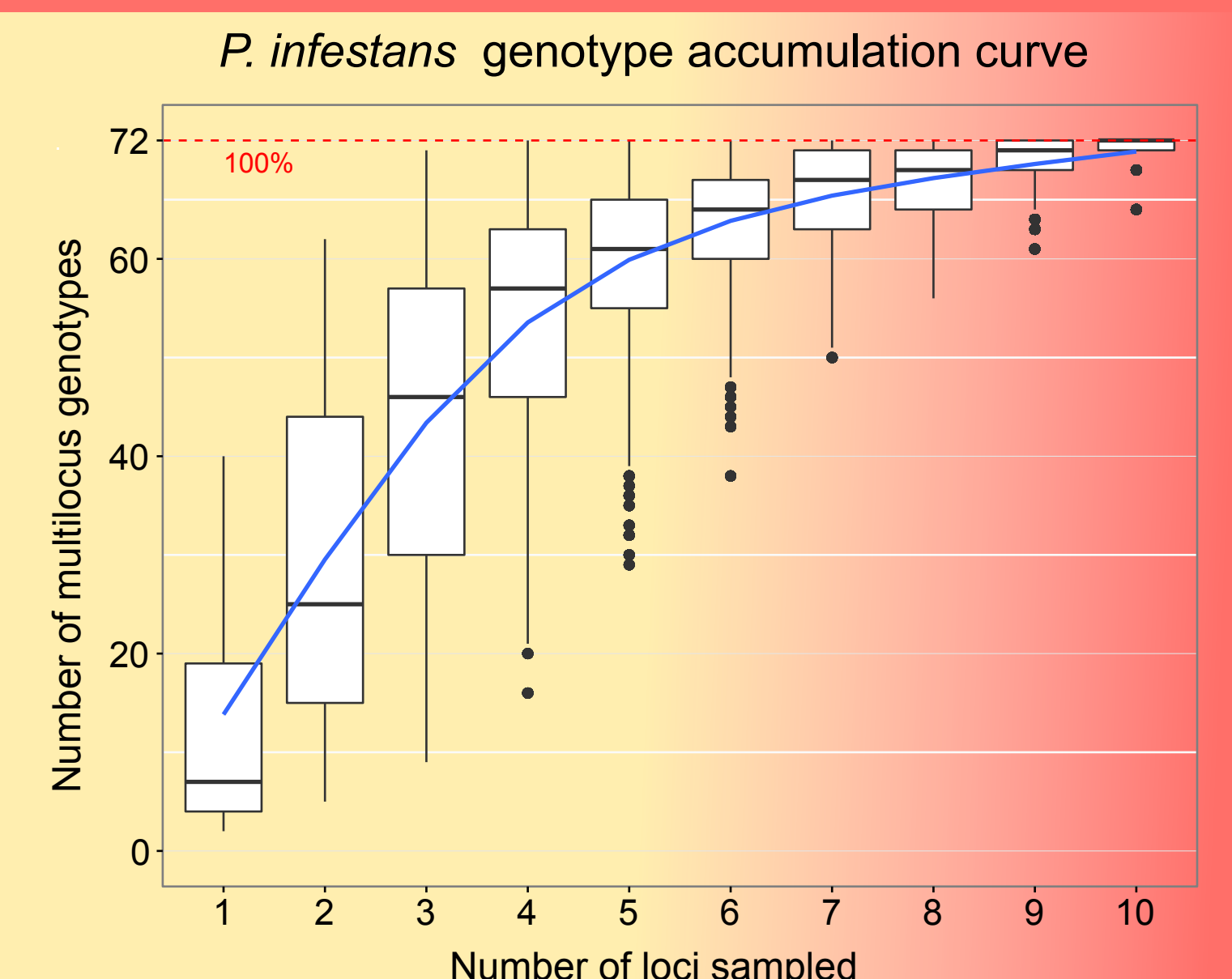
Minimum spanning networks are a useful tool for assessing the population structure of clonal organisms by showing the shortest connections between multilocus genotypes (G3-a, S4-a). New to poppr v2 is the ability to create minimum spanning networks showing connections between MLLs (G3-a, S4-a).

## Index of Association for SNPs

The index of association (and its standardized form) is a measure of multilocus linkage disequilibrium, which can serve as a metric for clonality. It was initially implemented for traditional markers such as AFLP or SSR, but with high-throughput sequencing technologies like GBS, assessing linkage across thousands of loci became computationally intensive. We have created a unique method for quickly assessing the clonality of populations by sampling SNPs within a sliding window or randomly. Fig. G4 was generated by randomly sampling 500 SNPs 1000 times over each population showing that the CA samples have a strong signal of clonal reproduction.

## SSR Data

*Phytophthora infestans*: 86 genotypes, 11 SSRs  
(Goss et al., 2014)



## Acknowledgements

This work was supported in part by US Department of Agriculture (USDA) Agricultural Research Service Grant 5358-22000-039-00D, USDA National Institute of Food and Agriculture Grant 2011-68004-30154, USDA APHIS, the USDA-ARS Floriculture Nursery Initiative, and the USDA-Forest Service Forest Health Monitoring Program. GBS data were first filtered by Brian Knaus with the package `vcfR` before running the analysis (Knaus & Grünwald, 2016). You can find the scripts/tutorials at [https://github.com/knausb/vcfR\\_class](https://github.com/knausb/vcfR_class)

Kamvar ZN, Brooks JC and Grünwald NJ (2015) Novel R tools for analysis of genome-wide population genetic data with emphasis on clonality. *Front. Genet.* 6:208. doi: 10.3389/fgene.2015.00208

Goss, Erica M., et al. "The Irish potato famine pathogen *Phytophthora infestans* originated in central Mexico rather than the Andes." *Proceedings of the National Academy of Sciences* 111.24 (2014): 8791-8796. doi: 10.1073/pnas.1401884111

Knaus, Brian J., and Niklaus J. Grünwald. In press. `vcfR`: a package to manipulate and visualize variant call format data in R. *Molecular Ecology Resources*. doi: 10.1111/1755-0998.12549.