

A possible way to identifying mixed infected samples

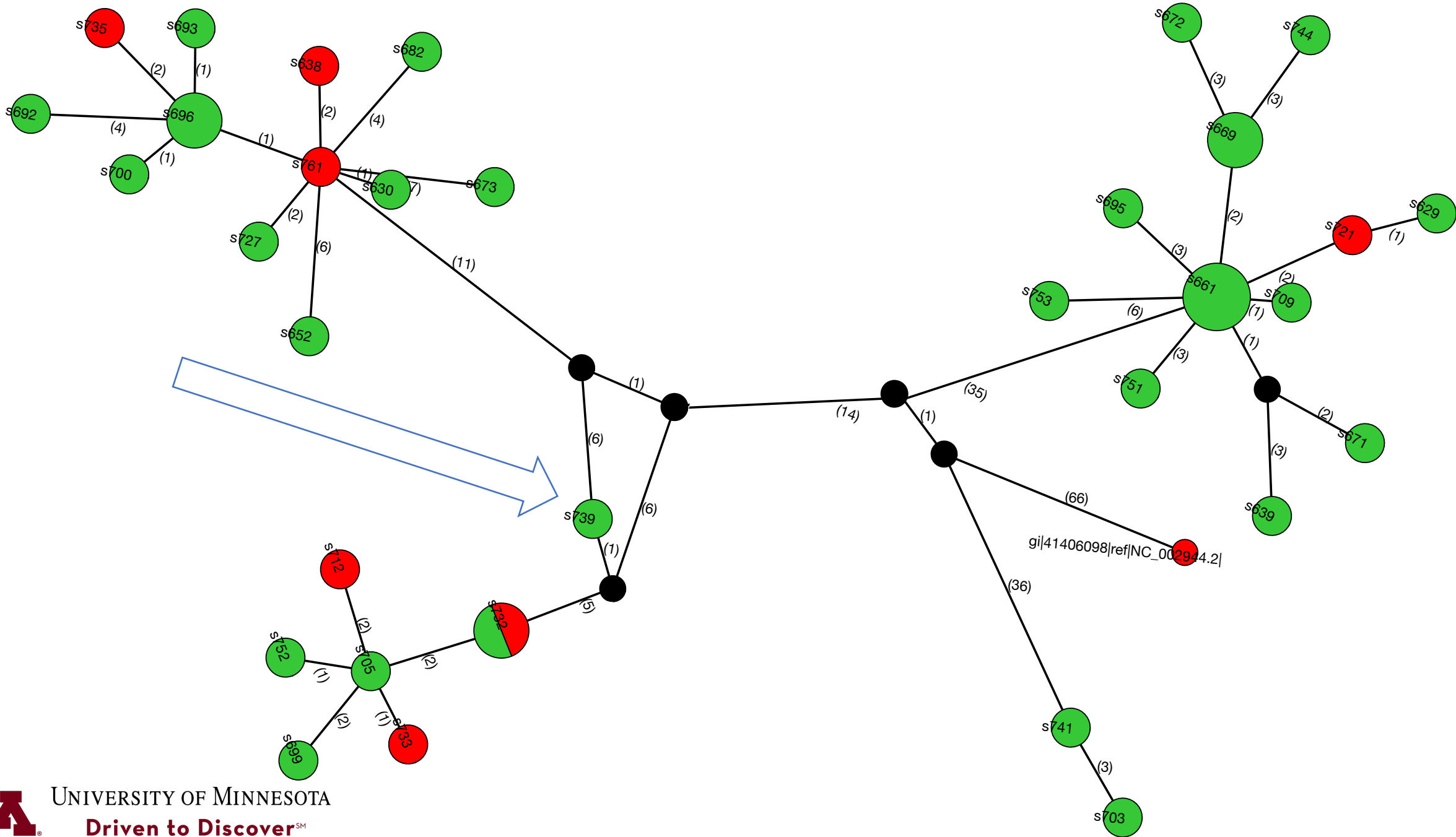
Ph. D candidate Yuanyuan Wang

Computational Biology and Bioinformatics

University of Minnesota, Twin Cities

1. Multiple infections need to be identified

- Why do we care about multiple infections?
 - No recognition of mixture sample will result in **wrong** phylogenetic trees
 - Multiple infected individuals associated **increased transmission intensity** in the local environment. (Manske M, Miotto O, et al., 2012). Identifying those individuals maybe crucial to disease control.



2. Definitions of multiple infections

- From WGS data, how to tell if the animal has multiple infection?

1) [Biology]

Heterozygous loci:

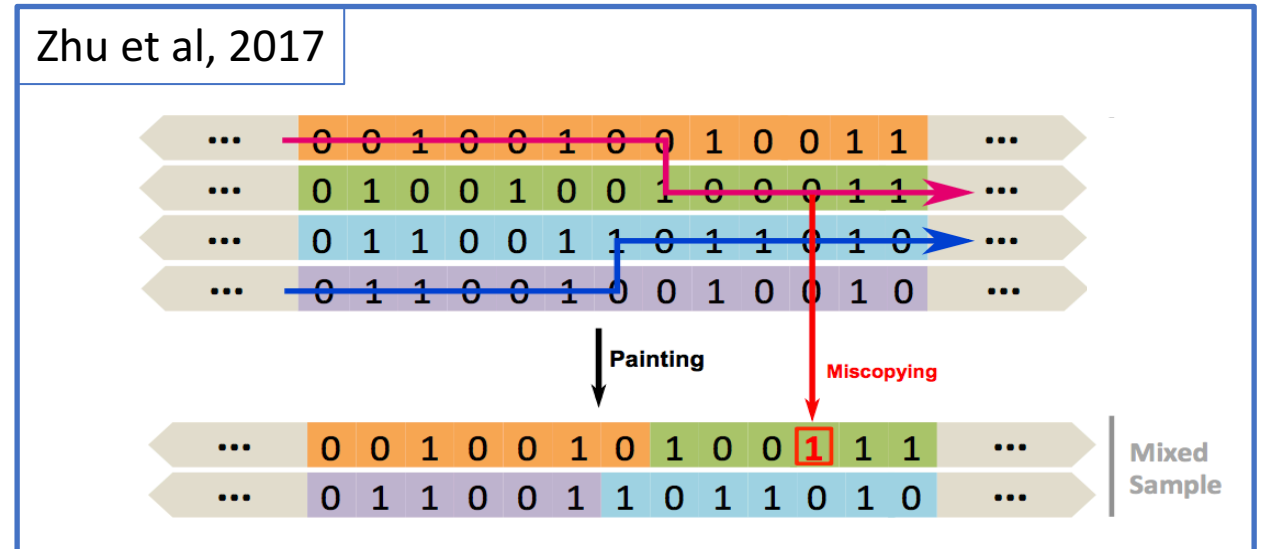
- when a loci that has multiple alternative alleles. (not suitable for single-colony isolates)

2) [Transmission]

Patterns/Haplotypes:

- when a sequence (sample) shares haplotypes/ evolutionary path with other sequences.

MN dataset	Position on chromosome			
	20	68	250	296
Reference allele	A	G	C	T
Minor allele #1	C	C	T	A
Minor allele #2	G	T	G	G



3. A quick way to screening for mixture samples using allele frequency plots

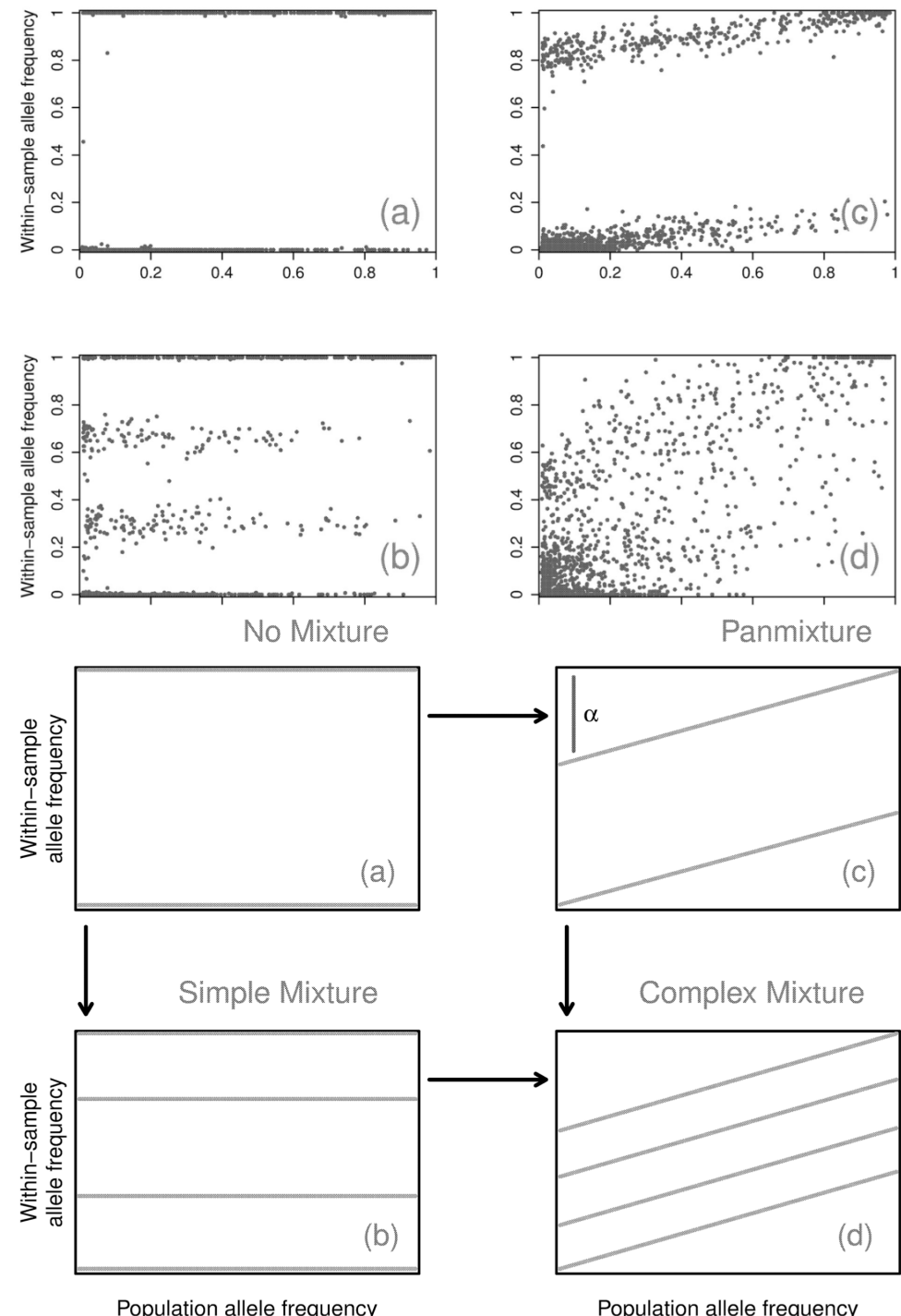
RESEARCH ARTICLE

Inferring Strain Mixture within Clinical *Plasmodium falciparum* Isolates from Genomic Sequence Data

John D. O'Brien^{1*}, Zamin Iqbal², Jason Wendler³, Lucas Amenga-Etego^{2,4}

¹ Mathematics Department, Bowdoin College, Brunswick, Maine, United States of America, ² Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford, Oxfordshire, United Kingdom, ³ Pacific Northwest National Laboratory, Richland, Washington, United States of America, ⁴ Navrongo Health Research Centre, Navrongo, Upper East Region, Ghana

*jobrien@bowdoin.edu



3. A quick way to screening for mixture samples using allele frequency plots

RESEARCH ARTICLE

Inferring Strain Mixture within Clinical *Plasmodium falciparum* Isolates from Genomic Sequence Data

John D. O'Brien^{1*}, Zamin Iqbal², Jason Wendler³, Lucas Amenga-Etego^{2,4}

¹ Mathematics Department, Bowdoin College, Brunswick, Maine, United States of America, ² Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford, Oxfordshire, United Kingdom, ³ Pacific Northwest National Laboratory, Richland, Washington, United States of America, ⁴ Navrongo Health Research Centre, Navrongo, Upper East Region, Ghana

$$\mathbb{P}(\mathcal{Q}, \mathcal{P}, \mathcal{W}, \alpha, \nu, K | \mathcal{D}_i) \propto \mathbb{P}(\mathcal{D}_i | \mathcal{Q}, \mathcal{P}, \nu, K) \cdot \mathbb{P}(\mathcal{Q} | \mathcal{P}, \nu, K, \mathcal{W}, \alpha) \cdot \mathbb{P}(\mathcal{P}) \cdot \mathbb{P}(\nu) \cdot \mathbb{P}(\mathcal{W} | K) \cdot \mathbb{P}(K) \cdot \mathbb{P}(\alpha).$$

Table 1. Parameters and definitions for the model and its description.

Parameter	Definition
N	Number of samples
M	Number of SNPs
K	Number of strains
$i = 1, \dots, N$	Index for samples
$j = 1, \dots, M$	Index for SNPs
$r = 1, \dots, 2^K$	Index for bands / strain mixtures
p_j	(Non-reference) allele frequency for SNP j
$\mathcal{P} = [p_j]$	The PLAF for all SNPs
$\mathcal{Q} = [q_{ij}]$	Within-sample allele frequency for SNP j in sample i
α	Degree of panmixia within a sample, panmixia coefficient
$\mathcal{S} = [s_1, \dots, s_K]$	Strains in a sample
$\mathcal{W} = [w_1, \dots, w_K]$	Strain proportions in a sample
λ_r	Band proportions within sample
ν	Variation parameter for Beta-binomial
WSAF	Within-sample allele frequency
PLAF	Population-level allele frequency

doi:10.1371/journal.pcbi



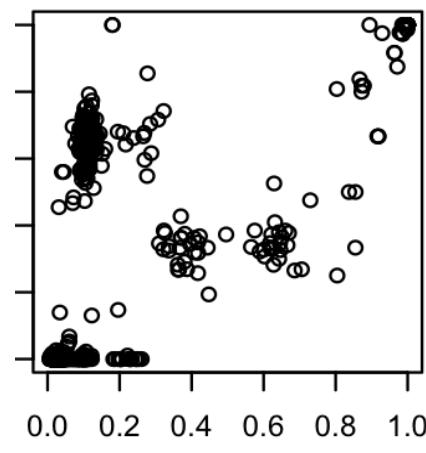
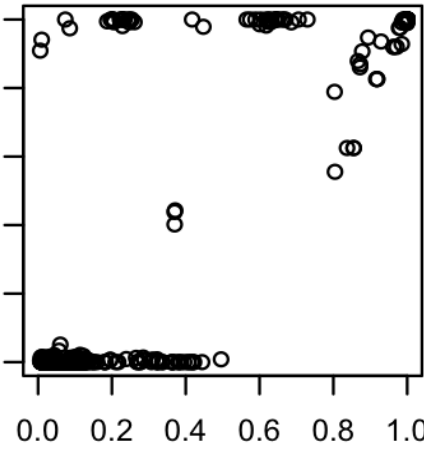
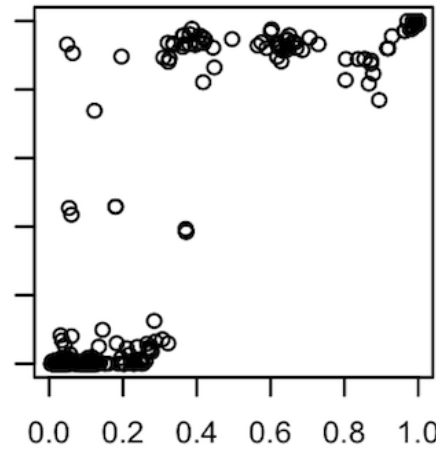
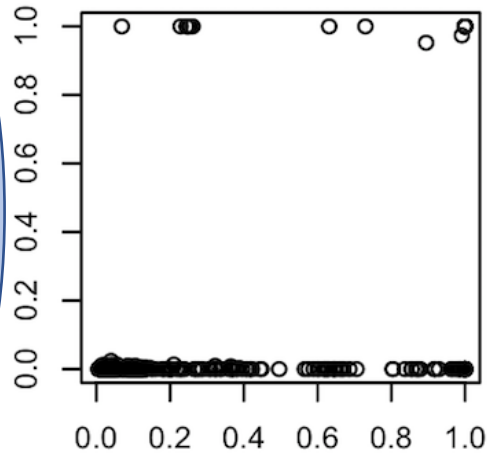
sample: s642

sample: s653

sample: s809

sample: s788

within-sam allele freq



S farm

Pop-level allele freq

Pop-level allele freq

Pop-level allele freq

Pop-level allele freq

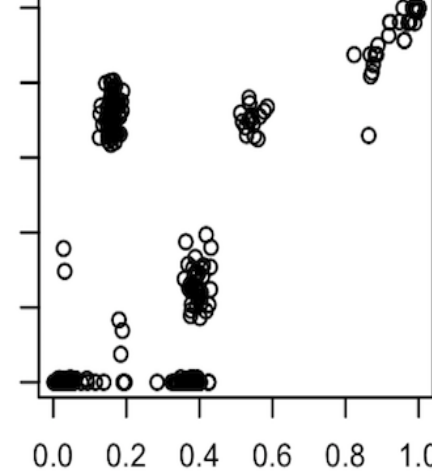
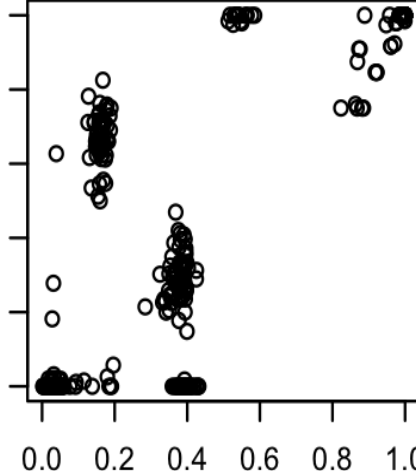
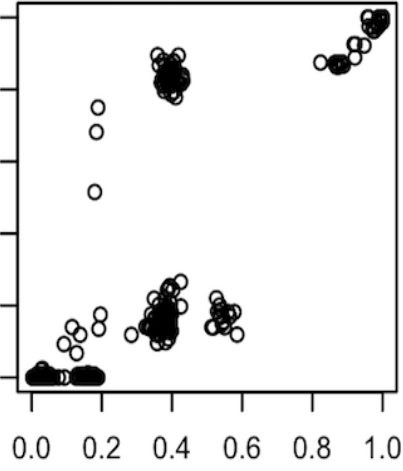
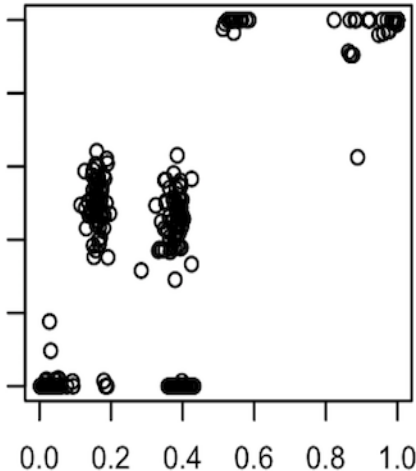
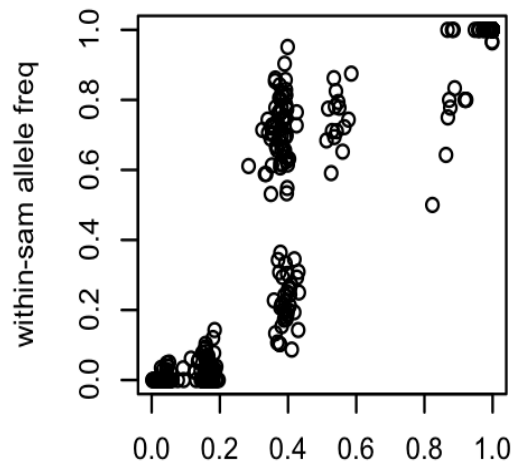
sample: s761

sample: s739

sample: s698

sample: s743

sample: s732



Q farm

Pop-level allele freq

Pop-level allele freq

Pop-level allele freq

Pop-level allele freq

Pop-level allele freq