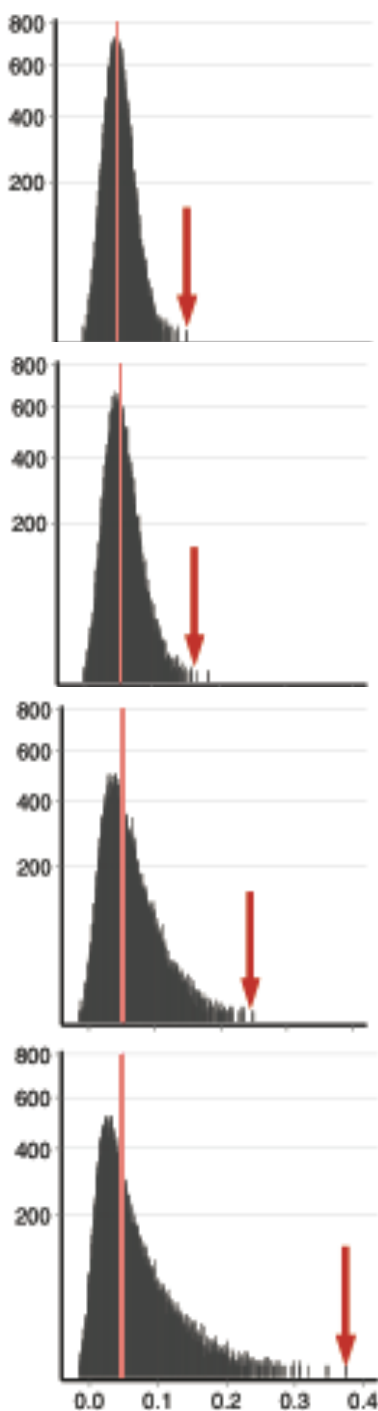


Reducing false positives in F_{ST} outlier tests with OutFLANK



Katie E. Lotterhos^{1,2}

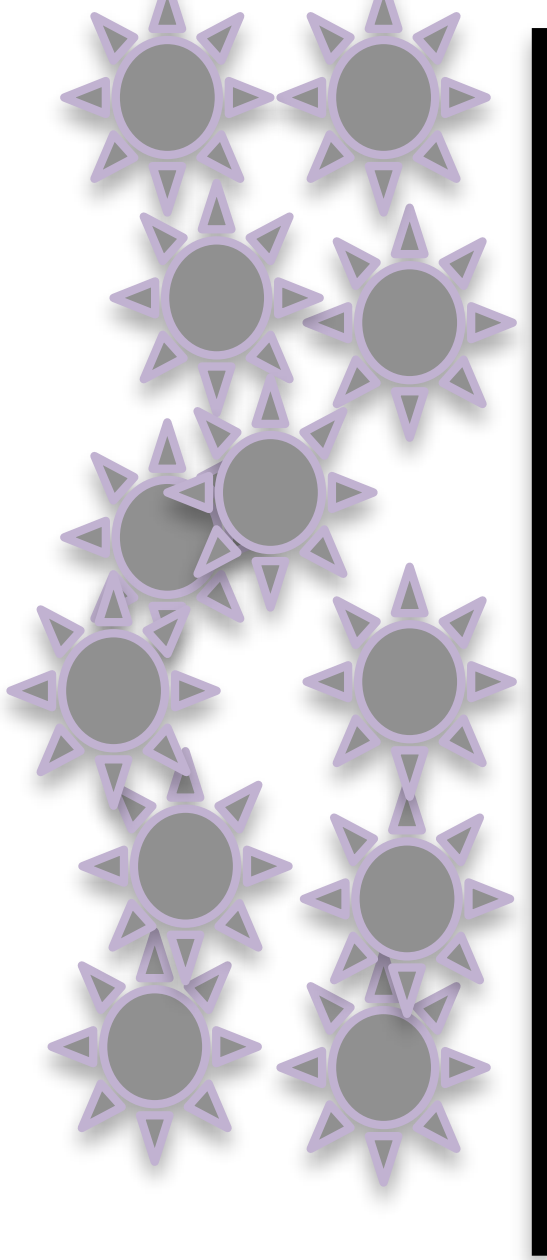
Michael C. Whitlock¹

¹University of British Columbia

²Wake Forest University

Question:

What is the genomic basis of local adaptation in heterogeneous environments?



N

Cold

S

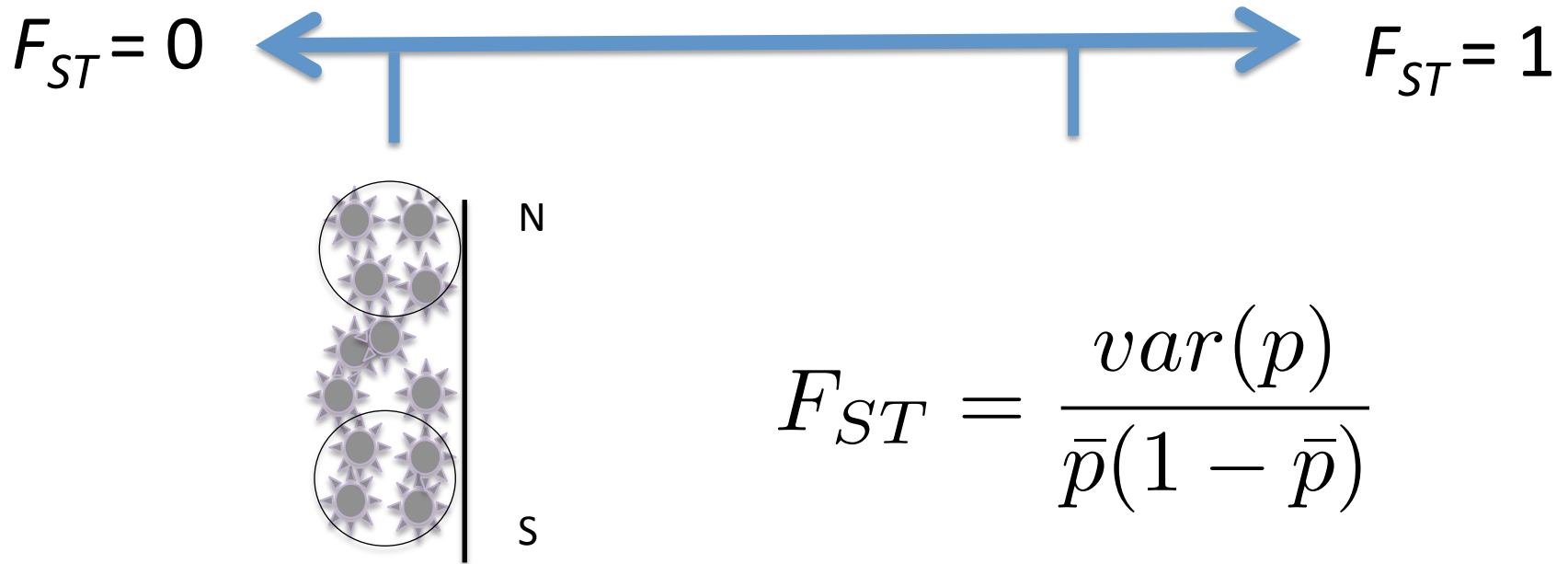
Warm

F_{ST} measures how different allele frequencies are between populations



$$F_{ST} = \frac{\text{var}(p)}{\bar{p}(1 - \bar{p})}$$

F_{ST} measures how different allele frequencies are between populations

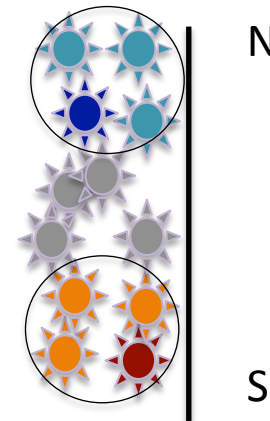


Non-adapted locus

F_{ST} measures how different allele frequencies are between populations



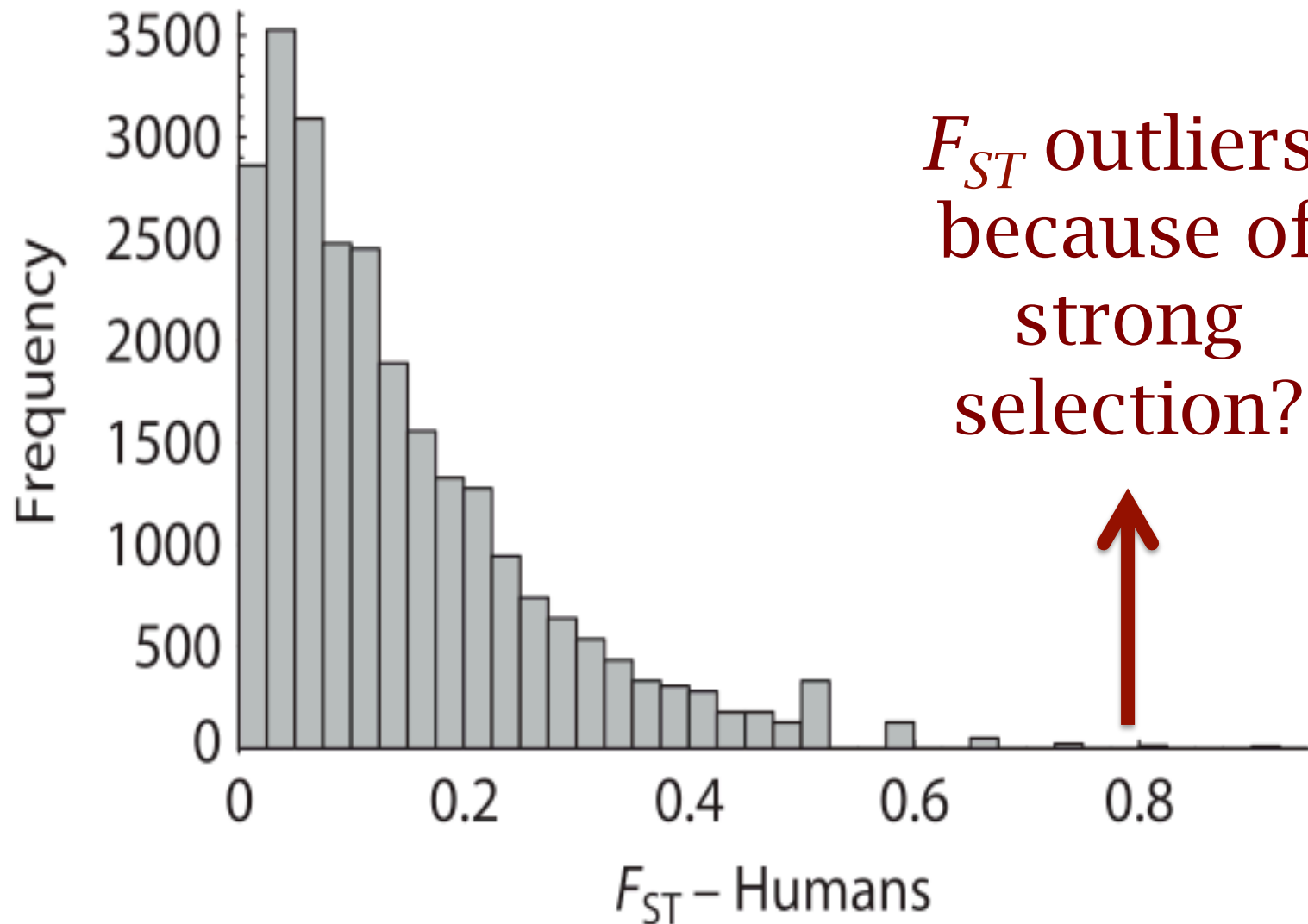
$$F_{ST} = \frac{var(p)}{\bar{p}(1 - \bar{p})}$$



Adapted locus

$$F_{ST} = 0$$

$$F_{ST} = 1$$



Problem:

How to separate the
effects of selection from
effects of an unknown
demographic history?

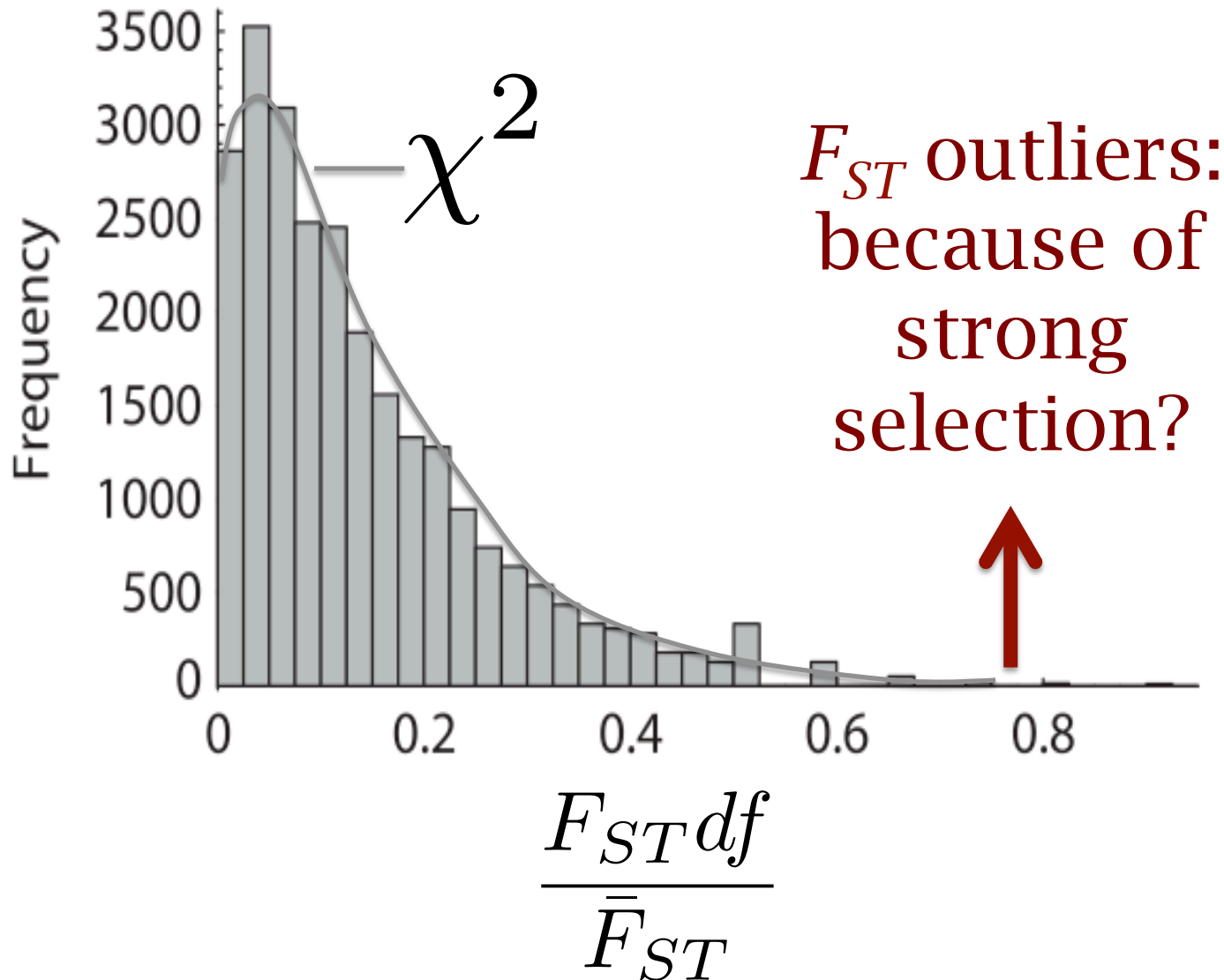
The good news:

The genome is filled with
approximately neutral
genes

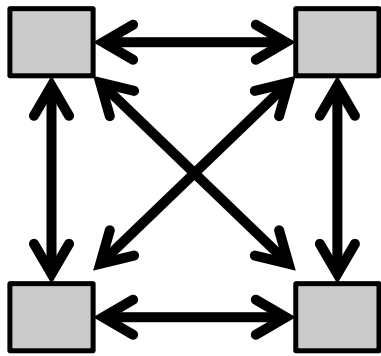
The bad news:

We don't know which
genes are neutral, or
precisely how they behave

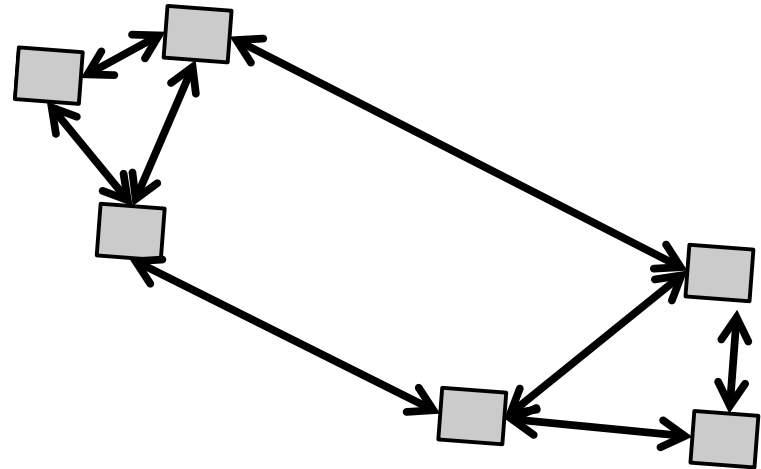
Lewontin-Krakauer Test



Difficulty: LK assumed sampled populations are equally independent of each other



True:
Island Model

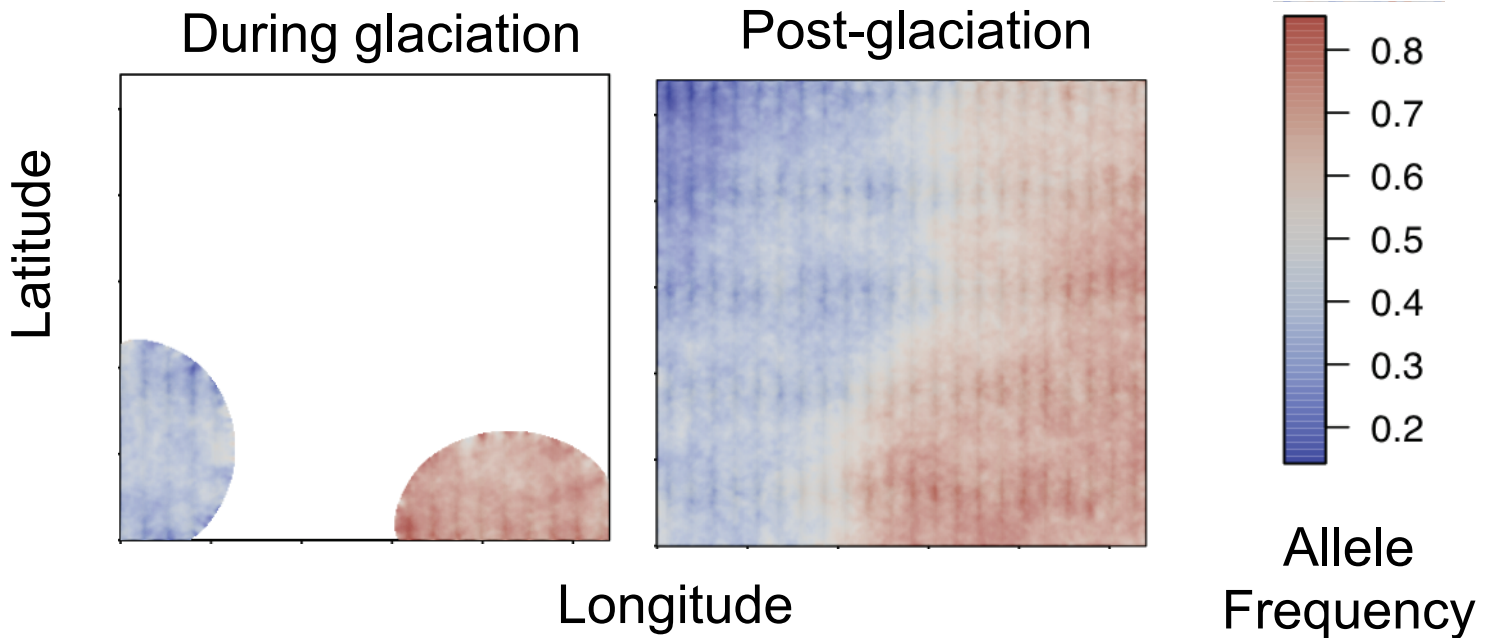


Not True:
Isolation by distance

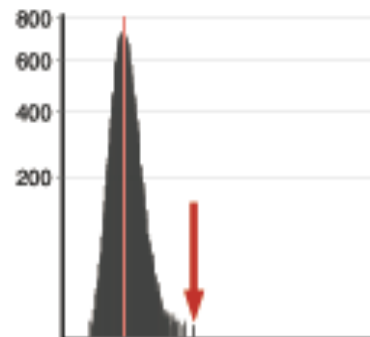
How sensitive are
 F_{ST} outlier tests to
demographic
history?

LandSHARC:

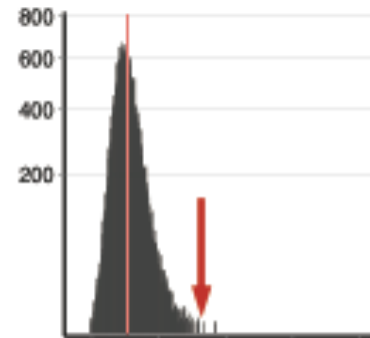
Landscape Simulator for Haploid Alleles in Realistic Climates



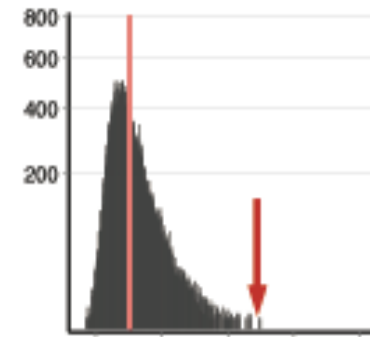
F_{ST}
distributions:
neutral loci



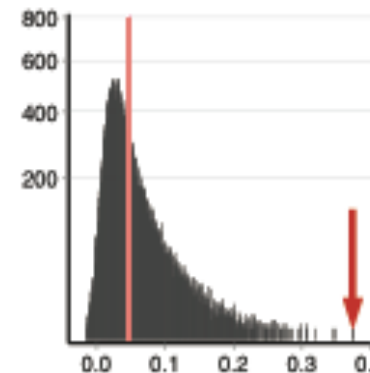
IM
Island Model
(at equilibrium)



IBD
Isolation by distance
(at equilibrium)



1R
Expansion from one refuge
(non-equilibrium)

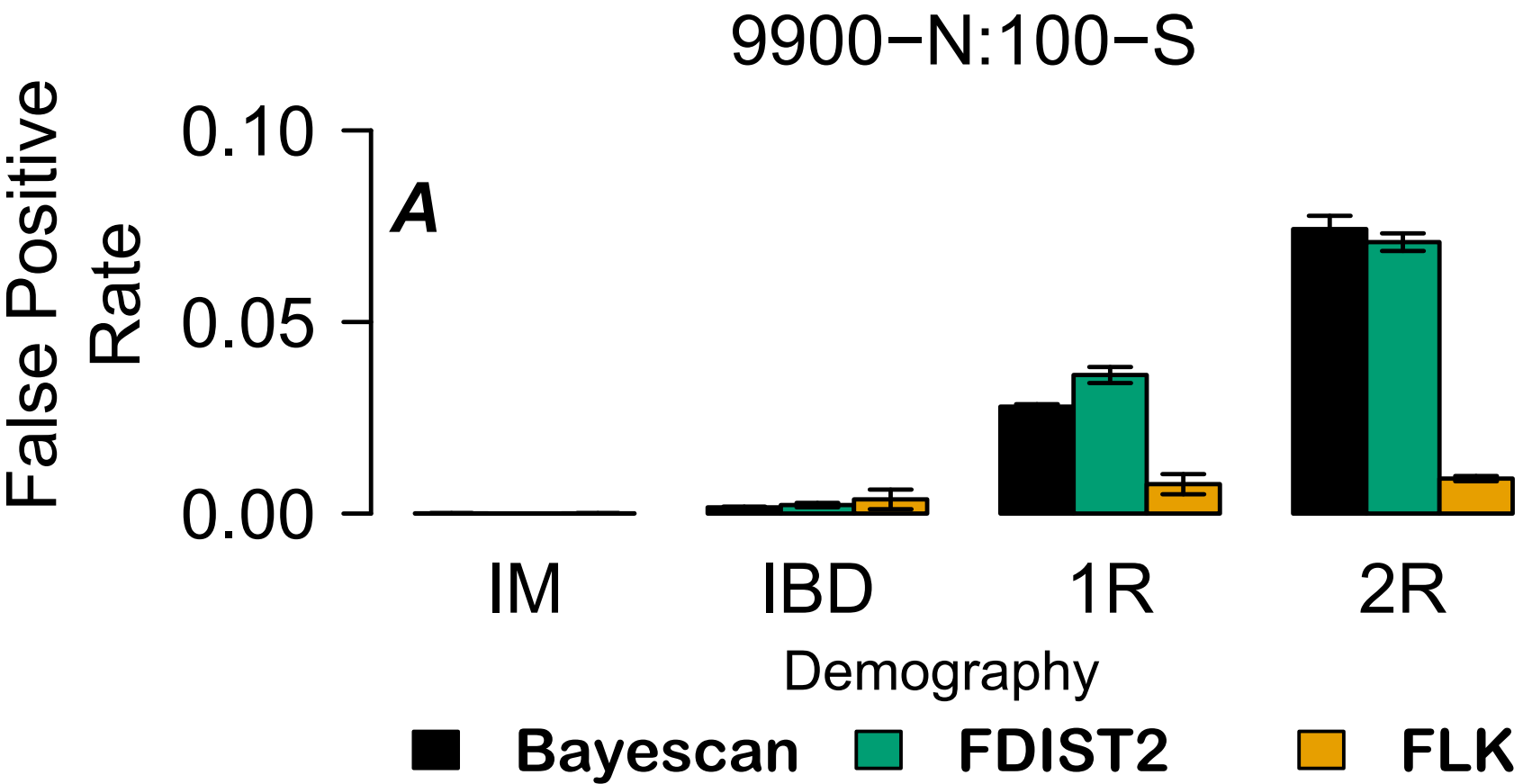


2R
Expansion from two refugia
(non-equilibrium)

False Positive Rates

$(\text{False Positive Neutral}) / (\text{Total Neutral})$

Want: $1/1000$



The bad news:

A false positive rate of 1%
is still too high

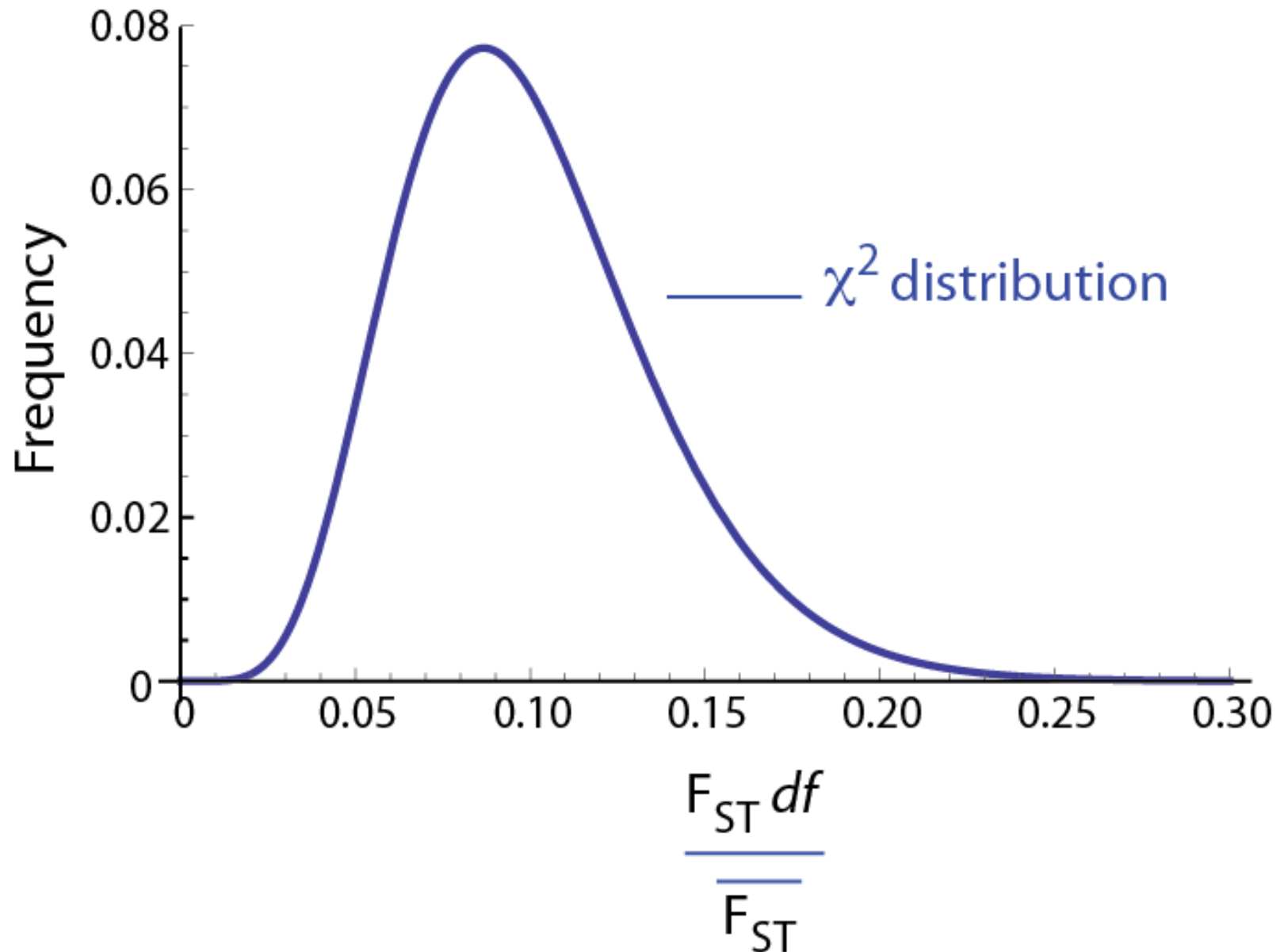
Error rates depend on
neutral parameterization

The good news:

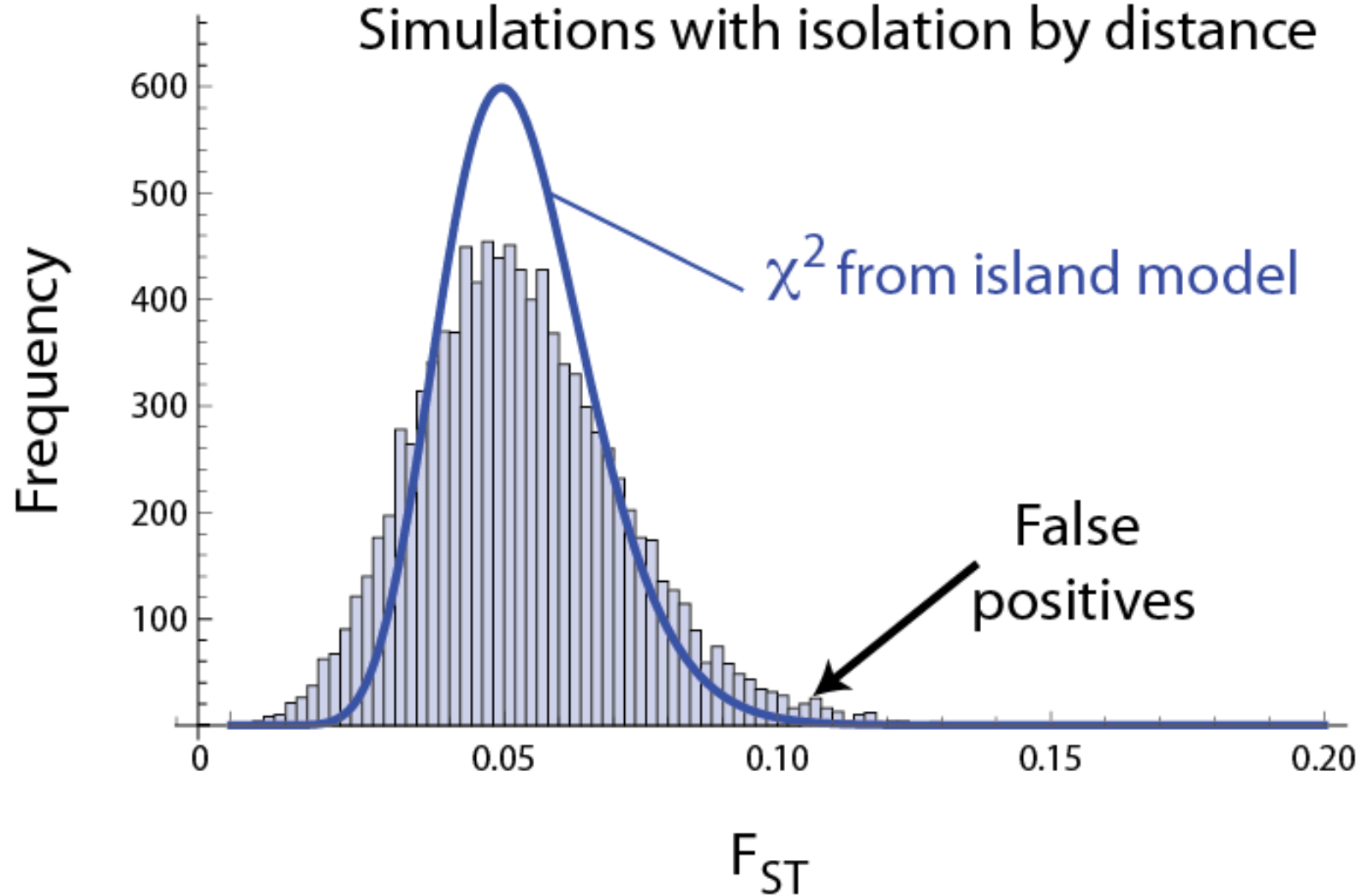
OutFLANK

...doesn't rely on having a
neutral set

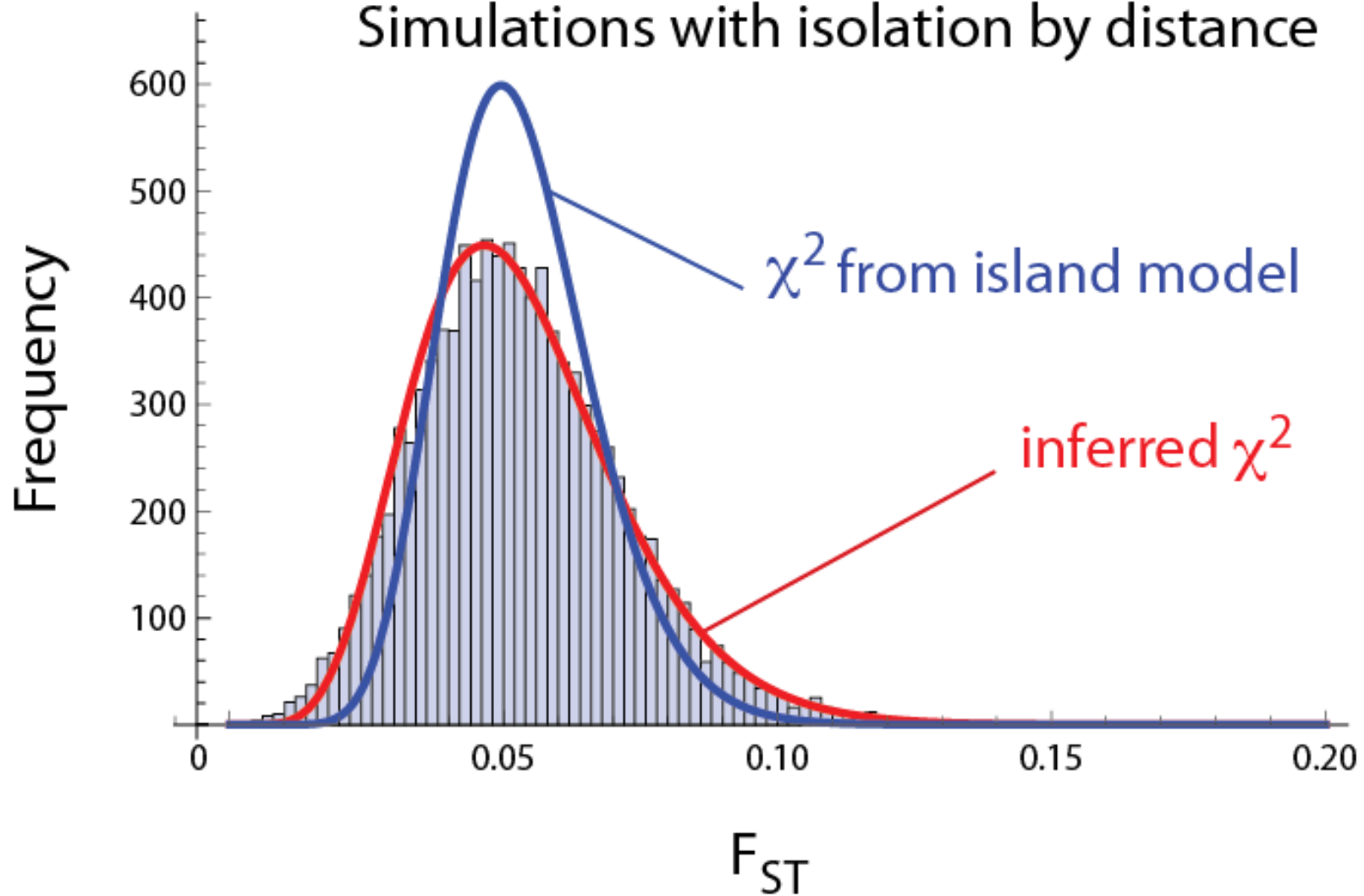
Let's revisit the Lewontin-Krakauer test



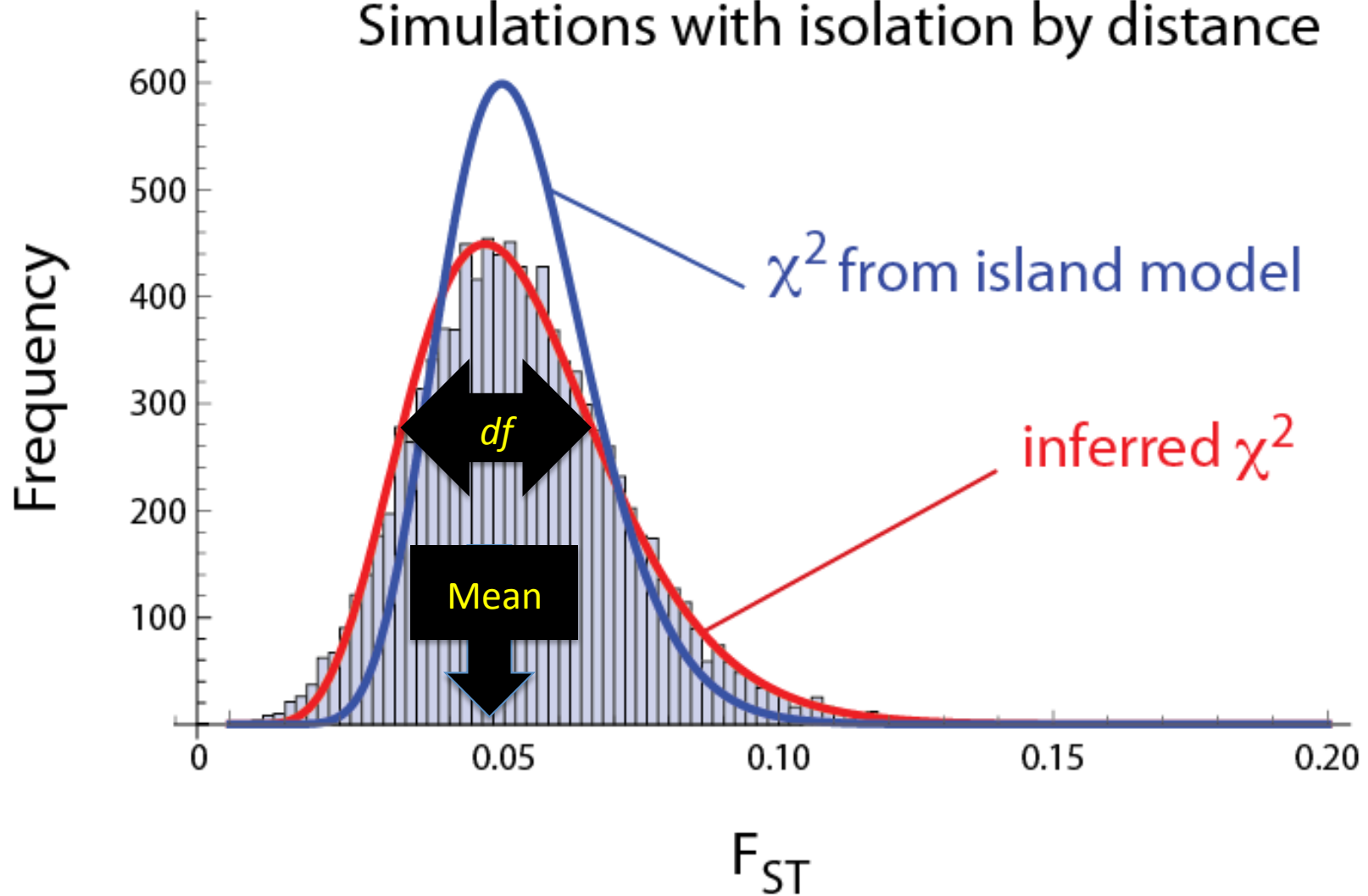
Simulations with isolation by distance



Simulations with isolation by distance

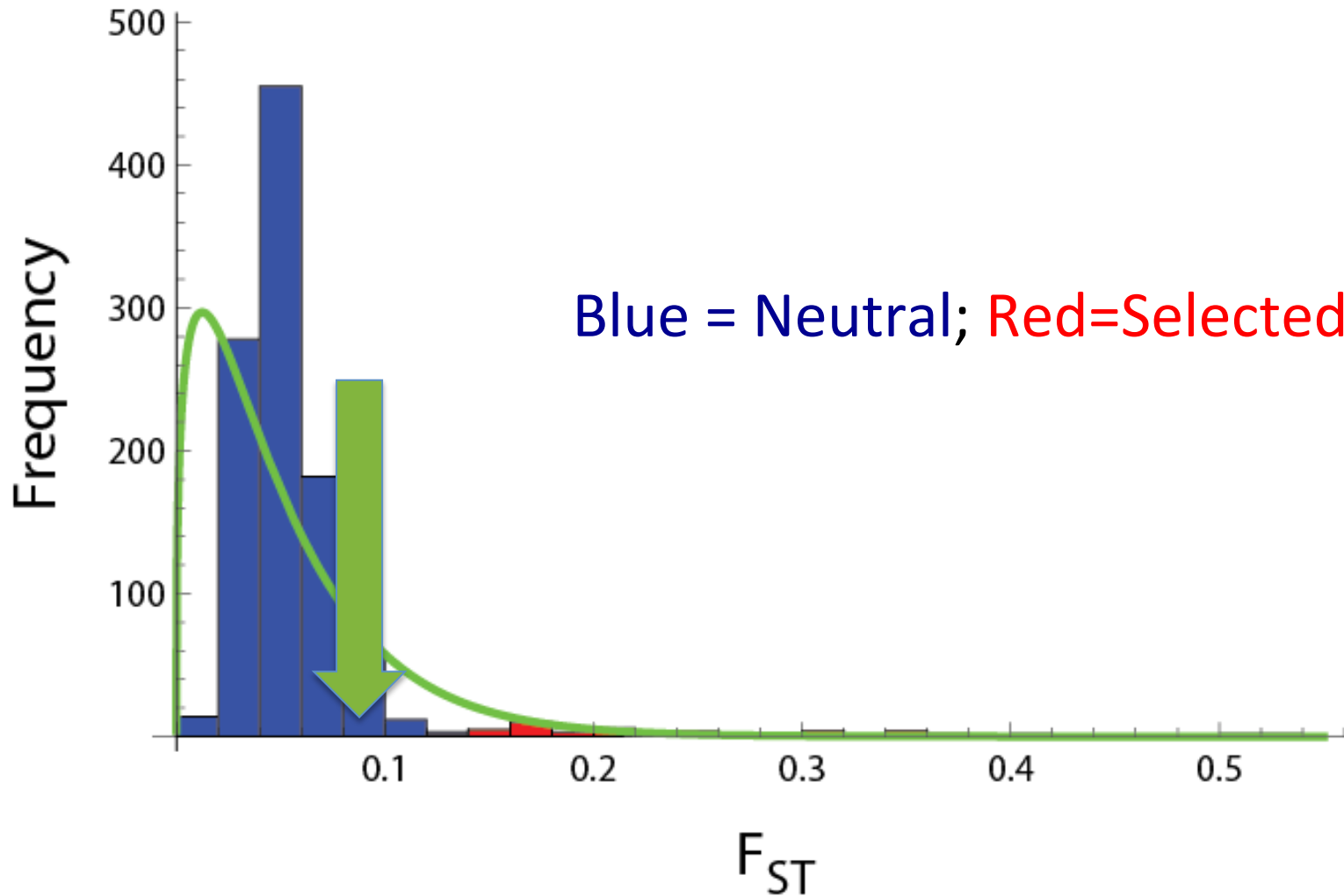


Simulations with isolation by distance

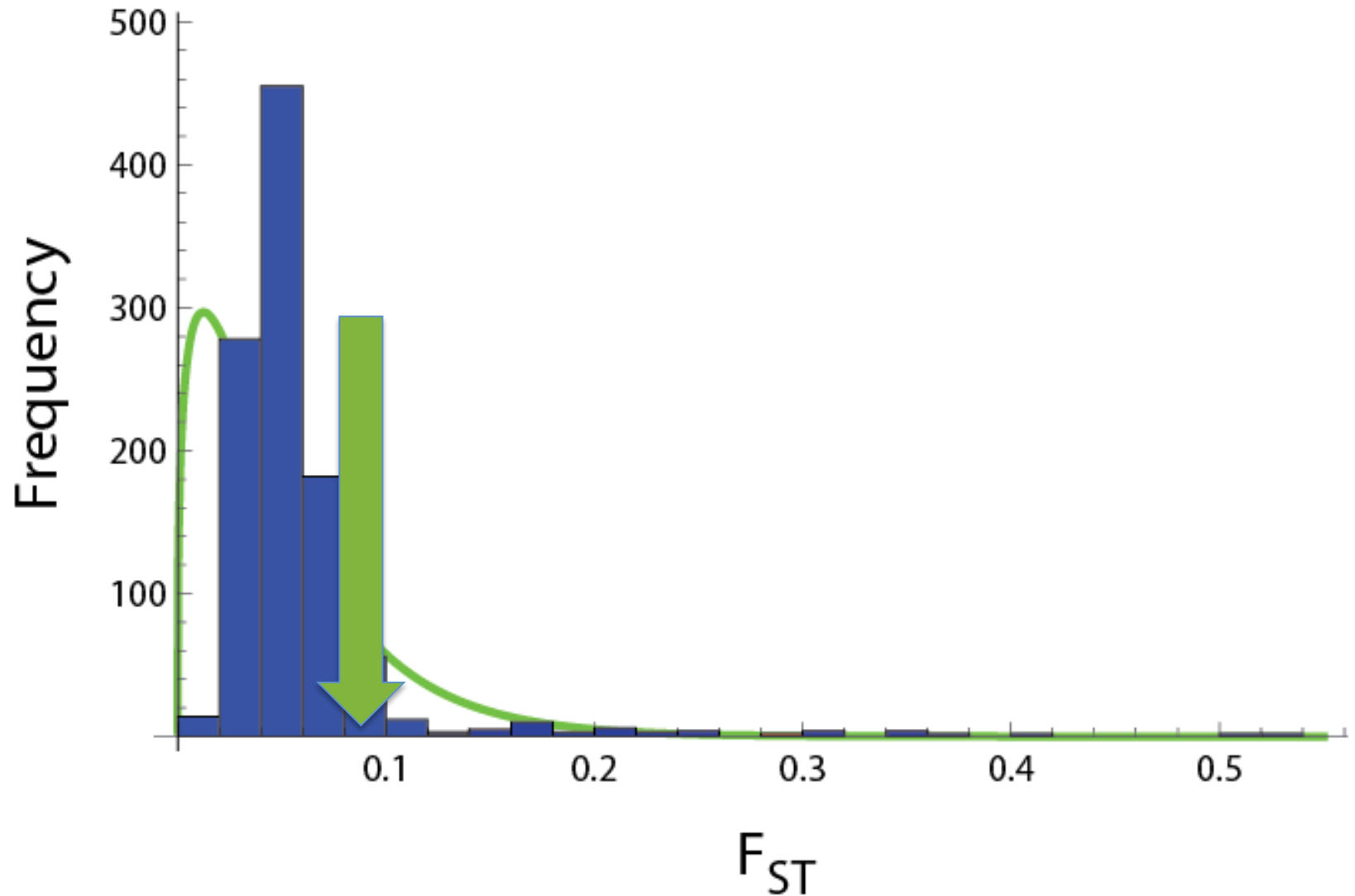


Finding
neutral mean F_{ST}
and
appropriate df
for χ^2

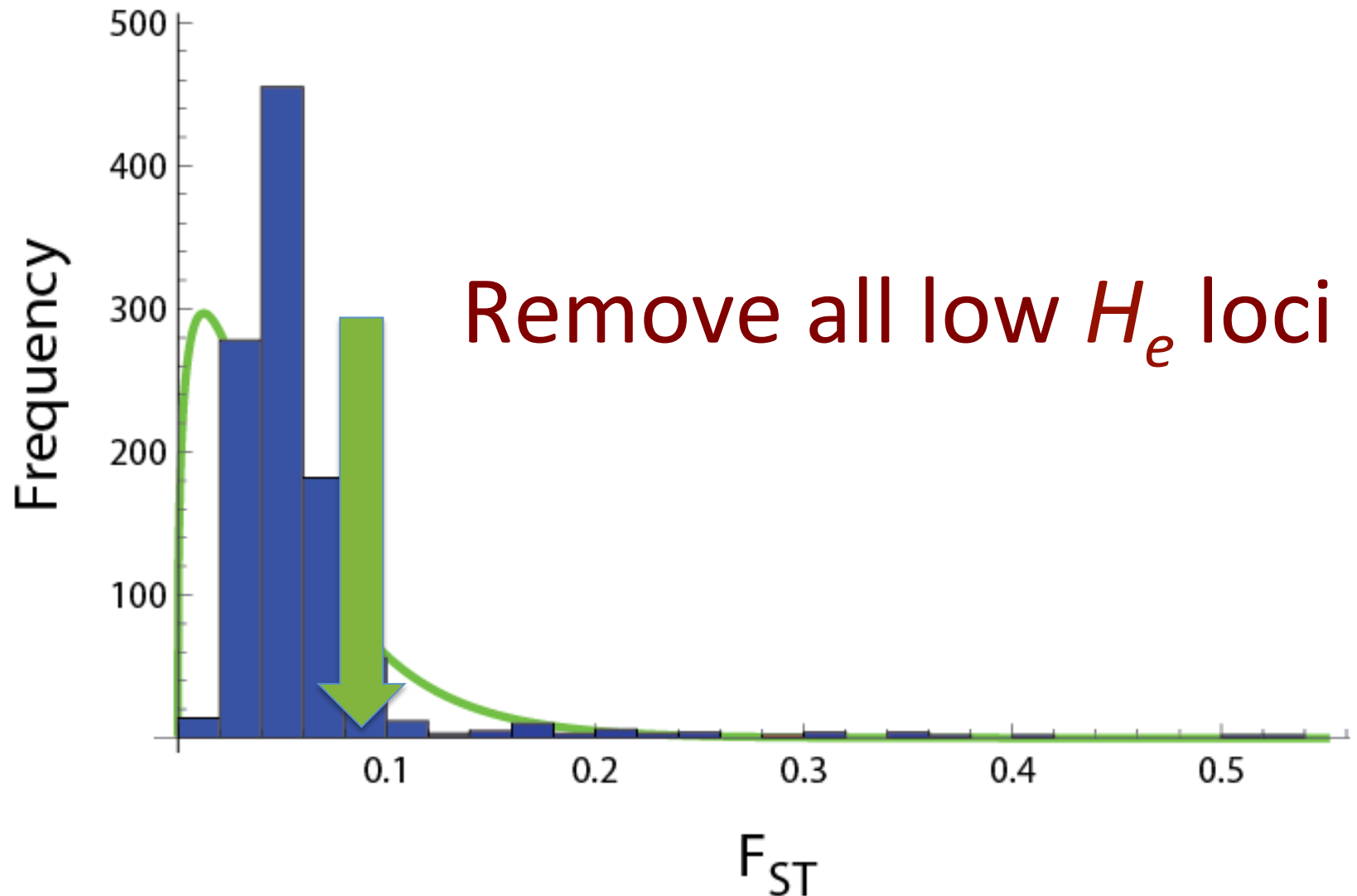
Finding neutral mean F_{ST} and appropriate df for χ^2



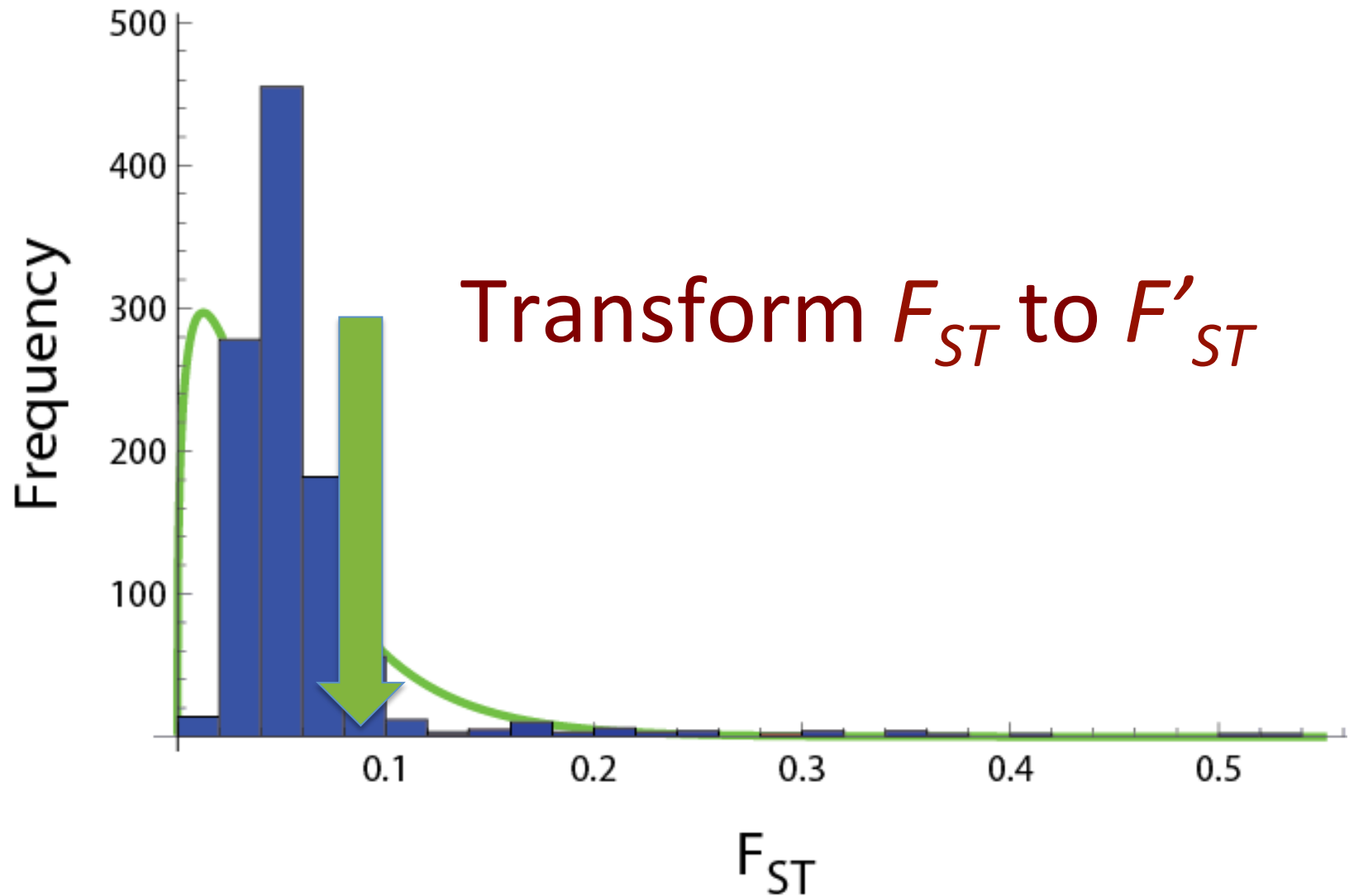
Finding neutral mean F_{ST} and appropriate df for χ^2



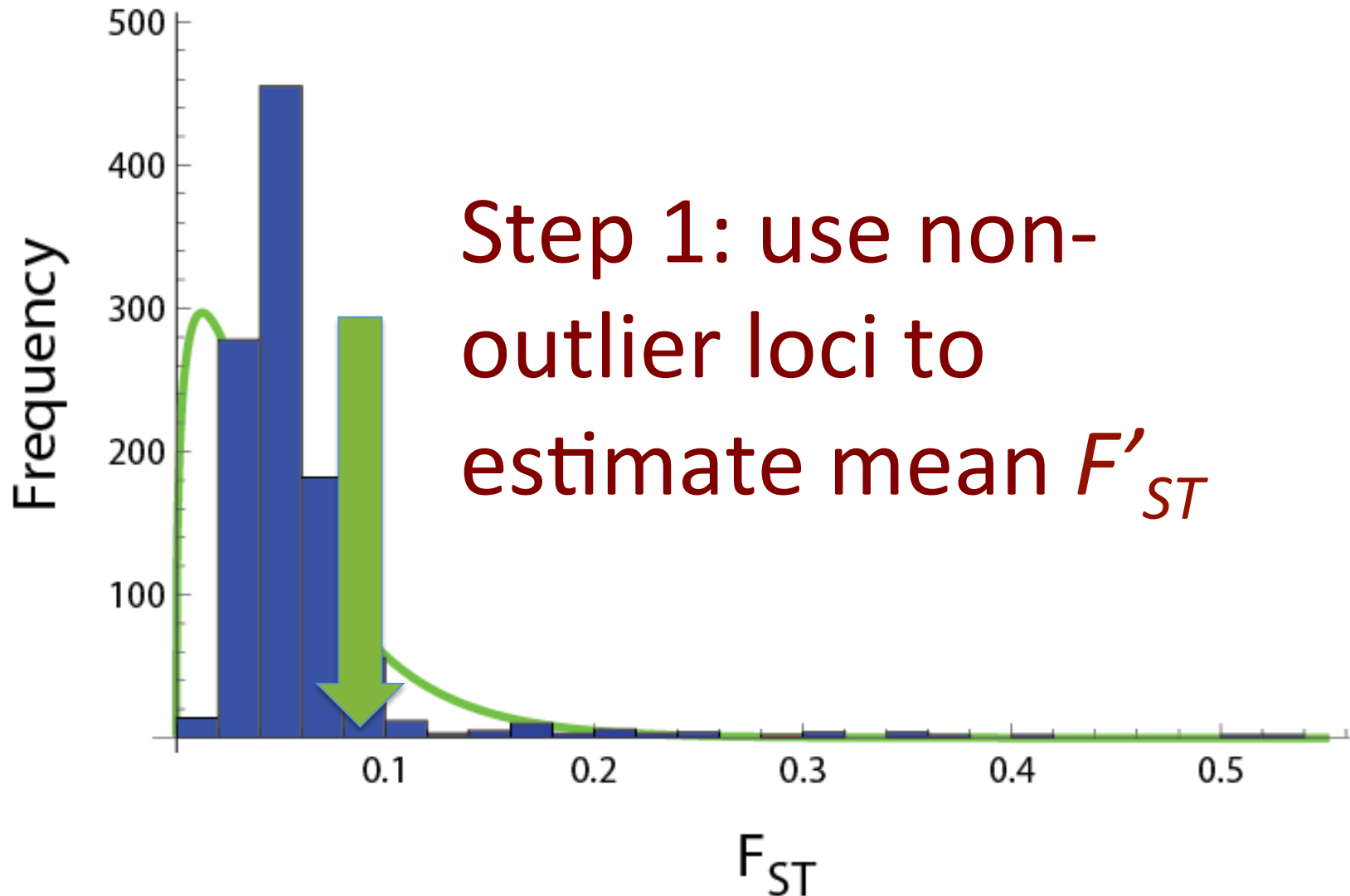
Finding neutral mean F_{ST} and appropriate df for χ^2



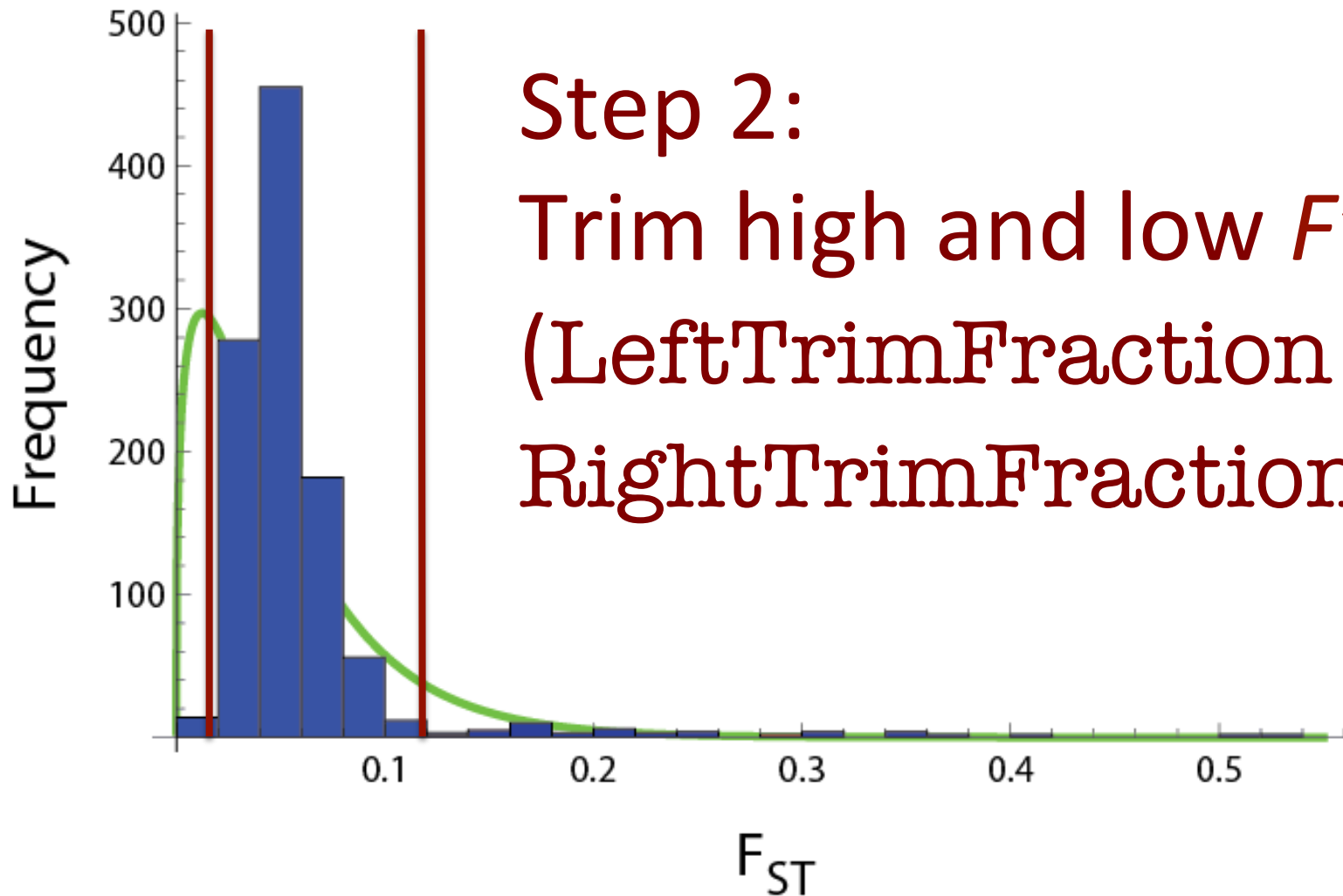
Finding neutral mean F_{ST} and appropriate df for χ^2



Finding neutral mean F_{ST} and appropriate df for χ^2

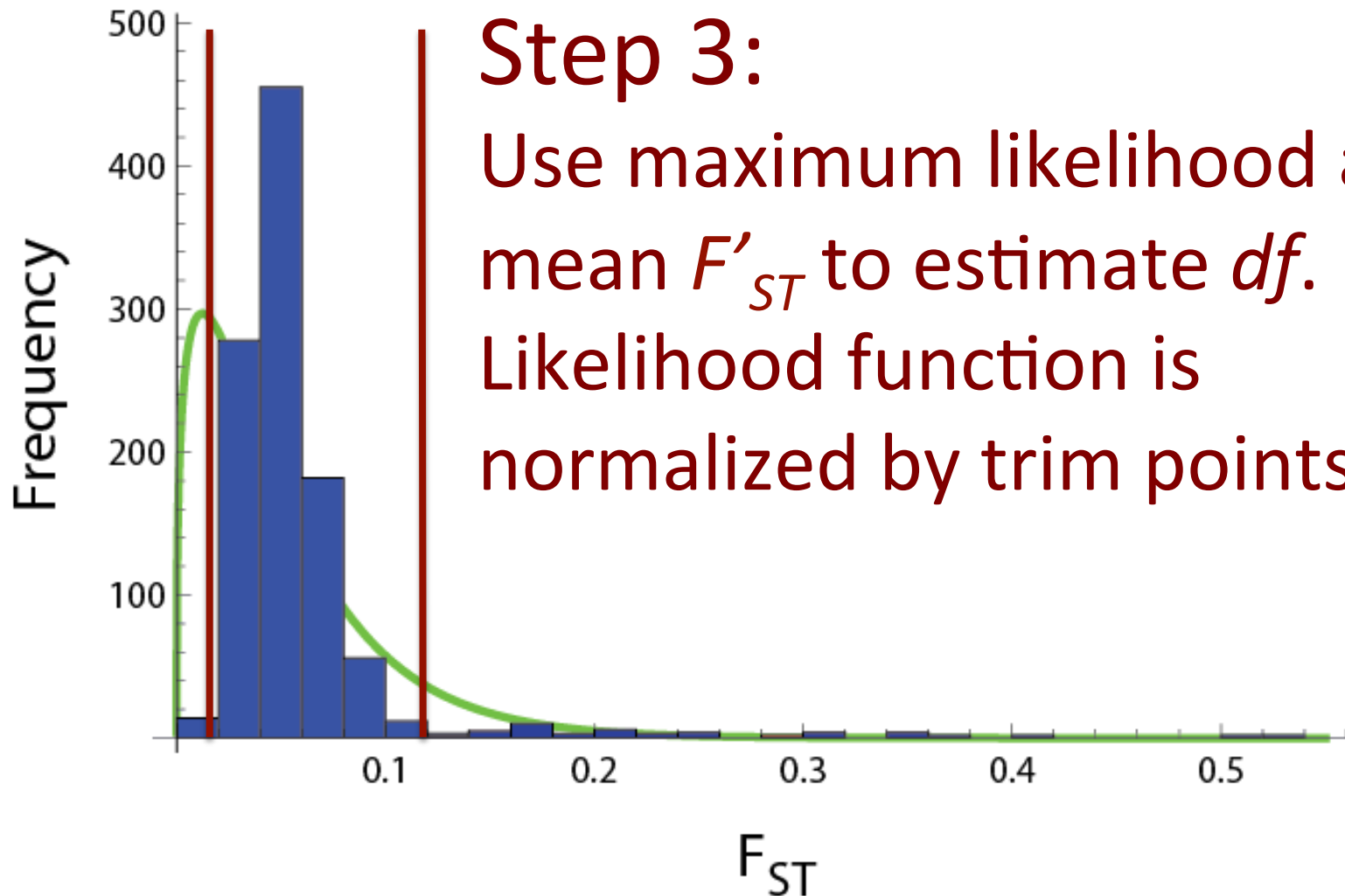


Finding neutral mean F_{ST} and
appropriate df for χ^2

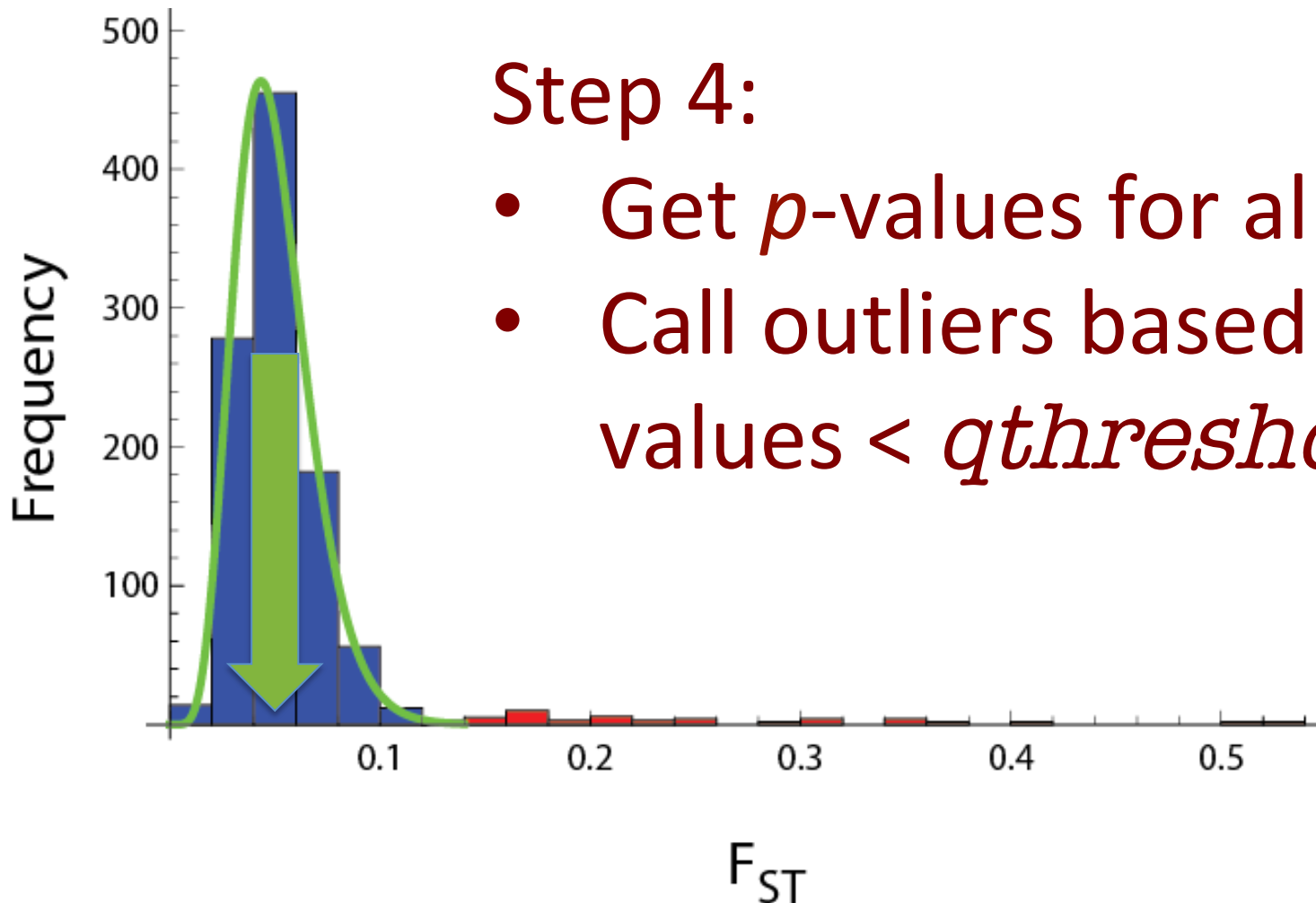


Step 2:
Trim high and low F'_{ST} s
(LeftTrimFraction
RightTrimFraction)

Finding neutral mean F_{ST} and appropriate df for χ^2



Finding neutral mean F_{ST} and appropriate df for χ^2



Step 4:

- Get p -values for all loci
- Call outliers based on q -values $< qthreshold$

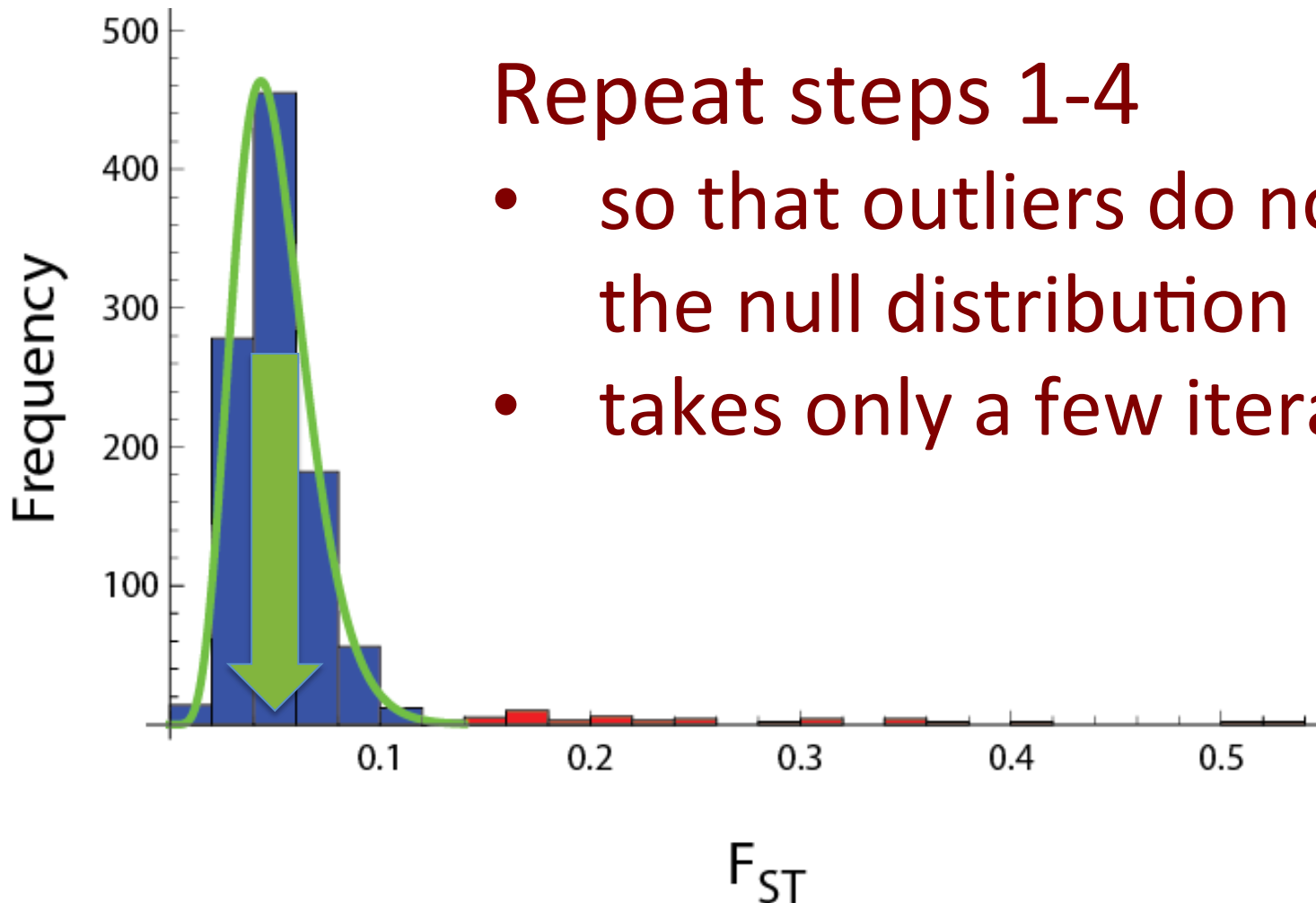
P-value vs. *q*-value

- The *p*-value in OutFLANK is calculated from the fit of the chi-square distribution
- Because there are 1000's of SNPs in the dataset, we need to correct for multiple tests
- OutFLANK controls for **False Discovery Rate (FDR)** with *q-values* using the method of Storey and Tibsharani (2003)

$$FDR = \frac{\textit{Number of false positives}}{\textit{Total number of positive tests}}$$

- *q* < 0.05 means that 5% of positive results are expected to be false positives

Finding neutral mean F_{ST} and appropriate df for χ^2



Repeat steps 1-4

- so that outliers do not affect the null distribution
- takes only a few iterations

How well does OutFLANK work?

Comparison to current methods

False Positive
Rate

0.10
0.05
0.00

9900-N:100-S

< 1 in 1000

IM

IBD

1R

2R

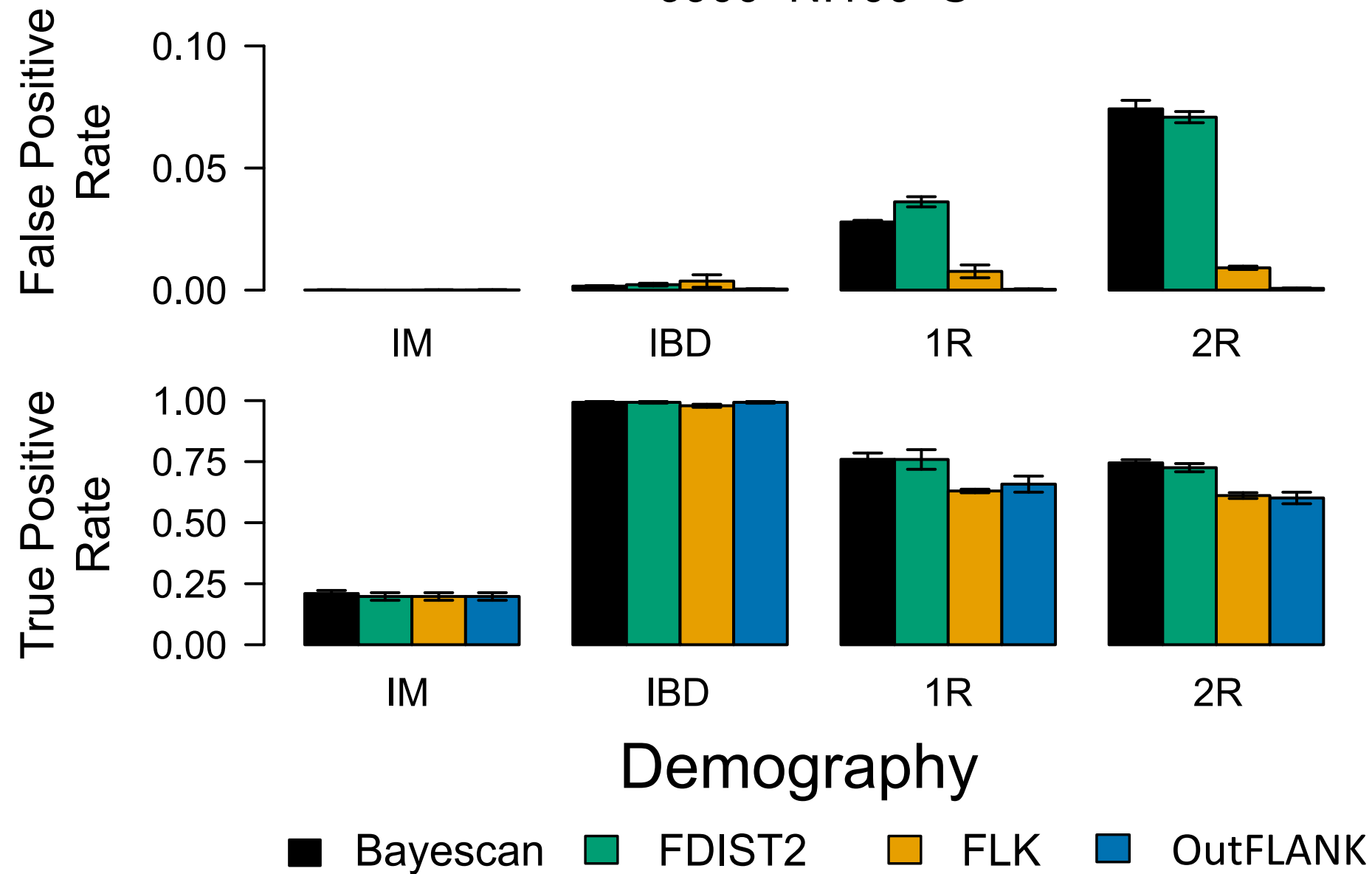
Demography

■ Bayescan ■ FDIST2 ■ FLK ■ OutFLANK

Does the new
method have
power?

$(\text{True positive}) / (\text{Total number under Selection})$

9900-N:100-S



Caution

- Power of OutFLANK increases when more populations are sampled and more individuals per population are sampled
- OutFLANK also needs a large number of loci (>1000 SNPs, see discussion in manuscript)

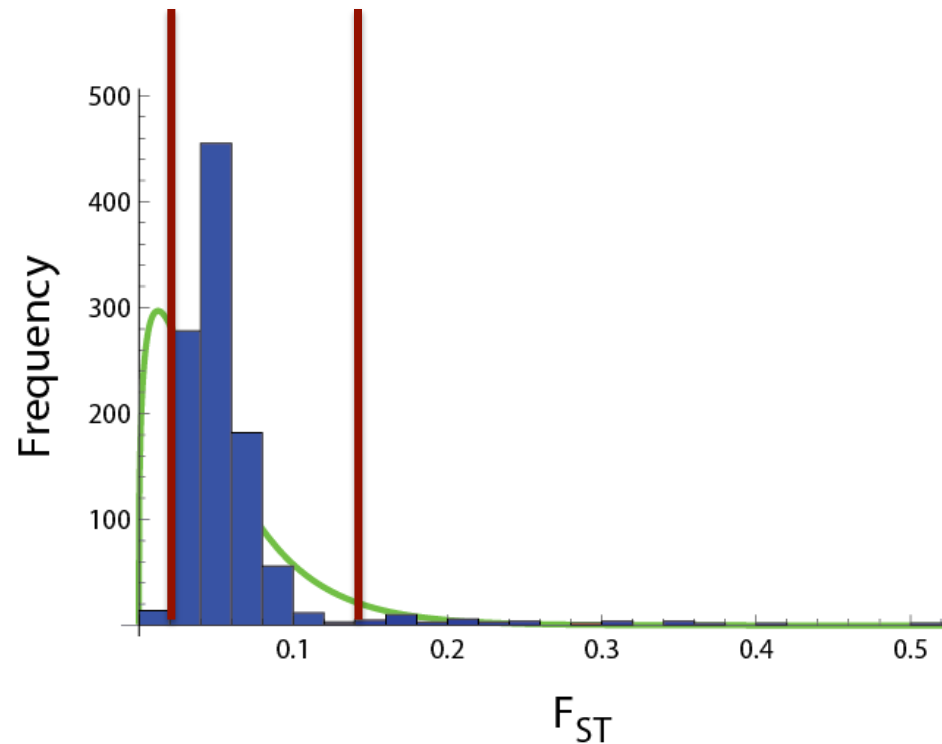
No. individuals per population	5 populations	10 populations	20 populations	40 populations
5	0	.09	.52	.84
10	.10	.56	.82	.94
20	.37	.75	.90	.95
40	.55	.81	.94	.97

Steps to running OutFLANK

1. Prepare a dataframe for the OutFLANK() function
 - From a formatted SNP dataset using MakeDiploidFSTMat()
 - (From your own data using functions provided with the package (see Section “For Advanced Users”))
2. Check for SNPs of low sample size or that have F_{ST} values differentially affected by sample size correction
3. Run OutFLANK
4. Plot results

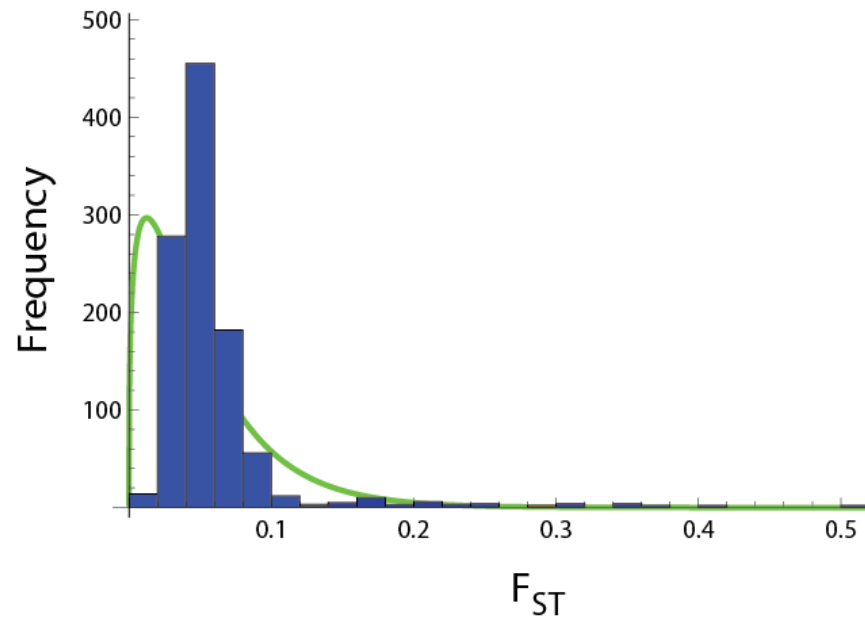
The OutFLANK() function

- **FstDataFrame**
- **LeftTrimFraction**=0.05
- **RightTrimFraction**=0.05
- **Hmin**=0.1 (loci with low H_e do not follow chi-square assumption)
- **NumberOfSamples** (Number of populations)
- **qthreshold**=0.05 (desired false discovery rate)

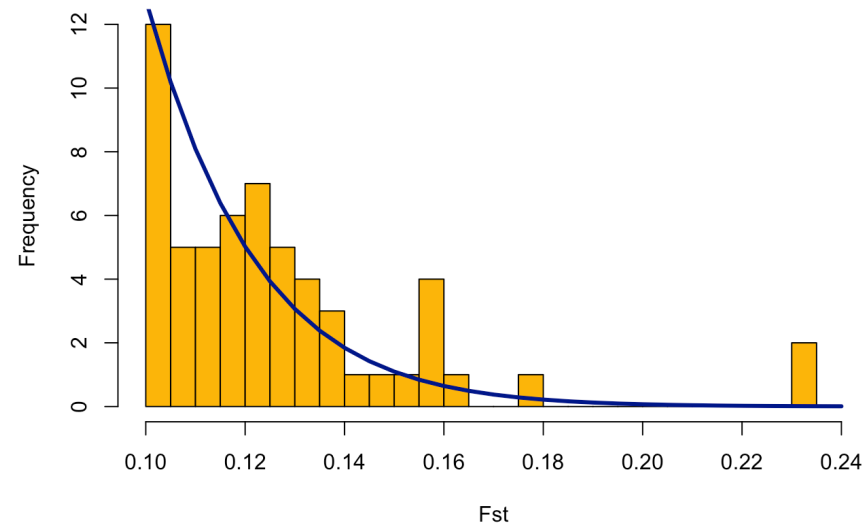


Hints

- **qthreshold** is used to call outliers. The trim points must be located within any potential outliers. If **qthreshold** is set too high, OutFLANK will return an error.
- The algorithm may fail if **RightTrimFraction** is set too high
- Use the plotting function to check the fit of the inferred distribution to the right tail
- OutFLANK will not fit the left tail of the F_{ST} distribution well



Fst without sample size correction



Other questions

- Can I use OutFLANK on pool-seq data? **Not Recommended**
- Does pool-seq data fit a chi-square distribution (the big assumption of the method)?
 - We understand how sampling variance from finite populations and finite individuals per population affect the variance of the Fst distribution, and we know it still fits a chi square for the cases we've simulated.
 - However, we are uncertain how the random sampling of chromosomes for pool-seq data would affect the Fst distribution.
 - Other issues arise because we are not using a sample size correction in outflank (so we assume loci have roughly equal sample sizes), and if individuals contribute unequally to pools then this could be a huge violation of this assumption.
- So, in short, if you wanted to do it I would suggest the following steps:
 - 1) Use the haploid Fst estimator that comes with the distribution instead of the default in preparing the matrix for the outflank function
 - 2) Show that your Fst distribution fits the chi-square output by outflank
 - 3) Do some simulations that show that outflank performs well for your study system and sampling design

Questions

- What if some populations have fewer individuals than others? **THAT'S OK – only a problem if some *loci* have lower sample sizes**

Notes based on student questions

- Confusion about `q_threshold` and the `qvalue` used for decision making
- a need to illustrate how increasing `q_threshold` can make small changes in the tail and the resulting `q-values` of the outliers
- a need to plot `p-value` histograms and `q-q` plots