

Short INDELS: genetic markers
for adaptive divergence

Original aspects of the short INDELs project

- Divergent natural selection vs neutral processes
- Species with high diversity
- Systems with imperfect genomes can still contain useful functional information

INDEL-SNP comparisons

1. Outlier sharing
2. Clustering of (different types) markers
3. Derived allele frequencies (in progress and for now simply minor)
4. Distributions of cline parameters

1. Outlier sharing

Total number of SNP: 11225

Total number of INDEL: 1752

Proportion of SNP with significant clines.

CZA left: 0.5317595
CZA right: 0.4457016
CZB left: 0.3277506
CZB right: 0.4244989
CZD left: 0.4473942
CZD right: 0.4823163

Proportion of INDEL with significant clines.

0.5296804
0.4549087
0.3413242
0.4092466
0.4737443
0.4834475

Proportions of SNP outliers that are shared.

CZA left and right: 0.6160714
CZB left and right: 0.5178571
CZD left and right: 0.6339286
CZA and CZB: 0.359375
CZA and CZD: 0.4107143
CZB and CZD: 0.484375

Proportions of INDEL outliers that are shared.

0.7058824
0.4705882
0.6470588
0.3529412
0.4117647
0.4411765

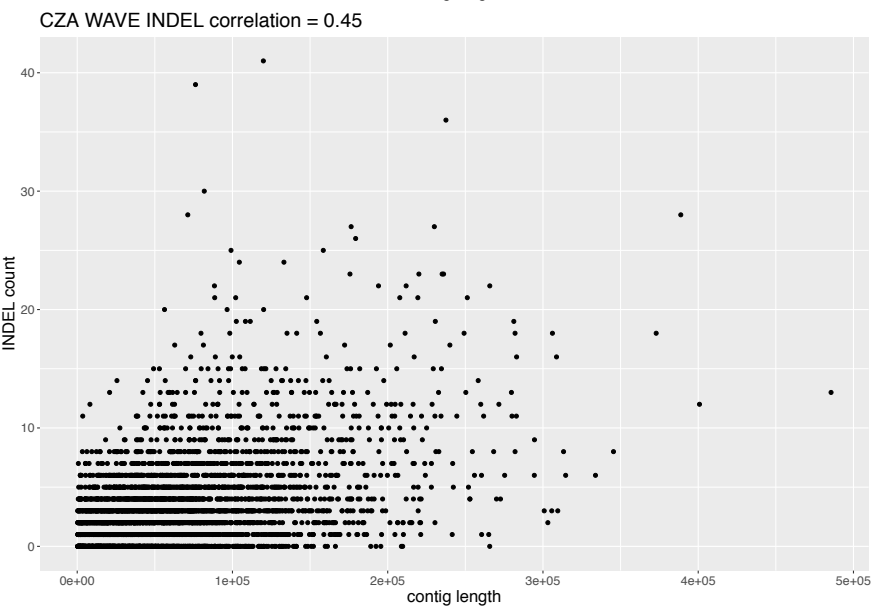
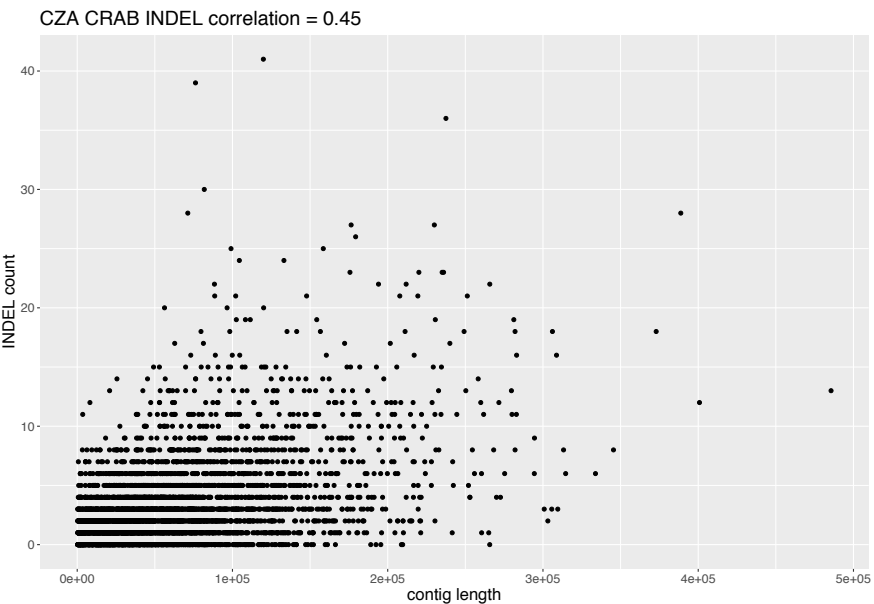
Number of SNP outliers found in 1 hybrid zone(s): 142
Number of SNP outliers found in 2 hybrid zone(s): 66
Number of SNP outliers found in 3 hybrid zone(s): 29
Number of SNP outliers found in 4 hybrid zone(s): 27
Number of SNP outliers found in 5 hybrid zone(s): 25
Number of SNP outliers found in 6 hybrid zone(s): 13

24
7
7
5
1
3

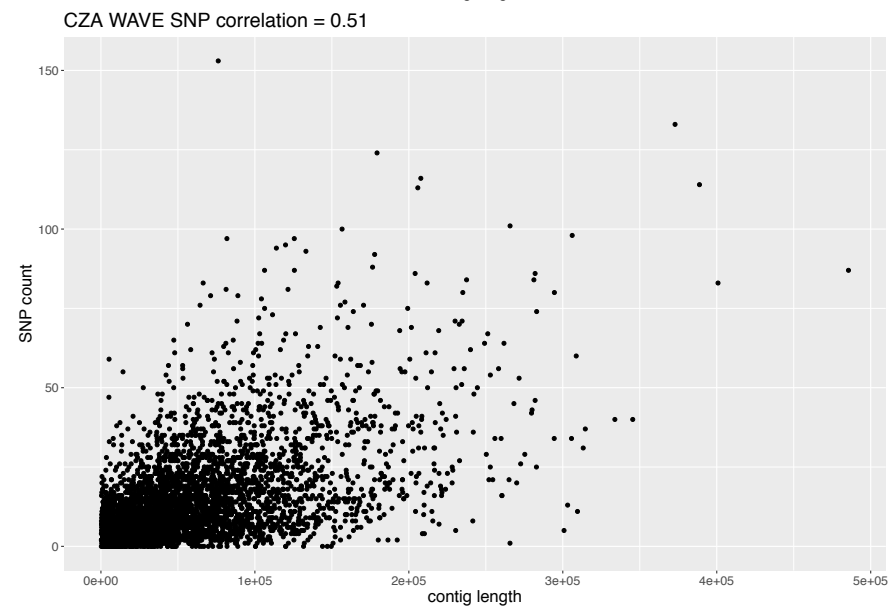
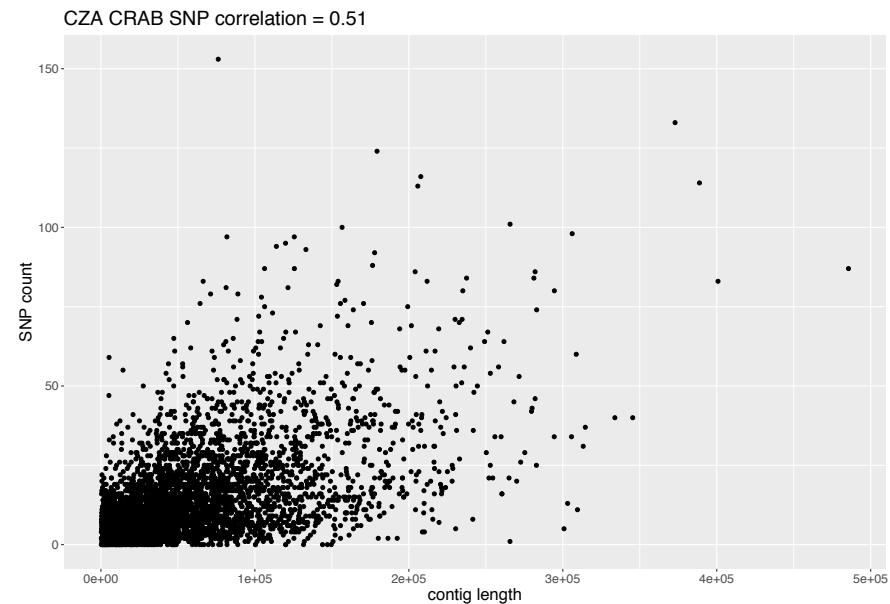
Prop. of SNP outliers in inversions found in 1 zone(s): 0.556
Prop. of SNP outliers in inversions found in 2 zone(s): 0.636
Prop. of SNP outliers in inversions found in 3 zone(s): 0.862
Prop. of SNP outliers in inversions found in 4 zone(s): 0.889
Prop. of SNP outliers in inversions found in 5 zone(s): 0.92
Prop. of SNP outliers in inversions found in 6 zone(s): 1

0.625
0.57
0.86
1
1
1

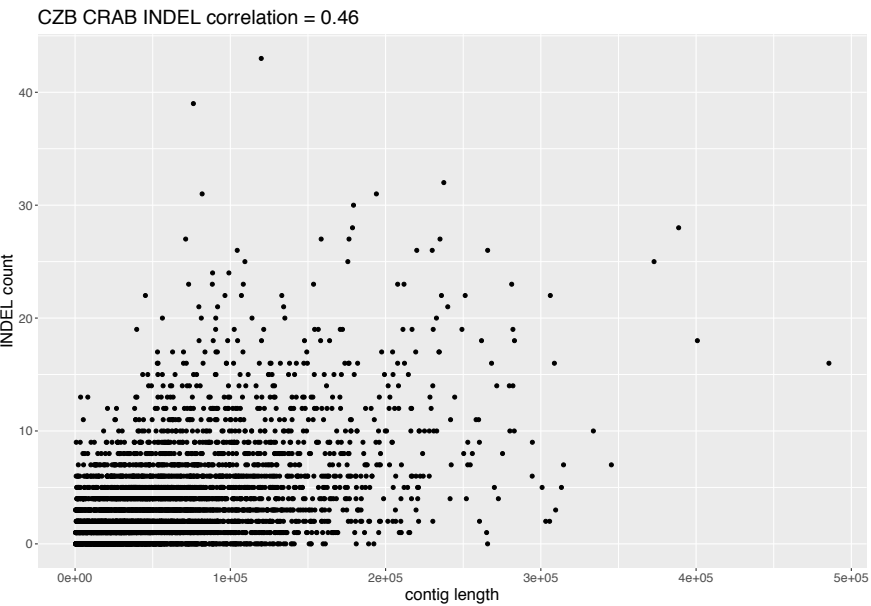
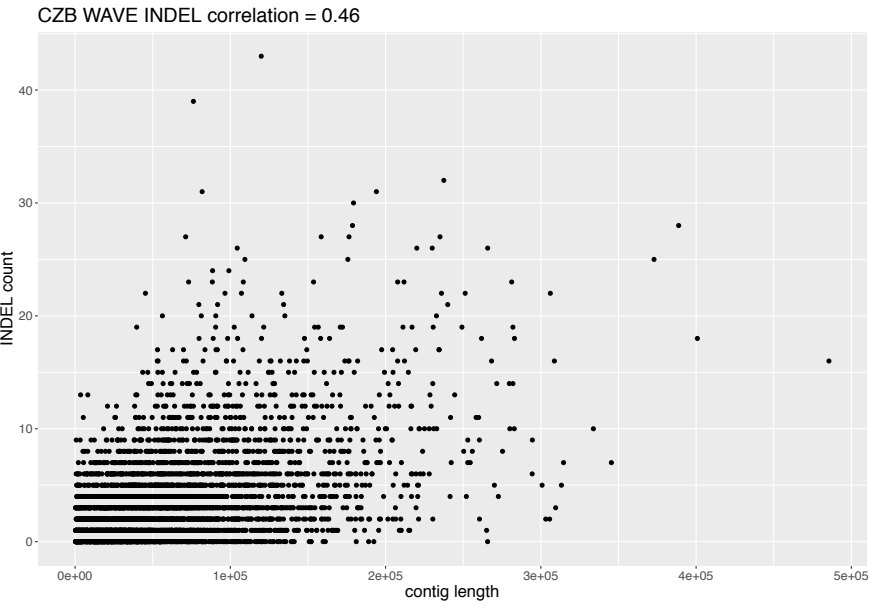
2. Clustering of (different types) markers



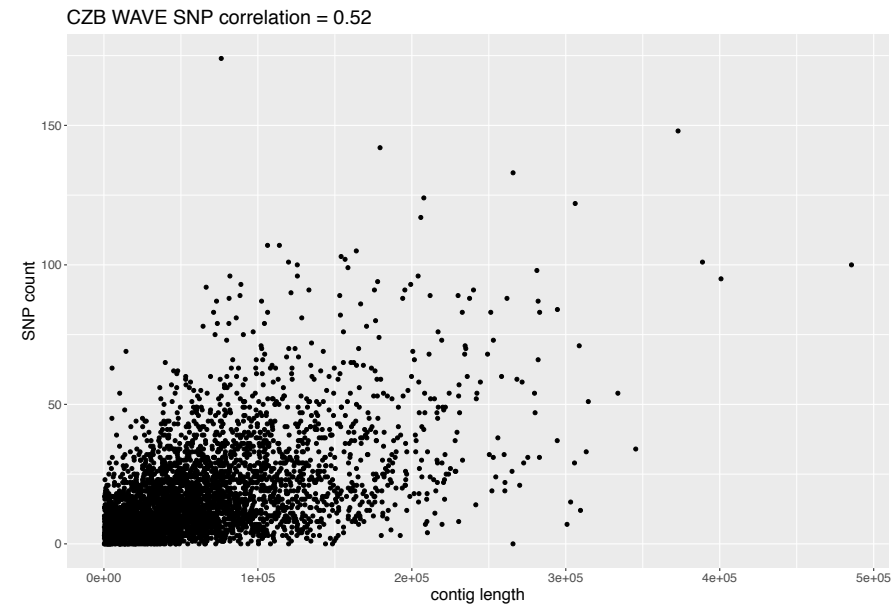
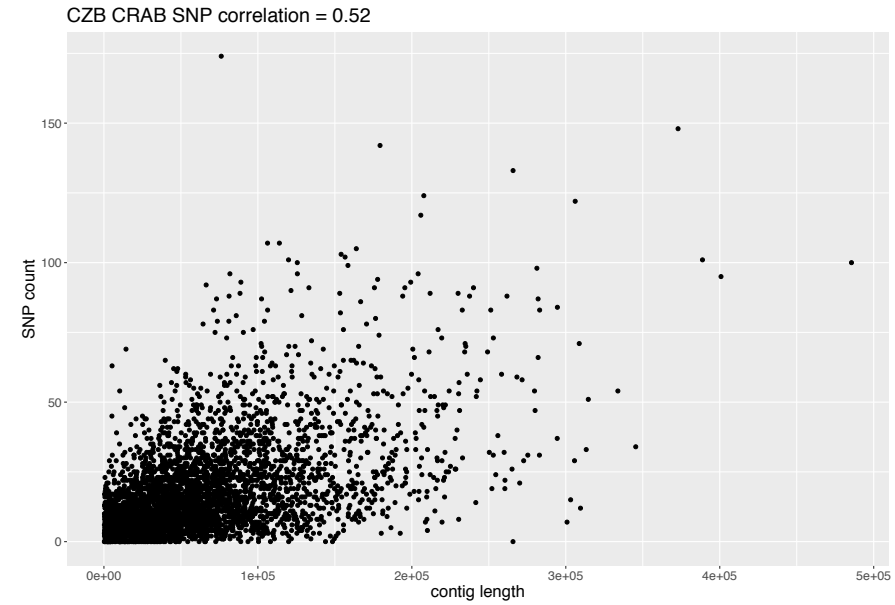
CZA: INDELs and SNPs after filtering but before cline analysis.



2. Clustering of (different types) markers

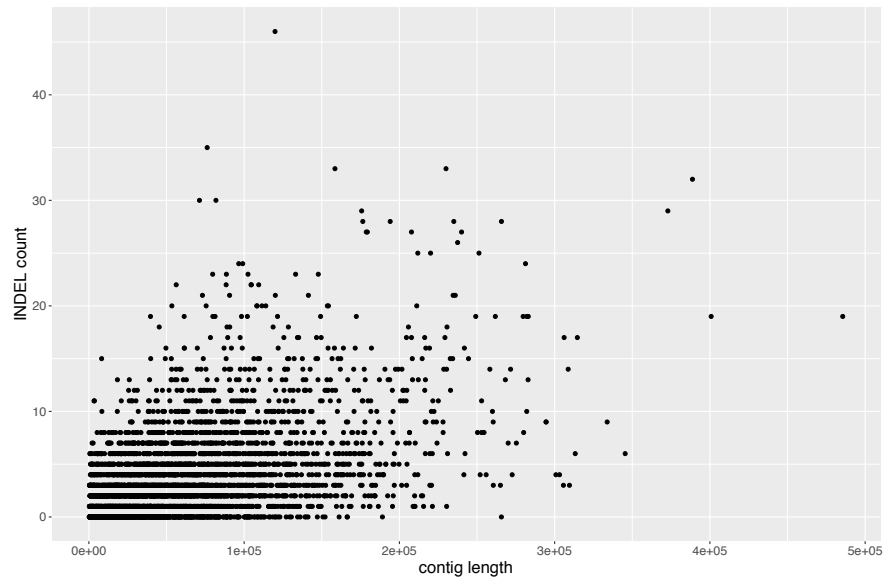


CZB: INDELs and SNPs after filtering but before cline analysis.

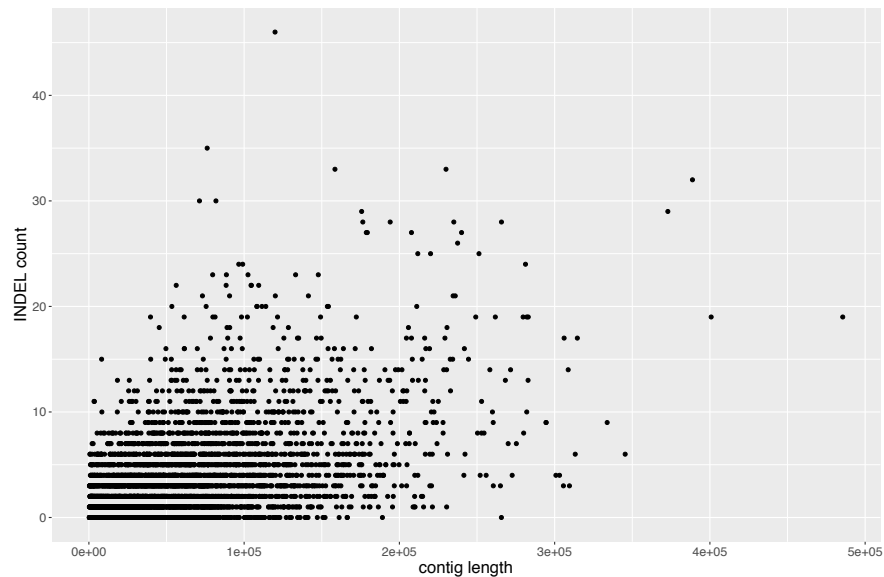


2. Clustering of (different types) markers

CZD CRAB INDEL correlation = 0.47

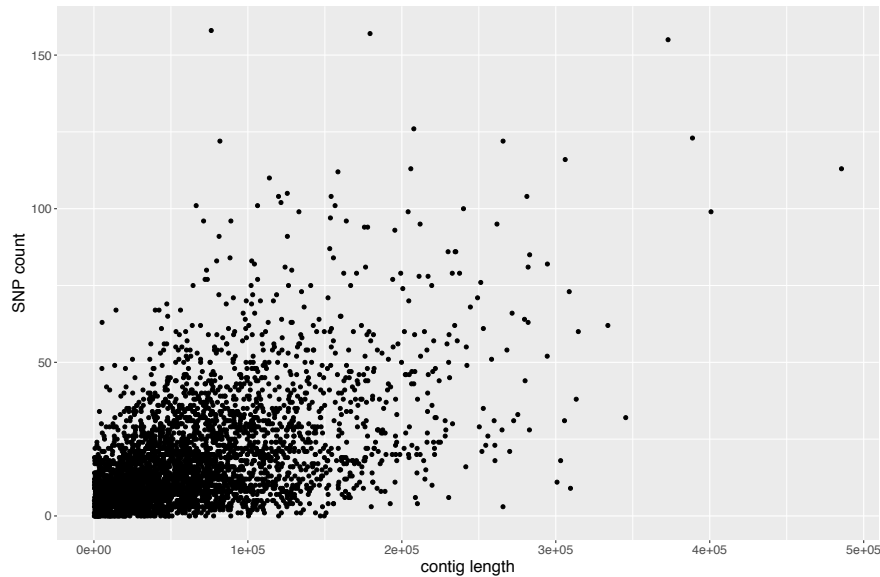


CZD WAVE INDEL correlation = 0.47

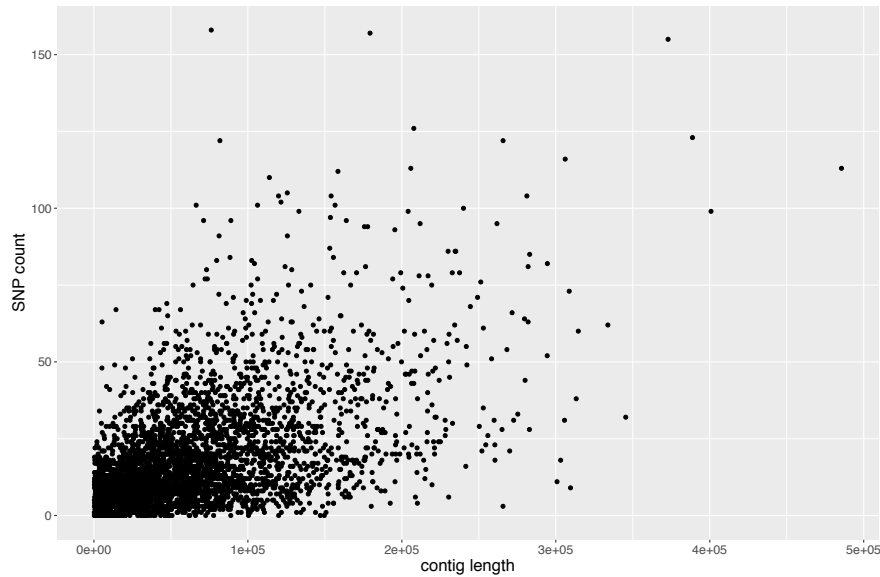


CZD: INDELs and SNPs after filtering but before cline analysis.

CZD CRAB SNP correlation = 0.52



CZD WAVE SNP correlation = 0.52



2. Clustering of (different types) markers

- INDELs and SNPs after filtering and cline analysis.
- All six hybrid zones combined.

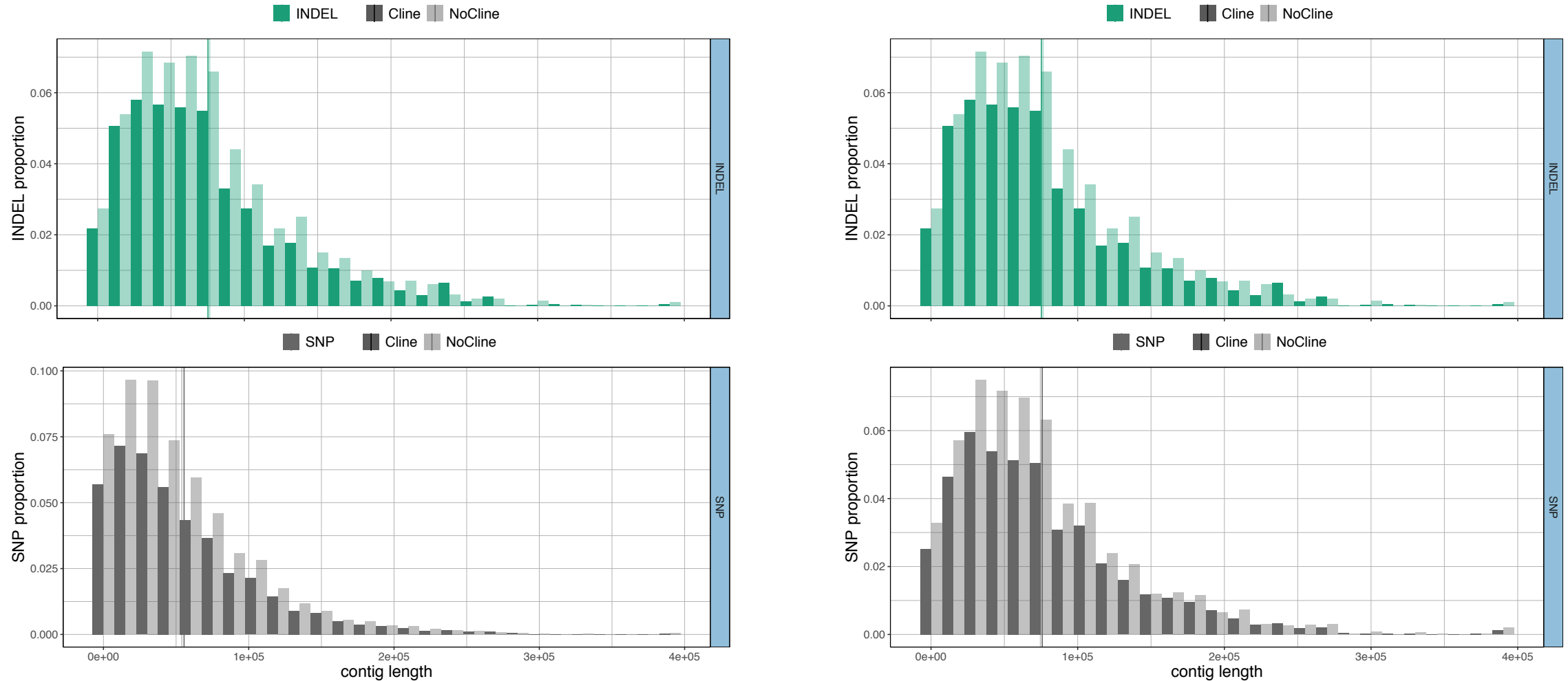


Figure 3a. Marker proportions over contig length. Proportion = $\text{count} / \sum \text{count per marker type}$ and bin width = 15000 base pairs. Clinal variants are dark coloured and non-clinal variants are light coloured. Left: SNP call using SAMtools and INDEL call using GATK. Right: both INDELs and SNPs were called with GATK.

3. Derived allele frequencies - INDELs

- Ancestral state was inferred from called genotypes:
 - Reference allele = ancestral allele = ref_anc
compressa is homo for the reference allele (0)
 - Alternative allele = ancestral allele = alt_anc
compressa is homo for the alternative allele (2)
 - Unknown ancestry = het
compressa is het (1)

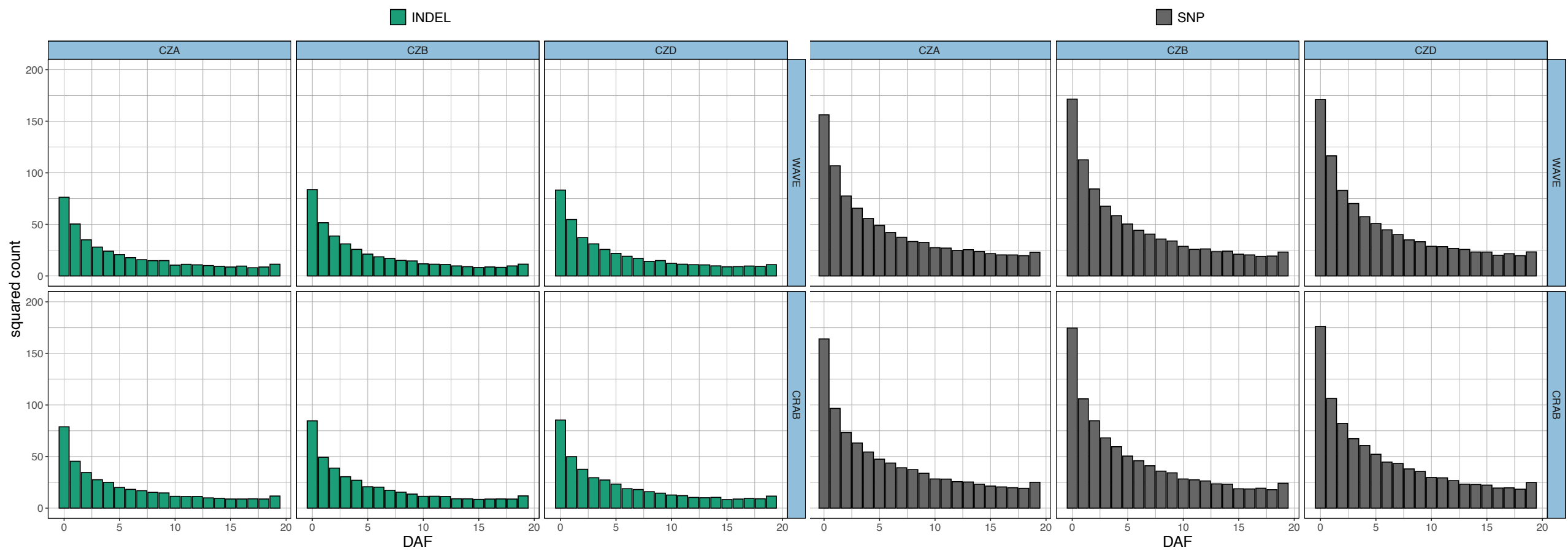
For the rest of the results, kept only variants highlighted in green.

*Table 1. Count of INDELs and SNPs for each combination of possible allelic states given one outgroup (*L. compressa*) with two samples (NE and W). There are two combinations in which the allelic state is concordant in both samples (in green), eight in which the allelic state can only be retrieved from one sample (in yellow) and finally, five in which the allelic state cannot be inferred (in red).*

NE_Lcomp	W_Lcomp	INDEL	SNP
alt_anc	alt_anc	5305	27188
alt_anc	het	528	3543
alt_anc	NA	245	1097
alt_anc	ref_anc	511	2195
het	alt_anc	2231	12439
het	het	1577	9691
het	NA	151	627
het	ref_anc	3120	17198
NA	alt_anc	158	765
NA	het	33	267
NA	ref_anc	449	1831
ref_anc	alt_anc	693	3292
ref_anc	het	1422	7462
ref_anc	NA	1003	3675
ref_anc	ref_anc	38884	178715

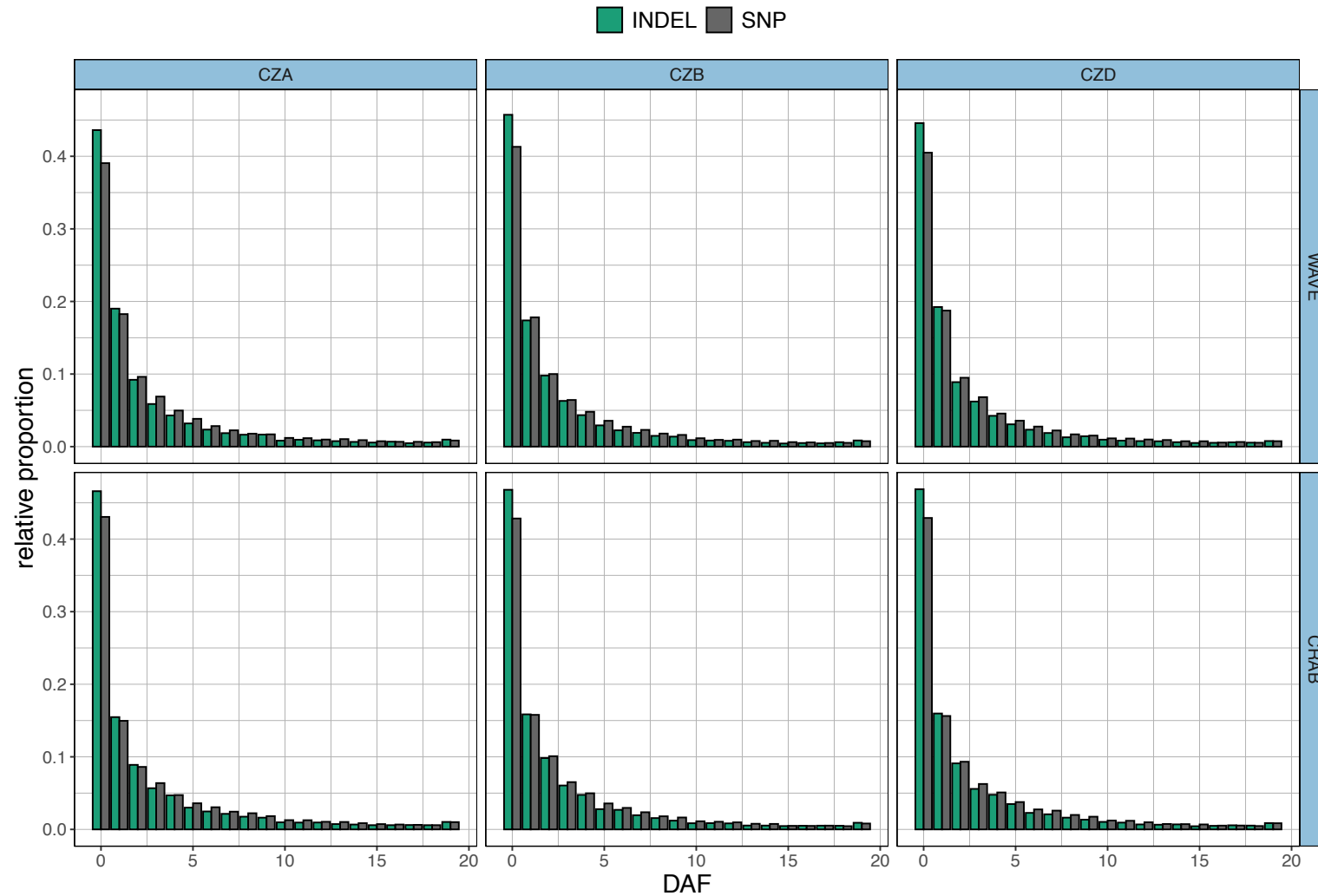
3. Derived allele frequencies (GATK call)

- Square root of count of INDELs and SNPs after filtering but before cline analysis.



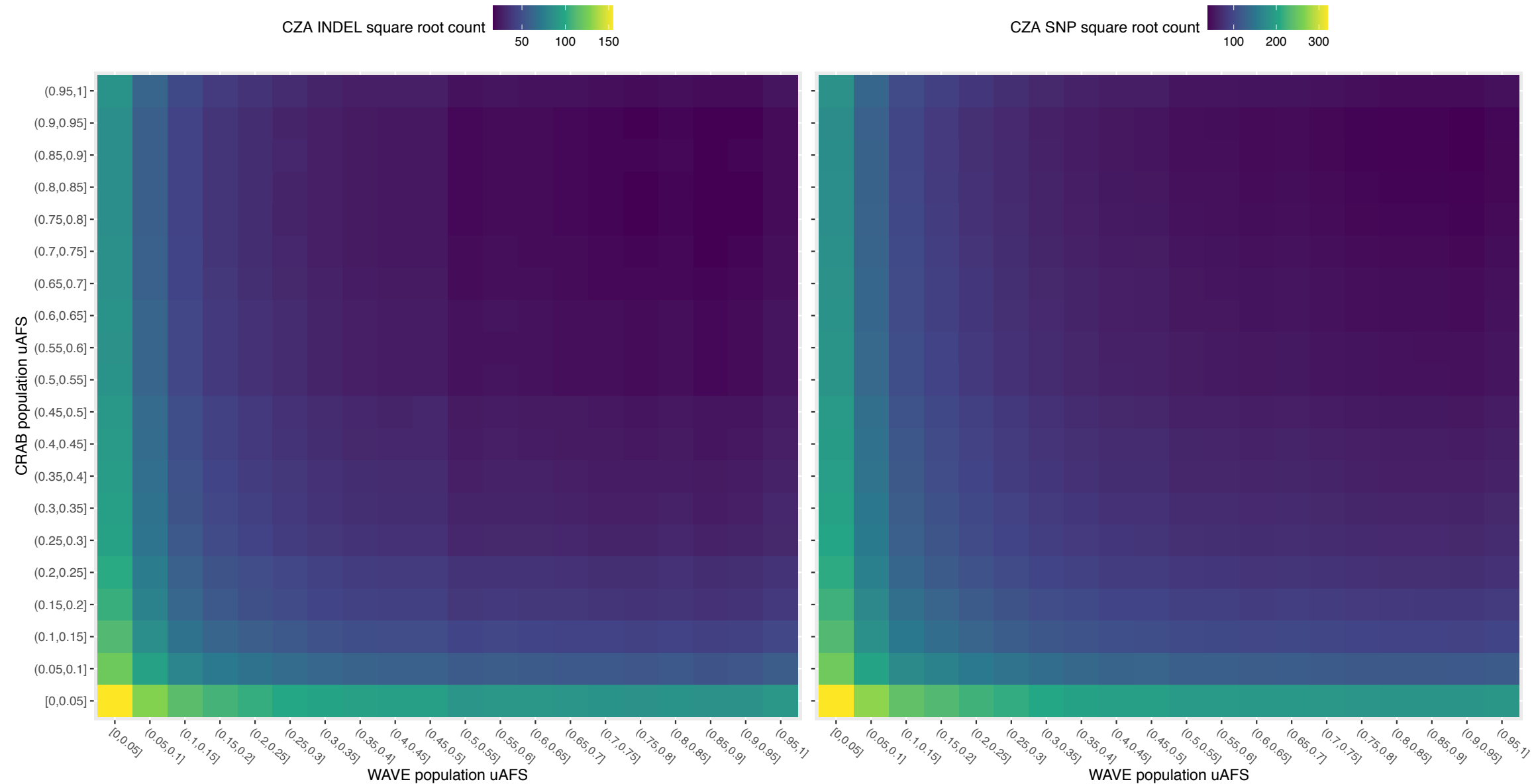
3. Derived allele frequencies (GATK call)

- Relative proportion of INDELs and SNPs after filtering but before cline analysis.



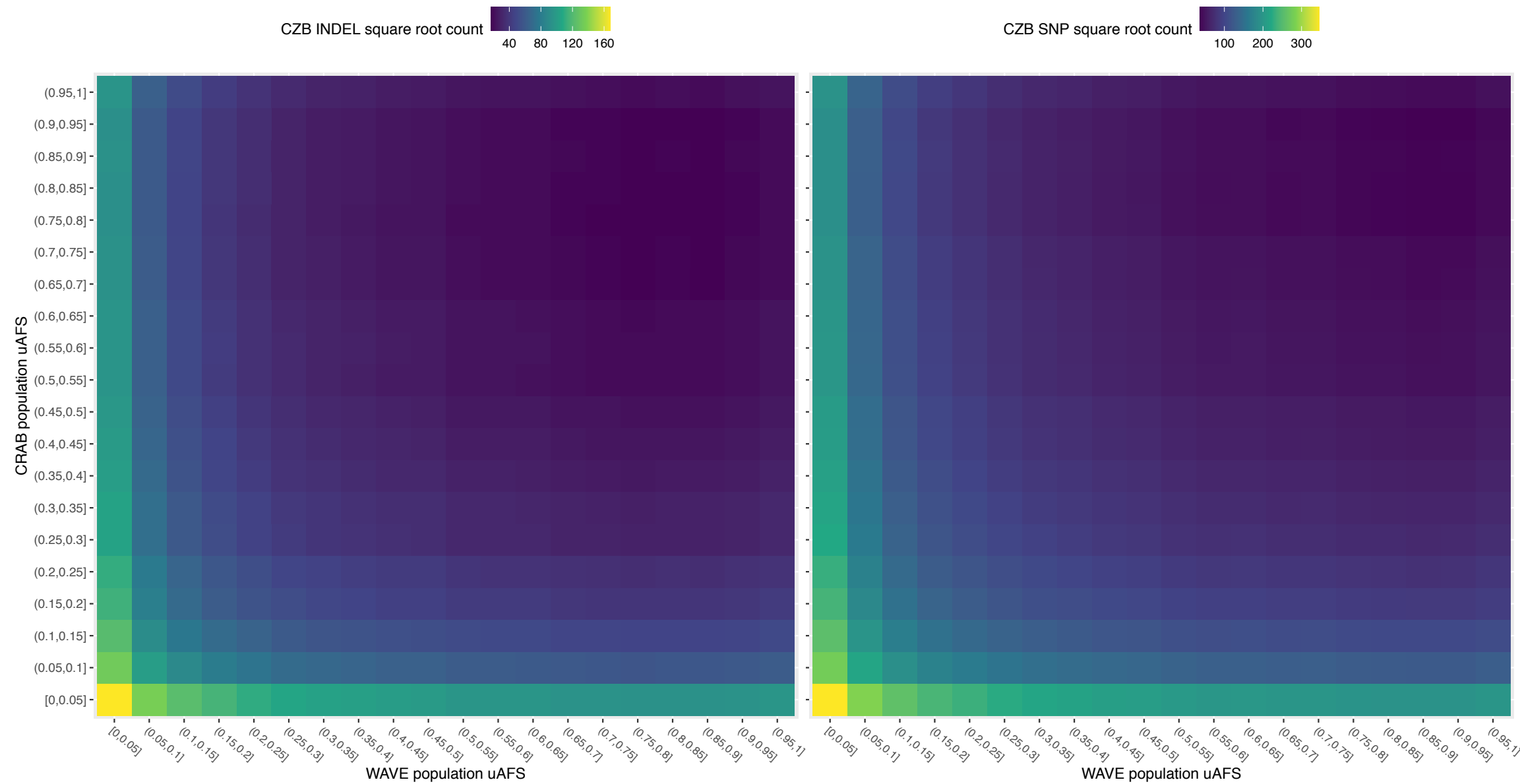
3. Derived allele frequencies (GATK call)

CZA CRAB-WAVE joint allele frequency spectra: square root count.



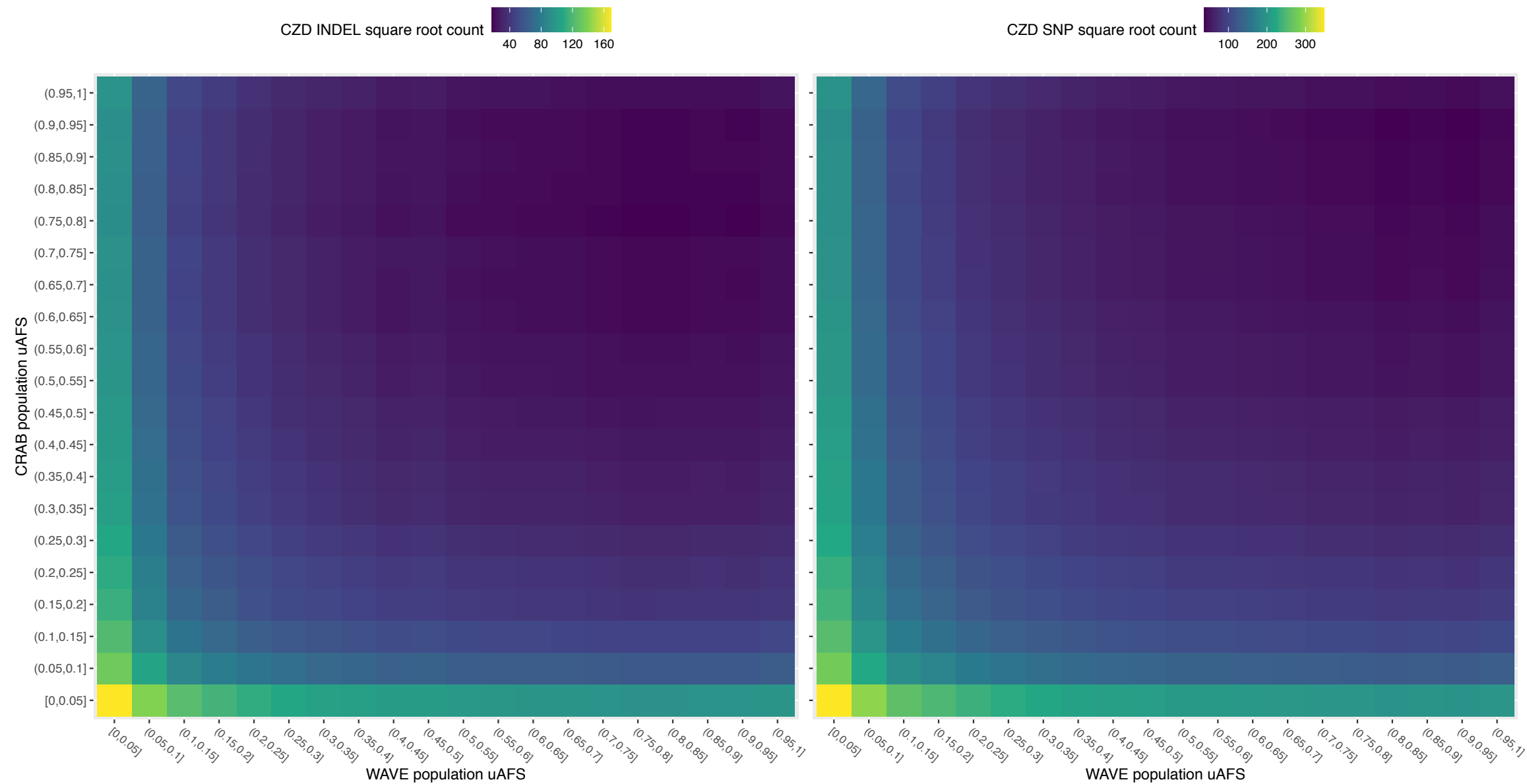
3. Derived allele frequencies (GATK call)

CZB CRAB-WAVE joint allele frequency spectra: square root count.



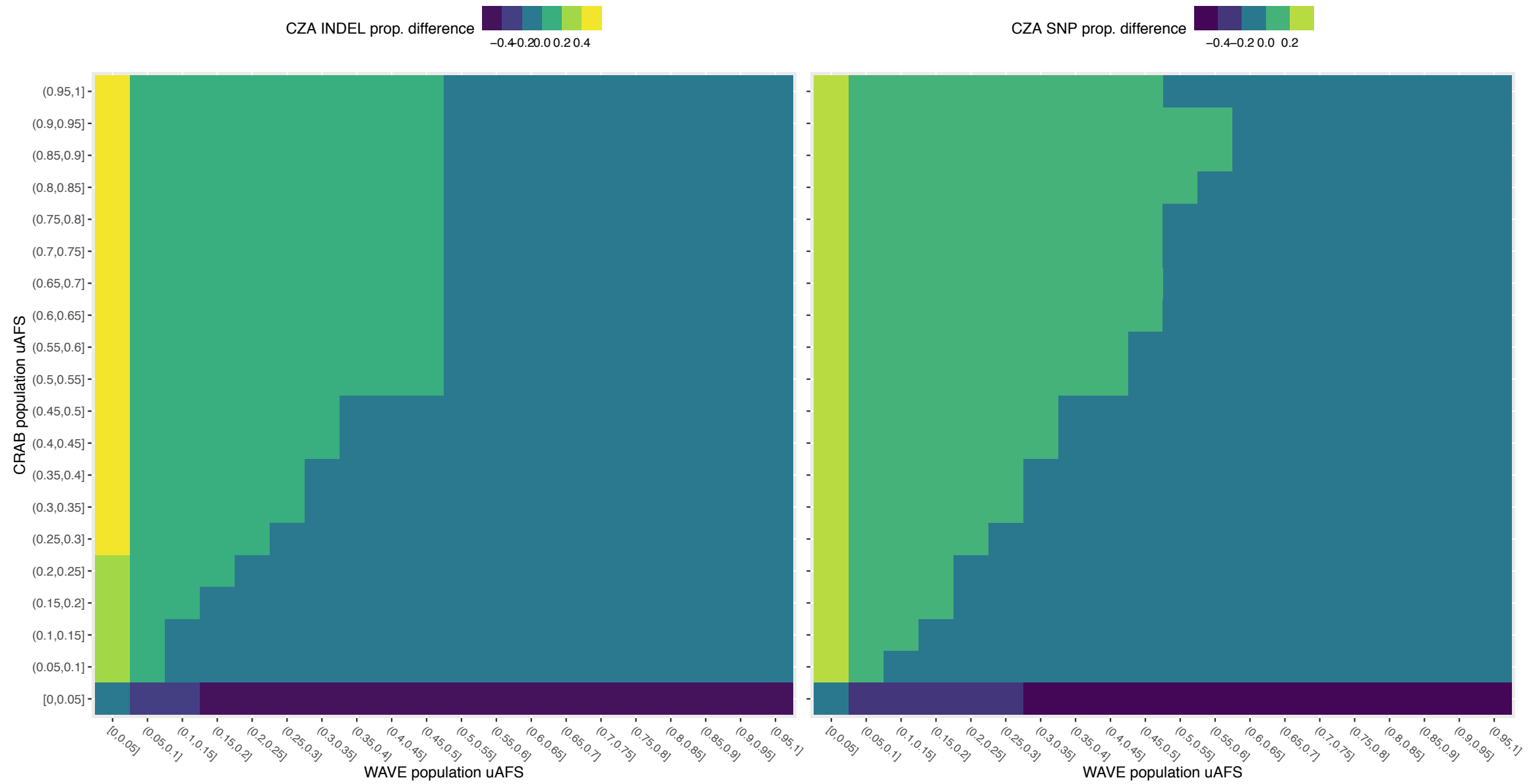
3. Derived allele frequencies (GATK call)

CZD CRAB-WAVE joint allele frequency spectra: square root count.



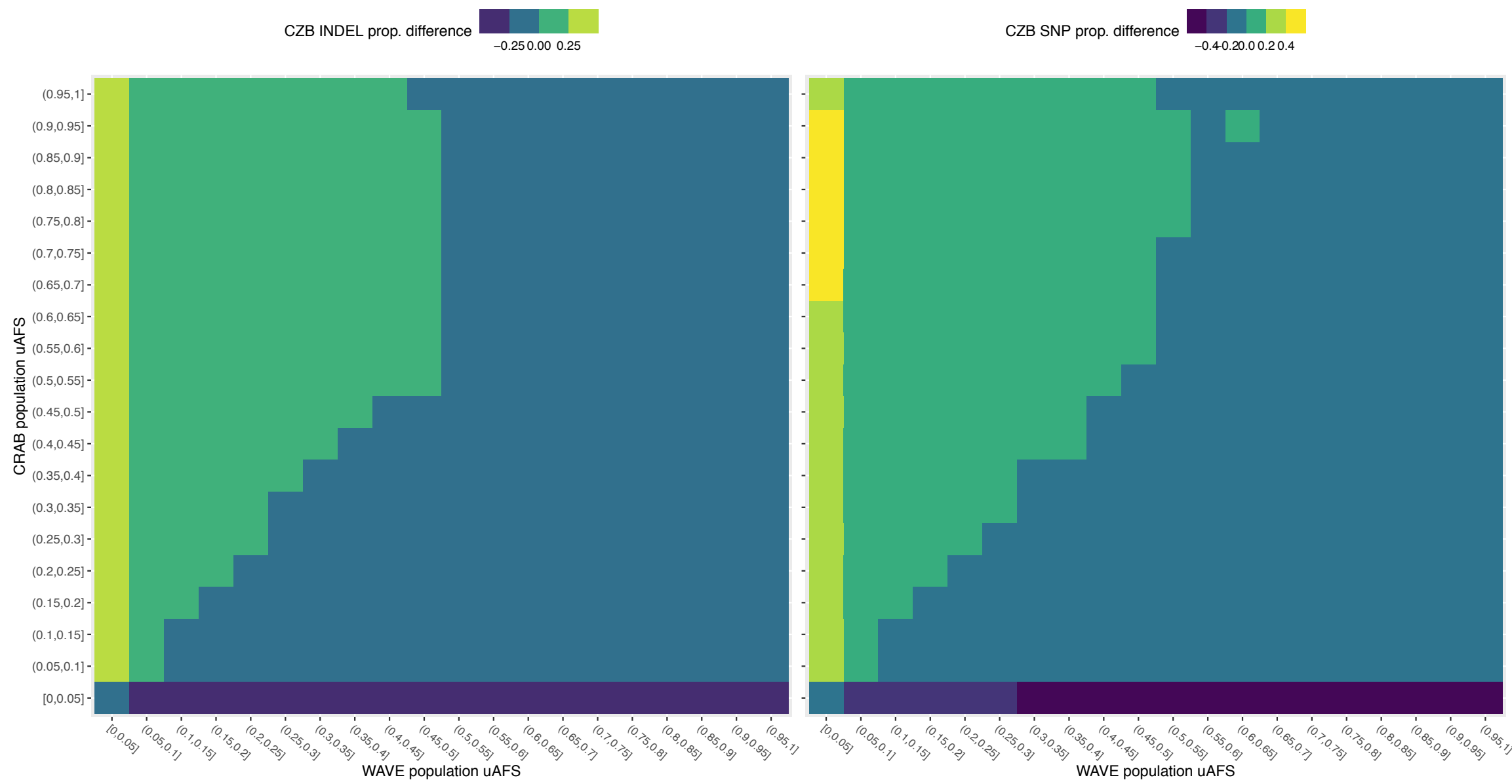
3. Derived allele frequencies (GATK call)

CZA CRAB-WAVE joint allele frequency spectra: difference in relative proportions.



3. Derived allele frequencies (GATK call)

CZB CRAB-WAVE joint allele frequency spectra: difference in relative proportions.



3. Derived allele frequencies (GATK call)

CZD CRAB-WAVE joint allele frequency spectra: difference in relative proportions.

