

Copy number variants underlie the major selective sweeps in insecticide resistance genes in *Anopheles arabiensis* from Tanzania.

Eric R. Lucas, Sanjay C. Nagi, Bilali Kabula, Arjen E. Van't Hof, Emily J. Rippon, Dimitra Pipini, Nicholas J. Harding, Naomi A. Dyer, Chris S. Clarkson, Alistair Miles, Martin J. Donnelly, David Weetman

Electronic Supplementary Material
Supplementary figures and tables

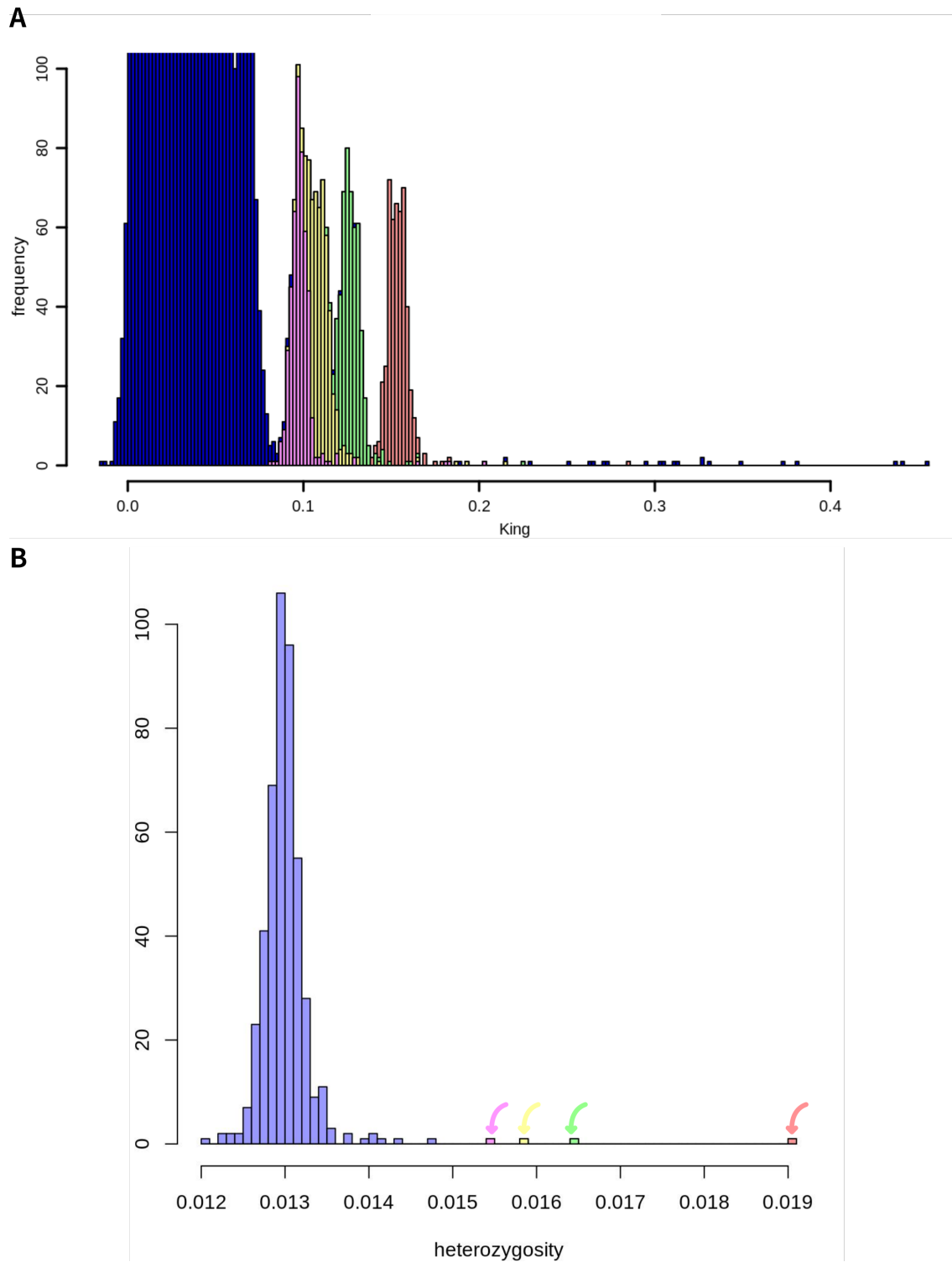


Fig. S1: A Histogram (zoomed in on the x axis to show right-hand tail) of all pair-wise KING relatedness values across all samples in our dataset. Four samples had universally high relatedness values to all samples. The relatedness values attributable to these four samples have been respectively coloured in pink, yellow, green and red on the histogram. **B** Histogram of heterozygosity values across all samples, showing that the same four samples (coloured and indicated with arrows) have elevated heterozygosity. As a result, these four samples had artificially augmented KING values.

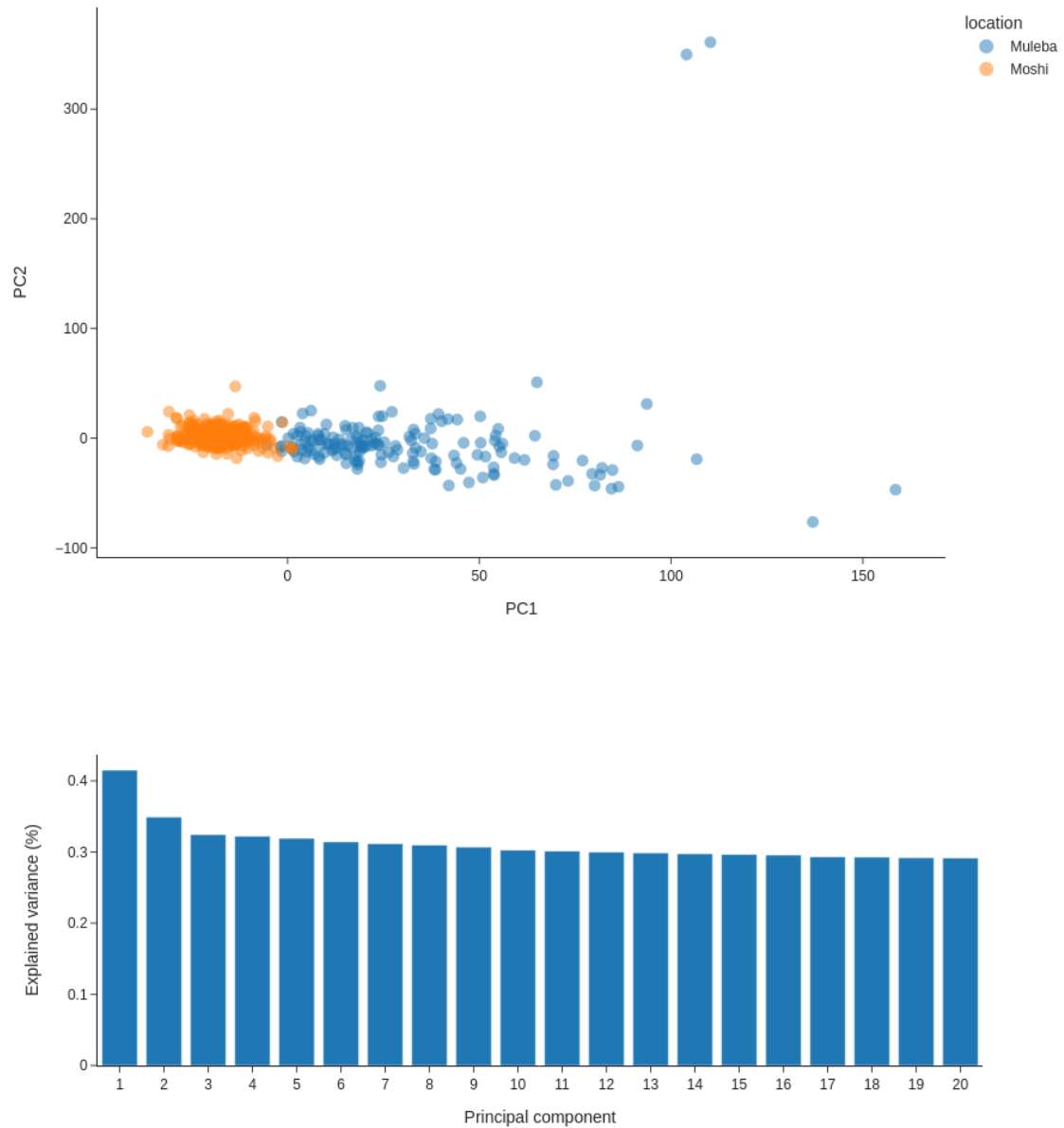


Fig. S2: PCA (using quality-filtered biallelic SNPs from genomic region 3L:15,000,000-41,000,000, euchromatic and free of chromosomal inversions). Top panel shows clustering of samples by region. Bottom panel show variance explained by the first 10 PCs, indicating that PCs 3 onwards explain similar levels of variance and are thus likely only capturing noise.

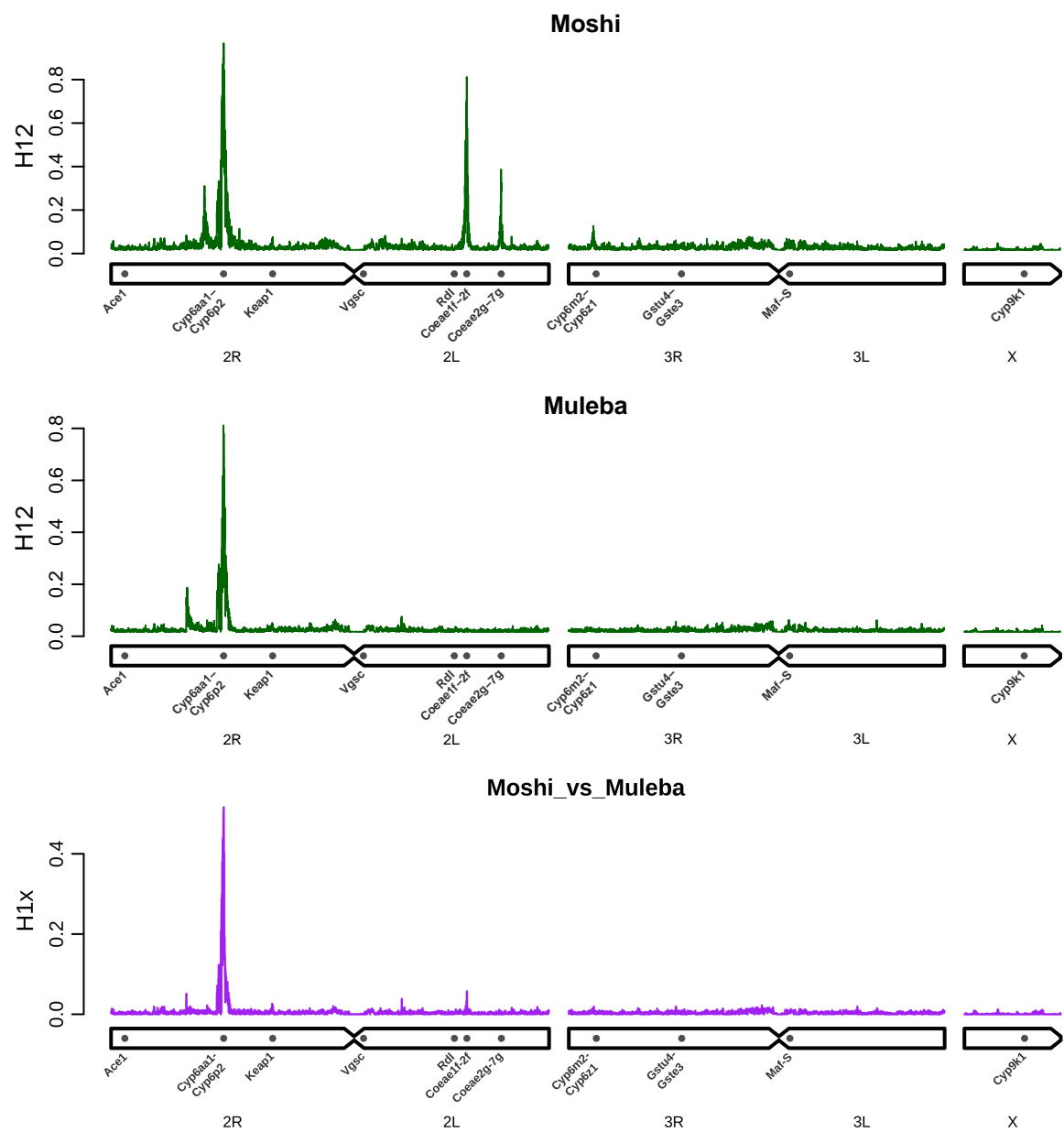


Fig. S3: Selection scans showing genome-wide H12 signal in Moshi (top) and Muleba (middle), as well as shared signals of selection (H1x) between the two sites (bottom).

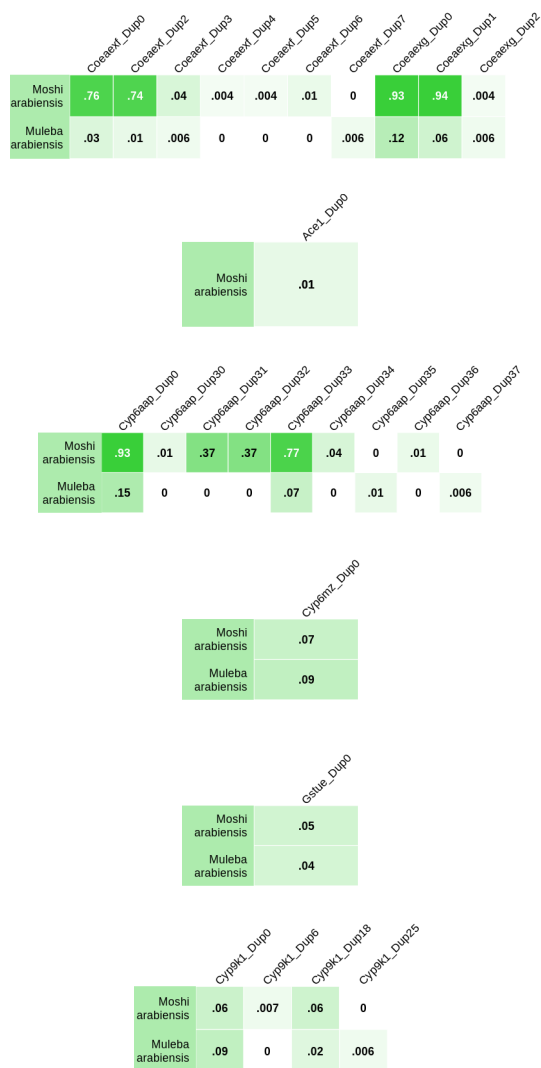


Fig. S4: Frequency (proportion of samples carrying at least one copy) of known CNV alleles detected using diagnostic reads around the *Coeaexf* cluster, the *Coeaexg* cluster (combined into a single subfigure), *Ace1*, the *Cyp6aa* / *Cyp6p* cluster, the *Cyp6m* / *Cyp6z* cluster, the *Gste* cluster and *Cyp9k1*. Only CNV alleles with frequency > 0% are shown. Cell darkness provided as a visual aid for the magnitude of the value in each cell. The genomic coordinates of each CNV allele can be found in Supplementary Data S@. In each cluster, the “Dup0” column indicates the presence of increased copy number in any of the genes in the cluster. Where this is larger than the sum of known alleles, it suggests the presence of uncharacterised CNV alleles. “Del” alleles in *Ace1* represent secondary deletions within the *Ace1*-Dup1 CNV.

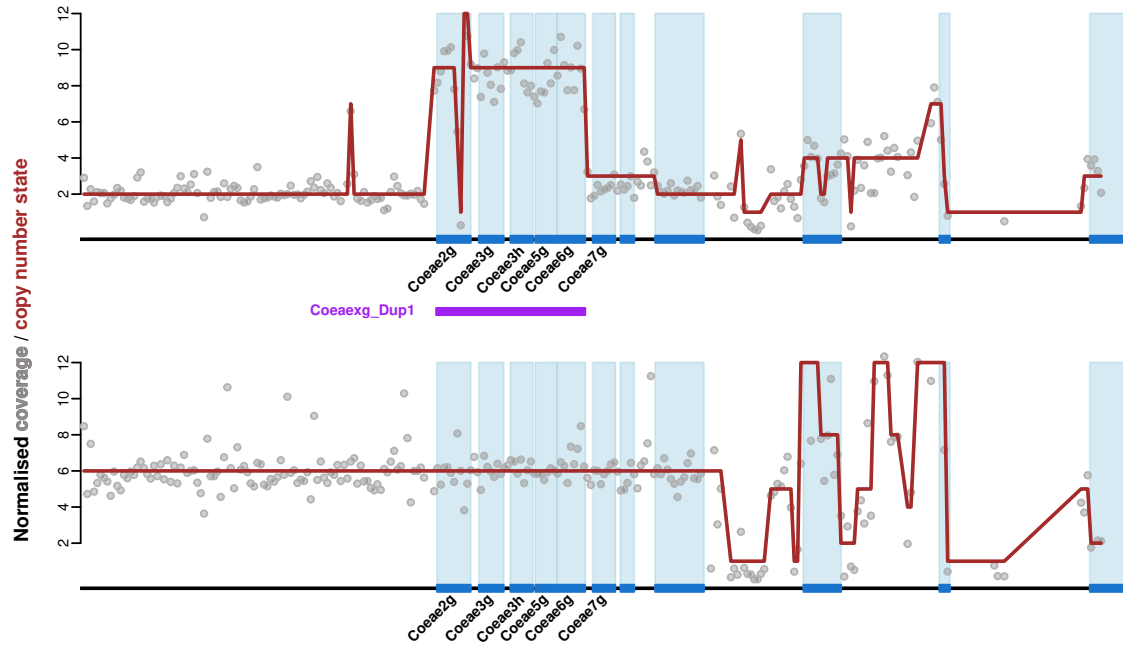


Fig. S5: Example coverage traces for CNVs in *Coeaexg*. Grey points show raw coverage (in 300 bp windows) normalised to a copy number of two (normal diploid copy number); brown line shows the output of the Hidden Markov Model (HMM) through the coverage data, indicating the predicted copy number state in each window. The genomic region to the right of the *Coeaexg* cluster has erratic coverage, suggesting a repeat region or poor genome assembly. CNVs in *Coeaexg* are evident from the HMM being consistently above the normal value of 2 across the region. The top plot shows an example sample carrying the main CNV allele found in our dataset of *An. arabiensis* from Tanzania (*Coeaexg_Dup1*, region covered by the CNV shown by purple bar). The bottom plot shows an example sample of *An. coluzzii* from Korle-Bu (Ghana), where the CNV extends far to the left, and into the region of erratic coverage to the right, thus meaning that we could not identify discordant reads that could tag individual CNV alleles in this population.

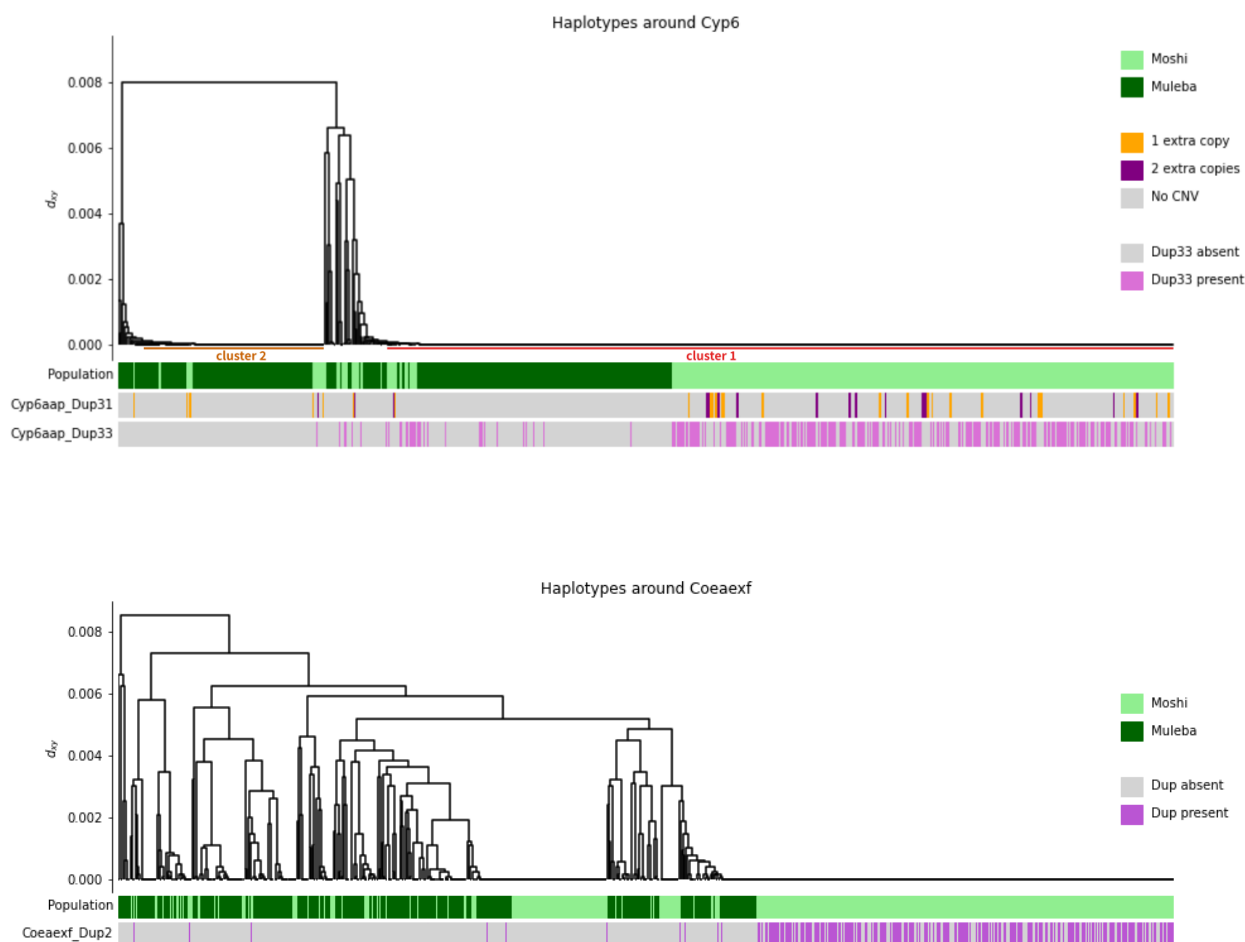


Fig. S6: Haplotype clustering of the genomic region around *Cyp6aa/Cyp6p* (top panel) showing that nearly all haplotypes belong to one of two selective sweeps. The less common cluster (cluster2) is predominantly found in Muleba and is not associated with CNVs. The more common cluster (cluster1), is found in both regions. Both *Cyp6aap_Dup31* and *Cyp6aap_Dup33* form a subset of haplotypes in cluster1. For *Cyp6aap_Dup33*, it was possible to assign presence (mauve) or absence (grey) of the CNV for each haplotype. For *Cyp6aap_Dup31*, it was only possible to determine whether the mosquito to which the haplotype belongs had a single extra copy (yellow), two (purple) or none (grey). A single extra copy indicates the sample is heterozygous for the CNV, and thus haplotypes labelled in yellow may not themselves carry the CNV. Similarly, in the *Coeaexf* region (bottom panel), haplotypes bearing the CNV allele *Coeaexf_Dup2* represented a subset of haplotypes from a large swept cluster. Clustering was performed using 500 SNPs in each region, and CNV alleles were phased by identifying SNPs that were highly correlated with their presence / absence.. Full workings to reproduce this analysis can be found at https://github.com/vigg-lstm/GAARD_east/blob/main/CNV_analysis/sweeps.

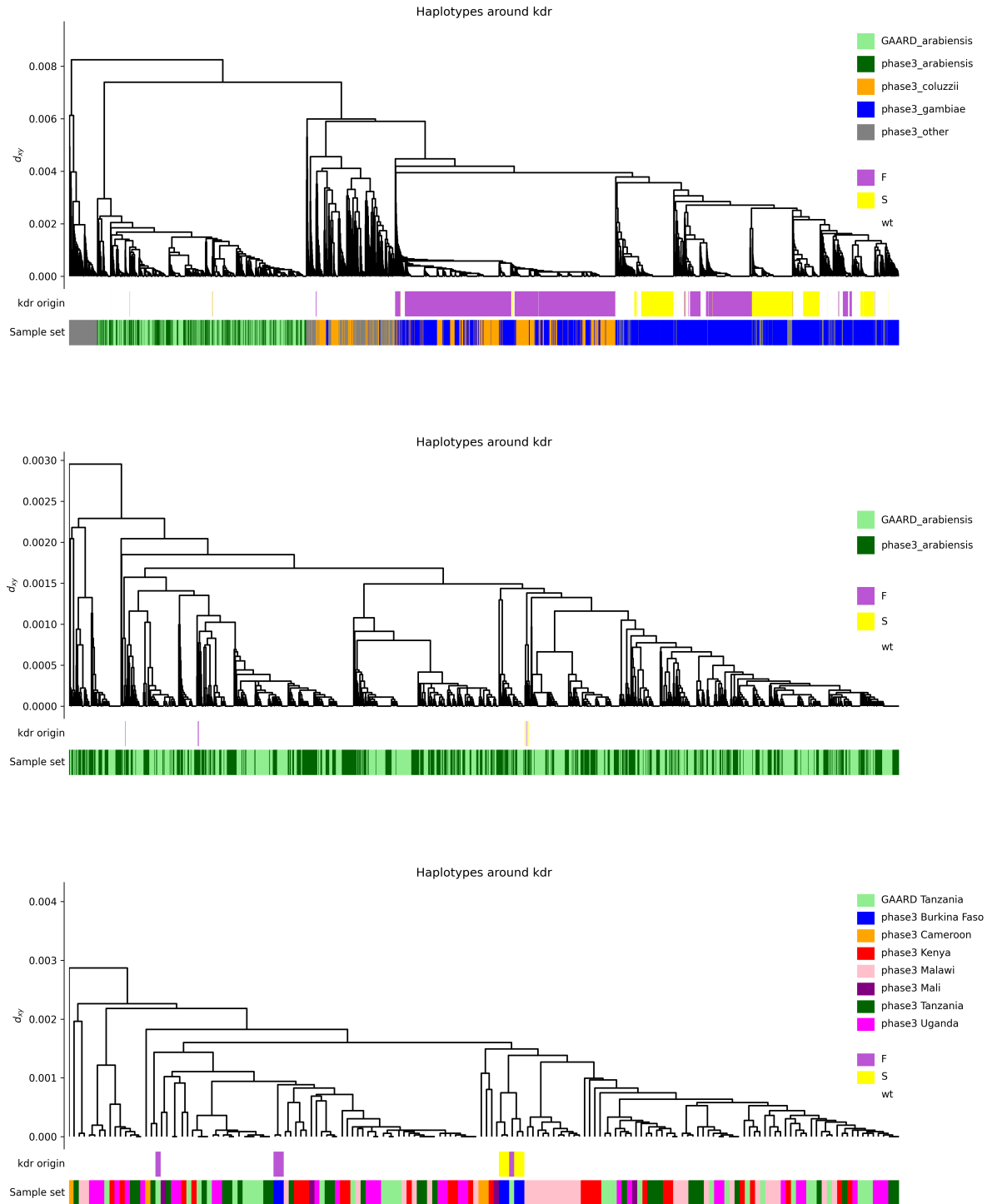


Fig. S7: Clustering of haplotypes around the *Vgsc* genomic region (2L:xxx-xxx, as in clarkson ref) reveals a diversity of *Vgsc*-995 origins in *An. arabiensis*, none of which are introgressed from *An. gambiae* or *An. coluzzii*. Combining our data with all haplotypes from phase 3 of Ag1000G (top) shows *An. arabiensis* haplotypes forming their own cluster, distinct from other species. When keeping only *An. arabiensis* haplotypes (middle), three different *Vgsc*-995F clusters are seen, despite only four such haplotypes existing in the dataset. The bottom plot shows all eight *Vgsc*-995 mutant haplotypes (two from our Tanzanian data, six from phase3 samples from Burkina Faso) and a random sub-sample of wild-type *An. arabiensis*, allowing a closer view of sample set labels for interpretation. The two Tanzanian *Vgsc*-995F haplotypes appear to be independent origins, one of which clusters more closely with haplotypes from Burkina Faso than with other Tanzanian haplotypes.

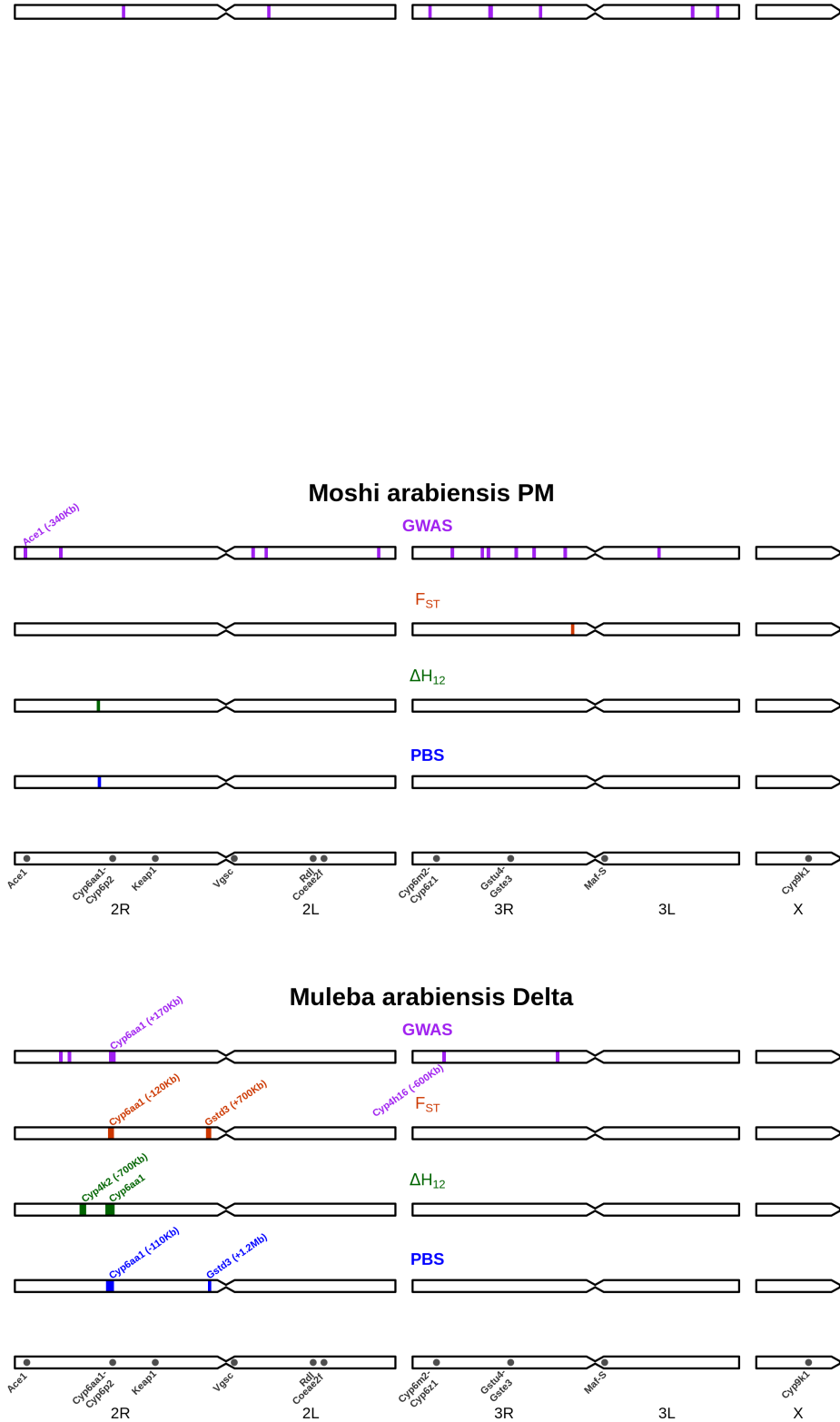


Fig. S8: Genomic regions implicated in insecticide resistance by each of our four approaches. For the global GWAS method, these are 100,000 bp windows which contained at least 10 of the top 1000 significant SNPs. For F_{ST} , these are significant peaks which contained at least one haplotype significantly positively associated with resistance (Supplementary Data S2). For ΔH_{12} and PBS, these are significant positive peaks (ie: indicating stronger signals of selection in resistant compared to susceptible samples). Regions are annotated with genes discussed in the manuscript as possibly causing the signal. Genomic distances in brackets indicate the distance of the peak either to the left (-) or right (+) of the gene in question.