

Demography-informed selection signatures reveal different roles of two polymorphic inversions in local adaptation of *Anopheles gambiae*

--Manuscript Draft--

Manuscript Number:	PGENETICS-D-15-00120
Full Title:	Demography-informed selection signatures reveal different roles of two polymorphic inversions in local adaptation of <i>Anopheles gambiae</i>
Short Title:	Demography-informed selection signatures in inversions
Article Type:	Research Article
Section/Category:	Evolution
Keywords:	Inversions, selection, demography, <i>Anopheles gambiae</i>
Abstract:	Chromosomal inversions are important structural changes that may facilitate divergent selection when they capture co-adaptive loci. However, identifying selection targets within inversions can be challenging. The high degree of differentiation between heterokaryotypes as well as the differences in demographic histories of collinear regions compared to inverted ones reduce the power of detection with traditional outlier scans. Here, we developed a new approach that uses discriminant functions to classify loci that are under selection (or drift) informed from inversion-specific expectations for selected and neutral patterns as characterized by a diversity of summary statistics. We demonstrate the approach with an analysis of RAD sequencing data we collected in a classic Dipteran species with polymorphic inversion clines - <i>Anopheles gambiae</i> , a malaria vector species from sub-Saharan Africa. Contrary to minimal geographic structure among populations in collinear regions, individuals are clustered by SNPs from two polymorphic inversions with the 2La and 2Rb arrangements predominating in dry habitats compared with prevalence of the 2L+a and 2R+b arrangements in wet habitats. Despite of both being adaptive introgressions, estimates of the demographic histories, and hence expectations for patterns of neutral and selected within the two inversions differ. Specifically, the origin of 2La predated the divergence of two species. Nevertheless, with our approach we were able to detect putative regions subject to selection within both chromosomal inversions, with much higher power than traditional FST-outlier analysis. Moreover, the arrangement 2L+a that associates with wet habitat exhibits much more selection signatures than the arrangement 2La that associates with dry habitat. We discuss the implications of these results with respect to studies of rapid adaptation in these malaria vectors, and in particular, the insights our newly developed approach offers for identifying not only potential targets of selection but also the lineage that has undergone adaptive change.
Additional Information:	
Question	Response
Data Availability	Yes - all data are fully available without restriction
PLOS journals require authors to make all data underlying the findings described in their manuscript fully available, without restriction and from the time of publication, with only rare exceptions to address legal and ethical concerns (see the PLOS Data Policy and FAQ for further details). When submitting a manuscript, authors must provide a Data Availability Statement that describes where the data underlying their manuscript can be found.	
Your answers to the following constitute your statement about data availability and will be included with the article in the event of publication. Please note that	

<p>simply stating 'data available on request from the author' is not acceptable. If, however, your data are only available upon request from the author(s), you must answer "No" to the first question below, and explain your exceptional situation in the text box provided.</p> <p>Do the authors confirm that all data underlying the findings described in their manuscript are fully available without restriction?</p>	
<p>Please describe where your data may be found, writing in full sentences. Your answers should be entered into the box below and will be published in the form you provide them, if your manuscript is accepted. If you are copying our sample text below, please ensure you replace any instances of XXX with the appropriate details.</p>	<p>All relevant data are within the paper and its Supporting Information files. All customized scripts, setting files for programs and genomic data will be deposited on Dryad.</p>
<p>If your data are all contained within the paper and/or Supporting Information files, please state this in your answer below. For example, "All relevant data are within the paper and its Supporting Information files."</p> <p>If your data are held or will be held in a public repository, include URLs, accession numbers or DOIs. For example, "All XXX files are available from the XXX database (accession number(s) XXX, XXX)." If this information will only be available after acceptance, please indicate this by ticking the box below.</p> <p>If neither of these applies but you are able to provide details of access elsewhere, with or without limitations, please do so in the box below. For example:</p> <p>"Data are available from the XXX Institutional Data Access / Ethics Committee for researchers who meet the criteria for access to confidential data."</p> <p>"Data are from the XXX study whose authors may be contacted at XXX."</p>	
<p>* typeset</p>	
<p>Additional data availability information:</p>	<p>Tick here if the URLs/accession numbers/DOIs will be available only after acceptance of the manuscript for publication so that we can ensure their inclusion before publication.</p>

Dear *PloS Genetics* journal,

We are delighted to submit the manuscript “**Demography-informed selection signatures reveal different roles of two polymorphic inversions in local adaptation of *Anopheles gambiae***” as possible publication in *PloS Genetics*.

With the increasing availability of physical maps among different species, the important role of chromosomal inversions in maintaining adaptive divergence has been demonstrated in many systems. Because chromosomal inversions suppress recombination between heterokaryotypes, co-adaptive genotypes within inversions that confer local adaptations can be protected. However, a method that identifies selection signature of loci/region within inversions in an empirical system of adaptive divergence is lacking.

In this study, we develop a new approach for locating a selection signature when population divergence (or speciation) is promoted by inversions, a genomic architecture that limits the power of traditional tests for divergent selection (i.e., F_{ST} -outlier tests). Genomic outlier scans either have an inherent assumption for a simple demographic model which might not capture the characteristics of the species' history or treat the genome as a homogenous pool evolving neutrally under a single common demographic history. In addition, it depends on an expectation of the background level of divergence generated by genetic drift to be distinguishable from that of selected scenarios. Therefore, these are improper methods to be applied for regions with reduced recombination (i.e., high levels of background divergence) or evolved under different demographic histories than the rest of the genome, such as inversions. The approach we designed circumvents these problems, with the intent to lower false discovery rates as well as increase detection power. First, separate demographic histories are inferred for inversion and collinear regions to generate region-specific background levels of differentiation. Second, genetic measures other than differentiations, such as haplotype and genetic diversity, are also included in the analyses to increase the power of distinguishing adaptive alleles from neutral ones. Third, discriminant functions are used to predict selection instead of an outlier approach by using training datasets from neutral versus selection simulations.

We applied this approach to an analysis of two common inversions in a severe malaria vector species, *Anopheles gambiae*, which are associated with aridity adaptations. We found the signature of selection traces predominantly in the wet-adapted 2L+^a arrangement rather than the dry-adapted 2La arrangement, which was previously overlooked in studies of 2La inversions. We think this study is worthy of publishing in *PloS Genetics* because it is the first attempt to identify selection signature within inversions in an empirical system using inversion-specific coalescent expectations. It has a wide application as it directly tests whether inversions evolved under Kirkpatrick and Barton (2006)'s local adaptation theory as well as how many loci/regions were involved.

Best regards,

Qixin He and L. Lacey Knowles

Department of Ecology and Evolutionary Biology, University of Michigan

1 **Title: Demography-informed selection signatures reveal different roles of two polymorphic
2 inversions in local adaptation of *Anopheles gambiae***

3

4 **Running Title: Demography-informed selection signatures in inversions**

5

6 **Type:** Article

7

8 **Authors:** Qixin He and L. Lacey Knowles

9

10 **Contact information:**

11 Qixin He, Ph. D. candidate

12 Email: heqixin@umich.edu

13 Affiliation: Department of Ecology & Evolutionary Biology, 1109 Geddes Ave., Museum of
14 Zoology, University of Michigan, Ann Arbor MI 48109-1079. phone (734) 763-7943, fax (734)
15 763-4080

16

17 L. Lacey Knowles, Professor and Curator

18 Email: knowlesl@umich.edu

19 Affiliation: Department of Ecology & Evolutionary Biology, 1109 Geddes Ave., Museum of
20 Zoology, University of Michigan, Ann Arbor MI 48109-1079. phone (734) 763-7943, fax (734)
21 763-4080,

22 **Corresponding Author:** Qixin He.

23 **Manuscript information:** 4 Figs + 1 Table; Supplementary, 6 Figs + 2 Tables + 1 Text file

24 **Data Archival Location:** Dryad upon acceptance

25 **Keywords:** Inversions, selection, demography, *Anopheles gambiae*

26

27

1 **Abstract**

2 Chromosomal inversions are important structural changes that may facilitate divergent selection
3 when they capture co-adaptive loci. However, identifying selection targets within inversions can
4 be challenging. The high degree of differentiation between heterokaryotypes as well as the
5 differences in demographic histories of collinear regions compared to inverted ones reduce the
6 power of detection with traditional outlier scans. Here, we developed a new approach that uses
7 discriminant functions to classify loci that are under selection (or drift) informed from inversion-
8 specific expectations for selected and neutral patterns as characterized by a diversity of summary
9 statistics. We demonstrate the approach with an analysis of RAD sequencing data we collected in
10 a classic Dipteron species with polymorphic inversion clines - *Anopheles gambiae*, a malaria
11 vector species from sub-Saharan Africa. Contrary to minimal geographic structure among
12 populations in collinear regions, individuals are clustered by SNPs from two polymorphic
13 inversions with the 2La and 2Rb arrangements predominating in dry habitats compared with
14 prevalence of the 2L+^a and 2R+^b arrangements in wet habitats. Despite of both being adaptive
15 introgressions, estimates of the demographic histories, and hence expectations for patterns of
16 neutral and selected within the two inversions differ. Specifically, the origin of 2La predated the
17 divergence of two species. Nevertheless, with our approach we were able to detect putative
18 regions subject to selection within both chromosomal inversions, with much higher power than
19 traditional F_{ST} -outlier analysis. Moreover, the arrangement 2L+^a that associates with wet habitat
20 exhibits much more selection signatures than the arrangement 2La that associates with dry
21 habitat. We discuss the implications of these results with respect to studies of rapid adaptation in
22 these malaria vectors, and in particular, the insights our newly developed approach offers for
23 identifying not only potential targets of selection but also the lineage that has undergone adaptive
24 change.

25 **Author summary**

26 Evidence supporting the prevalent role of inversions in facilitating adaptive divergence is now
27 common. Nonetheless, identifying the loci within inversions that confer local adaptation has
28 been difficult. Here, we present a new approach to identifying targets of selection, and
29 demonstrate the approach with an analysis of two common inversions in a malaria vector species,
30 *Anopheles gambiae*. Specifically, discriminant functions were built to detect regions under

1 selection based on the comparison of the observed data to expectations generated from
2 simulations of neutral versus selection under inversion specific demographic histories. With this
3 new approach we found the signature of selection traces predominantly to the 2L+^a arrangement
4 rather than the introgressed 2La arrangement, suggesting that the same region has involved in
5 two episodes of adaptations - specifically to wet and dry habitats, respectively. Our approach can
6 be widely applied for chromosomal regions with different histories, as well as non-model species
7 with no genome reference or sparse genomic markers, and represents a significant improvement
8 over traditional tests like F_{ST} -outlier analyses to detect targets of selection.

9

1 **Introduction**

2 Detecting the signature of natural selection from molecular data has been a central focus of
3 geneticists, not only because of the deeper understanding of molecular evolution it brings, but
4 also its potential in revealing important functional information [1,2]. With the advances in
5 sequencing technologies, a plethora of new methods have been developed for identifying
6 selected sites and regions of the genomes. These approaches capture different attributes the
7 signature of selection might leave. These include tests for patterns of exceptionally long
8 haplotypes [3-5], surges in linkage disequilibrium (LD) [6,7], and skewed site frequency
9 spectrums [8-10]. However, when the goal is to identify selection under spatially divergent
10 selection among populations, a predominant approach is to rely on F_{ST} -outlier tests to scan for
11 regions with exceptionally high divergence between ecologically divergent populations [11-15].
12 The appeal of the approach also extends from its broad applicability, especially with the lack of
13 genomic resources for detailed estimates of LD or haplotypes in most taxa. However, the power
14 of such methods becomes inherently limited when adaptive mutations occur in regions with
15 reduced recombination, such as with chromosomal inversions. This becomes particular
16 problematic for studying adaptive divergence given the important role inversions play in
17 maintaining co-adaptive genotypes [16,17].

18 Identifying selection targets within alternative inversions can be challenging because of two
19 characteristics that distinguish them from collinear regions (i.e., genomic regions without
20 inversions). First, in genomic regions with inversions there is an overall high divergence between
21 the inverted and non-inverted (i.e., standard) chromosomal arrangement [e.g., 18], which itself
22 diminishes the power to detect selection [19,20]. Second, as a consequence of the effect that
23 inversions have on the mixing of alleles among individuals (i.e., alleles captured by inversions
24 are protected from mixing with alleles on standard chromosomes), the genome is a mosaic, with
25 the amount of gene flow differing across regions depending on whether it includes an inversion.
26 This means that tests that rely upon a single neutral parameterization (either with simulations
27 under island-model used in FDIST2 [21] or coancestry matrix in FLK [13] will necessarily be
28 based on a mis-specification for expected patterns of divergence, which can then further
29 exacerbate the difficulties with detecting targets of selection.

1 In this study, we develop a new approach for locating a selection signature when population
2 divergence (or speciation) is promoted by inversions, a genomic architecture that limits the
3 power of traditional tests for selection, as described above. By first estimating region specific
4 demographic histories, we are able to generate region specific neutral expectations for
5 demography-adjusted selection tests [see also 22]. We then built a discriminant function using
6 combinations of summary statistics from sequences simulated under neutral and selected
7 scenarios, which can then be used to assign empirical sites into neutral or selection classes based
8 on predictions from the discriminant function. We apply the newly developed approach to
9 *Anopheles gambiae*, which like other Dipeteran species, has a long history of research on
10 inversion polymorphisms [23]. As a widespread species with large population size, it typically
11 lacks significant population structure, except for the geographic structure observed at genomic
12 regions characterized by seven commonly segregating inversions on the chromosome 2 [24-26].
13 Here we focus on two large inversions, 2La and 2Rb (21.4Mb and 7 Mb, respectively), that
14 exhibit stable clines, but which defy detection of putatively selected regions/genes within the
15 inversions because of high LD and F_{ST} across the inversions [27]. The high frequency of 2La in
16 dry geographic areas (savannas) compared to wet areas (forests) [28], a predictable cycle in its
17 frequency during dry and wet seasons, and higher resistance to dessication in experiments [29],
18 identify its role in adaptive divergence, as do similar trends observed in 2Rb [30]. In addition to
19 the general difficulties with identifying the targets of selection within inverted regions discussed
20 above, another potential complication is associated with the origins of the two inversions. Both
21 2La and 2Rb are examples of adaptive introgression from another species in the same species
22 complex, *An. arabiensis* [31], which is sympatric with *An. gambiae* in arid savanna areas [32,33].
23 This suggests that the region contained in the inversions might be old enough to be highly
24 differentiated from the standard chromosome in *An. gambiae*, which exacerbates the
25 aforementioned high divergence problem between heterokaryotypes.

26 Nevertheless, despite these challenges, as our analyses demonstrate, we are able to not only
27 identify genomic areas associated with targets of selection, but we are also able to make
28 statistical statements about what lineage underwent selective divergence with respect to the
29 ancestral state. We highlight how these new findings have interesting implications for how we
30 think about rapid divergence in the malaria vector *An. gambiae*. We also discuss the applicability
31 of our procedure for tests of selection in other taxa more generally, including those where

1 specific genomic regions have vastly different history than the rest of the genome because of the
2 mosaic nature.

3

4 **Results**

5 We collected genomic data in 259 *An. gambiae* individuals and 8 *An. arabiensis* from six sites
6 along a gradient of wet to dry habitats in Cameroon (see supplementary text S1 for specimen
7 identification results; Table S1 for sampling information). From individually-barcoded double
8 digest genomic Radseq libraries [34], and two lanes of 100bp paired-end sequencing on Illumina
9 HiSeq2000 platform, 25,966 loci (i.e., RADtags from the genomic preparation) were mapped
10 onto Chromosome 2, 3 and the X, after filtering for ambiguously mapped reads and loci with low
11 coverage per sample or low presence across samples. Although we recognize that this represents
12 a small proportion of the genome (about 1%), the goal of this manuscript is to demonstrate the
13 promise of the approach for detecting approximate targets of selection, as opposed to providing a
14 full analysis of the proportion of sites under selection within inversions contributing to adaptive
15 divergence, which would require additional sequencing that is beyond the scope of this study.

16 Our new approach of targeting specific selection within inversion involves several procedures
17 (summarized in Fig. 1a). First we estimated divergence time and introgression rate between the
18 two species from collinear regions (Fig. 1b) using site frequency spectrum [35]. We then
19 estimated inversion specific parameters, such as gene flux rate between alternative inversions
20 and the age of inversion mutations (Fig. 1c). This framework was then used for tests of selection
21 against neutral expectations that are region specific based upon discriminant functions derived
22 from summary statistics estimated from empirical versus simulated data under the inferred
23 demographic history.

24 ***Establishing a demographic null model for tests of selection***

25 Principal Components Analysis (PCA) and Discriminant Analysis of Principal Components
26 (DAPC) analyses of SNP data showed a lack of population genetic structure in the species in the
27 collinear regions. Specifically, no distinctive geographic structure among the six sampled
28 populations was apparent from the PCA (Fig. S3) and one group ($K = 1$) received the highest

1 support in the DAPC analysis (see supplementary text for details about tests of geographic
2 structure). Therefore individuals were pooled across populations to estimate a demographic
3 history that could be used to establish a null expectation for patterns of genomic divergence
4 under a model of drift.

5 Specific demographic histories were inferred for different genomic regions from the region-
6 specific site frequency spectrum (SFS) under a composite-likelihood approach implemented in
7 fastsimcoal2 (for details see methods) [35]. Estimations showed that the two species diverged
8 fairly recently around 87K years ago (~1.4Ne generations ago, T_{div} in Table 1; assuming 12
9 generations a year [26]) with small but constant genetic exchange (on the order of 1E-7 to 1E-8,
10 which is about exchanging 0.01 to 0.2 individuals per generation; m in Table 1).

11 In contrast to the collinear regions, SNPs from 2La (2L: 20524058-42165532) and 2Rb (2R:
12 19023925-26758676) were distributed into one of three clusters in the PCAs (first PC explaining
13 20.3% and 11.9% of the total variance in 2La and 2Rb, respectively; Fig. S4 a,d). The results
14 from the DAPC also supported $K = 3$ as the most likely number of genetic clusters (Fig. S4b, e).
15 Individuals that form the three clusters identified from these analyses correspond to the three
16 genotypes associated with the inverted genomic region: namely, the inverted homokaryotypes
17 (referred to as *I/I* hereafter), the heterokaryotypes (*I/S*), and the standard (i.e., non-inverted)
18 homokaryotypes (*S/S*) based on comparison with molecular karyotyping results (Fig. S4c, f).
19 Separate demographic histories were therefore inferred for the different genomic arrangements
20 (i.e., inverted versus standard chromosomal regions) for both the 2La and 2Rb regions, with a
21 recombination parameter to accommodate occasional gene flux between alternative karyotypes
22 (Fig. 1c).

23 Inversions were modeled as an introgressed between *An. arabiensis* and *An. gambiae* after the
24 species' divergence time (see demographic model in Fig. 1c). The age of inversion mutation and
25 their introgression time will strongly influence the baseline divergence expected for detecting
26 selection. Interestingly, the coalescent time of *S* and *I* (T_{IS} , Fig. 1c) of 2La was around 3Ne
27 generations ago (Table 1), which is much more ancient than species divergence time (T_{div}). In
28 contrast, T_{IS} of 2Rb is similar to T_{div} . Despite its ancient origin (which predates the divergence of
29 *An. arabiensis* and *An. gambiae*), the introgression event of 2La was found to have occurred
30 more recently (~0.5Ne) than 2Rb (~0.8Ne) (Table 1).

1 ***Signature of selection within inversions***

2 A modified F_{ST} -outlier scan based on inversion specific demographic history was first performed.
3 Due to structural constraints of inversions, double recombination rates are higher in the center of
4 an inversion versus breaking points. We therefore divided inversion regions into 150kb segments
5 and adjusted recombination rates from our average estimates based on local 2Mb F_{ST} estimations
6 (for details see methods) and then simulated 1000 neutral cases per region. We found that
7 because the divergence of *I* and *S* of 2La was quite old, F_{ST} estimates between I and S for
8 neutrally evolved DNA segments tended to be very high (e.g., the 95 percentile and 99 percentile
9 of simulation estimates generally exceed 0.8; see Fig. 2a). Consequently, the power of
10 differentiating selection from neutral genes based on F_{ST} estimates alone is severely limited
11 within the 2La inversion (~0.3 and 0.1 with 0.05 and 0.01 false-positive rates respectively; Fig 4),
12 in contrast to collinear regions where most F_{ST} estimates are concentrated around 0 (Fig. 2c).
13 This is less exacerbated in the younger 2Rb (Fig. 2b; ~0.42 and 0.3 with 0.05 and 0.01 false-
14 positive rates respectively; Fig 4). Yet, the empirical F_{ST} distribution has much longer tail than
15 the collinear F_{ST} distribution (Fig. 2d), again highlighting the challenge with detecting putatively
16 selected regions given the elevated F_{ST} -values that reduces detection power with very high false-
17 discovery rates (Fig. 4; see Material and Methods for the distinction between false-positive rates
18 and false-discovery rates).

19 Since outlier analysis based solely on F_{ST} measures have limited power to differentiate selection
20 from drift (see Fig. 4), we designed a new approach to make use of a diversity of summary
21 statistics. 1000 simulations of 50kb sequences containing selected locus were ran for each of the
22 three scenarios: 1) neutral; 2) selection occurs on sites associated with *S* in *Anopheles gambiae*; 3)
23 selection occurs on sites associated with *I* in *An. gambiae* and *An. arabiensis* (Fig. 1a).
24 Discriminant functions (DAPC) were built to differentiate the three scenarios based on the values
25 observed across the summary statistics of simulated RADtags (Fig. 3a). This approach yielded a
26 significant gain in power (Fig. 4a-d) and a reduction in false-discovery rates (Fig. 4e-h). In order
27 to minimize random effects from individual RADtags, average summary statistics of random loci
28 sampled across the simulated 50kb regions at a similar density as the empirical data were
29 calculated to build the discriminant function which classify regions as containing selected locus
30 or not. The power of assigning regions correctly to three scenarios is therefore above 0.9 (Table

1 S2) with a comparable false-positive rates (~0.06) as the outlier analysis with a top 5% cut-off
2 (Fig. 4). Moreover, in contrast to the pervasive problem with high false-discovery rates based on
3 F_{ST} -values alone (i.e., falsely identified targets of selection), such cases significantly decreased
4 when multiple summary statistics are used (Fig. 4e-h; Table S2). Interestingly, heterozygosity (H)
5 and θ_π contributed more of the variation used in differentiating each scenario compared to F_{ST} .

6 After empirical estimates of summary statistics were transformed into discriminant scores,
7 regions were assigned into neutral or selection classes based on predictions from the discriminant
8 function estimated from the simulations (see Fig 3a). We found that on 2Rb, four regions
9 contained selected sites associated with I (red bars on Fig 3b), while four contain selection
10 associated with S (blue bars in Fig. 3b). These regions span some loci that were identified as F_{ST}
11 outliers, but not all of them. In 2La, we recovered a much higher number of regions identified as
12 being selected in S compared to those associated with I (Fig. 3).

13

1 **Discussion**

2 Our study is the first attempt to identify selection signature within inversions in an empirical
3 system using inversion-specific coalescent expectations. By modeling the unique origin and
4 introgression histories of each inversion compared to collinear regions, our approach allows for
5 identifying targets of selection through discriminant function analyses trained with a suite of
6 genetic measures. Compared to traditional F_{ST} -outlier analysis, simulation-trained discriminant
7 function classification analysis has a much higher detection power and reduced false discovery
8 rates (Fig. 4). When the F_{ST} -outlier analysis was applied to *Anopheles gambiae* populations,
9 older inversions (i.e., 2La) suffered more from reduced power of detection because of the
10 increasing baseline of divergence estimators for neutral sequences in older inversions (see Fig. 2,
11 4). In contrast, the simulation-trained discriminant function classification analysis achieved
12 similar detection power and false discovery rates for both 2La and 2Rb (Fig. 4). In addition to
13 the higher detection power, the new approach can also decipher the specific branch where
14 selection occurred. In particular, selection was found to have occurred predominantly on
15 *Anopheles gambiae* standard arrangement (2L+^a) rather than the introgressed arrangement (2La),
16 suggesting that the same region has involved in two episodes of adaptations - specifically to wet
17 and dry habitats, respectively. Our study highlighted the importance of adaptive introgression via
18 inversions as a mechanism for maintaining adaptive divergence in species with high gene flow
19 among populations inhabiting heterogeneous environments.

20

21 **Detecting selection in regions with reduced recombination rates**

22 The traditional indexes (F_{ST} , D_{XY}) in genomic scan to quantify inter-population differentiation are
23 still the most effective way to detect targets of divergent selection in most cases (e.g., Soria-
24 Carrasco et al. 2014). With the availability of population genomic data, such genomic scan
25 studies become a common practice (reviewed in [36]) and many studies reported regions of
26 elevated divergence, termed genomic islands (e.g., [37-39]; but see [40]. While common
27 practices do not accommodate regions with reduced recombination or different demographic
28 histories, we tried to circumvent the overall high divergence problem in inversions by generating
29 null expectations from neutral evolution in selection identification. However, this is still not ideal
30 because the detection power is highly diminished when variability of divergence estimators

1 among neutral sequences increases with the age of divergence (see overall low power and high
2 false discovery rates in F_{ST} -outlier results in Fig. 4, S5-6).

3 The simulation-trained discriminant function classification analysis in this study is similar to an
4 ABC parameter estimation process, in which simulations close to observations are retained, and
5 linear functions are built to estimate parameters from PCs transformed from summary statistics
6 [41,42]. The difference is that we are choosing models (neutral or under selection) instead of
7 estimating parameters because our simulations are fixed with specific parameters estimated from
8 reconstructed demographic histories. In addition, the increased power of applying multiple
9 summary statistics, converting them into PCs and transforming into discriminant functions is a
10 good way to detect selection when expectations for a single statistic do not differ sufficiently
11 under selection scenarios to provide discriminatory power. This approach has the great
12 advantage of jointly considering all summary statistics together across different
13 arrangements/species instead of relying on F_{ST} measures alone. Since RADtag sequences are
14 relatively short (~100bp), simulated sequences still show a great degree of variation in summary
15 statistics for each scenario so that discriminant functions do not differentiate them completely
16 (Fig. 4). Nevertheless, true positive rate are ~0.8 in assigning to the correct scenarios. However,
17 these rates can be improved significantly when we identify regions under selection using average
18 summary statistics across a 50kb window encompassing randomly distributed short Radtags.
19 Using regional average of statistics greatly reduced stochastic variability of each scenario so that
20 correct assignment rates are above 90% in all scenarios (Fig. 4, Table S2).

21 Due to the rapid decay of linkage disequilibrium in the *Anopheles gambiae* complex (r^2 close to
22 0 between SNPs 1Kb apart; see [43]), local changes in linkage disequilibrium inside inversions
23 are hard to detect, especially with our sparsely sampled genomic dataset. Yet, LD or extended
24 haplotype tests might be informative within individuals of the same arrangements to detect
25 selection signatures [e.g., 44] because although recombination between heterokaryotypes is
26 highly reduced, it is not reduced within same karyotypes (e.g., LDhat [45] did not find
27 significant changes in recombination rates within same karyotypes in our data). Nevertheless, our
28 demography-informed discriminant analysis is still powerful for tests of selection in genomic
29 regions with different evolution histories. These results highlight the general utility of this
30 approach in non-model species, especially with sparse sampling of the genome or no reference

1 genome for applying methods for detecting changes in LD within the same karyotype from
2 individually-barcoded whole genome sequences.

3 ***Age of inversions and its implications on adaptive introgression***

4 Our estimation of the older divergence of 2L+^a from 2La corroborates recent karyotype and
5 genomic phylogenies for the *Anopheles gambiae* complex [43,46], which predicted that 2L+^a
6 evolved in the ancestor of *An. gambiae* and *An. arabiensis* and alternative karyotypes got fixed
7 in the two species after their divergence. The coalescent time of 2L+^a and 2La is estimated to be
8 around 3Ne, which was predicted to be a good age to detect selection from estimating coalescent
9 time between alternative arrangements [47]. However, our F_{ST} -outlier analysis showed very
10 limited power in 2La region, with a slightly higher power in the younger 2Rb region (see Fig. 4).
11 The reason might rise from the fact that the introgression age of 2La is too short compared to its
12 origin time so that gene flux between heterokaryotypes is not frequent enough to reduce the
13 divergence between neutrally evolving regions inside inversions. 2Rb, in which the background
14 divergence is not significantly older than the time of introgression, showed clearer patterns of
15 regions with significantly higher divergence.

16 With the new approach, we gained equal power of detecting selection in both of the inversions
17 and identified more loci/regions under selection in 2La than 2Rb (Fig. 3). More interestingly, the
18 signature overwhelmingly showed more prevalent selection in sites associated with *S* rather than
19 *I* in 2La, which has never been identified before. This might not be surprising given 2La's origin
20 history. Past molecular studies have shown that 2La is the ancestral arrangement [48], from
21 which 2L+^a arose and got fixed in *Anopheles gambiae* [31]. The new phylogenomic study on the
22 seven species in the complex even challenged the long held belief of direction of introgression,
23 proposing the possibility that 2La might have been a retained polymorphism in *An. gambiae*
24 and introgressed into *An. arabiensis* [43]. Although our study does not provide evidence for the
25 direction of introgression, the fact that the younger 2L+^a arrangement shows stronger selection
26 signal highlights the important fitness advantage associated in wet forest environments, which
27 was previously overlooked in studies of 2La inversions[49]. Therefore, our data suggest that the
28 alternative arrangements have each facilitated either wet or dry habitat adaptation at different
29 stages of the species history. Additional genomic study will be needed to pinpoint the exact
30 genes that are under selection given such inferences are beyond our dataset given the relatively
31 low density of Radtag markers (on average ~10kb between adjacent markers).

1 ***Adaptation from mosaic genotypes***

2 With the increasing availability of physical maps among different species, the important role of
3 chromosomal inversions in maintaining adaptive divergence has been demonstrated in many
4 systems, such as controlling flowering differences in *Mimulus guttatus* between different
5 ecotypes [50] and wing patters that form Batesian mimicry in *Heliconius numata* [51]. The
6 unique aspect of adaptation via polymorphic inversions in Sub-Saharan mosquito species
7 *Anopheles gambiae* is the prevalent introgression and sharing of inversions among sibling
8 species [31,49], which has posed challenges in recovering phylogenetic relationships within the
9 *Anopheles gambiae* complex [31,52-54].

10 Our study showed how different sets of adaptive loci for different habitat (e.g., dry vs. wet) can
11 be maintained in alternative rearrangements throughout a widespread species. The fact that the
12 species complex does not have complete reproductive isolation and that they "borrow" pre-
13 adapted inversions from each other while exploring new environments provides an interesting
14 example of how adaptation leads to mosaic genotypes instead of new species, especially for
15 species with big populations and high connectivity. This mode of adaptation coincides with
16 recent theories that predict that mechanisms that reduce or suppress recombination, or increase
17 linkage between co-adaptive and maladapted genotypes, will be advantageous because the gene
18 complex can avoid being swamped when gene flow persists between populations that are under
19 divergent selection [16,55-57].

20

21

22

1 **Material and Methods**

2 ***Sample collection and DNA extraction***

3 Mosquitoes were collected indoors at each site using either aspirators or insecticide spray and
4 preserved individually in 0.5ml tubes containing 100% ethanol. Morphologies of each sample
5 were examined according to Gillies and Meillon [58] and Gillies and Coetze [59] under
6 dissecting microscope before DNA extraction with QIAamp DNA Mini Kit, whose yield ranged
7 from 10-200 ng per sample.

8 ***Molecular identification of species and karyotypes***

9 The species status of each sample (i.e., *gambiae*, *coluzzi*, or *arabiensis*) was determined by a
10 PCR-RFLP method following Fanello et al. [60]. Briefly, the species was identified by the
11 difference in the number of bands and/or fragment length after *HhaI* digestion of part of the
12 intergenic spacer (IGS) of the ribosomal DNA PCR products. The presence of inversions was
13 determined by PCR of unique breakpoint regions of alternative arrangements. For 2La, primers
14 were chosen to amplify a 492 bp region of 2La distal breakpoint and a 207 bp product from 2L+^a
15 proximal breakpoint [28]. For 2Rb, three primers amplify a 429 bp fragment on 2Rb breakpoint
16 and a 630bp fragment on 2R+^b breakpoint [61]. If only one of the two PCR bands is present on a
17 gel electrophoresis, then the sample is considered to be homokaryotype of one arrangement;
18 alternatively, if both bands are present with similar brightness, the sample is considered to be
19 heterokaryotype.

20 ***ddRAD library preparation and sequence analysis***

21 Genomic DNA from each sample was individually barcoded and used in a reduced complexity
22 library for Illumina sequencing using a double digestion Restriction Associated DNA sequencing
23 procedure (ddRADseq; for details see [34]). Briefly, DNA was digested with the two most
24 frequent restriction enzymes, MluCI and MseI, to maximize the number of unique short
25 fragments. The digested products were then ligated by part of an Illumina adaptor sequence and a
26 unique barcode. Ligation products were pooled among samples and size-selected between 340
27 and 420 base pairs (excluding adaptor lengths) using a Pippin Prep (Sage Science) machine. The
28 targeted-size ligation products were amplified by iProof™ High-Fidelity DNA Polymerase
29 (BIO-RAD) with 12 cycles. The library was sequenced in two lanes on the Illumina HiSeq2000

1 platform to generate paired-end 100 base pair reads. Sequences were identified to each sample
2 based on the barcodes. Only reads with an average quality score of at least 30 (Phred) and an
3 unambiguous barcode and restriction cut site were retained.

4 After filtering, sequences were mapped to the *Anopheles gambiae* reference genome AGam30
5 [62] using BWA-MEM algorithm in bwa with default settings [63]; mappings with quality scores
6 above 10 using SAMTOOLS [64] were retained. SNPs were called from mapped contigs and
7 genotypes were assigned using a maximum-likelihood statistical model [65,66] with the STACKS
8 v1.03 pipeline [67]; default settings were used except where noted below. Specifically, loci
9 (termed as “stacks” in the program) were identified from genomic locations with mappings of at
10 least 5 copies of RAD sequences in each individual using the PSTACKS program to ensure
11 credible calling of heterozygous SNPs in an individual [67]. A catalog of loci was built with the
12 CSTACKS program from the PSTACKS output files across individuals to check the presence or
13 absence of a particular locus at a genomic location. We retained loci that were present in at least
14 50% of all the samples and no more than two haplotypes per locus within each sample. All
15 customized scripts, setting files for programs and genomic data are available on Dryad under doi:
16 xxxxxxxxxxxx.

17 ***Population genetic structure of collinear and inverted regions***

18 Geographic structure of *An. gambiae* populations were examined by measuring population
19 divergence and performing principal component analysis (PCA). We performed separate
20 analyses on collinear regions and inverted regions of each chromosome. Weir and Cockerham’s
21 F_{ST} (1984) and nucleotide diversity (π) were estimated on a per-site basis as well as for a window
22 of 150kb (50kb steps were used to slide along the chromosome) using the POPULATIONS program
23 in the STACKS pipeline [67]. SNPs were thinned to be at least 1000bp apart on the genome and
24 imported into adegenet 1.4-1 package [68] in R [69] for PCA analyses. Only SNPs that are
25 present in all populations and at least 80% of all individuals were included in the study. Missing
26 genotypes were represented by the average value for the PCA analyses. Based on PCA results,
27 discriminant analysis of principal components (DAPC; [70]) were run to determine genetic
28 clusters within the species without prior assumptions on the model of population subdivision.
29 DAPC runs K-means clustering on the transformed PCs to identify groups of individuals that
30 maximize between-group genetic variation while minimizing within-group variation. The best

1 supported number of clusters is then determined by the comparison of model likelihoods through
2 Bayesian Information Criterion (BIC) similar to the program STRUCTURE [71,72]. The agreement
3 between genetic clusters inferred from inverted regions and molecular karyotyping assignment
4 was also checked to assess the reliability of PCR identification methods.

5 ***Demographic history of collinear and inverted regions***

6 Inference of demographic history implemented in FASTSIMCOAL2 [35] calculates composite-
7 likelihood of joint-SFS across populations under user-specified demographic scenario by
8 parameterized simulations sampled from priors and optimizes the parameter estimation through a
9 conditional maximization algorithm (ECM). The derived and ancestral states of the SNPs were
10 inferred from the comparison of four species in the *An. gambiae* complex: *An. gambiae* s.s., *An.*
11 *arabiensis*, *An. quadriannulatus*, and *An. merus*. Whole genome scaffolds of the latter three
12 species were mapped to *gambiae* genome using MUMMER 3.23 [73,74] and unique alignments
13 were kept. Majority state of the diallelic SNP among four species was considered ancestral and
14 multi-states SNPs were filtered. In order to maximize the number of SNPs included and ensure
15 reasonable running time for each scenario, we excluded individuals with less complete
16 sequencing coverage and subsampled SNPs in each case to infer region specific demographic
17 histories. Point estimates were obtained from the run with the best maximum likelihood out of 50
18 realizations. Confidence intervals of parameter estimates were then obtained by 100 parametric
19 bootstrapping runs from the point estimates.

20 ***Coalescent history between A. gambiae and A. arabiensis in collinear regions***

21 We first estimated the divergence time (T_{div}), introgression rate (m) and recent population
22 expansion (N_{cur} , T_{exp}) (Fig. 1b) of *An. gambiae* and *An. arabiensis* using joint-SFS built from
23 Chromosome 3 as proxies for collinear regions (X chromosome has a different effective
24 population size and different selection regime). In order to include more SNPs with missing data
25 while obtaining an accurate estimation of the site frequency in the population, all SNPs were
26 subsampled to 40 copies in *An. gambiae* and 6 in *An. arabiensis*. One variable SNP per RADtag
27 were chosen to build the site frequency spectrum to ensure that SNPs are not linked. We fixed
28 the population size (N_e) of *An. gambiae* to estimate other free parameters (see also [35]) (Fig. 1b,
29 c); this parameter was set to $\sim N_e = 750,000$ using a mutation rate of 3.5E-9 per base per
30 generation (estimation from Drosophila resequencing, [75]) given that *An. gambiae* was

1 estimated to have a nucleotide diversity (π) of $0.01024 \pm 4.0\text{E-}5$ from all Radtags. Due to
2 pervasive introgression between the *An. arabiensis* and *An. gambiae* (Besansky et al. 2003), we
3 estimated the rate of gene flow between the two species after the divergence from their common
4 ancestor (Fig. 1b).

5 ***Coalescent history between A. gambiae and A. arabiensis in regions with inversion***
6 ***polymorphisms***

7 We estimated the divergence time between standard chromosomes and inverted chromosomes
8 (T_{IS}), time of introgression for inversions (T_{int}), and recombination rates between alternative
9 rearrangements (r) using joint-SFS built from 2La or 2Rb regions while fixing other parameters
10 that were estimated from the collinear region. SNPs were subsampled to 20 copies in *I* and *S*, and
11 6 in *An. arabiensis*. We fixed the parameters that have been estimated from collinear region
12 models and focused on inversion specific parameters because collinear region has a larger SNP
13 dataset. *I* and *S* of *An. gambiae* were treated as separate populations, but connected by a severely
14 reduced recombination (r) (Fig. 1c). The same introgression rate estimated from collinear regions
15 was applied to that between the inverted copy *I* in *An. gambiae* and *An. Arabiensis*, while the
16 introgression rate between *S* in *An. gambiae* and *I* in *An. Arabiensis* was set to zero.

17 ***Selection signature in inversion regions***

18 Based on reconstructed demographic histories for inversion regions, we obtained neutral
19 expectations of population genomic measures through simulations. These measures can then be
20 compared against empirical data to detect selection. We first applied F_{ST} -outlier analyses similar
21 to traditional approaches. One difficulty lies in the heterogeneity of recombination rates inside
22 inversions between heterokaryotypes (i.e., recombination rates are higher in the center and
23 decrease sharply towards the breaking point region [76]), whereas the recombination rates
24 estimated from SFS demographic modeling was an average. In order to make neutral simulations
25 more realistic, we adjusted the recombination rate to be higher in the center region and lower on
26 the two sides. First, empirical 2Mb F_{ST} -windows in 150kb-step were calculated and fed into a
27 smooth spline function in R to get fitted values for each 150kb segment. Recombination rate was
28 adjusted to the value that generated the empirical mean F_{ST} -value for the segment. 1000
29 demographic simulations were carried out using estimated parameters for each segment to
30 generate 100bp DNA sequences. The population divergence measures, F_{ST} , between *I* and *S* were

1 estimated for the sequences. Lastly, empirical divergence measures for each locus were
2 compared against the range of values from simulations to identify outliers.

3 Our second approach utilized sets of summary statistics to detect selection through simulation-
4 trained discriminant function classification analysis. Three scenarios were run for 1000
5 replications under the current demographic model using MSMS [77] : a) pure neutral evolution; b)
6 selection occurred on the branch of *An. arabiensis* and inverted chromosomes (I); c) selection
7 occurred on the branch of *An. gambiae* and continued in standard chromosomes. Selection
8 started from when two species diverged with selection coefficient ranging from 0.01 to 0.0001
9 (Fig. 2c). Selected locus is located in the center of a 50kb-long simulated region. We tested two
10 ways of building discriminant functions: 1) based on summary statistics of one neutral locus that
11 is close to the selected locus; 2) based on an average of summary statistics across several neutral
12 loci of a region that contains selected locus. For the first approach, we sampled 100bp long
13 sequences that are located 5Kb away from the selected locus (an average distance between
14 empirical adjacent Radtags in our data). 9 summary statistics, including heterozygosity (H), θ_n of
15 each population, and population pairwise- F_{ST} , were calculated for each simulated sequence. For
16 the second approach, we sampled random sets of 100bp short sequences across the entire 50kb
17 according empirical Radtag distributions and calculated an average of these summary statistics.
18 In both cases, discriminant functions using principle components transformed from summary
19 statistics (DAPC; [70]) were built based on the simulated training sets to differentiate the three
20 scenarios. We then used the discriminant function to predict which scenario each empirical
21 locus/region belonged to based on their estimated summary statistics.

22 We compared the performance of the outlier test with the discriminant function classification
23 analysis by evaluating their false-positive rate, detection power (true positive rate) and false
24 discovery rates. We define these rates, according to Lotterhos and Whitlock [20], as follows:
25 false-positive rate is ratio of the number of significant neutral loci and the total number of neutral
26 loci; the detection power is the ratio of identified selected loci and the total number of selected
27 loci; the false-discovery rate is the number of significant neutral loci divided by the total number
28 of significant loci. Three levels of selected loci were tested in simulated datasets to calculate
29 power and false-discovery rates. Specifically, out of 1000 otherwise neutral loci, 1%, 5%, or 10%
30 are simulated under selection. In our first approach, cut-off values for loci to be considered

1 significant were 95% and 99% of neutral distributions. Hence, we set the false-positive rates to
2 be compared in our first approach as 5% and 1% (see Fig. 4). For the second approach, the cut-
3 off values to be classified as an outlier (i.e., false positive rate) are not arbitrarily decided.
4 Instead, the percentage of correct assignment from discriminant functions represent how likely
5 different scenarios can be differentiated. Hence, the error rate in assigning neutral cases to
6 selection cases is the false-positive rates.

7

8 **Acknowledgments**

9 We appreciate useful discussions with Anthony J. Cornel, Yooksook Jin, Huateng Huang, Mark
10 Kirkpatrick, and Mark Christie. We would also like to thank Seraphin Menzepoh and Kevin
11 Njabo for tremendous help in the field work. Financial support was provided by a NSF Doctoral
12 Dissertation Improvement Grant (DEB-1210359) to Q. He and L. L. Knowles.

13

14

1 **References**

- 3 1. Nielsen R, Hellmann I, Hubisz M, Bustamante C, Clark AG (2007) Recent and ongoing selection in the
4 human genome. *Nature Reviews Genetics* 8: 857-868.
- 5 2. Hohenlohe PA, Phillips PC, Cresko WA (2010) Using population genomics to detect selection in
6 natural populations: key concepts and methodological considerations. *International Journal of
7 Plant Sciences* 171: 1059-1071.
- 8 3. Sabeti PC, Reich DE, Higgins JM, Levine HZP, Richter DJ, et al. (2002) Detecting recent positive
9 selection in the human genome from haplotype structure. *Nature* 419: 832-837.
- 10 4. Sabeti PC, Varilly P, Fry B, Lohmueller J, Hostetter E, et al. (2007) Genome-wide detection and
11 characterization of positive selection in human populations. *Nature* 449: 913-918.
- 12 5. Voight BF, Kudaravalli S, Wen XQ, Pritchard JK (2006) A map of recent positive selection in the
13 human genome. *PLoS Biology* 4: 446-458.
- 14 6. Kim Y, Nielsen R (2004) Linkage disequilibrium as a signature of selective sweeps. *Genetics* 167:
15 1513-1524.
- 16 7. Wang ET, Kodama G, Baidi P, Moysis RK (2006) Global landscape of recent inferred Darwinian
17 selection for *Homo sapiens*. *Proceedings of the National Academy of Sciences of the United
18 States of America* 103: 135-140.
- 19 8. Carlson CS, Thomas DJ, Eberle MA, Swanson JE, Livingston RJ, et al. (2005) Genomic regions
20 exhibiting positive selection identified from dense genotype data. *Genome research* 15: 1553-
21 1565.
- 22 9. Nielsen R, Williamson S, Kim Y, Hubisz MJ, Clark AG, et al. (2005) Genomic scans for selective
23 sweeps using SNP data. *Genome Research* 15: 1566-1575.
- 24 10. Ronen R, Udpa N, Halperin E, Bafna V (2013) Learning natural selection from the site frequency
25 spectrum. *Genetics* 195: 181-193.
- 26 11. Antao T, Lopes A, Lopes RJ, Beja-Pereira A, Luikart G (2008) LOSITAN: A workbench to detect
27 molecular adaptation based on a F(st)-outlier method. *BMC Bioinformatics* 9.
- 28 12. Beaumont MA, Balding DJ (2004) Identifying adaptive genetic divergence among populations from
29 genome scans. *Molecular Ecology* 13: 969-980.
- 30 13. Bonhomme M, Chevalet C, Servin B, Boitard S, Abdallah J, et al. (2010) Detecting selection in
31 population trees: The Lewontin and Krakauer test extended. *Genetics* 186: 241-U406.
- 32 14. Foll M, Gaggiotti O (2008) A genome-scan method to identify selected loci appropriate for both
33 dominant and codominant markers: a Bayesian perspective. *Genetics* 180: 977-993.
- 34 15. Gunther T, Coop G (2013) Robust identification of local adaptation from allele frequencies. *Genetics*
35 195: 205-220.
- 36 16. Kirkpatrick M, Barton N (2006) Chromosome inversions, local adaptation and speciation. *Genetics*
37 173: 419-434.
- 38 17. Yeaman S (2013) Genomic rearrangements and the evolution of clusters of locally adaptive loci.
39 *Proceedings of the National Academy of Sciences* 110: E1743-E1751.
- 40 18. Cheng C, White BJ, Kamdem C, Mockaitis K, Costantini C, et al. (2012) Ecological Genomics of
41 *Anopheles gambiae* Along a Latitudinal Cline: A Population-Resequencing Approach. *Genetics*
42 190: 1417-1432.
- 43 19. Beaumont MA (2005) Adaptation and speciation: what can F_{ST} tell us? *Trends in Ecology &
44 Evolution* 20: 435-440.
- 45 20. Lotterhos KE, Whitlock MC (2014) Evaluation of demographic history and neutral parameterization
46 on the performance of FST outlier tests. *Molecular ecology* 23: 2178-2192.
- 47 21. Beaumont MA, Nichols RA (1996) Evaluating loci for use in the genetic analysis of population
48 structure. *Proceedings of the Royal Society of London Series B: Biological Sciences* 263: 1619-
49 1626.

- 1 22. Rafajlović M, Klassmann A, Eriksson A, Wiehe T, Mehlig B (2014) Demography-adjusted tests of
2 neutrality based on genome-wide SNP data. *Theoretical Population Biology*.
- 3 23. Coluzzi M, Sabatini A, Petrarca V, Dideco MA (1979) Chromosomal differentiation and adaptation to
4 human environments in the *Anopheles gambiae* complex. *Transactions of the Royal Society of
5 Tropical Medicine and Hygiene* 73: 483-497.
- 6 24. Czeher C, Labbo R, Vieville G, Arzika I, Bogueau H, et al. (2010) Population Genetic Structure of
7 *Anopheles gambiae* and *Anopheles arabiensis* in Niger. *Journal of Medical Entomology* 47: 355-
8 366.
- 9 25. Lanzaro GC, Toure YT, Carnahan J, Zheng LB, Dolo G, et al. (1998) Complexities in the genetic
10 structure of *Anopheles gambiae* populations in west Africa as revealed by microsatellite DNA
11 analysis. *Proceedings of the National Academy of Sciences of the United States of America* 95:
12 14260-14265.
- 13 26. Lehmann T, Hawley WA, Grebert H, Collins FH (1998) The effective population size of *Anopheles
14 gambiae* in Kenya: Implications for population structure. *Molecular Biology and Evolution* 15:
15 264-276.
- 16 27. White BJ, Hahn MW, Pombi M, Cassone BJ, Lobo NF, et al. (2007) Localization of candidate regions
17 maintaining a common polymorphic inversion (2La) in *Anopheles gambiae*. *PLoS Genet* 3: e217.
- 18 28. White BJ, Santolamazza F, Kamau L, Pombi M, Grushko O, et al. (2007) Molecular karyotyping of
19 the 2LA inversion in *Anopheles gambiae*. *American Journal of Tropical Medicine and Hygiene*
20 76: 334-339.
- 21 29. Gray EM, Rocca KAC, Costantini C, Besansky NJ (2009) Inversion 2La is associated with enhanced
22 desiccation resistance in *Anopheles gambiae*. *Malaria Journal* 8.
- 23 30. Simard F, Ayala D, Kamdem G, Pombi M, Etouna J, et al. (2009) Ecological niche partitioning
24 between *Anopheles gambiae* molecular forms in Cameroon: the ecological side of speciation.
25 *BMC Ecology* 9: 17.
- 26 31. Besansky NJ, Krzywinski J, Lehmann T, Simard F, Kern M, et al. (2003) Semipermeable species
27 boundaries between *Anopheles gambiae* and *Anopheles arabiensis*: Evidence from multilocus
28 DNA sequence variation. *Proceedings of the National Academy of Sciences of the United States
29 of America* 100: 10818-10823.
- 30 32. Neafsey D, Lawniczak M, Park D, Redmond S, Coulibaly M, et al. (2010) SNP genotyping defines
31 complex gene-flow boundaries among African malaria vector mosquitoes. *Science* 330: 514-517.
- 32 33. White BJ, Cheng C, Sangaré D, Lobo NF, Collins FH, et al. (2009) The population genomics of trans-
33 specific inversion polymorphisms in *Anopheles gambiae*. *Genetics* 183: 275-288.
- 34 34. Peterson BK, Weber JN, Kay EH, Fisher HS, Hoekstra HE (2012) Double digest RADseq: an
35 inexpensive method for de novo SNP discovery and genotyping in model and non-model species.
36 *PLoS one* 7: e37135.
- 37 35. Excoffier L, Dupanloup I, Huerta-Sánchez E, Sousa VC, Foll M (2013) Robust demographic
38 inference from genomic and SNP data. *PLoS genetics* 9: e1003905.
- 39 36. Nosil P, Feder JL (2012) Genomic divergence during speciation: causes and consequences
40 Introduction. *Royal Society Philosophical Transactions Biological Sciences* 367: 332-342.
- 41 37. Harr B (2006) Genomic islands of differentiation between house mouse subspecies. *Genome Research*
42 16: 730-737.
- 43 38. Nadeau NJ, Whibley A, Jones RT, Davey JW, Dasmahapatra KK, et al. (2012) Genomic islands of
44 divergence in hybridizing *Heliconius* butterflies identified by large-scale targeted sequencing.
45 *Philosophical Transactions of the Royal Society B: Biological Sciences* 367: 343-353.
- 46 39. Turner TL, Hahn MW, Nuzhdin SV (2005) Genomic islands of speciation in *Anopheles gambiae*.
47 *Plos Biology* 3: 1572-1578.
- 48 40. Cruickshank TE, Hahn MW (2014) Reanalysis suggests that genomic islands of speciation are due to
49 reduced diversity, not reduced gene flow. *Molecular ecology*.
- 50 41. Beaumont MA, Zhang WY, Balding DJ (2002) Approximate Bayesian computation in population
51 genetics. *Genetics* 162: 2025-2035.

- 1 42. Wegmann D, Leuenberger C, Excoffier L (2009) Efficient approximate Bayesian computation
2 coupled with Markov chain Monte Carlo without likelihood. *Genetics* 182: 1207-1218.
- 3 43. Fontaine MC, Pease JB, Steele A, Waterhouse RM, Neafsey DE, et al. (2015) Extensive introgression
4 in a malaria vector species complex revealed by phylogenomics. *Science* 347.
- 5 44. Lang M, Murat S, Clark AG, Gouppil G, Blais C, et al. (2012) Mutations in the neverland gene turned
6 *Drosophila pachea* into an obligate specialist species. *Science* 337: 1658-1661.
- 7 45. Auton A, McVean G (2007) Recombination rate estimation in the presence of hotspots. *Genome
8 research* 17: 1219-1227.
- 9 46. Kamali M, Xia A, Tu Z, Sharakhov IV (2012) A new chromosomal phylogeny supports the repeated
10 origin of vectorial capacity in malaria mosquitoes of the *Anopheles gambiae* complex. *PLoS
11 pathogens* 8: e1002960.
- 12 47. Guerrero RF, Rousset F, Kirkpatrick M (2012) Coalescent patterns for chromosomal inversions in
13 divergent populations. *Philosophical Transactions of the Royal Society B: Biological Sciences*
14 367: 430-438.
- 15 48. Sharakhov IV, White BJ, Sharakhova MV, Kayondo J, Lobo NF, et al. (2006) Breakpoint structure
16 reveals the unique origin of an interspecific chromosomal inversion (2La) in the *Anopheles
17 gambiae* complex. *Proceedings of the National Academy of Sciences of the United States of
18 America* 103: 6258-6262.
- 19 49. Coluzzi M, Sabatini A, della Torre A, Di Deco MA, Petrarca V (2002) A polytene chromosome
20 analysis of the *Anopheles gambiae* species complex. *Science* 298: 1415-1418.
- 21 50. Lowry DB, Willis JH (2010) A widespread chromosomal inversion polymorphism contributes to a
22 major life-history transition, local adaptation, and reproductive isolation. *Plos Biology* 8: 14.
- 23 51. Joron M, Frezal L, Jones RT, Chamberlain NL, Lee SF, et al. (2011) Chromosomal rearrangements
24 maintain a polymorphic supergene controlling butterfly mimicry. *Nature* 477: 203-206.
- 25 52. Besansky NJ, Powell JR, Caccone A, Hamm DM, Scott JA, et al. (1994) Molecular phylogeny of the
26 *Anopheles gambiae* complex suggests genetic introgression between principal malaria vectors.
27 *Proceedings of the National Academy of Sciences of the United States of America* 91: 6885-6888.
- 28 53. Bhutkar A, Gelbart WM, Smith TF (2007) Inferring genome-scale rearrangement phylogeny and
29 ancestral gene order: a *Drosophila* case study. *Genome Biology* 8.
- 30 54. White BJ, Collins FH, Besansky NJ (2011) Evolution of *Anopheles gambiae* in relation to humans
31 and malaria. *Annual Review of Ecology, Evolution, and Systematics* 42: 111-132.
- 32 55. Aeschbacher S, Bürger R (2014) The effect of linkage on establishment and survival of locally
33 beneficial mutations. *Genetics* 197: 317-336.
- 34 56. Barton N (1995) A general model for the evolution of recombination. *Genetical research* 65: 123-144.
- 35 57. Yeaman S, Whitlock MC (2011) The genetic architecture of adaptation under migration-selection
36 balance. *Evolution* 65: 1897-1911.
- 37 58. Gillies MT, Meillon Bd (1968) The Anophelinae of Africa South Or the Sahara (Ethiopian
38 Zoogeographical Region): South African Institute of Medical Research Johannesburg.
- 39 59. Gillies M, Coetzee M (1987) A Supplement to the Anophelinae of Africa South of the Sahara.
40 Publications of the South African Institute for Medical Research 55: 1-143.
- 41 60. Fanello C, Santolamazza F, Della Torre A (2002) Simultaneous identification of species and
42 molecular forms of the *Anopheles gambiae* complex by PCR-RFLP. *Medical and veterinary
43 entomology* 16: 461-464.
- 44 61. Lobo NF, Sangare DM, Reger AA, Reidenbach KR, Bretz DA, et al. (2010) Breakpoint structure of
the *Anopheles gambiae* 2Rb chromosomal inversion. *Malaria Journal* 9.
- 46 62. Holt RA, Subramanian GM, Halpern A, Sutton GG, Charlab R, et al. (2002) The genome sequence of
47 the malaria mosquito *Anopheles gambiae*. *Science* 298: 129-+.
- 48 63. Li H, Durbin R (2009) Fast and accurate short read alignment with Burrows-Wheeler transform.
49 *Bioinformatics* 25: 1754-1760.
- 50 64. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, et al. (2009) The sequence alignment/map format
51 and SAMtools. *Bioinformatics* 25: 2078-2079.

- 1 65. Catchen JM, Amores A, Hohenlohe P, Cresko W, Postlethwait JH (2011) Stacks: building and
2 genotyping loci de novo from short-read sequences. *G3: Genes, Genomes, Genetics* 1: 171-182.
- 3 66. Hohenlohe PA, Catchen J, Cresko WA (2012) Population genomic analysis of model and nonmodel
4 organisms using sequenced RAD tags. *Data Production and Analysis in Population Genomics:*
5 Springer. pp. 235-260.
- 6 67. Catchen J, Hohenlohe PA, Bassham S, Amores A, Cresko WA (2013) Stacks: an analysis tool set for
7 population genomics. *Molecular ecology* 22: 3124-3140.
- 8 68. Jombart T (2008) adegenet: a R package for the multivariate analysis of genetic markers.
9 *Bioinformatics* 24: 1403-1405.
- 10 69. R Core Team (2012) R: A Language and Environment for Statistical Computing. Vienna, Austria: R
11 Foundation for Statistical Computing.
- 12 70. Jombart T, Devillard S, Balloux F (2010) Discriminant analysis of principal components: a new
13 method for the analysis of genetically structured populations. *BMC genetics* 11: 94.
- 14 71. Falush D, Stephens M, Pritchard JK (2003) Inference of population structure using multilocus
15 genotype data: linked loci and correlated allele frequencies. *Genetics* 164: 1567-1587.
- 16 72. Pritchard JK, Stephens M, Donnelly P (2000) Inference of population structure using multilocus
17 genotype data. *Genetics* 155: 945-959.
- 18 73. Delcher AL, Kasif S, Fleischmann RD, Peterson J, White O, et al. (1999) Alignment of whole
19 genomes. *Nucleic Acids Research* 27: 2369-2376.
- 20 74. Delcher AL, Phillippy A, Carlton J, Salzberg SL (2002) Fast algorithms for large-scale genome
21 alignment and comparison. *Nucleic acids research* 30: 2478-2483.
- 22 75. Keightley PD, Trivedi U, Thomson M, Oliver F, Kumar S, et al. (2009) Analysis of the genome
23 sequences of three *Drosophila melanogaster* spontaneous mutation accumulation lines. *Genome*
24 *research: gr.* 091231.091109.
- 25 76. Navarro A, Betran E, Barbadilla A, Ruiz A (1997) Recombination and gene flux caused by gene
26 conversion and crossing over in inversion heterokaryotypes. *Genetics* 146: 695-709.
- 27 77. Ewing G, Hermission J (2010) MSMS: a coalescent simulation program including recombination,
28 demographic structure and selection at a single locus. *Bioinformatics* 26: 2064-2065.

29

30

31

1 **Figure Legends**

2 Figure 1. Schematic illustration of study design and population demographic scenarios of
3 *Anopheles gambiae* and *An. arabiensis*. a) procedures involved in the detection of targets of
4 selection; b) in collinear regions, *An. gambiae* and *An. arabiensis* diverged from a common
5 ancestor at time T_{div} with a population size N_e . They experienced recent population expansion at
6 time T_{exp} to the current population size N_{cur} . The two species have constant gene flow since
7 divergence; c) in regions with alternative arrangements, the arrangement, S , split with the
8 alternative arrangement, I , at time T_{IS} ; at time T_{int} , introgression occurred so that I split into two
9 lineages, *An. gambiae* and *An. Arabiensis*, respectively. Alternative arrangements within *An.*
10 *gambiae* have reduced recombination rate (r), while within the same arrangements I , the
11 introgression rate (m) is the same as that in collinear regions. Parameters that maintain the same
12 in inverted and collinear regions are shown in gray, while inversion specific parameters are
13 shown in black.

14

15 Figure 2. Outlier analysis of inverted region scan. a) and b), dots represent F_{ST} measures of each
16 RADtag locus between S and I chromosomes in *An. gambiae* populations along the region.
17 Shades of yellow show quantiles of 25%, 50%, 75%, 95%, 99% respectively, of simulated values
18 of divergence measures under reconstructed demographic histories with region-adjusted
19 recombination rate. c) and d) empirical distributions of F_{ST} between individuals with S and I
20 chromosomes on inverted region (green) or collinear region (red). a) and c), 2La regions; b) and
21 d), 2Rb regions;

22

23 Figure 3. Candidate loci and regions under selection. All the dots represent RADtag loci that are
24 classified to have experienced selection in the S lineage (red) or I lineage (blue). Y-axis shows
25 the posterior probability of such classifications. Dots with solid color are the ones with larger
26 than 0.9 in posterior probabilities of assignment. Bars represent regions that have been classified
27 as either experienced selection in the S lineage (red) or I lineage (blue). Width of each bar
28 corresponds to 50kb in our analysis.

29 Figure 4. Comparison of detection power and false-discovery rates between two selection
30 detection methods. Upper panel, 2La region; lower panel, 2Rb region. a), c), e), and g) are the
31 results for the scenario in which selection occurs on the branch of S . b, d), f), and h) are the
32 results for selection on the branch of I . The x -axis presents cases where the proportion of selected
33 loci consists of 1%, 5%, and 10% out of all simulated loci (i.e., rest of the loci are simulated
34 under neutral scenario). All selected loci are generated with $s = 0.01$. See Fig. S5&6 for cases
35 where $s = 0.001$ & 0.0001 . See methods for the details of how rates are calculated. The inset
36 shows the false-positive rates for each method.

37

38

1 **Supporting Information Captions**

- 2 S1 Text. Detailed results on population genetic structures of collinear and inverted regions.
- 3 Figure S1. Sampling locations and species composition of *Anopheles gambiae* species complex.
4 The area of each pie chart correspond to the sample size. Map color from blue to red stands for
5 humid to dry areas.
- 6 Figure S2. Principle component analyses using SNPs in different collinear genomic regions.
7 Color of the dots represent different populations. Red dots are individuals of *An. arabiensis*.
- 8 Figure S3. Principle component analyses of *An. gambiae* using SNPs in different collinear
9 genomic regions. Color of the dots represent different populations.
- 10 Figure S4. Principle component analyses of *Anopheles gambiae* using SNPs from 2La and 2Rb.
11 Left and Right panels are the result for 2La and 2Rb, respectively. Top panel is the result for
12 PCA clustering of individuals from different populations. Middle panel finds the best number of
13 clusters based on BIC scores. The bottom panel shows how divergent each cluster is from each
14 other on the discriminant function space.
- 15 Figure S5. Comparison of detection power and false-discovery rates between two selection
16 detection methods. Upper panel, 2La region; lower panel, 2Rb region. a), c), e), and g) are the
17 results for the scenario in which selection occurs on the branch of *S*. b, d), f), and h) are the
18 results for selection on the branch of *I*. The x-axis presents cases where the proportion of selected
19 loci consists of 1%, 5%, and 10% out of all simulated loci (i.e., rest of the loci are simulated
20 under neutral scenario). All selected loci are generated with $s = 0.001$. See methods for the
21 details of how rates are calculated. The inset shows the false-positive rates for each method.
- 22 Figure S6. Comparison of detection power and false-discovery rates between two selection
23 detection methods. Upper panel, 2La region; lower panel, 2Rb region. a), c), e), and g) are the
24 results for the scenario in which selection occurs on the branch of *S*. b, d), f), and h) are the
25 results for selection on the branch of *I*. The x-axis presents cases where the proportion of selected
26 loci consists of 1%, 5%, and 10% out of all simulated loci (i.e., rest of the loci are simulated
27 under neutral scenario). All selected loci are generated with $s = 0.0001$. See methods for the
28 details of how rates are calculated. The inset shows the false-positive rates for each method.
- 29 S1 Table. Collection sites, coordinates and sampling sizes of *Anopheles gambiae*.
- 30 S2 Table. Power of differentiating selection form drift using average summary statistics of loci in
31 a 5kb segment with/without selected locus inside.
- 32
- 33

1 **Tables**

2 Table 1. Estimations of population genetic and demographic parameters using region-specific
3 SFS implemented in FASTSIMCOAL2. All parameters were estimated relative to a fixed prior of
4 *An. gambiae* ($Ne = 750,000$; see Material and Methods for details). See Fig. 1b, c for
5 explanations of parameters. Note that the estimations of time are in the unit of generations.

Regions	Parameters	Point estimates	Relative to Ne	95% Confidence Interval	
Collinear	Na_{cur}	2,101,300	2.80	1,126,650	33,035,200
	Ng_{cur}	5,374,400	7.17	4,578,413	10,547,442
	Na	338,900	0.45	304,350	406,750
	m	1.75E-07		1.52E-08	2.70E-07
	Ta_{exp}	186,800	0.25	136,900	228,700
	Tg_{exp}	278,500	0.37	256,600	290,200
	T_{div}	1,052,500	1.40	824,600	1,240,400
2La	r	7.99E-08		9.25E-08	1.08E-07
	T_{int}	340,200	0.45	320,700	352,000
	T_{IS}	2,353,800	3.14	2236100	2568100
2Rb	r	4.79E-07		4.64E-07	5.15E-07
	T_{int}	591,800	0.79	570,100	625,200
	T_{IS}	1091300	1.46	916,600	1,140,500

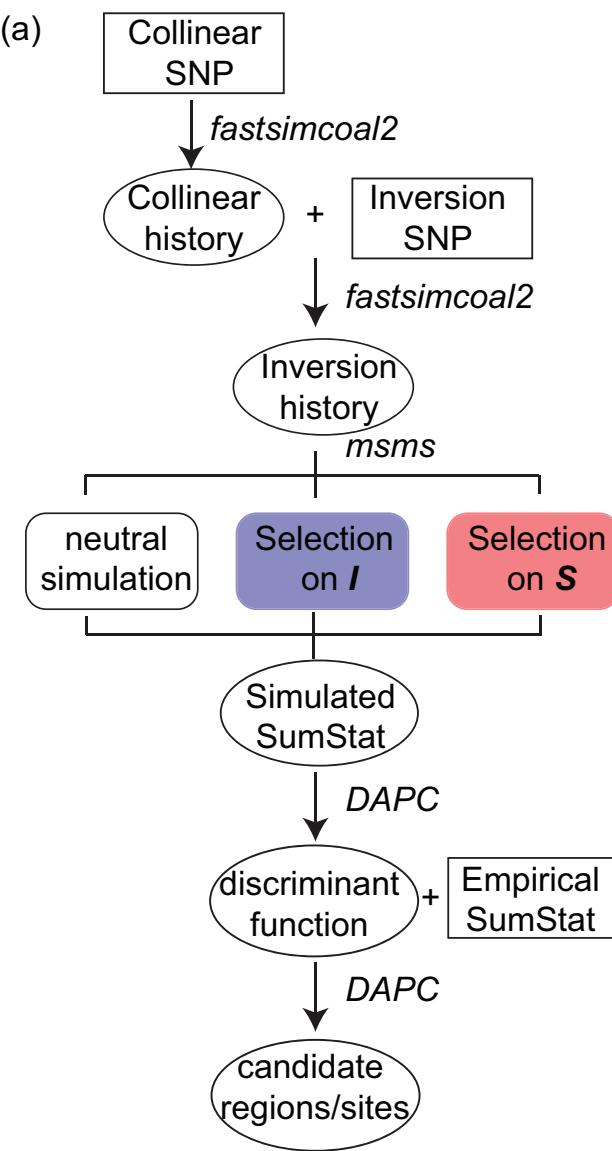
6

7

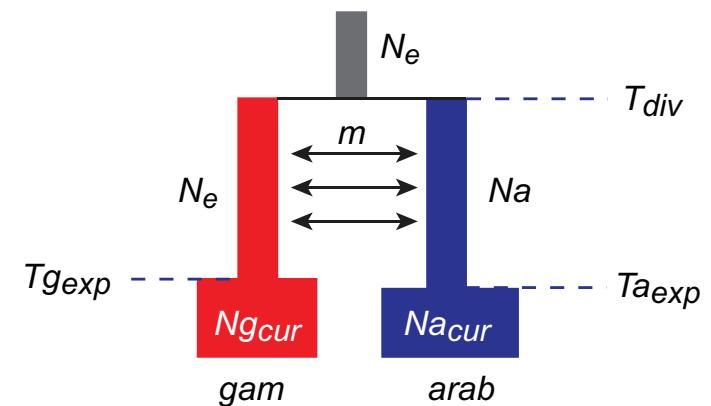
8

Fig. 1

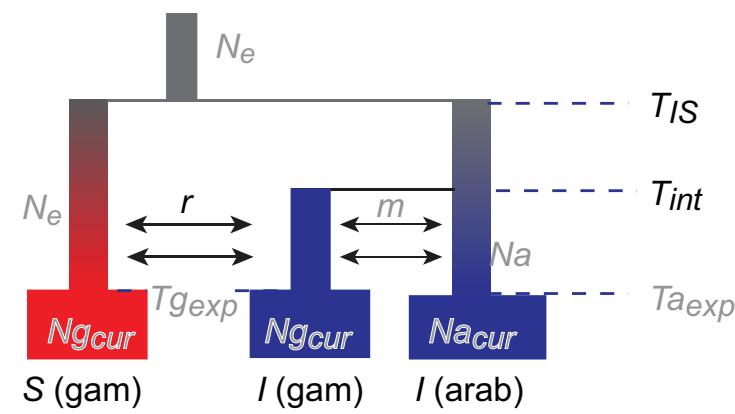
[Click here to download Figure: Fig1.eps](#)



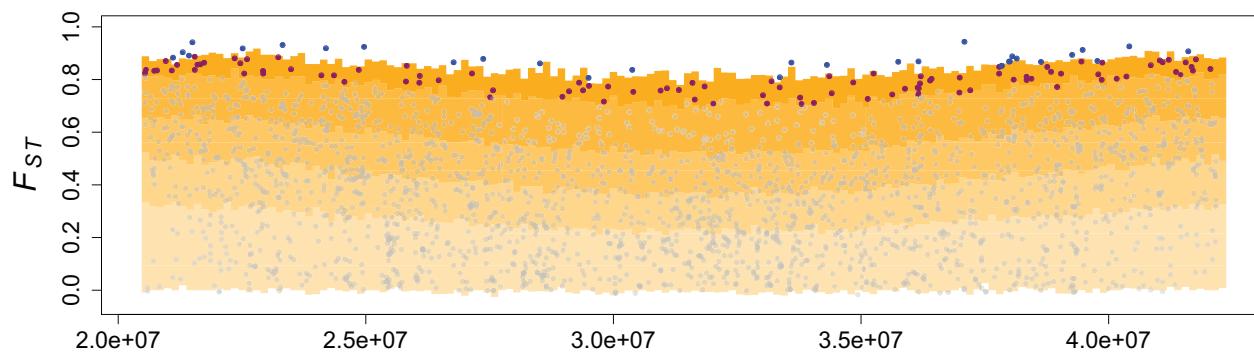
(b) Collinear region demographic history



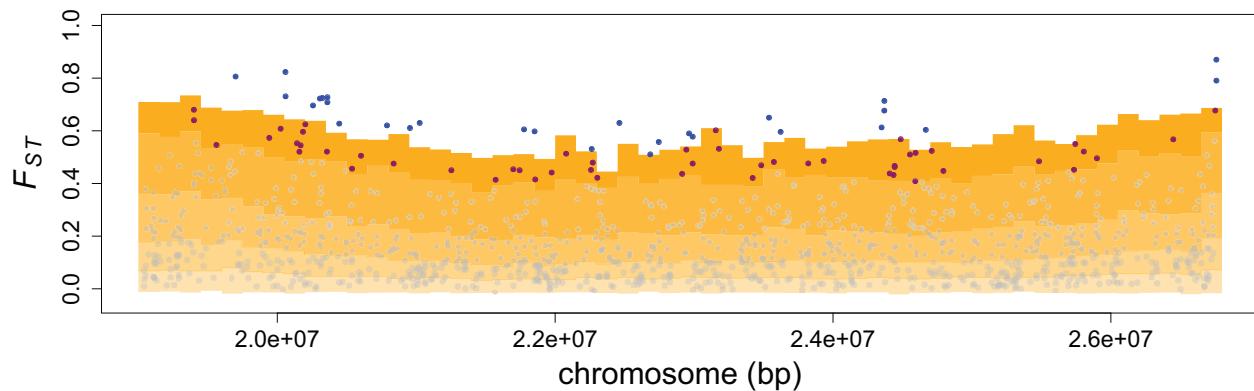
(c) Inversion region demographic history



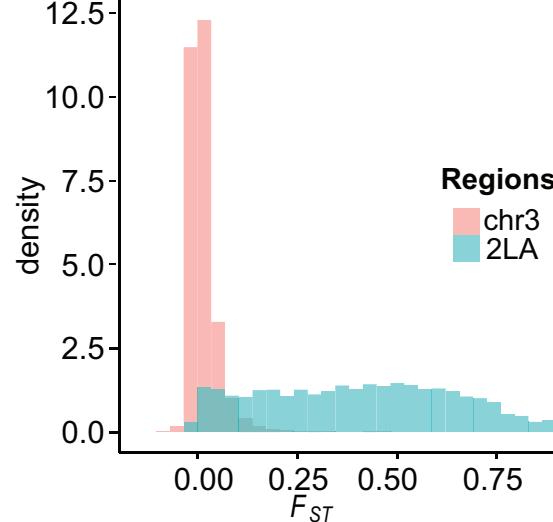
(a)

Chromosome 2La region

(b)

Chromosome 2Rb region

(c)



(d)

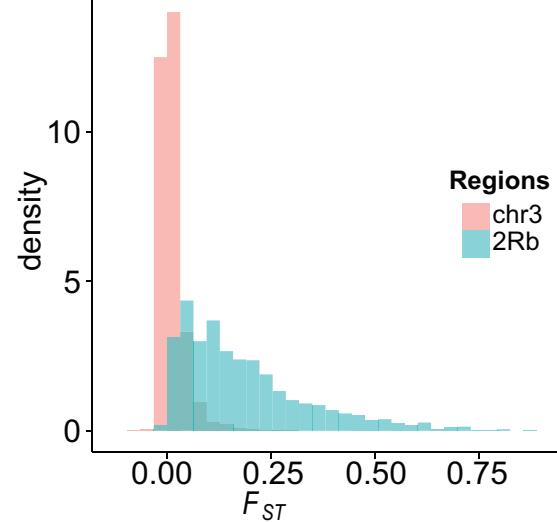
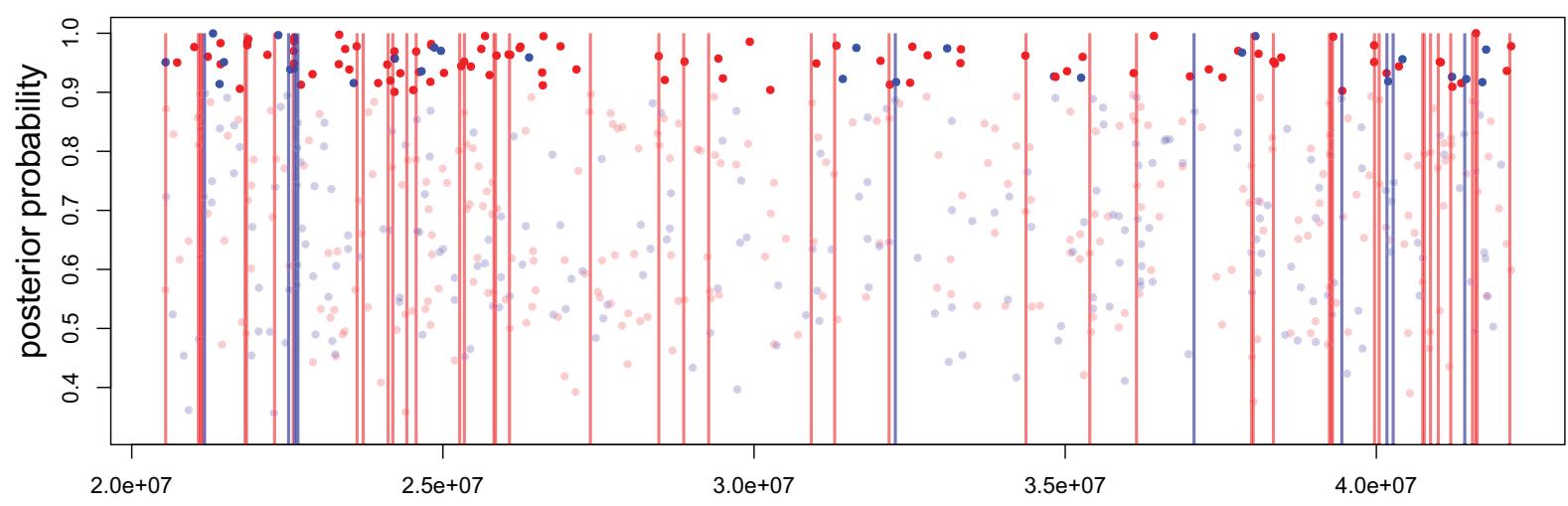


Fig.3

[Click here to download Figure: Fig3.eps](#)

selection on 2La



selection on 2Rb

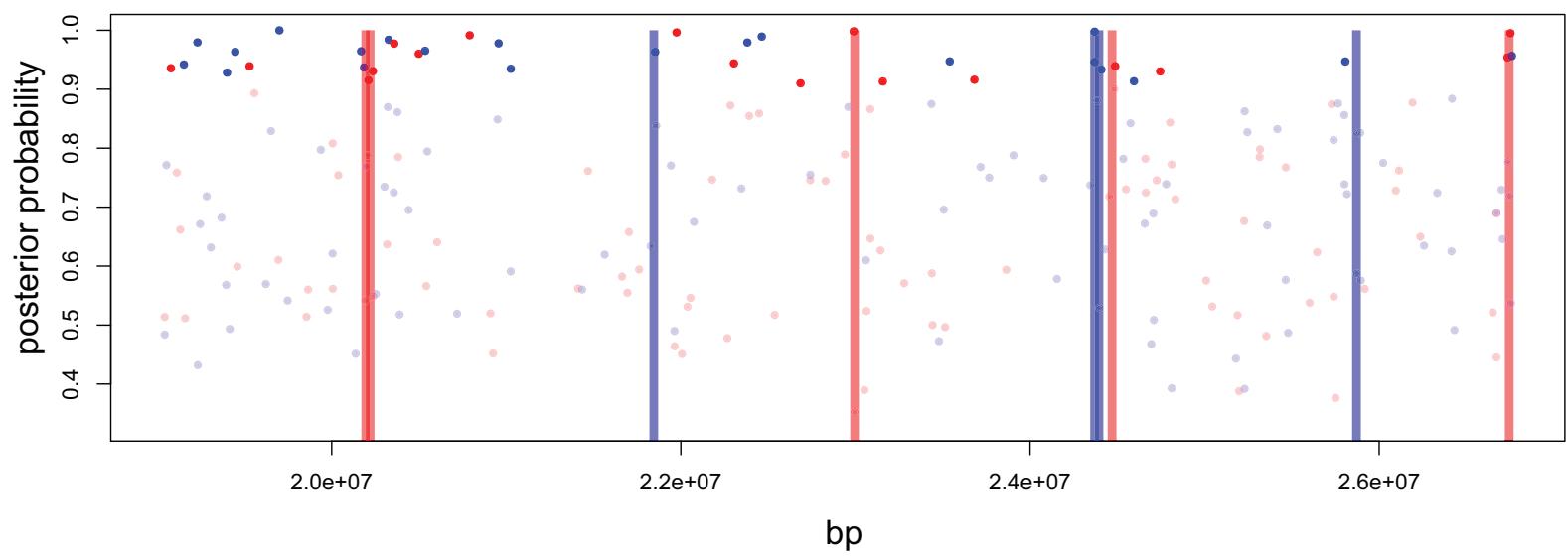
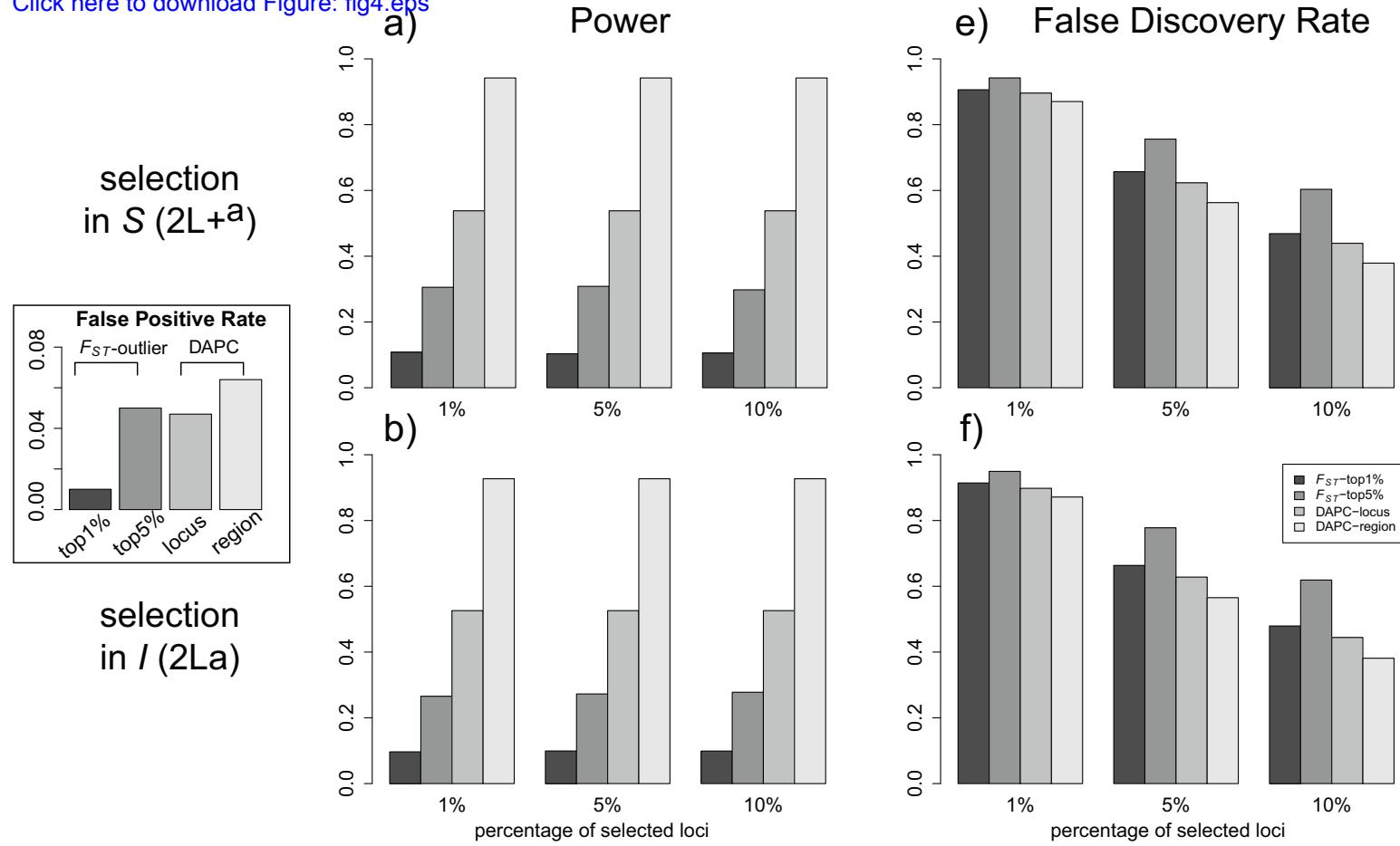


Fig.4

[Click here to download Figure: fig4.eps](#)

2La



2Rb

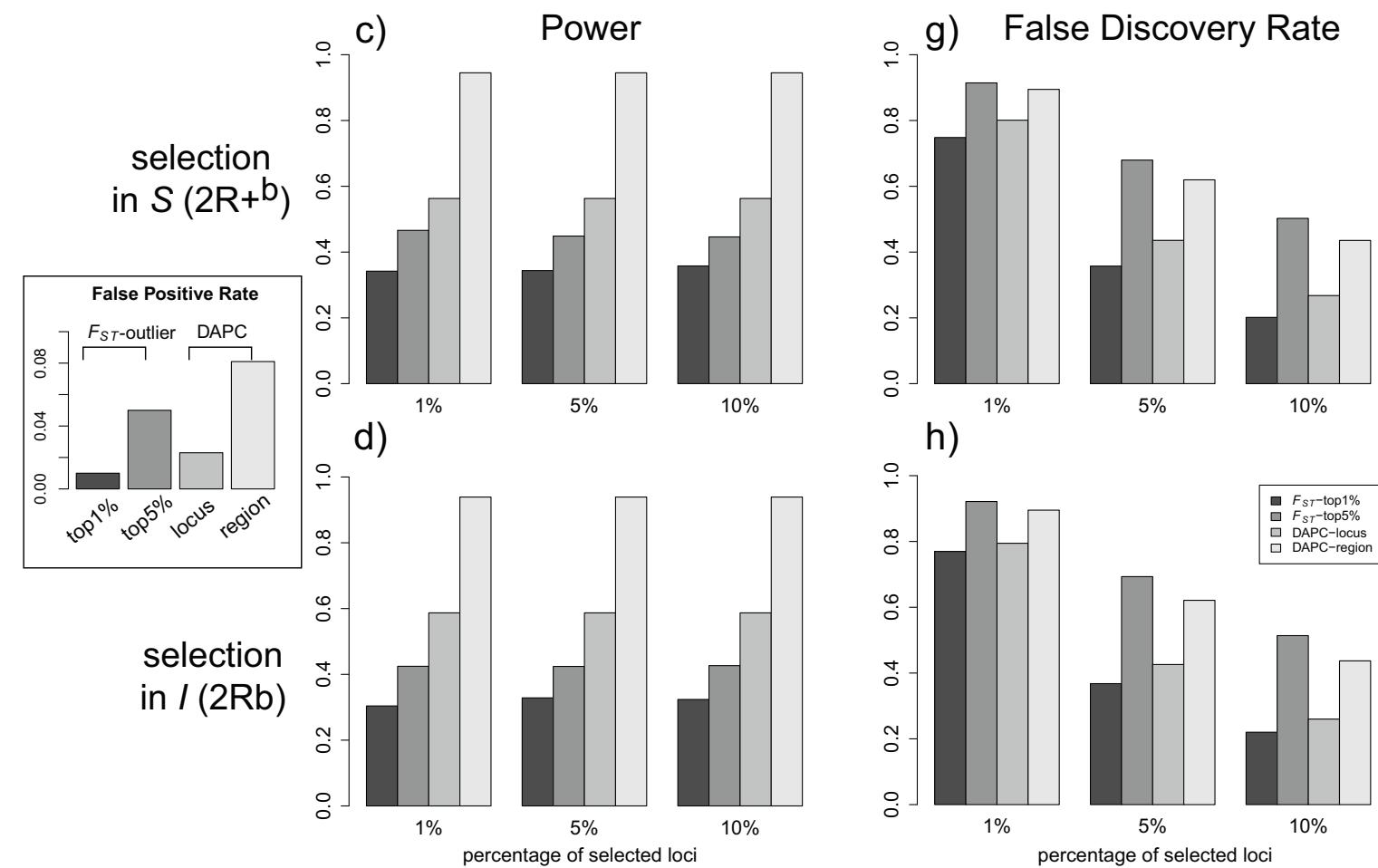


Fig. S1

[Click here to download Figure: FigS1.eps](#)

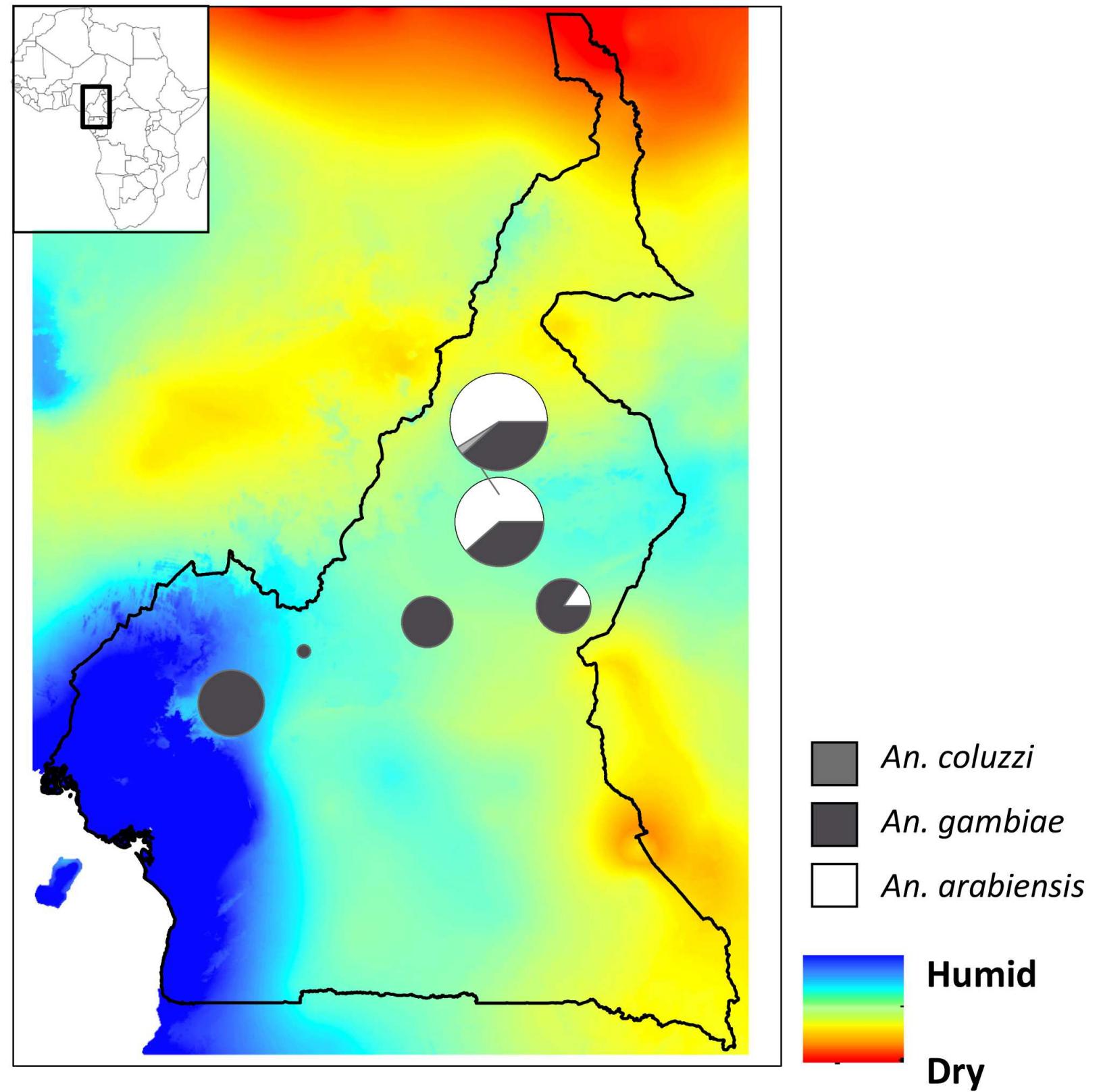
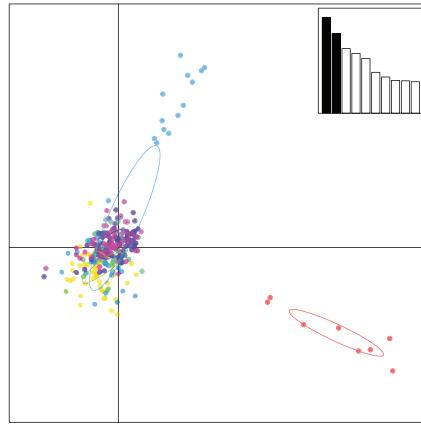


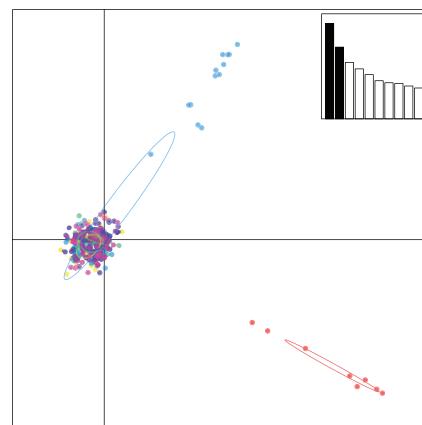
Fig. S2

[Click here to download Figure: FigS2.eps](#)

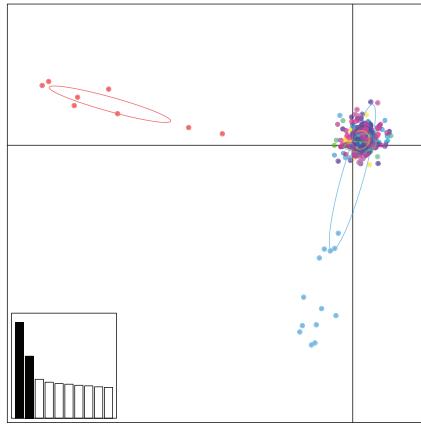
(a) 2L-collinear



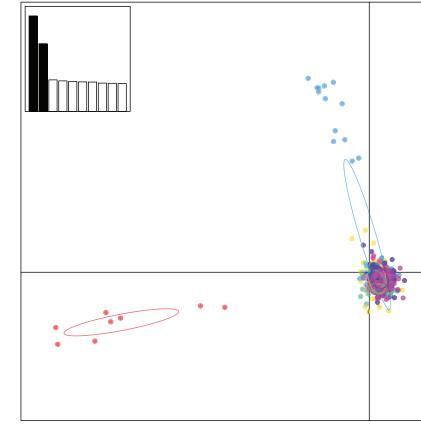
(b) 2R-collinear



(c) 3L



(d) 3R



(e) X

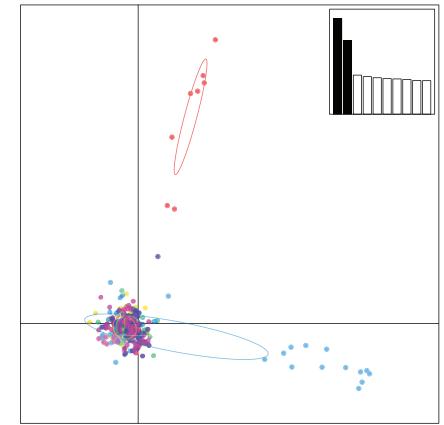
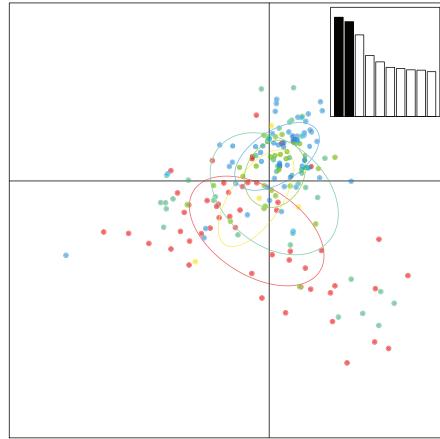


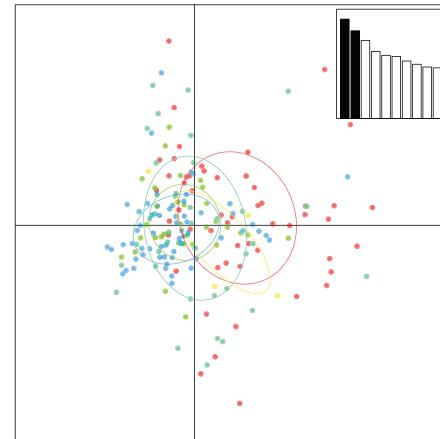
Fig. S3

[Click here to download Figure: FigS3.eps](#)

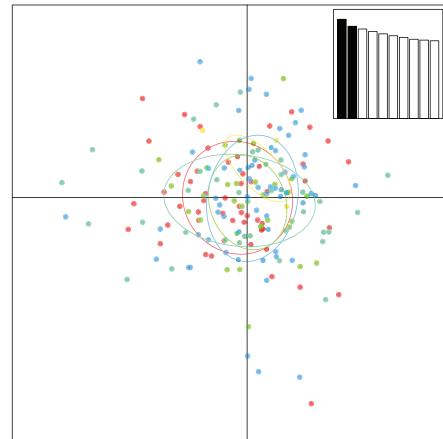
(a) 2L-collinear



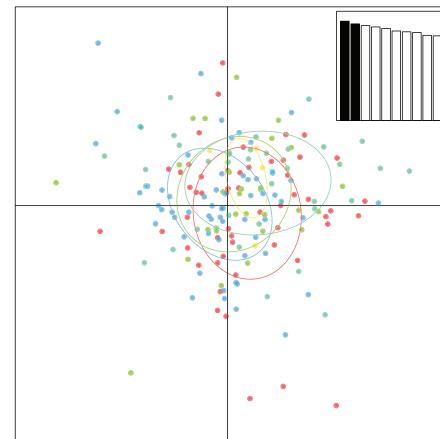
(b) 2R-collinear



(c) 3L



(d) 3R



(e) X

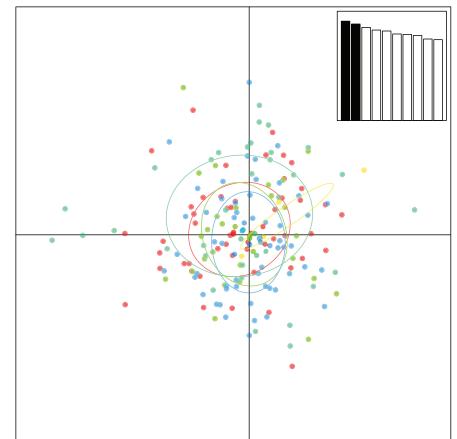
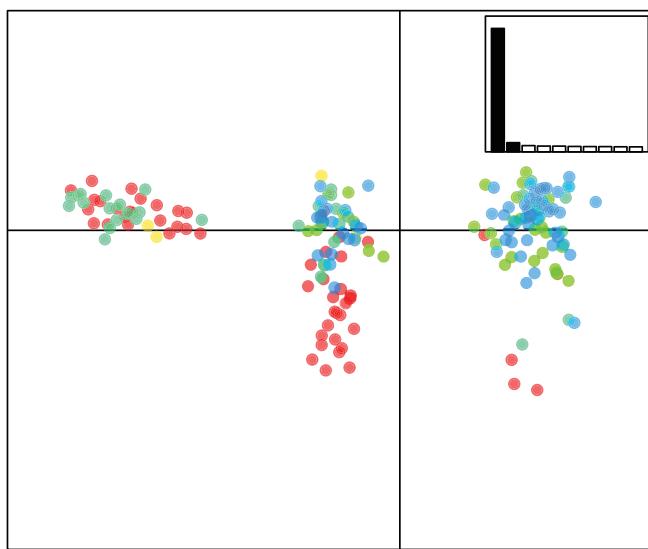


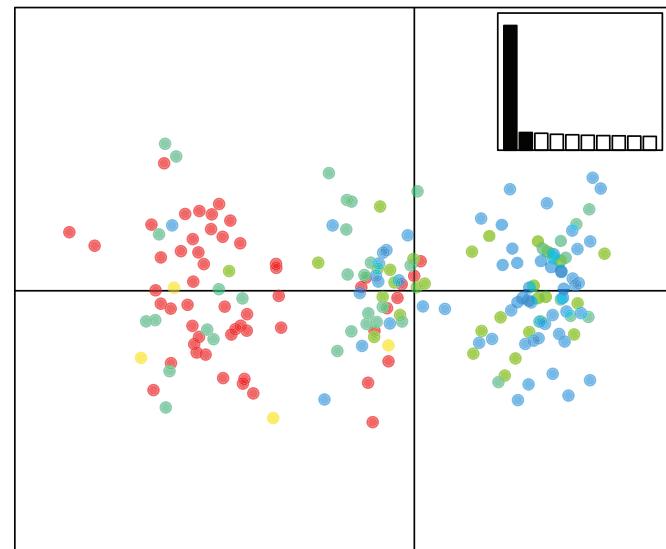
Fig. S4

[Click here to download Figure: FigS4.eps](#)

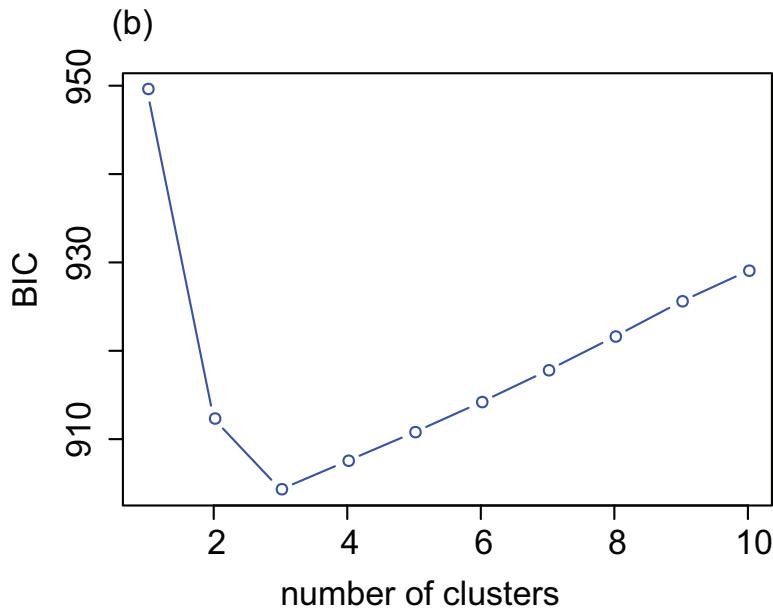
(a)



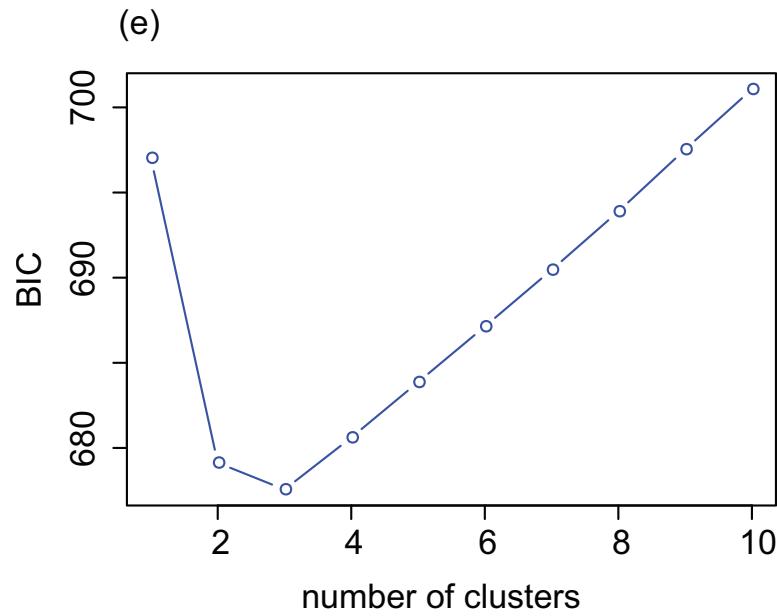
(d)



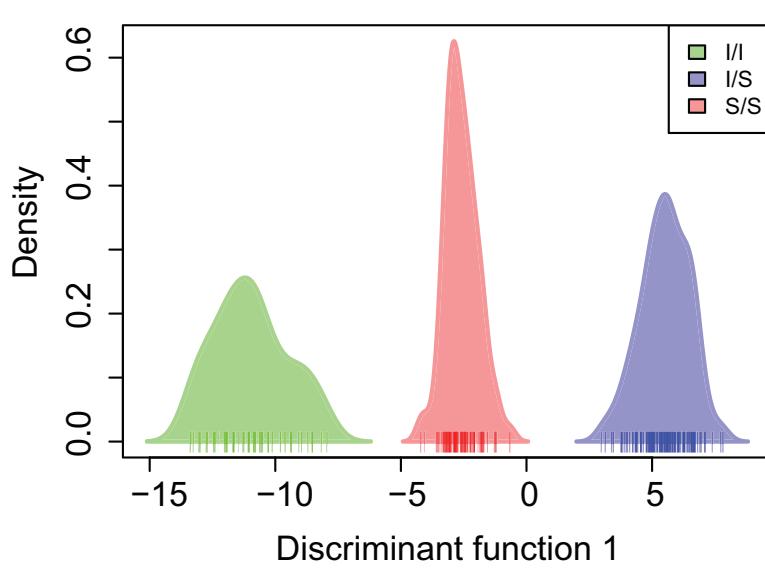
(b)



(e)



(c)



(f)

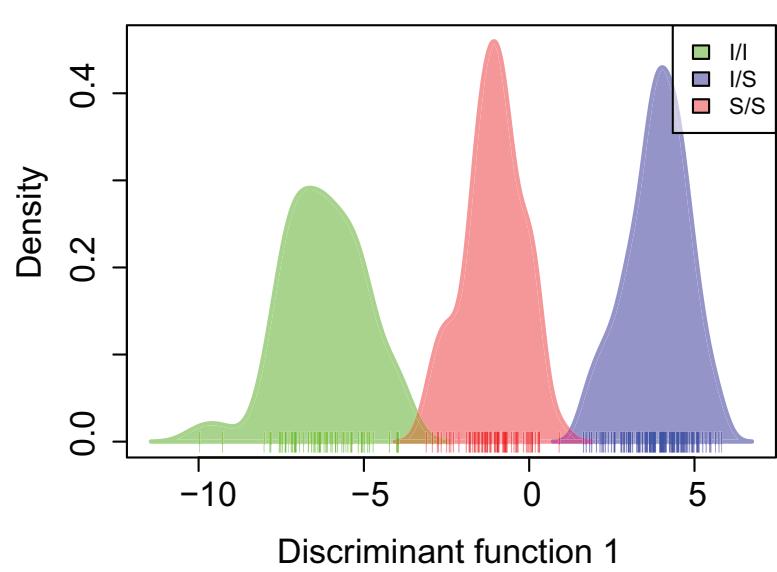
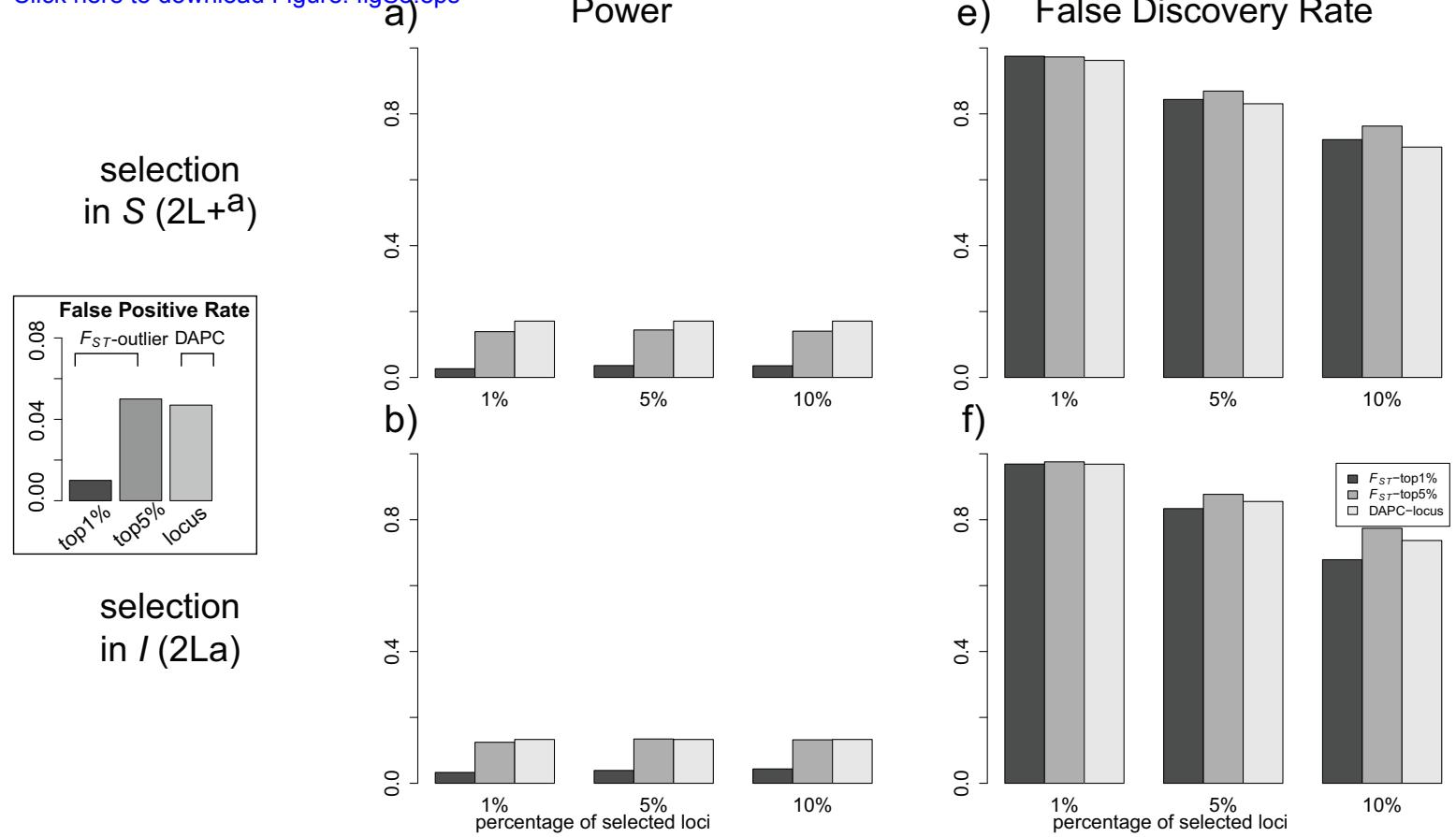


Fig. S5

[Click here to download Figure: figS5.eps](#)

2La



2Rb

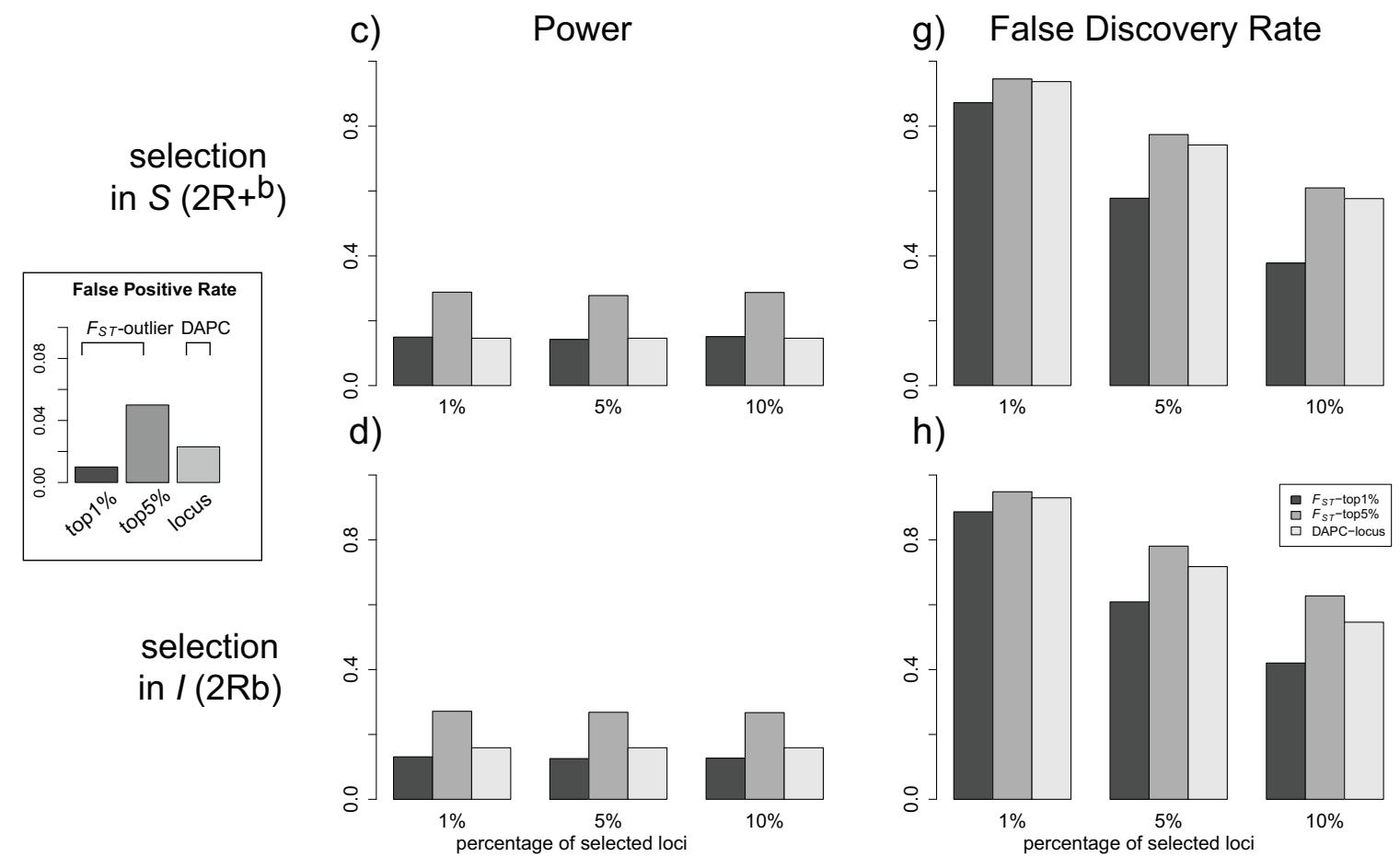


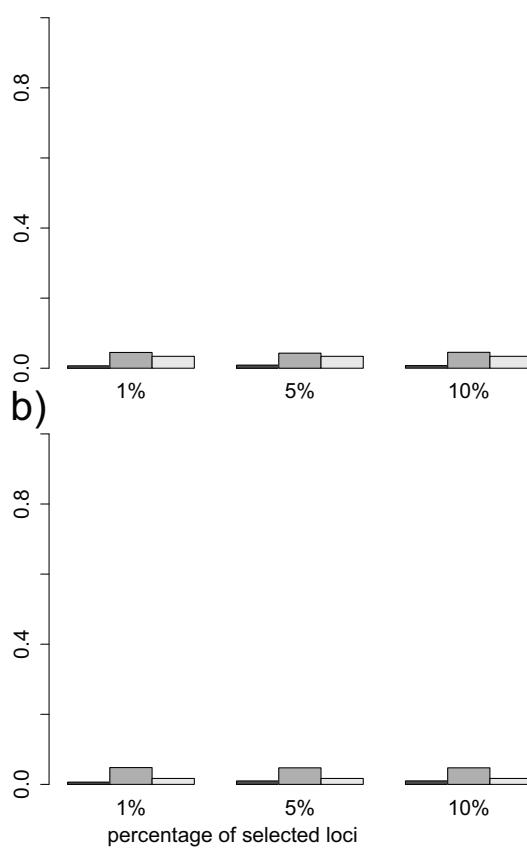
Fig. S6

[Click here to download Figure: figS6.eps](#)

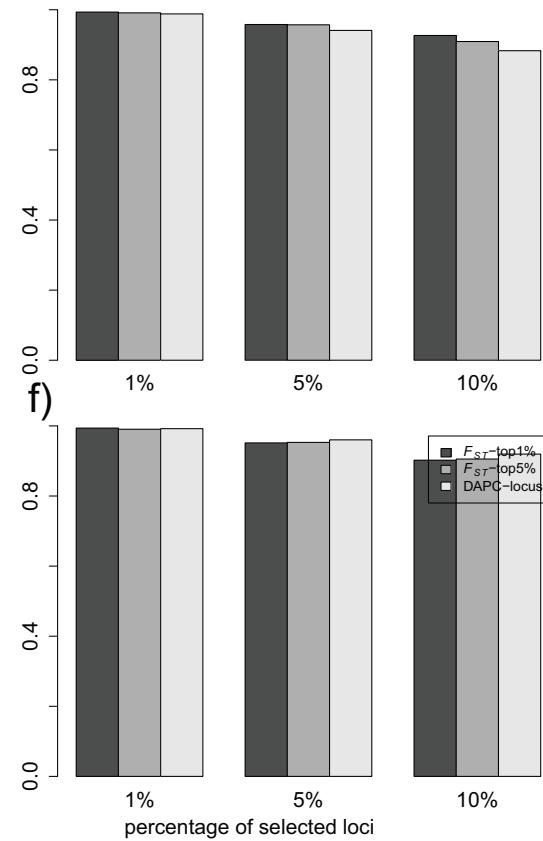
2La

Power

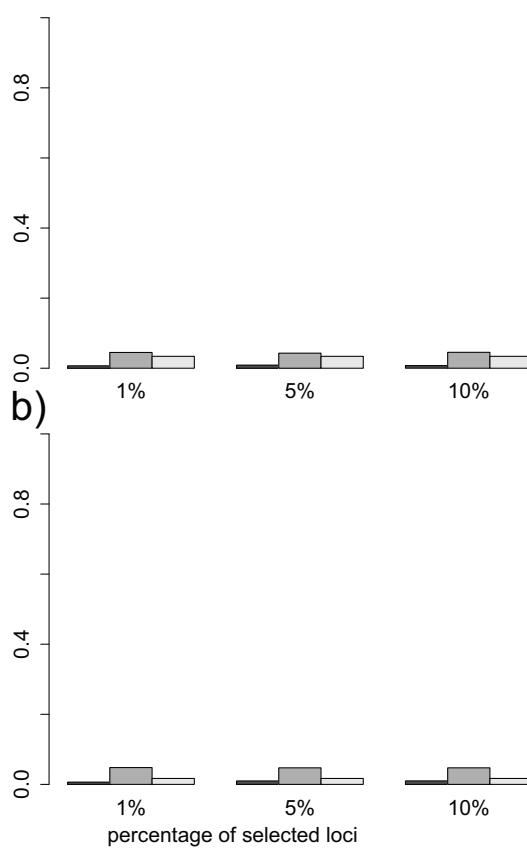
selection
in *S* (2L+^a)



e) False Discovery Rate



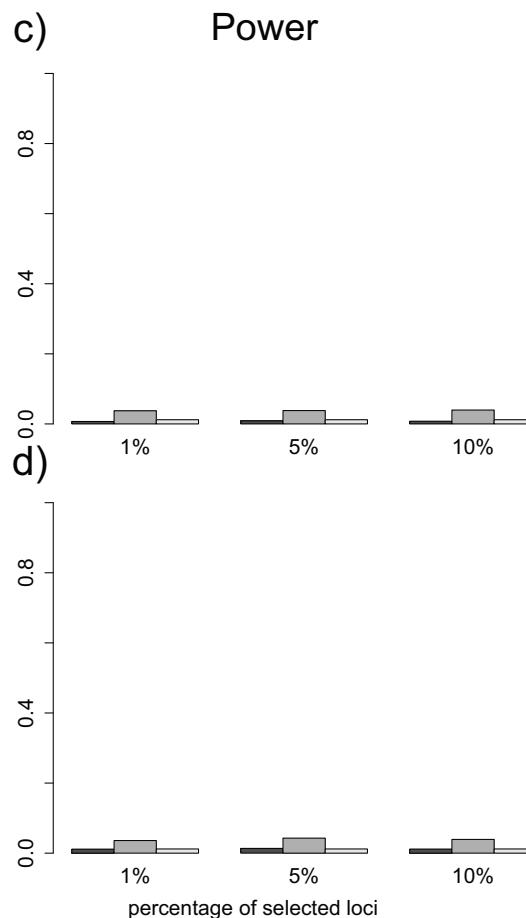
selection
in *I* (2La)



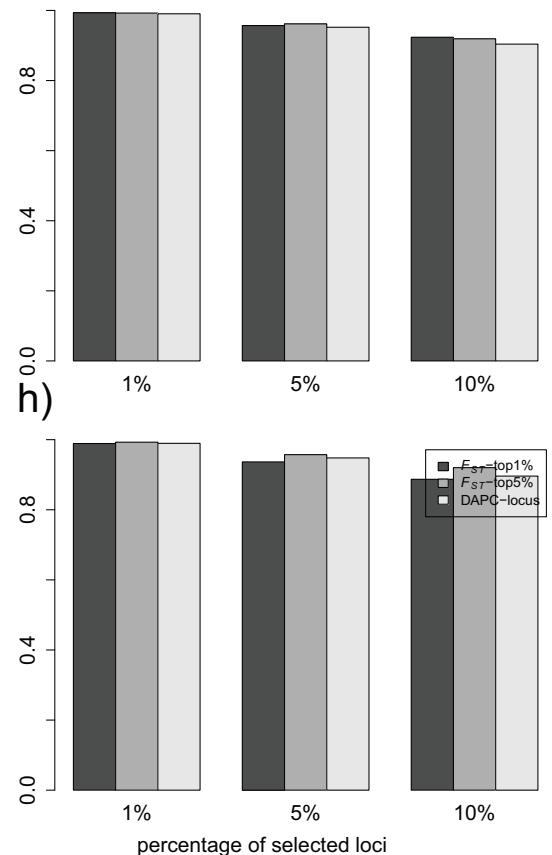
2Rb

Power

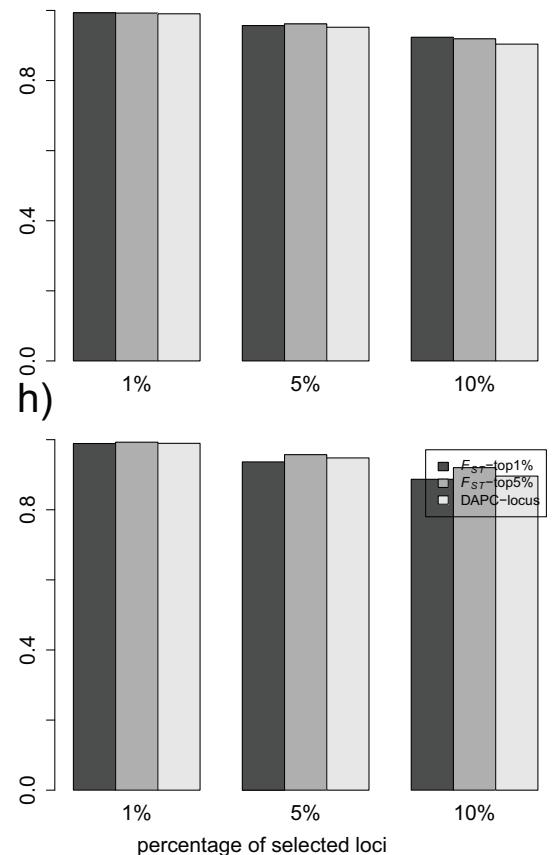
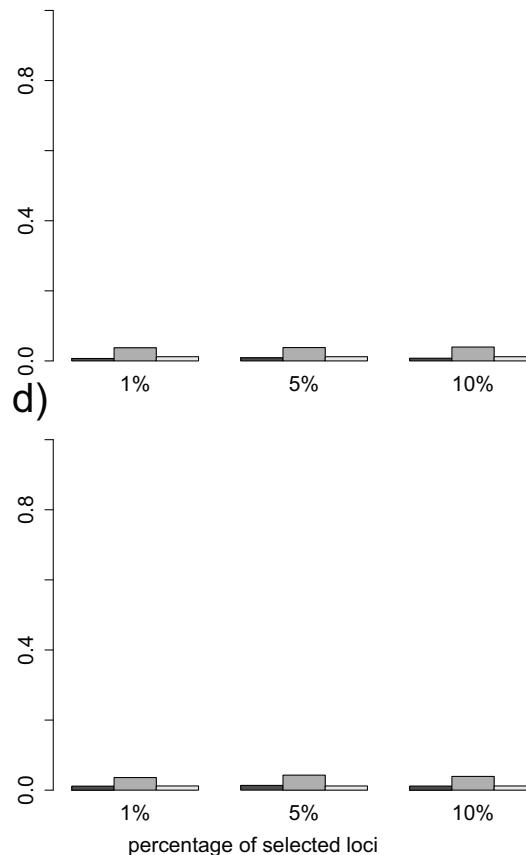
selection
in *S* (2R+^b)



g) False Discovery Rate



selection
in *I* (2Rb)



Supporting Information

[Click here to download Supporting Information: QXH&LLK_2014_SuppleText.docx](#)