



HAL
open science

Detecting sex-linked genes using genotyped individuals sampled in natural populations

Jos Käfer, Nicolas Lartillot, Gabriel a B Marais, Franck Picard

► **To cite this version:**

Jos Käfer, Nicolas Lartillot, Gabriel a B Marais, Franck Picard. Detecting sex-linked genes using genotyped individuals sampled in natural populations. *Genetics*, 2021, 218 (2), 10.1093/genetics/iyab053 . hal-03250679

HAL Id: hal-03250679

<https://hal.science/hal-03250679>

Submitted on 4 Jun 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Detecting sex-linked genes using genotyped individuals sampled in natural populations

Jos Käfer^{*,1}, Nicolas Lartillot^{*}, Gabriel A.B. Marais^{*,†,2} and Franck Picard^{*,2}

^{*}Université de Lyon; Université Lyon 1; CNRS; UMR 5558; Laboratoire de Biométrie et Biologie Evolutive; F-69622, Villeurbanne, France, [†]Current address: LEAF- Linking Landscape, Environment, Agriculture and Food, Instituto Superior de Agronomia, Universidade de Lisboa, Portugal

ABSTRACT We propose a method, SDpop, able to infer sex-linkage caused by recombination suppression typical of sex chromosomes. The method is based on the modeling of the allele and genotype frequencies of individuals of known sex in natural populations. It is implemented in a hierarchical probabilistic framework, accounting for different sources of error. It allows statistical testing for the presence or absence of sex chromosomes, and detection of sex-linked genes based on the posterior probabilities in the model. Furthermore, for gametologous sequences, the haplotype and level of nucleotide polymorphism of each copy can be inferred, as well as the divergence between them. We test the method using simulated data, as well as data from both a relatively recent and an old sex chromosome system (the plant *Silene latifolia* and humans), and show that, for most cases, robust predictions are obtained with 5 to 10 individuals per sex.

KEYWORDS Sex chromosomes; population genomics; probabilistic inference; hierarchical model

1 Introduction

Sex chromosomes, which are found in many species with genetic sex determination, are key elements of the biology of the sexes (e.g. between-sexes differences in development, morphology, physiology) and are consequently studied by different branches of biology (e.g. genetics, genomics, developmental biology, physiology, evolutionary biology, medical research, agronomy). While sex chromosomes share striking features even between independently evolved systems, there is also much diversity (reviewed in [Bachtrog et al. 2014](#)).

Sex chromosomes are thought to evolve from a pair of autosomes ([Lahn and Page 1999](#); [Charlesworth et al. 2005](#)). One sex is characterized by a pair of sex chromosomes for which recombination is suppressed over part of the length, while the other sex has two identical chromosomes that recombine normally. Thus, the sex chromosome which is limited to one of the two sexes contains a part that never recombines, and evolves independently from the homologous part on the other sex chromosome. These

systems are termed according to the sex that is heterogametic; when the males are the heterogametic sex, the sex chromosomes are named X and Y, and when the females are heterogametic, they are named Z and W; Y and W are the Y or W chromosomes in these systems, respectively. The recombining part of the sex chromosomes is called the pseudo-autosomal region.

The complete lack of homologous recombination in a part of a chromosome has dramatic consequences for its evolution ([Charlesworth et al. 2005](#); [Lynch 2007](#)): it accumulates deleterious mutation (Muller's ratchet) and transposable elements, which deteriorate gene function and finally chromosome integrity. Once most genes have lost their function, parts of the chromosome can be lost without negative effects on fitness. One typical outcome of this evolutionary process is a large X or Z chromosome and a small Y or W chromosome with only few functional genes, as in humans. However, in many sex chromosome systems, the Y or W chromosome does not differ significantly in size from its homologous counterpart, and recombination is suppressed in only a small region of these chromosomes ([Charlesworth 2016](#); [Muyle et al. 2017](#)). As more sex chromosome systems are being studied in phylogenetically distant lineages, more questions arise on their similarities and differences (e.g. [Charlesworth 2016](#); [Muyle et al. 2017](#); [Ponnikas et al. 2018](#)).

Even though sequencing technologies have been improving constantly, sex chromosomes have remained difficult to study for a long time. In particular, Y or W chromosome sequences

doi: 10.1534/genetics.XXX.XXXXXX

Manuscript compiled: Friday 29th January, 2021

²These authors contributed equally to this work.

¹Corresponding author: Jos Käfer, Laboratoire de Biométrie et Biologie Evolutive; CNRS; UMR 5558; Université Lyon 1; Bat. Grégor Mendel, 43 bd du 11 novembre 1918, 69622 Villeurbanne cedex, E-mail: jos.käfer@univ-lyon1.fr.

1 used to be scarce for two reasons (Tomaszkiewicz *et al.* 2017).
2 First, genome projects have been typically using whole-genome
3 shotgun (WGS), in which DNA is fragmented and sequenced
4 using short read technologies and then re-assembled in contigs.
5 The assembly of the non-recombining regions of the Y or W
6 chromosome, which are full of repeats, using short reads is
7 virtually impossible (Hughes and Rozen 2012). Second, in the
8 heterogametic sex, read coverage for the sex chromosomes is
9 half the coverage for the autosomes, making the assembly of
10 both sex chromosomes more difficult, while sequencing of the
11 homogametic sex only yields one of the sex chromosomes.

12 Recently, the development of new methods to study sex
13 chromosomes has boosted their study. The first Y sequence
14 of humans and other species (e.g. papaya) have been obtained
15 with specific methods based on bacterial artificial chromosomes
16 (BACs), which require intensive laboratory work (Skaletsky *et al.*
17 2003; Hughes and Rozen 2012; Wang *et al.* 2012; Bellott *et al.* 2014).
18 Although costs have decreased and methods have become more
19 standardized (e.g. Bellott *et al.* 2018), they are still substantial
20 for any Y or W chromosome larger than a few megabases, and
21 their application remains limited to a few model species. More
22 efficient and cheaper methods have thus been developed to po-
23 tentially detect and study sex chromosomes in many non-model
24 species (reviewed in Muyle *et al.* 2017; Palmer *et al.* 2019).

25 A first alternative approach was to compare male and a fe-
26 male genomes obtained using short reads (WGS). Here, absence
27 in the female genome is enough to identify the Y or W con-
28 tigs (Carvalho and Clark 2013; Cortez *et al.* 2014). Then, the X
29 or Z contigs can be identified by computing the read depth in
30 male and female: the read depth ratio for the heterogametic
31 vs. homogametic sex is expected to be 1 for autosomes and 0.5
32 for the X or Z, as the heterogametic sex has one copy and the
33 homogametic sex two (Vicoso and Bachtrog 2011). Variants of
34 this approach have been widely used (e.g. Vicoso *et al.* 2013a,b;
35 Hall *et al.* 2013; Li *et al.* 2018). However, this approach works
36 only with strongly diverged, old sex chromosomes, in which
37 the reads map uniquely to one of the chromosomes but not to
38 both. In more recently evolved, weakly diverged systems, the
39 sex chromosomes still share a high level of sequence similar-
40 ity, and reads cannot be unambiguously assigned to one of the
41 chromosomes with alignment methods.

42 To study weakly divergent sex chromosomes, a second cate-
43 gory of empirical methods uses segregation patterns in genotyp-
44 ing data from natural populations. An excess of heterozygotes
45 is expected in the heterogametic sex compared to the homog-
46 ametic sex for the sex-linked genes. Such patterns could be iden-
47 tified using male versus female F_{ST} (e.g. Zhou *et al.* 2017) or
48 differences in heterozygote frequencies (e.g. Picq *et al.* 2014; Kirk-
49 patrick and Guerrero 2014). GWAS using sex as the studied
50 trait is another option which will work based on different allele
51 frequencies in the sexes, but a precise model of the expected as-
52 sociation is preferable (Galichon *et al.* 2012). Indeed, the patterns
53 these empirical approaches look for can be caused by other pro-
54 cesses than sex linkage, such as sex-antagonistic selection on the
55 pseudo-autosomal part of the sex chromosome (Qiu *et al.* 2013;
56 Kirkpatrick and Guerrero 2014). Such methods thus have high
57 false positive rates and either must use conservative cut-offs,
58 which significantly reduces their sensitivity, or are applicable
59 only to species with well-assembled genomes (e.g. Picq *et al.*
60 2014; Mathew *et al.* 2014). An empirical approach somewhat
61 intermediate between the former and the latter categories is to
62 identify sequences specific to the Y or W chromosome in the

63 genomic reads (RAD-tags, k-mers), but this only gives a first
64 characterization of the sex chromosome system, and further se-
65 quencing is needed to identify the X or Z copies and study the
66 sex-linked regions (Scharmann *et al.* 2017; Torres *et al.* 2018).

67 Recently, the SEX-DEtector statistical framework, relying on
68 the modeling of the transmission of alleles that account for the
69 observed genotypes, was introduced (Muyle *et al.* 2016). The
70 method has been shown to yield high power and high sensitivity
71 with as few as five offspring of each sex, and has successfully
72 been used on a number of species (Muyle *et al.* 2018; Veltsos *et al.*
73 2019; Martin *et al.* 2019; Prentout *et al.* 2020; Badouin *et al.* 2020;
74 Fruchard *et al.* 2020). It can be used for several purposes, using
75 a single dataset: testing whether a species has sex chromosomes
76 and of which kind (XY or ZW), identifying sex-linked genes,
77 reconstructing the haplotypes of the gametolog copies, and es-
78 timating expression of each of the copies if RNAseq data are
79 used. However, its requirement to produce a controlled cross
80 for sequencing limits its use to species than can be easily bred or
81 cultivated in controlled conditions, and hinders its application
82 to species with long generation times.

83 A method that can be applied to sequencing data from in-
84 dividuals sampled in natural populations, DETSEX, was intro-
85 duced by Gautier (2014). While it offers some attractive features,
86 like automatic sexing of some individuals in the sample, it relies
87 on ploidy levels to detect sex-linked sequences, and can thus
88 only be used for sufficiently divergent sex chromosomes. The
89 method has, to our knowledge, not been applied in any study
90 (apart from the original publication).

91 We here introduce a new hierarchical framework based on
92 the modeling of genotype and allele frequencies in a popula-
93 tion, according to autosomal and sex-linked segregation types.
94 As SEX-DEtector statistically analyses Mendelian segregation
95 and its deviations due to sex-linkage, the new method does so
96 based on the Hardy-Weinberg equilibrium and its deviations.
97 The method we present here is “SEX-DEtector in a population”,
98 and is thus termed “SDpop”. It can be applied to any sam-
99 ple collected from natural populations, provided the sex of the
100 individuals can be determined morphologically. Despite the
101 similarity of its goals with those of SEX-DEtector (notably, ap-
102 plicability to sex chromosome systems of any age, statistical
103 testing, likelihood-based inferences, prediction of gametolog
104 haplotypes), it is an entirely different model due to the data
105 it handles, that can be obtained by directly sampling in natu-
106 ral populations. Thus, the underlying genetic expectations are
107 different, as well as the population genetic predictions it can
108 produce.

109 Materials and Methods

110 Model

111 **Terminology** A sex-linked gene (i.e. a gene that is present on
112 the non-recombining region of the sex chromosomes) has two
113 independent copies. These copies will accumulate mutations,
114 be subject to selection and drift, and accumulate allele fixations
115 independently from one another. These copies thus diverge after
116 the suppression of recombination, in a way that is reminiscent of
117 the divergence of paralogous genes after gene duplication; the
118 two copies of a sex-linked gene are thus termed “gametologs”
119 (Garcia-Moreno and Mindell 2000). Eventually, due to neo- or
120 sub-functionalization, or the accumulation of deleterious muta-
121 tions, one of the copies can get lost, most often the copy present
122 on the Y or W chromosome. The remaining gene, on the X or Z,
123 is termed “X-” or “Z-hemizygous”.

1 Gametology and hemizyosity cause differences in the rela- 60
2 tion between genotype and allele frequencies with respect to 61
3 autosomal segregation. More precisely, these differences occur 62
4 in one sex, i.e. the heterogametic one. If one considers a sex- 63
5 ual, random mating population of diploid individuals, most 64
6 autosomal genes will be at Hardy-Weinberg equilibrium. This 65
7 results in equilibria between allele and genotype frequencies 66
8 that are different from the Hardy-Weinberg equilibrium in the
9 heterogametic sex. These can be described assuming panmixia,
10 and are the basis of the present model. For gametologs, two
11 independent copies are present in the heterogametic sex, while
12 for hemizygous genes, only one copy is present; for both sex-
13 linked segregation types, equilibria are thus different between
14 the sexes. For non-sex-linked genes, the presence of two inde-
15 pendent copies can also occur for both sexes under paralogy,
16 and the presence of only one copy for the haploid mitochondrial
17 and chloroplastic genes, or when only one allele is expressed in
18 transcriptome data. We thus distinguish five possible segrega-
19 tion types in a diploid population: diploid autosomy (hereafter
20 just called "autosomy"), haploidy, paralogy, hemizyosity, and
21 gametology.

22 We model genotypes that are obtained by mapping sequenc-
23 ing data on a reference sequence. As the organisms are supposed
24 to be diploid, the input data at all positions are diploid geno-
25 types, and deviations from diploidy will be detected afterwards.
26 We consider polymorphisms consisting of two alleles; these
27 can be single nucleotide polymorphisms (SNPs) or structural
28 variants (indels and length polymorphisms). For paralogs and
29 gametologs, which are the result of the mapping of two non-
30 recombining copies, we assume that cases in which both copies
31 are polymorphic at the same position and for the same two alle-
32 les are very rare and can be neglected. Indeed, such cases would
33 arise for sites that had polymorphisms that existed at the time
34 when recombination suppression (or gene duplication) occurred,
35 and under a neutral model, these would become rapidly fixed
36 (in less than $4N_e$ generations).

37 **Genotype frequencies** We adopt a hierarchical probabilistic
38 framework in which the distribution of alleles is modeled given
39 a segregation type for each polymorphism. A technical presen-
40 tation and derivations are given in the Appendix; here, we
41 describe the general principles of the model.

42 Sex-linkage produces different genotype-allele equilibria for
43 each sex. We consider sites with two alleles, a and b , and three
44 possible genotypes $g' \in \{aa, ab, bb\}$. The genotype frequencies
45 are denoted $p(g')$ if they are equal for both sexes, and indexed
46 with σ and φ symbols when different. In the following, the allele
47 frequency f is the frequency of allele a , unless otherwise stated.

48 For **autosomal segregation**, we can simply write

$$49 \quad p(aa) = f^2 \quad ; \quad p(ab) = 2f(1-f) \quad ; \quad p(bb) = (1-f)^2,$$

50 which is the Hardy-Weinberg equilibrium.

51 **Haploid genes** correspond to haploid genotypes, but geno-
52 typing will have given apparently diploid genotypes. No sex-
53 specific difference is expected, so

$$54 \quad p(aa) = f \quad ; \quad p(ab) = 0 \quad ; \quad p(bb) = 1-f.$$

55 **Paralogy** is caused by the mapping of the reads of two more
56 or less recently duplicated genes on one locus in the reference.
57 Thus, paralogous genes have tetraploid genotypes, which again
58 are considered diploid by the genotyper. There is no recombina-
59 tion between the copies, that thus evolve independently. Only

60 biallelic sites are modeled here; for simplicity, we assume that
61 one of the copies is fixed for one of the alleles. Thus, if one copy
62 is fixed for allele b , the tetraploid genotypes are $aabb$, $abbb$ and
63 $bbbb$; the former two will yield the diploid genotype ab . Again,
64 the allele frequency f is the frequency of allele a in the polymor-
65 phic copy. No difference between the sexes is expected, and we
66 thus obtain the diploid genotype frequencies

$$67 \quad p(aa) = 0 \quad ; \quad p(ab) = f^2 + 2f(1-f) \quad ; \quad p(bb) = (1-f)^2.$$

68 A choice has to be made about which allele to consider fixed,
69 as this is not always clear *a priori*. The previous segregation
70 types, autosomal and haploid, are symmetrical with respect this
71 choice, i.e., the expected genotype frequencies are strictly identi-
72 cal whether the frequency of allele a was used or the frequency
73 of allele b . The paralogous segregation type is asymmetrical with
74 respect to this choice: considering allele a or allele b as fixed does
75 not lead to identical genotype frequencies. The details about the
76 calculations are specified in the Appendix.

77 For **hemizyously segregating genes**, the members of one
78 sex are haploid while the others are diploid. Thus, in one sex,
79 the Hardy-Weinberg equilibrium is expected, while in the other,
80 the expectations are the same as for the haploid segregation type.
81 In the case of an XY system, expectations are thus

$$82 \quad p_{\varphi}(aa) = f^2 \quad ; \quad p_{\varphi}(ab) = 2f(1-f) \quad ; \quad p_{\varphi}(bb) = (1-f)^2$$

$$83 \quad p_{\sigma}(aa) = f \quad ; \quad p_{\sigma}(ab) = 0 \quad ; \quad p_{\sigma}(bb) = 1-f.$$

84 **Gametologous segregation** is characterized by the presence
85 of two independent copies in one sex, and one copy in the other.
86 As for the paralogous case, we assume that an allele is fixed in at
87 least one of the copies. We have to distinguish two cases: either
88 the X copy is polymorphic, or the Y copy. The fraction of the sites
89 in gametologous genes for which the X copy is polymorphic is
90 ρ , which is a parameter of the model. In either case (X and Y
91 polymorphism), we have to choose which allele we consider
92 fixed on one of the copies, thus leading to four different allele-
93 genotype equilibria. Here, one case of X-polymorphism and one
94 case of Y-polymorphism are described; details about the four
95 equilibria are given in the Appendix.

96 In the case of polymorphism on the X, female genotypes
97 are modeled by the Hardy-Weinberg equilibrium, while male
98 genotypes show a specific deviation. If we consider allele b to
99 be fixed on the Y chromosome, and define f as the frequency of
100 allele a on the X chromosome, we get

$$101 \quad p_{\varphi}(aa) = f^2 \quad ; \quad p_{\varphi}(ab) = 2f(1-f) \quad ; \quad p_{\varphi}(bb) = (1-f)^2$$

$$102 \quad p_{\sigma}(aa) = 0 \quad ; \quad p_{\sigma}(ab) = f \quad ; \quad p_{\sigma}(bb) = 1-f.$$

103 For a polymorphism on the Y chromosome, considering allele
104 b to be fixed on the X-chromosome and defining f to be the
105 frequency of allele a on the Y chromosome yields

$$106 \quad p_{\varphi}(aa) = 0 \quad ; \quad p_{\varphi}(ab) = 0 \quad ; \quad p_{\varphi}(bb) = 1$$

$$107 \quad p_{\sigma}(aa) = 0 \quad ; \quad p_{\sigma}(ab) = f \quad ; \quad p_{\sigma}(bb) = 1-f.$$

108 The X-hemizygous and XY segregation described here can easily
109 be converted to Z-hemizygous and ZW segregation (see Ap-
110 pendix for details).
111
112
113

Hierarchical structure of the model and likelihood function SDpop relies on a hierarchical probabilistic model for posterior inference of segregation types at each polymorphic site. The input of SDpop consists in the observed genotype (g) frequencies at each biallelic site. For a given individual, the observed genotype can differ from the (unobserved) true genotype g' with probability given by an error model e , thus specifying the conditional probability $p(g|g')$. Then we introduce the segregation type, modeled a latent variable S with multinomial prior distribution of parameter vector π . Finally, the conditional probabilities of the true genotypes under each segregation type, $p(g')$, are obtained through the population genetic equations described above, using the empirical allele frequencies as a plugin estimator of the true unknown allele frequencies in the population.

Conditional on the segregation type, the probability of the observed genotype is obtained by summing over all possible true genotypes:

$$p(g|S) = \sum_{g'} p(g'|S) p(g|g')$$

In turn, the marginal probability of the observed genotype is obtained by summing over all segregation types:

$$p(g) = \sum_S p(g|S) \pi(S)$$

Finally, taking the product over all sites gives the likelihood, as a function of the parameters of the model (π, ρ, e).

SDpop can investigate several models of sexual systems: absence of sex-linkage, sex-linkage of the XY type, sex-linkage of the ZW type, and sex-linkage of both types. Apart from the sex-linkage, models can be run with or without haploid and paralogous segregation types. The XY and ZW model include both hemizygous and gametologous segregation types.

Parameter estimation and model comparison The parameters of the model are estimated by maximum likelihood, using the Expectation-Maximization (EM) algorithm. The likelihood score of the model (when parameter value convergence has been attained) is then used to calculate the Bayesian Information Criterion (BIC), with the number of polymorphic sites in the data as the number of observations, which allows comparison of models with different numbers of parameters. For a given model, the segregation type at each polymorphic site is inferred using empirical Bayes posterior probabilities, $p(S|g)$. For the gene or contig, we calculate the geometric mean of the site likelihoods (cf Nelson 2017), which are readjusted to sum to one. As the fraction π of sex-linked SNPs can be very small, we do not use it as a prior, but use equal priors instead. This procedure yields contig-wise scores that can be interpreted as probabilities, and one can use a threshold value to assign contigs to segregation types. All the technical derivations of SDpop are provided in the Appendix.

Haplotypes and population genetics inferences For the sequences inferred as gametologous, several statistics can be calculated from the optimized values. The posterior probabilities of the segregation subtypes indicate which one of the copies (X or Y) is polymorphic, and which one of the alleles is fixed in one of the copies. We can thus infer the frequency of the alleles in the X and Y copies (\hat{f}_X and \hat{f}_Y), which is the mean of the empirical allele frequencies, weighted by the posterior probability for each subtype. These allow to reconstruct the X and Y haplotypes, by

using only alleles that are fixed (or nearly so) in the X and Y copies.

Furthermore, the inferred allele frequencies allow to calculate the nucleotide diversity of the X and Y copies (π_X and π_Y) and divergence between both gametolog copies (D_{XY}) as the means of the diversity and the divergence over the whole contig:

$$\pi_X = \langle 2\hat{f}_X (1 - \hat{f}_X) \rangle; \quad \pi_Y = \langle 2\hat{f}_Y (1 - \hat{f}_Y) \rangle;$$

$$D_{XY} = \langle \hat{f}_X (1 - \hat{f}_Y) + \hat{f}_Y (1 - \hat{f}_X) \rangle$$

Implementation and availability SDpop is written in C with some C++ functionalities. It is available on <https://gitlab.in2p3.fr/sex-det-family/sdpop>, which also contains a user manual and several helper programs.

Simulations

Sets of 10000 contigs of n individuals per sex (i.e., $2n$ individuals in total) were generated using *ms* (Hudson 2002), that simulates samples drawn from a population under the neutral Wright-Fisher model of genetic variation. Samples of autosomal contigs were generated by simulating $4n$ haploid sequences and combining them arbitrarily into $2n$ diploid samples. The parameter $\theta = 4N_e\mu = 4N_euL$ gives the average number of segregating sites. Here, μ is the overall mutation rate for the sequence with length L and per base mutation rate u . The heterozygosity rate (per site) is thus $H = \frac{\theta/L}{1+\theta/L}$, which is equal to the level of polymorphism π . Here, we performed simulations with sequence length $L = 500$ base pairs and a level of polymorphism $\pi = 0.002$, we thus used the parameter $\theta = \frac{\pi L}{1-\pi} = 1.002$.

Gametologous contigs were generated from a sample of $3n$ X-linked sequences and n Y-linked sequences, that were simulated by supposing two populations that split at time t from a population with size N_e into one population of size $\frac{3}{4}N_e$ (X-linked sequences) and one of size $\frac{1}{4}N_e$ (Y-linked sequences). In *ms*, the time t is scaled by $4N_e$, such that, according to coalescent theory, $t = 1$ is the average time to fixation for a neutral mutation. X-hemizygous sequences were simulated similarly, except that no Y-linked sequences were simulated; male “diploid” genotypes were obtained from the haploid X sequence for each individual. We used a constant number of contigs (10000), with fractions ranging from 0.1% to 10% of sex-linked contigs. All of these sex-linked contigs were XY gametologs, or half of those were transformed into X-hemizygous contigs.

Errors were introduced after the coalescent simulations, such that each individual at each position had a probability e to have an erroneous genotype. Two types of errors were introduced: either homozygous genotypes were rendered heterozygous, or heterozygous genotypes homozygous, both with the same rate. As there are typically many more homozygous than heterozygous genotypes, many more errors occur on monomorphic sites than on polymorphic sites. SDpop estimates parameters using polymorphic sites only, and in this subset of all sites, the fraction of sites with errors is higher than in the whole genome, which includes monomorphic sites as well. Thus, the dataset of polymorphic sites on which SDpop runs is enriched in sites with errors; we refer to the error rate among polymorphic sites as the “effective error rate”, e_e . As an example, for a series of simulations with 2, 5, 10, 20, 50, and 100 individuals per sex, which were all conducted with an error rate of 0.0001 and a level of polymorphism (π) of 0.002, e_e is 0.015, 0.011, 0.0088, 0.0068, 0.0045, and 0.0030, respectively.

1 The four models of SDpop (without sex chromosomes, with
 2 XY chromosomes, with ZW chromosomes, and with both XY and
 3 ZW chromosomes) were run on each simulation. The best model
 4 was chosen based on the BIC values. The contigs were classified
 5 as sex-linked or not based on SDpop's posterior probabilities of
 6 the contigs, using a posterior probability threshold of 0.8, and
 7 this classification was used to count the number of true posi-
 8 tives (TP), false positives (FP), and false negatives (FN). Power
 9 and precision of the classification of contigs as sex-linked or not
 10 is quantified by the true positive rate ($TPR = TP / (TP + FN)$)
 11 and the positive predictive value ($PPV = TP / (TP + FP)$). Nu-
 12 cleotide polymorphism and divergence were calculated using
 13 the observed allele frequencies before errors were added.

14 *Silene latifolia*

15 We here use data from Muyle *et al.* (2020) comprising transcrip-
 16 tome sequencing of 34 individuals of *S. latifolia* collected all
 17 over Europe. These data were mapped on a *de novo* transcrip-
 18 tome assembled for the SEX-DETECTOR analysis by Muyle *et al.*
 19 (2016), using the same mapping (BWA, Li and Durbin 2009)
 20 and genotyping (reads2snp, Gayral *et al.* 2013) procedures. We
 21 assessed population structure by performing a Principal Component
 22 Analysis (PCA) of the genotyping data using PLINK
 23 version 1.9 (Chang *et al.* 2015).

24 We calculated true and false positives (TP, FP) as well as false
 25 negatives (FN) of SDpop's assignment based on two datasets.
 26 First, we used SEX-DETECTOR's results published by Muyle *et al.*
 27 (2016), that were inferred using sequencing data of a family
 28 (mother, father, and five offspring of each sex). We considered a
 29 gene as XY gametologous or X-hemizygous whenever the poste-
 30 rior probability of the corresponding segregation type exceeded
 31 0.8, not using SEX-DETECTOR's filtering for the presence of at
 32 least one SNP without genotyping errors in the gene. To infer
 33 positives in SDpop's inferences of sex-linkage, we used a similar
 34 criterion as for SEX-DETECTOR, i.e. the posterior probabilities
 35 of the XY and the X-hemizygous segregation types should be
 36 higher than 0.8. All other contigs were classed as negatives.

37 Furthermore, we placed the contigs on the genetic map of
 38 Papadopulos *et al.* (2015), by identifying the scaffolds from their
 39 genome assembly that had the best blast hits (blastn with stan-
 40 dard parameters; Altschul *et al.* 1990). We considered as sex-
 41 linked the genes located in the non-recombining region of the
 42 X chromosome (cf Krasovec *et al.* 2020). It was not possible
 43 to distinguish between XY gametologous and X-hemizygous
 44 genes, thus we considered genes as sex-linked in SDpop's out-
 45 put whenever the sum of the posterior probabilities for these
 46 two segregation types exceeded 0.8.

47 Human data

48 Exome-targeted sequence data for Finnish individuals (Auton
 49 *et al.* 2015) were downloaded from the 1000 genomes project
 50 website. To facilitate file handling in the analysis pipeline, only
 51 individuals for which all reads were present in two fastq files
 52 (forward and reverse) were retained, i.e. a total of 66 individ-
 53 uals (44 females and 22 males). These reads were mapped on
 54 the human genome reference GRCh37 after removal of the Y
 55 sequence (the reference used includes the primary assembly, i.e.
 56 chromosomal plus unlocalized and unplaced contigs, and the
 57 rCRS mitochondrial sequence, Human herpesvirus 4 type 1 and
 58 the concatenated decoy sequences) using BWA (Li and Durbin
 59 2009) with standard parameters. The individual alignments to
 60 the reference were merged, and the genotypes were called using

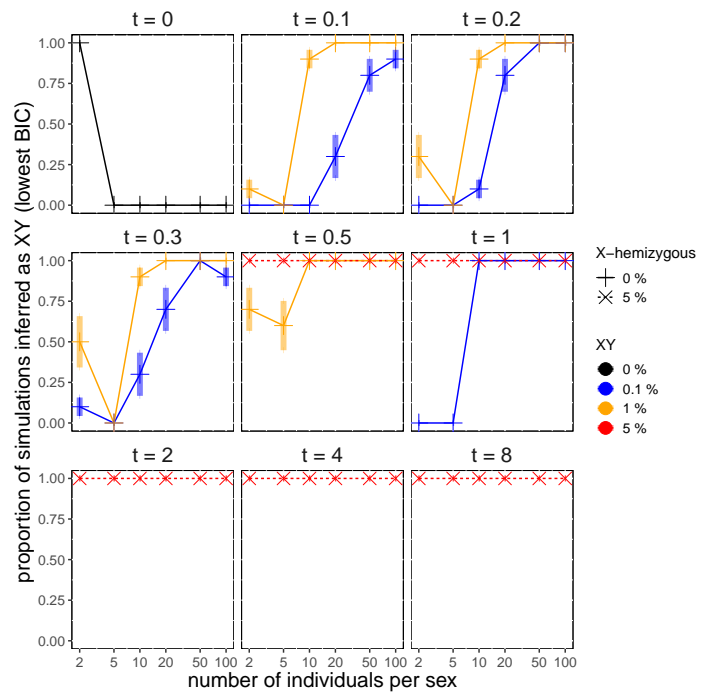


Figure 1 Model choice by SDpop on simulated data. The proportion of simulations for which the XY model had the lowest BIC is indicated; each combination of simulation parameter values was repeated 10 times from a random seed. Vertical bars indicate the expected variance based on the binomial distribution. Different panels represent result for different values of the time since recombination suppression t ; the percentage of simulated X-hemizygous genes is indicated by the line types and symbols (solid lines with "+" symbols indicate no X-hemizygous genes; dashed lines with "x" symbols 5% of X-hemizygous genes); the colors indicate the percentage of XY gametologous genes (black 0%, blue 0.1%, orange 1% and red 5%). These simulations were carried out with error rate $e = 0.0001$; results with $e = 0.001$ are shown in Figure S1.

61 bcftools (mpileup & call), and the targeted exons (file provided
 62 by the 1000 genome project) were extracted using bedtools. Ex-
 63 ons were grouped by gene using the exon list for the reference
 64 GRCh37 retrieved through Ensembl's biomart tool.

65 Results

66 Tests on simulations

67 Model choice: detecting the presence of sex chromosomes

68 We used simulated data to test the range of validity of SDpop
 69 with a controlled population genetic background. We used vary-
 70 ing numbers of individuals, different fractions of sex-linked
 71 sequences in the genome and varying times since recombina-
 72 tion suppression, excluding biologically implausible scenarios (i.e.
 73 a very small fraction of sex-linked sequences and long times
 74 since recombination suppression, or the inverse). The results are
 75 shown in Figure 1 and Figure S1.

76 The behavior of the model when only 2 individuals per sex
 77 are used is somewhat erratic, and leads to a high rate of type
 78 I errors (Figure 1), which is perhaps not surprising given the
 79 limited information that is present in a sample of only 4 individ-
 80 uals. With 5 or more individuals per sex, no type I errors were

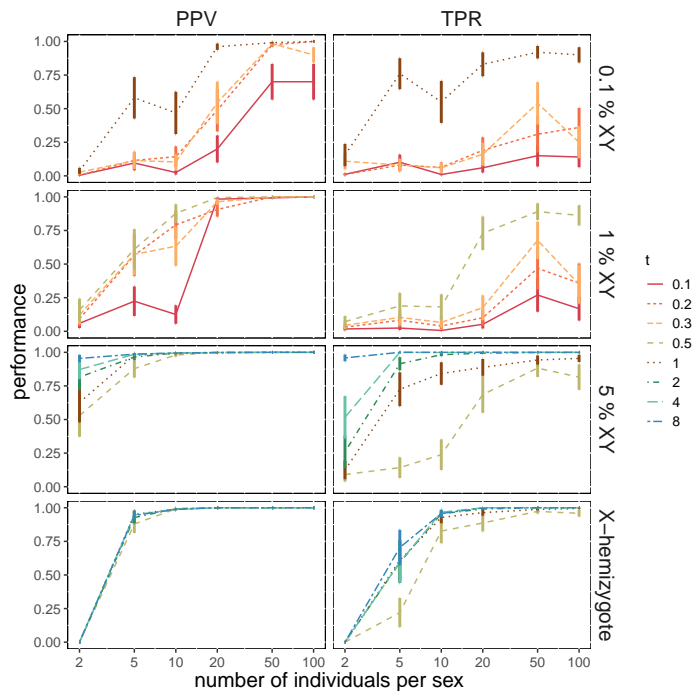


Figure 2 Precision (Positive Predictive Value, left) and power (True Positive Rate, right) of the detection of sex-linked contigs in simulated data, using a threshold for the posterior probability of 0.8. First three rows: XY gametologs, grouped by the proportion of simulated gametologs in the genome (0.1%, 1%, 5%). Bottom graphs: X-hemizygous genes, for which the simulated proportion in the genome was 5%. The color and line scales indicate the simulated time since recombination suppression t . Each point is the average of 100 simulations, with the bars representing the standard error. For all cases shown here, the simulated error rate was 0.0001; for $e = 0.001$, see Figure S2.

1 observed, and the power to detect sex-linkage increases with the
 2 proportion of sex-linked sequences and the time since recombina-
 3 tion suppression, as expected (Figure 1). Indeed, SDpop relies
 4 on polymorphic sites that show evidence for sex linkage: for
 5 gametologs that have stopped recombining $4N_e$ generations ago
 6 ($t = 1$; this is the average time to fixation in a neutral model), we
 7 expect and observe that most “true” polymorphisms (i.e. those
 8 without errors) in sex-linked genes are polymorphic in only one
 9 of the gametologs, or have different alleles fixed on both game-
 10 tologs. However, for recombination suppression much less than
 11 $4N_e$ generations ago, most polymorphisms in gametologs are
 12 either derived from ancestral polymorphism (a case which is not
 13 explicitly modelled), or due to recent mutations with low alter-
 14 native allele frequency in one of the copies, making detection of
 15 sex-linkage much harder. Furthermore, the number of sex-linked
 16 genes is typically low in such situations. Our simulations indi-
 17 cate that even with a time since recombination suppression as
 18 low as $0.1 \times 4N_e$ generations and 0.1% of sex-linked sequences,
 19 the method can nevertheless select the appropriate model in most
 20 cases when 50 or more individuals per sex are used (Figure
 21 S1).

22 **Assignment of genes to segregation types** Contigs are consid-
 23 ered sex-linked when their posterior probability to be either XY

24 or X-hemizygous exceeds a threshold value. We use the thresh-
 25 old value of 0.8 throughout the manuscript, but users can choose
 26 other values depending on the balance between false positives
 27 and false negatives they consider acceptable. We measure the
 28 precision of this assignment, i.e. the fraction of contigs assigned
 29 as sex-linked that were indeed simulated as sex-linked contigs,
 30 as quantified by the positive predictive value (PPV), and the
 31 power of the assignment, i.e. the fraction of contigs that were
 32 simulated as sex-linked that we are able to detect, as quantified
 33 by the true positive rate (TPR). Results are shown in Figure 2
 34 and Figure S2.

35 For XY gametolog pairs, precision and power increase with
 36 time since recombination suppression and the size of the non-
 37 recombining region. When the time since recombination sup-
 38 pression exceeds $4N_e$ generations (i.e. $t \geq 1$) and the non-
 39 recombining region is sufficiently large (i.e. more than 1% of
 40 the genome), the precision is larger than 95% and the power
 41 larger than 70% with as few as 5 individuals per sex (Figure
 42 2). Even with relatively recent recombination suppression ($2N_e$
 43 generations ago, i.e. $t = 0.5$), the method has a precision close
 44 to 100% and a power close to 70% with 20 individuals per sex.
 45 For shorter time since recombination suppression, both decrease
 46 rapidly. Indeed, in these cases, X- and Y-linked SNPs will still
 47 have similar frequencies, and many individuals are needed to
 48 test whether the observed allele frequency differences are due
 49 to sampling or not. For almost all cases, the precision is higher
 50 than the power, meaning that the type I error (false positives) is
 51 lower than the type II error (false negatives).

52 The time since recombination suppression has no clear effect
 53 on the detection of X-hemizygous sequences, understandably,
 54 as it is only the nucleotide polymorphism in X-linked sequences
 55 that creates a signal for detection, and this level of polymor-
 56 phism only depends on N_e . Of course, when the error rate
 57 increases (Figure S2), power and precision decrease, as expected.

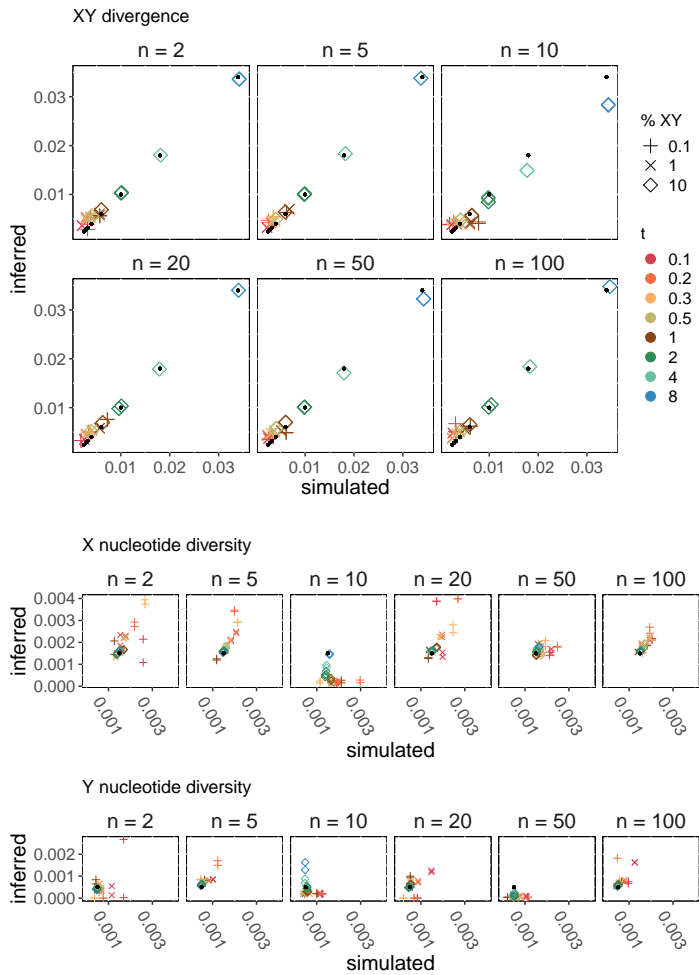


Figure 3 Population genetic inferences of SDpop. The values of nucleotide diversity and divergence calculated directly from the simulation results are compared to the values inferred from SDpop’s output. Comparisons are based on SDpop’s assignment of the genes (i.e. all genes with a posterior probability > 0.8 were used). Top: gametolog divergence (D_{XY}); middle: X nucleotide diversity (π_X); bottom: Y nucleotide diversity (π_Y). Facets are separated by the number of individuals per sex used (n). Color indicates the time since recombination suppression t , and symbols the simulated proportion of gametologs $\%XY$. The black points indicate the theoretical values (D_{XY} : one for each t ; π_X and π_Y : one value for all runs). Here, $e = 0.0001$; for higher error rates, see Figure S3.

indeed, genotyping and sequencing errors have a larger influence here. Importantly, for $t < 1$, as for π_X and π_Y the effect of random variation is large, as expected due to stochastic effects (Takahata and Nei 1985).

Application to *Silene latifolia*

The plant *Silene latifolia* has a pair of well-differentiated X and Y sex chromosomes, with the oldest stratum being 11 My old (Krasovec et al. 2018). The sex chromosomes, and especially the Y chromosome, are large, and previous studies identified about 1000 gametologous and 300 X-hemizygous genes, either using a highly inbred line or a controlled cross (Papadopulos et al. 2015; Muyle et al. 2016). Muyle et al. (2020) provided RNAseq data of 34 plants collected across Europe, that we mapped on the transcriptome used for testing SEX-DETECTOR, a *de novo* assembly based on RNAseq data from male and female plants (Muyle et al. 2016). Mapping and genotyping were performed using the same pipelines as for SEX-DETECTOR.

Prior to the analysis with SDpop, we inspected population structure using a Principal Component Analysis (PCA) of the genetic variation, as shown in Figure S4. As expected, there is substantial structure as the samples were collected from all over Europe (Muyle et al. 2020) and we chose to perform an analysis with SDpop both on the total dataset as well as a subsample of six females and six males, which we termed the “central cluster” (Figure 4).

The four models of SDpop were run on two datasets, one with all individuals and one with twelve individuals from the central cluster. Using all individuals, the dataset consisted of 2118574 SNPs in 29654 contigs; the dataset for the central cluster contained 1106628 SNPs in 26114 contigs. As shown in table 1, the XY model had lower BIC values than the ZW model or the model without sex chromosomes, but a model with both XY and ZW segregation performed slightly better, although only marginally so when SDpop was run on the central cluster. It is possible that the results are influenced by sex-biased expression, as these are RNAseq data: female-biased expression might lead to a higher degree of heterozygosity in females than in males, and thus create ZW-like patterns; likewise, male-biased expression could lead to patterns resembling Z-hemizygosity. It should be noted that a model with both XY and ZW chromosomes in the same individuals has never been demonstrated, so this model is biologically unrealistic. It could be used to test for segregation distortion, e.g. as reported by Martin et al. (2019).

The optimised parameters of the XY model of SDpop indicate that between 1.7% (all individuals) and 2.4% (central cluster) of the SNPs is gametologous, $7.1 \times 10^{-4}\%$ to 1.5% X-hemizygous, and 94% autosomal. These percentages also roughly correspond to the proportion of genes detected as XY gametologous and X-hemizygous (Table 1).

We compared the results obtained here with those of SEX-DETECTOR (Muyle et al. 2016), and with the genes positioned on a genetic map (Papadopulos et al. 2015), as shown in Table 2. SEX-DETECTOR distinguishes between X-hemizygous and XY gametologous genes, but such a distinction was not possible based on the genetic map. There were few false positives, yielding a high Positive Predictive Value (PPV), both for X-hemizygous and XY gametologous genes. There are however many false negatives, in particular concerning X-hemizygous genes. Identifying such genes is notably difficult (Bergero and Charlesworth 2011; Crowson et al. 2017), as there needs to be sufficient polymorphism on the X copies to create a signal.

Population genetic inferences An original feature of SDpop is that for genes inferred as gametologs, the parameter estimates and posterior probabilities can be used to estimate the allele frequencies on the X and Y copies. We use these to calculate the level of diversity in X and Y sequences, as well as their divergence. In Figure 3, the estimates of nucleotide diversity and divergence calculated directly from the simulations, before errors were added, are compared to the estimates based on SDpop’s output.

D_{XY} is expected to increase with time since recombination suppression, and for $t \geq 1$, the values based on SDpop’s inference correlate well with the simulated values. For more recent recombination suppression, SDpop slightly overestimates D_{XY} ;

model	all individuals (34)					central cluster (12)				
	BIC	number of genes				BIC	number of genes			
		Xh	XY	Zh	ZW		Xh	XY	Zh	ZW
no	4.3772×10^7	-	-	-	-	1.2608×10^7	-	-	-	-
XY	4.3366×10^7	98	442	-	-	1.2545×10^7	41	340	-	-
ZW	4.3700×10^7	-	-	36	23	1.2606×10^7	-	-	39	8
both	4.3329×10^7	92	392	22	8	1.2545×10^7	38	312	23	1

Table 1 Summary of model results based on *Silene latifolia* data. Indicated are the Bayesian Information Criterion (BIC), as well as the number of genes inferred as hemizygote (Xh: X-hemizygote; Zh: Z-hemizygote) and gametologs (XY, ZW) using a probability threshold of 0.8.

		all individuals						central cluster					
		N	TP	FP	FN	TPR	PPV	N	TP	FP	FN	TPR	PPV
SEX-DEtector	XY	15534	407	8	706	0.37	0.98	15240	327	3	762	0.30	0.99
	Xh		17	2	139	0.11	0.89		7	1	142	0.047	0.88
Genetic map	XY & Xh	3693	117	5	233	0.33	0.96	3500	74	2	265	0.22	0.98

Table 2 Performance of SDpop on *Silene latifolia* data compared to other methods, SEX-DEtector and a genetic map. Positives in the SEX-DEtector analysis (Muyle *et al.* 2016) were contigs with a posterior probability higher than 0.8, and inferences of SDpop were based on this same criterion. Positives in the genetic map were genes inferred to be in the non-recombining region of the X chromosome (Papadopulos *et al.* 2015); SDpop’s inferences were based on the sum of the posterior probabilities of the X-hemizygote and XY segregation type, that should be higher than 0.8. *N* is the total number of genes in the comparisons (i.e. those present in SDpop’s output and in SEX-DEtector’s output, or on the genetic map); TP, FP, FN the number of true and false positives and false negatives, respectively; TPR and PPV are the True Positive Rate and Positive Predictive Value.

As shown in Figures 4 and S5, when placing the inferences on a genetic map, the X chromosome clearly differs from the rest of the genome, and there is a clear distinction between the non-recombining and the pseudo-autosomal part (cf Krasovec *et al.* 2020). The inferred divergence between X and Y copies is higher than the values for dS found by Papadopulos *et al.* (2015), although the variation, and the differences between stratum I (0 – 40 cM) and stratum II (45 – 63 cM) are similar. The nucleotide diversity of the Y copies seems to be somewhat higher in stratum I than the diversity of the X copies, while the inverse is true in stratum II. However, these calculations would need to be performed on a dataset of samples from a panmictic population to allow interpretation in biological terms.

Performance on human data

The human genome has a strongly heteromorphic XY chromosome pair, with the Y chromosome being much smaller than the X. This difference is due to the strong degeneration the Y chromosome has been subject to: it has lost many genes. Recombination suppression has occurred several times in the lineage leading to humans, leading to distinct strata characterized by different levels of degeneration (Lahn and Page 1999). Only in the youngest stratum have genes retained both X and Y copies, while in the older strata, Y copies have most often been lost. Note that recombination was suppressed much earlier than $4N_e$ generations ago. The most recent stratum is estimated to have stopped recombining about 30×10^6 years ago (Ross *et al.* 2005), while a gross estimate of $4N_e$ would be 10^6 years (the human N_e has varied greatly, but is in the order of magnitude of 10^4 (Auton *et al.* 2015), and the generation time is about 25 years); *t*

would thus be around 30, which is much larger than in the cases we’ve simulated.

Using five, ten or twenty individuals per sex, the model in SDpop with the lowest BIC always was the XY model. The human genome has a low level of polymorphism, and as a consequence, SDpop’s gene-wise inferences are mostly based on a few SNPs (using ten individuals per sex, 15372 genes had SNPs, with a median value of three). SDpop clearly identifies the XY chromosome pair in human sequencing data (Figure 5; Figure S6). Here, the reference consists of the 22 chromosomes and the X (excluding the Y). The larger part of the X chromosome is detected as X-hemizygous, as most of the genes on the X have lost their Y copy. Genes on the extremities of the X chromosome have a high probability to be autosomal, which is again expected as these regions are pseudo-autosomal and do recombine between X and Y. Only one small region with XY gametologs is detected, near the left pseudo-autosomal region, at the position of the youngest stratum with XY gametologs. The genes in the older strata, for which the Y copies have been lost, are detected as X-hemizygous by SDpop. The Y also has several genes without X homologs, that probably resulted from transpositions or translocations from other autosomes, the so-called ampliconic genes (Skaletsky *et al.* 2003). Indeed, the autosomal gene DAZL, known to have given rise to the Y-ampliconic family DAZ, had a very high probability to be sex-linked (> 0.99). Other genes are also known to be homologous with sequences on the sex chromosomes (Galichon *et al.* 2012), and two of these (PPP1R12B and TPTE2) harbor the majority of sex-linked SNPs detected outside chromosome X.

In total, when using ten individuals per sex, SDpop inferred

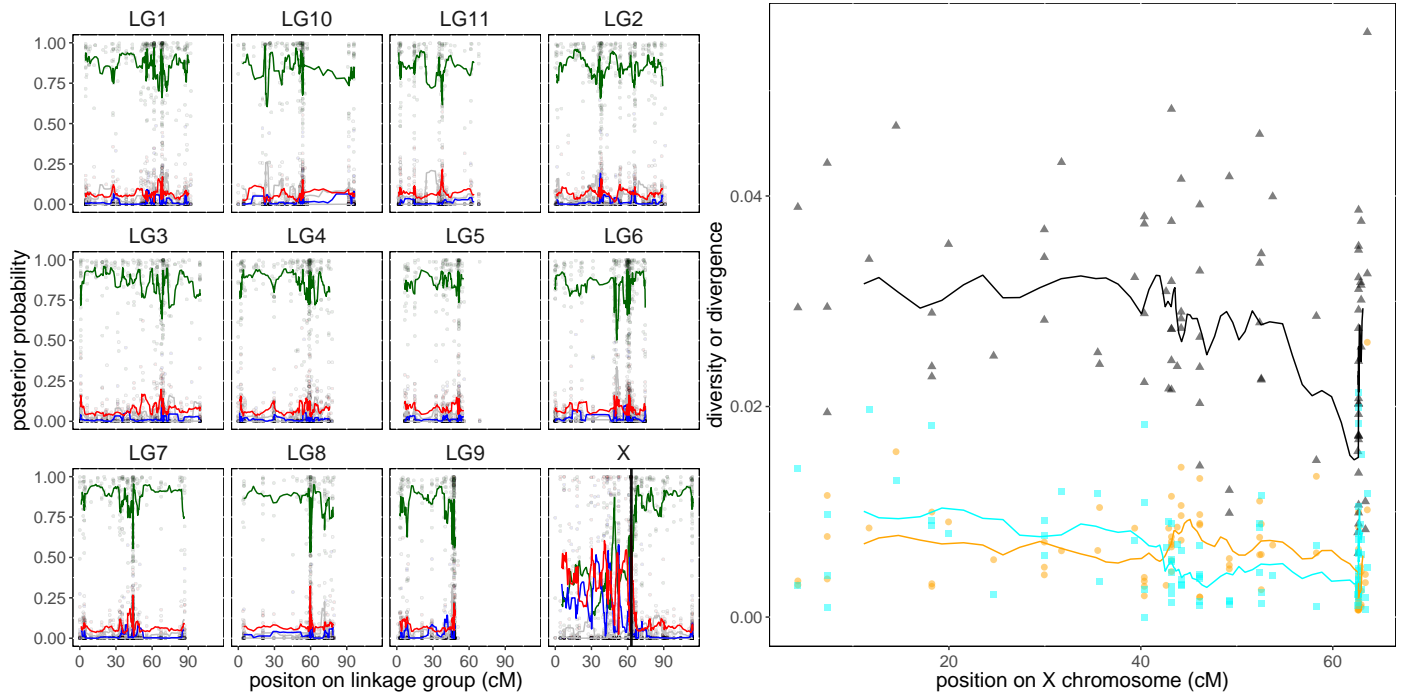


Figure 4 SDpop’s inferences of sex-linkage in *Silene latifolia*, placed on the genetic map of Papadopulos *et al.* (2015). Left panels: posterior probabilities for all placed contigs: autosomal segregation in green, x-hemizyosity in blue, and XY gametology in red; the (uninformative) haploid and paralogous segregation types are indicated in grey. Lines represent running averages, using sliding windows of 10 contigs. The “fuzzy boundary” between the non-recombining region and the pseudoautosomal region on the X chromosome (Krasovec *et al.* 2020) is indicated by the horizontal line. Right panel: predicted divergence (black triangles) and nucleotide diversity of X and Y copies (orange circles and cyan squares) based on SDpop’s output. The lines are the running averages over 10 genes. Figure S5 shows the results obtained with the 12 individuals from the “central cluster”.

1 221 true positives, 77 false positives, and 248 false negatives,
 2 yielding a True Positive Rate of 0.47 and a Positive Predictive
 3 Value of 0.74. As we’ve argued, the fact that these values are
 4 lower than the predictions from simulations is due to the low
 5 level of polymorphism in the human genome, and mapping
 6 errors caused by the dynamics of gene families and an old sex
 7 chromosome system.

8 As Y-specific genes are often present in a few copies, which
 9 allows gene conversion to rescue these sequences that suffer
 10 from the lack of recombination, Y-specific diversity and XY-
 11 divergence are not correctly estimated from the outputs of SD-
 12 pop. Indeed, we observe that the parameter ρ , the proportion
 13 of X-polymorphism in XY gametologs, is 0.12 when using 20
 14 individuals per sex, indicating that 88% of polymorphism in XY
 15 gametologs is due to Y polymorphism. This could be due to the
 16 mapping of several copies of the Y gene to the same position on
 17 the X chromosome.

18 Discussion

19 SDpop is a probabilistic framework for the detection of sex-
 20 linked sequences, that relies on the modeling of the expected
 21 equilibrium between allele and genotype frequencies under sex-
 22 linked segregation. It combines the principles that are at the
 23 basis of several methods for detection of sex-linkage, such as
 24 increased frequencies of heterozygotes or allele frequency de-
 25 viations, in a specific framework. As such, it requires fewer
 26 individuals than methods that are based on allele frequencies
 27 or genotype frequencies alone: e.g., GWAS usually needs > 50

individuals, as do studies using tests for heterozygote frequen- 28
 cies (Picq *et al.* 2014). Expectedly, SDpop’s power depends on 29
 the size and age of the non-recombining region, as would be the 30
 case for any method. 31

32 The likelihood-based framework of SDpop allows compar- 32
 ing models with and without sex linkage. It is thus possible, 33
 with a moderate sequencing effort, to determine the sex chro- 34
 mosome system and to obtain the sex-linked sequences for any 35
 species whose individuals can be sexed and sampled in the field. 36
 The approach can be used on any kind of individual-based se- 37
 quencing data, such as RNA-seq, DNA-reseq and RAD-seq. The 38
 functionality to calculate gene-level posterior probabilities is of 39
 course only useful for gene-based sequencing, such as RNA-seq 40
 and exome capture. For DNA-reseq, per-site posterior probabili- 41
 ties could be aggregated for small scaffolds, or by splitting the 42
 chromosomes into windows of fixed size. 43

44 The underlying population genetics model uses the classical 44
 assumptions of the Hardy-Weinberg principle, notably random 45
 mating between the sexes in an infinitely large population. The 46
 model proposed here will thus perform best when used on a 47
 sample of individuals taken from a single, panmictic population. 48
 As shown in the application to *Silene latifolia*, population struc- 49
 ture will weaken the performance of the model. In this case, the 50
 individuals were sampled from different populations compris- 51
 ing both females and males from each population, and we’ve 52
 shown that the influence on SDpop’s performance is mainly a 53
 loss of power (i.e., an increased proportion of false negatives). 54
 If, on the contrary, females and males are sampled from separ- 55
 ate populations, this will lead to type I errors (false positives), 56

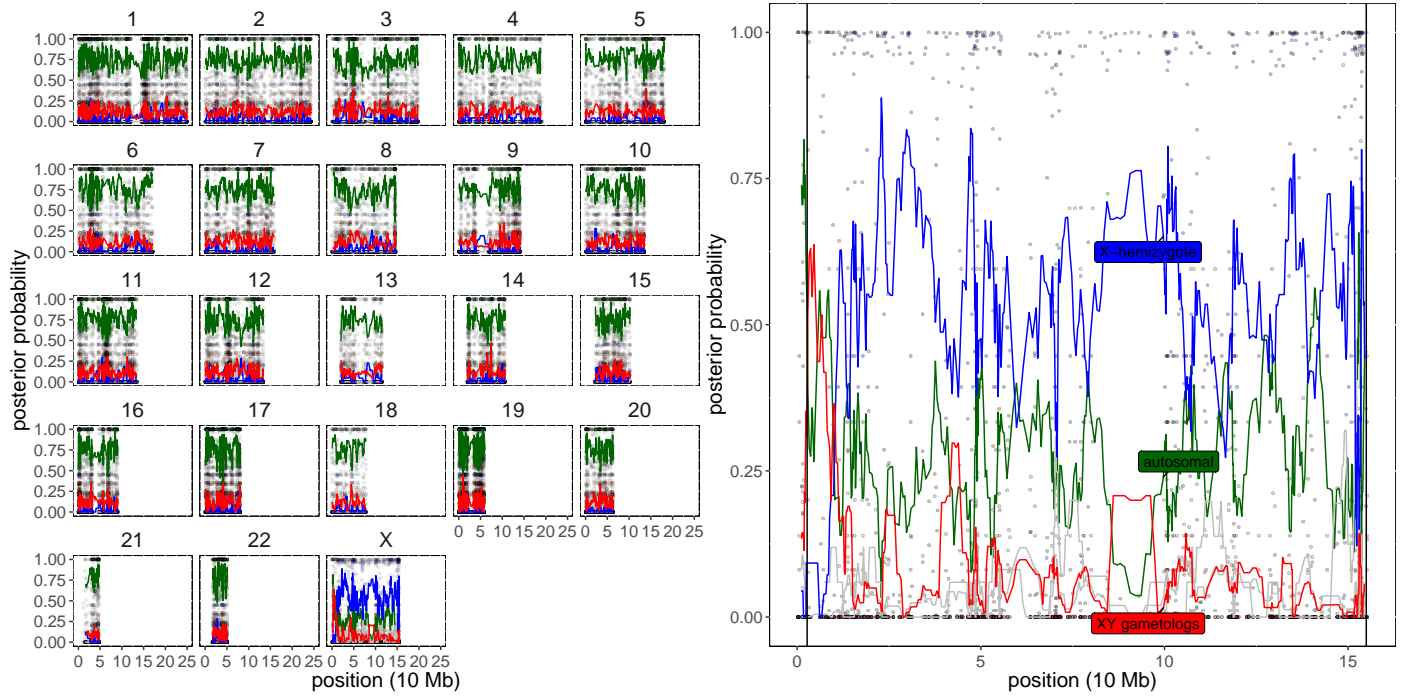


Figure 5 Test of SDpop’s performance on the human exome-targeted sequencing data from the 1000 genome project, using 10 individuals per sex. The exon-level posterior probabilities for autosomal (green), X-hemizygote (blue) and XY (red) segregation are shown, as well as their average in sliding windows of 10 genes; haploid and paralogous posterior probabilities (which are uninformative as they do not allow to distinguish between sex-linkage or not) are indicated in grey. The right panel shows the results on the X chromosome, with the vertical lines delimiting the pseudo-autosomal regions. Results obtained with 5 and 20 individuals per sex are shown in Figure S6.

1 as the population structure mimics the deviations from Hardy-
 2 Weinberg equilibrium that are expected for sex linkage. Even
 3 if all individuals are from one population, a recently migrated
 4 haplotype found in some individuals could lead to false infer-
 5 ences, if this haplotype happened to be overrepresented in the
 6 individuals of one sex present in the sample. It is however
 7 straightforward to check for such population structure and to
 8 exclude potentially problematic individuals.

9 We model different kinds of sources of error in SDpop. First,
 10 we explicitly model haploid and paralogous sequences. These
 11 are more or less frequently encountered in NGS data; haploid
 12 sequences could result from contamination (e.g. mitochondria,
 13 chloroplasts, or bacteria), monoallelic expression in RNAseq
 14 data, or redundancy in the reference (i.e. different alleles of a
 15 gene are split into several contigs). Paralogous sequences could
 16 result from recent paralogs that were not recognized as such
 17 in the reference genome or transcriptome, or from contamina-
 18 tion between samples. Furthermore, SDpop incorporates an
 19 error parameter, to account for other sequencing or genotyping
 20 errors. These sources of errors are modeled solely to improve
 21 the detection of sex-linkage: failing to take them into account
 22 would increase the chance that haploid sequences are inferred
 23 as X-hemizygotes, and paralog sequences as XY gametologs,
 24 while sequencing and genotyping errors would penalize these
 25 sex-linked segregation types more than the inference of auto-
 26 somal segregation. SDpop should not be used when detection
 27 of haploid or paralogous sequences is the goal of a study, as its
 28 performance for these goals has not been evaluated, and other
 29 tools might yield better results. However, when the goal is to
 30 detect sex-linked sequences, we recommend not to filter haploid

or paralogous sequences from a dataset prior to the application
 of SDpop, as such filtering might remove the hemizygous and
 gametologous sequences as well.

SDpop’s sex-linked segregation types include both gametolo-
 gous segregation and hemizygote segregation. In the first case,
 both gametologous copies are present, while in the second, there
 is no information about the sequence from the Y or W chro-
 some. X-hemizygous loci thus correspond to loss of the Y copy
 of a gene, presumably through Y degeneration. However, ap-
 parent X-hemizygosity can also be caused by artifacts that are
 more or less difficult to control. First, the X and Y copies might
 be incorporated as distinct genes in the genome or transcrip-
 tome assembly. In species with known and well described sex
 chromosomes, such as humans, the Y assembly could simply be
 excluded from mapping, as we did here. In species with more
 recently evolved or less well described sex chromosomes, one
 should thus preferably use the homogametic sex for preparing
 a mapping reference. The case that is harder to solve arises
 when Y sequences are too divergent to map on the X copy in
 the reference; thus, when some sex-linked genes have high XY
 divergence values, one should be aware of the fact that some
 genes that have been inferred as X-hemizygous might actually
 be gametologs.

As shown here and by others (e.g. [Bergero and Charlesworth 2011](#); [Crowson et al. 2017](#)), XY gametologs are much easier to
 detect than X-hemizygous genes, as, first, XY gametologs will
 have more SNPs than X-hemizygous genes, and second, the
 information contained in a fixed XY SNP is much less ambiguous
 than for a X-hemizygous SNP. These are additional reasons to
 try to reduce the number of artifactual X-hemizygous as much

1 as possible, in order to obtain more reliable inferences.

2 Importantly, there is further information to be obtained from
3 XY gametologs. We use the inferred allele frequencies on the X
4 and the Y copies to calculate the nucleotide diversity in the X
5 and Y copies, and their divergence. Thus, using the information
6 on which alleles are fixed in each copy, SDpop is also able to
7 reconstruct the haplotypes of the X and Y copies, even when the
8 input data have not previously been phased. We've shown that
9 the estimation of population genetic parameters comes quite
10 close to the simulated values, even though these estimations are
11 based on empirical allele frequencies (i.e., they do not take the
12 error rate into account). However, even if the estimate would
13 be perfectly unbiased, we expect that the variance of these pa-
14 rameters on a per-gene basis (cf Takahata and Nei 1985) is much
15 larger than the bias introduced by not taking the error rate into
16 account. Thus, obtaining reliable estimates of these parameters
17 (i.e., estimates that reflect population processes and not mere
18 stochasticity) requires many independent samples, and cannot
19 be addressed solely by modeling.

20 We use several approximations in SDpop. First, we use the
21 empirical allele frequencies instead of incorporating their es-
22 timation in our model. Second, we only allow one of the ga-
23 metologous (or paralogous) copies to have segregating alleles,
24 assuming one of the alleles is fixed in the other copy. Third, sites
25 are treated as unlinked when calculating the model likelihood.
26 At the contig-level, this leads to overestimating the number of
27 independent observations which would lead to inflating the
28 posterior probability contrasts between the segregation types.
29 For this reason, we use the geometric mean of site likelihoods
30 to calculate the contig-wise posterior probabilities. It would,
31 in principle, be possible to model these points exactly, but this
32 would come with a considerable cost, adding more parameters
33 and assumptions, while we do not expect this to yield a sig-
34 nificant increase in performance at the sample size SDpop is
35 intended for (i.e. 5 to 20 individuals per sex). We've shown here
36 that these approximations yield reliable results in our simulation
37 experiments.

38 We also applied the method to two real datasets, transcrip-
39 tome sequencing for the plant *Silene latifolia* (Muyle et al. 2020),
40 and human exome sequencing from the 1000 genomes project
41 (Auton et al. 2015). For *Silene latifolia*, a species with relatively
42 young sex chromosomes, we used data that show considerable
43 population structure, but nevertheless, SDpop is able to distin-
44 guish the non-recombining region of the X chromosome from
45 the rest of the genome, including the pseudoautosomal region of
46 the sex chromosome. The population structure, which is equally
47 large for females and males, causes the power to be consider-
48 ably reduced, but it importantly doesn't lead to an increase of
49 false positives. The human sex chromosomes are considered
50 old, and much of the Y chromosome has degenerated, so most
51 sex-linked genes are X-hemizygous genes. Although these are
52 more difficult to detect than gametologs, as discussed above,
53 SDpop clearly distinguishes the non-recombining and pseudo-
54 autosomal regions of the X chromosome (Figure 5). Importantly,
55 in these applications to real data, the power and sensitivity for
56 the detection of sex-linked genes is reduced compared to the per-
57 formance in simulations, because of the upstream data treatment,
58 notably the mapping on a reference. In the human data, the exis-
59 tence of gene families with copies on the sex chromosomes as
60 well as on autosomes results in the detection of sex-linkage on
61 autosomes. In *Silene latifolia*, a high-quality genome assembly
62 is not yet available, so mapping was performed on a *de novo*

transcriptome assembly (Muyle et al. 2016), and SDpop's output
was tested using a genetic map containing about twenty per-
cent of the genes. Both the transcriptome and the genetic map
might still contain some errors (e.g. chimeric contigs, mis-placed
scaffolds).

SDpop can be a useful tool to detect sex-linkage in both re-
cent and old sex chromosome systems. Note that we've also
successfully applied SDpop to RNAseq data of the shrub *Am-
borella trichopoda* (Käfer et al. 2020). The framework of SDpop
allowed to detect sex chromosomes of the ZW type, which is the
first report of sex chromosomes in this species. SDpop also pro-
vided a characterization of the non-recombining region, which is
about 4 Mb large and shows almost no gene loss. Thus, SDpop is
applicable to sex chromosomes that are less diverged than those
of *Silene latifolia*.

Although it has similarities with available methods, SDpop is
unique in the combination of input data it requires and the pre-
dictions it can produce. As a consequence, we cannot compare
its performance directly to any of previously published methods.
Our simulations and tests show that reasonable performance
can be achieved with as few as 5 individuals per sex. This is
close to the sequencing effort required for SEX-DETECTOR (Muyle
et al. 2016), which relies on a controlled cross to infer sex-linkage.
Given the fact that usage of SDpop puts less constraint on the in-
put data (kinship does not need to be known, and is assumed to
be absent), it might seem somewhat surprising that SDpop can
work with a similar number of individuals as can SEX-DETECTOR.
The reason for this is that not all sites are informative for SEX-
DETECTOR (e.g. when both parents are heterozygous) and are
discarded, whereas SDpop's likelihoods are calculated over all
polymorphic sites.

Apart from the main practical advantage that SDpop does
not require controlled breeding of the study organism, it has
a few more advantages compared to SEX-DETECTOR. First, as
SEX-DETECTOR considers parents and their F1 offspring, there
has been little recombination between the homologous chro-
mosomes. For the sex chromosomes, this implies that genes
from the pseudo-autosomal region are genetically linked to the
non-recombining region, and will have more or less distorted
segregation. This leads SEX-DETECTOR to overestimate the size
of the non-recombining region, especially when the pseudo-
autosomal region is large and the sex-linked region small (with
a large non-recombining and a small pseudo-autosomal region,
this effect is less important, as the sparse recombination events
of the sex chromosomes will be located in the small pseudo-
autosomal region, and linkage disequilibrium will decrease
rapidly). Thus, we expect SDpop to yield a more precise indica-
tion of the pseudo-autosomal boundary. Note that, in the case of
a very small and recently evolved non-recombining region that
will be difficult to detect by SDpop, SEX-DETECTOR's behavior of
overestimating the non-recombining region might be beneficial
as it increases the capacity to detect the sex chromosome; in
such cases, model choice could be done with SEX-DETECTOR on a
controlled cross, and delimitation of the non-recombining region
with SDpop on population data.

Second, an advantage of the use of population data in SDpop
is that estimates of population genetic parameters are possi-
ble. In a cross, there will be three X chromosomes and one Y,
so it will be impossible to distinguish fixed substitutions from
polymorphism on the Y chromosome.

DETSEX (Gautier 2014) uses a Bayesian framework modeling
samples collected in natural populations, to infer whether mark-

ers (SNPs) are sex-linked or not. Despite the obvious similarities with SDpop, there are several important differences. First, SDpop allows comparing of the total likelihoods of the models, and thus can be used as a statistical test for the presence or absence of sex chromosomes, while in DETSEX, the presence of sex chromosomes is assumed, but the individual's sex does not need to be known. Second, in DETSEX, X and Y-linked sequences are expected to map to different positions, which can be a safe assumption in well-studied species with an old sex chromosomes system (such as humans), but not for more recently evolved sex chromosomes.

Thus, SDpop fills a gap in the current panel of methodological approaches for the identification and study of sex chromosomes. It requires input data that have now become classical: short reads sequencing of genomes or transcriptomes of ten to twenty individuals, collected in any population. It uses standard vcf files as input, thus allowing integration in existing genotyping pipelines. Its probabilistic framework and its implementation in a widely used and efficient programming language (C/C++) furthermore allow future developments (including, but not limited to, inference of individual's sex, corrections for population structure).

Recommendations for the use of SDpop

SDpop is designed to be applicable to a wide range of organisms and using different types of sequencing data (e.g. RNAseq, DNA-reseq, RADseq) as input. Cleaning, mapping, genotyping and filtering of the data can thus be done in different ways to obtain input files (in the standard vcf format). The results, of course, will depend on the quality of the data and the bioinformatic pipelines used, and it's not possible to tailor a standard pipeline. Consider, for example, the case in which no genome of the species studied is available: one could have the choice either to map on a well-assembled genome of a closely related species, or to produce one's own reference genome, but this choice cannot be prescribed without any prior knowledge about the species.

The main requirement for SDpop is that the individuals, which should be sexed, are sampled from a panmictic population. Sampling from one local population seems the most appropriate strategy to ensure that sufficient gene flow has occurred between the individuals, and many population genetic tools are available to verify this (e.g., principal component analysis of genetic variation, clustering (Pritchard *et al.* 2000)).

Some output parameters of SDpop can be used to assess suitability of the data and the quality of the upstream data treatment. E.g., a high percentage of SNPs inferred as haploid indicates that rare alleles are more often found to be homozygous than expected, which could happen if the samples come from a highly spatially structured population, or if alternative alleles are often missed, either through mapping biases or through monoallelic expression in RNAseq data. Or, there could be a high percentage of paralogous SNPs, which could indicate polyploidy, a genome duplication with respect to the reference genome, and possibly other problems. Thus, users of SDpop should have a close look at the general mapping statistics, the genome-wide estimates of heterozygosity in all samples, and possible population structure (e.g. by performing a principal component analysis of genetic variation in the samples).

We recommend that for species with known sex chromosomes, the genome or transcriptome of the homogametic sex is used as a reference for mapping, in order to correctly identify

gametologous genes as such. However, SDpop can also be used to detect the sex chromosome system in species for which such knowledge is lacking. An obvious solution would be use two different references, one for each sex. In other cases, a high-quality assembly might be available for one sex only, which happens to be the heterogametic sex, and the assembly might contain both gametologous copies of some genes. These could be identified and removed from the reference; an example of such procedure is given in Käfer *et al.* (2020) where additional coverage data from DNA-resequencing was used to remove either Z- or W-specific parts of the reference genome.

As gametology, which yields the most powerful signal, is close to paralogy, we recommend that no prior filtering against paralogous sequences is performed. For this reason, we included a paralogous segregation type in SDpop, so the method is able to distinguish it from gametology. Of course, paralogy could also occur for gametologous genes (e.g. a gene duplication present on the sex chromosomes and an autosome), and such cases cannot be distinguished from standard paralogy in the current model. Thus, if many genes or SNPs are inferred as paralogous, the power to detect sex-linkage is reduced, and in that case, it could be worth to finetune the mapping algorithm.

We further recommend that the contig- or gene-wise posterior probabilities are used to identify sex-linked regions. Aggregating the site-wise likelihoods by calculating the geometric mean is a much more robust procedure than focusing on single sites, as the geometric mean gives more weight to sites that are informative (i.e. having a large difference in the likelihood for each of the segregation types). This would help researchers separating noise from signal. For transcriptome or exome data, where the unit of study is a gene, contig or exon, such contig-wise posterior probabilities are naturally calculated. When the data have larger scaffolds or even pseudo-molecules (chromosomes) as units, these could be cut into smaller windows for the calculation of contig-wise posterior probabilities. When the genotyped units only have one or a few SNPs (e.g. RADseq, GBS), we recommend using more individuals and higher thresholds for the inference of sex-linkage.

Acknowledgements

We thank Laurent Duret for help with the treatment of the 1000 genomes data, Aline Muyle for providing access to and advice on *Silene latifolia* data, François Gindraud for help with programming, and Sylvain Mousset and Thibault Latrille for discussions on population genetics issues. This work was performed using the computing facilities of the CC LBBE/PRABI, and was supported by funding from the Agence Nationale de la Recherche (grant number ANR-14-CE19-0021-01).

Author contributions

GM conceived the project and acquired funding; JK, NL, GM and FP conceived the model; JK, NL and FP developed the methodology and the formal analysis; JK developed the software, performed simulations and data analysis, and prepared the first draft manuscript; JK, GM and FP wrote the current version of the manuscript; all authors approve of the manuscript.

Literature Cited

Altschul, S. F., W. Gish, W. Miller, E. W. Myers, and D. J. Lipman, 1990 Basic local alignment search tool. *J. Mol. Biol.* **215**: 403–410.

- 1 Auton, A., L. D. Brooks, R. M. Durbin, E. P. Garrison, H. M. Kang,
2 *et al.*, 2015 A global reference for human genetic variation.
3 *Nature* **526**: 68–74.
- 4 Bachtrog, D., J. E. Mank, C. L. Peichel, M. Kirkpatrick, S. P. Otto,
5 *et al.*, 2014 Sex determination: why so many ways of doing it?
6 *PLoS Biol.* **12**: e1001899.
- 7 Badouin, H., A. Velt, F. Gindraud, T. Flutre, V. Dumas, *et al.*,
8 2020 The wild grape genome sequence provides insights into
9 the transition from dioecy to hermaphroditism during grape
10 domestication. *Genome Biol.* **21**: 223.
- 11 Bellott, D. W., T. J. Cho, J. F. Hughes, H. Skaletsky, and D. C. Page,
12 2018 Cost-effective high-throughput single-haplotype iterative
13 mapping and sequencing for complex genomic structures.
14 *Nat. Protoc.* **13**: 787–809.
- 15 Bellott, D. W., J. F. Hughes, H. Skaletsky, L. G. Brown, T. Pyn-
16 tikova, *et al.*, 2014 Mammalian Y chromosomes retain widely
17 expressed dosage-sensitive regulators. *Nature* **508**: 494–499.
- 18 Bergero, R. and D. Charlesworth, 2011 Preservation of the Y
19 transcriptome in a 10-million-year-old plant sex chromosome
20 system. *Curr. Biol.* **21**: 1470–1474.
- 21 Carvalho, A. B. and A. G. Clark, 2013 Efficient identification of Y
22 chromosome sequences in the human and *Drosophila* genomes.
23 *Genome Res.* **23**: 1894–1907.
- 24 Chang, C. C., C. C. Chow, L. C. Tellier, S. Vattikuti, S. M. Purcell,
25 *et al.*, 2015 Second-generation PLINK: rising to the challenge
26 of larger and richer datasets. *Gigascience* **4**: 7.
- 27 Charlesworth, D., 2016 Plant sex chromosomes. *Annu. Rev. Plant*
28 *Biol.* **67**: 397–420.
- 29 Charlesworth, D., B. Charlesworth, and G. Marais, 2005 Steps in
30 the evolution of heteromorphic sex chromosomes. *Heredity*
31 **95**: 118–128.
- 32 Cortez, D., R. Marin, D. Toledo-Flores, L. Froidevaux, A. Liechti,
33 *et al.*, 2014 Origins and functional evolution of Y chromosomes
34 across mammals. *Nature* **508**: 488–493.
- 35 Crowson, D., S. C. H. Barrett, and S. I. Wright, 2017 Purifying
36 and positive selection influence patterns of gene loss and gene
37 expression in the evolution of a plant sex chromosome system.
38 *Mol. Biol. Evol.* **34**: 1140–1154.
- 39 Fruchard, C., H. Badouin, D. Latrasse, R. S. Devani, A. Muyle,
40 *et al.*, 2020 Evidence for dosage compensation in *Coccinia gran-*
41 *dis*, a plant with a highly heteromorphic XY system. *Genes*
42 (Basel) **11**.
- 43 Galichon, P., L. Mesnard, A. Hertig, B. Stengel, and E. Ron-
44 deau, 2012 Unrecognized sequence homologies may confound
45 genome-wide association studies. *Nucleic Acids. Res.* **40**: 4774–
46 4782.
- 47 Garcia-Moreno, J. and D. P. Mindell, 2000 Rooting a phylogeny
48 with homologous genes on opposite sex chromosomes (gam-
49 etologs): a case study using avian CHD. *Mol. Biol. Evol.* **17**:
50 1826–1832.
- 51 Gautier, M., 2014 Using genotyping data to assign markers to
52 their chromosome type and to infer the sex of individuals: a
53 Bayesian model-based classifier. *Mol. Ecol. Resour.* **14**: 1141–
54 1159.
- 55 Gayral, P., J. Melo-Ferreira, S. Glemin, N. Bierne, M. Carneiro,
56 *et al.*, 2013 Reference-free population genomics from next-
57 generation transcriptome data and the vertebrate-invertebrate
58 gap. *PLoS. Genet.* **9**: e1003457.
- 59 Hall, A. B., Y. Qi, V. Timoshevskiy, M. V. Sharakhova, I. V.
60 Sharakhov, *et al.*, 2013 Six novel Y chromosome genes in
61 *Anopheles* mosquitoes discovered by independently sequenc-
62 ing males and females. *BMC Genomics.* **14**: 273.
- Hudson, R. R., 2002 Generating samples under a Wright-Fisher
neutral model of genetic variation. *Bioinformatics* **18**: 337–338.
- Hughes, J. F. and S. Rozen, 2012 Genomics and genetics of hu-
man and primate y chromosomes. *Annu. Rev. Genomics Hum.*
Genet. **13**: 83–108.
- Käfer, J., A. Bewick, A. Andres-Robin, G. Lapetoule, A. Harkess,
et al., 2020 A derived ZW chromosome system in *Amborella*
trichopoda, the sister species to all other extant flowering plants.
bioRxiv .
- Kirkpatrick, M. and R. F. Guerrero, 2014 Signatures of sex-
antagonistic selection on recombining sex chromosomes. *Ge-*
netics **197**: 531–541.
- Krasovec, M., M. Chester, K. Ridout, and D. A. Filatov, 2018 The
mutation rate and the age of the sex chromosomes in *Silene*
latifolia. *Curr. Biol.* **28**: 1832–1838.
- Krasovec, M., Y. Zhang, and D. A. Filatov, 2020 The location
of the pseudoautosomal boundary in *Silene latifolia*. *Genes*
(Basel). **11**.
- Lahn, B. T. and D. C. Page, 1999 Four evolutionary strata on the
human X chromosome. *Science* **286**: 964–967.
- Li, H. and R. Durbin, 2009 Fast and accurate short read align-
ment with Burrows-Wheeler transform. *Bioinformatics* **25**:
1754–1760.
- Li, S., M. Ajimura, Z. Chen, J. Liu, E. Chen, *et al.*, 2018 A new
approach for comprehensively describing heterogametic sex
chromosomes. *DNA Res.* **25**: 375–382.
- Lynch, M., 2007 *The Origins of Genome Architecture*. Sinauer Asso-
ciates, Inc.
- Martin, H., F. Carpentier, S. Gallina, C. Gode, E. Schmitt, *et al.*,
2019 Evolution of young sex chromosomes in two dioecious
sister plant species with distinct sex determination systems.
Genome Biol. Evol. **11**: 350–361.
- Mathew, L. S., M. Spannagl, A. Al-Malki, B. George, M. F. Torres,
et al., 2014 A first genetic map of year palm (*Phoenix dactylif-*
era) reveals long-range genome structure conservation in the
palms. *BMC Genomics* **15**: 285.
- Muyle, A., J. Käfer, N. Zemp, S. Mousset, F. Picard, *et al.*, 2016
SEX-DETECTOR: a probabilistic approach to study sex chromo-
somes in non-model organisms. *Genome Biol. Evol.* **8**: 2530–
2543.
- Muyle, A., H. Martin, N. Zemp, M. Mollion, S. Gallina, *et al.*,
2020 Dioecy is associated with high genetic diversity and
adaptation rates in the plant genus *Silene*. *Mol. Biol. Evol.* **0**:
O.
- Muyle, A., R. Shearn, and G. A. Marais, 2017 The evolution of sex
chromosomes and dosage compensation in plants. *Genome*
Biol. Evol. **9**: 627–645.
- Muyle, A., N. Zemp, C. Fruchard, R. Cegan, J. Vrana, *et al.*,
2018 Genomic imprinting mediates dosage compensation in a
young plant XY system. *Nat. Plants.* **4**: 677–680.
- Nelson, K. P., 2017 Assessing probabilistic inference by compar-
ing the generalized mean of the model and source probabilities.
Entropy **19**: 286.
- Palmer, D. H., T. F. Rogers, R. Dean, and A. E. Wright, 2019 How
to identify sex chromosomes and their turnover. *Mol. Ecol.* **28**:
4709–4724.
- Papadopulos, A. S., M. Chester, K. Ridout, and D. A. Filatov,
2015 Rapid Y degeneration and dosage compensation in plant
sex chromosomes. *Proc. Natl. Acad. Sci. U.S.A.* **112**: 13021–
13026.
- Picq, S., S. Santoni, T. Lacombe, M. Latreille, A. Weber, *et al.*, 2014
A small XY chromosomal region explains sex determination in

1 wild dioecious *V. vinifera* and the reversal to hermaphroditism
2 in domesticated grapevines. BMC. Plant Biol. **14**: 229.

3 Ponnikas, S., H. Sigeman, J. K. Abbott, and B. Hansson, 2018
4 Why do sex chromosomes stop recombining? Trends Genet.
5 **34**: 492–503.

6 Prentout, D., O. Razumova, B. Rhoné, H. Badouin, H. Henri,
7 *et al.*, 2020 An efficient RNA-seq-based segregation analysis
8 identifies the sex chromosomes of *Cannabis sativa*. Genome
9 Res. **30**: 164–172.

10 Pritchard, J. K., M. Stephens, and P. Donnelly, 2000 Inference of
11 population structure using multilocus genotype data. Genetics
12 **155**: 945–959.

13 Qiu, S., R. Bergero, and D. Charlesworth, 2013 Testing for the
14 footprint of sexually antagonistic polymorphisms in the pseu-
15 doautosomal region of a plant sex chromosome pair. Genetics
16 **194**: 663–672.

17 Ross, M. T., D. V. Grafham, A. J. Coffey, S. Scherer, K. McLay,
18 *et al.*, 2005 The DNA sequence of the human X chromosome.
19 Nature **434**: 325–337.

20 Scharmann, M., T. U. Grafe, F. Metali, and A. Widmer, 2017
21 Sex-determination and sex chromosomes are shared across
22 the radiation of dioecious *Nepenthes* pitcher plants. bioRxiv .

23 Skaletsky, H., T. Kuroda-Kawaguchi, P. J. Minx, H. S. Cordum,
24 L. Hillier, *et al.*, 2003 The male-specific region of the human Y
25 chromosome is a mosaic of discrete sequence classes. Nature
26 **423**: 825–837.

27 Takahata, N. and M. Nei, 1985 Gene genealogy and variance of
28 interpopulational nucleotide differences. Genetics **110**: 325–
29 344.

30 Tomaszkiwicz, M., P. Medvedev, and K. D. Makova, 2017 Y and
31 W Chromosome Assemblies: Approaches and Discoveries.
32 Trends Genet. **33**: 266–282.

33 Torres, M. F., L. S. Mathew, I. Ahmed, I. K. Al-Azwani,
34 R. Krueger, *et al.*, 2018 Genus-wide sequencing supports a
35 two-locus model for sex-determination in *Phoenix*. Nat. Com-
36 mun. **9**: 3969.

37 Veltsos, P., K. E. Ridout, M. A. Troups, S. C. Gonzalez-Martinez,
38 A. Muyle, *et al.*, 2019 Early sex-chromosome evolution in the
39 diploid dioecious plant *Mercurialis annua*. Genetics **212**: 815–
40 835.

41 Vicoso, B. and D. Bachtrog, 2011 Lack of global dosage compen-
42 sation in *Schistosoma mansoni*, a female-heterogametic parasite.
43 Genome Biol. Evol. **3**: 230–235.

44 Vicoso, B., J. J. Emerson, Y. Zektser, S. Mahajan, and D. Bachtrog,
45 2013a Comparative sex chromosome genomics in snakes: dif-
46 ferentiation, evolutionary strata, and lack of global dosage
47 compensation. PLoS Biol. **11**: e1001643.

48 Vicoso, B., V. B. Kaiser, and D. Bachtrog, 2013b Sex-biased gene
49 expression at homomorphic sex chromosomes in emus and its
50 implication for sex chromosome evolution. Proc. Natl. Acad.
51 Sci. U.S.A. **110**: 6453–6458.

52 Wang, J., J. K. Na, Q. Yu, A. R. Gschwend, J. Han, *et al.*, 2012
53 Sequencing papaya X and Yh chromosomes reveals molecular
54 basis of incipient sex chromosome evolution. Proc. Natl. Acad.
55 Sci. U.S.A. **109**: 13710–13715.

56 Zhou, Y., M. Massonnet, J. S. Sanjak, D. Cantu, and B. S. Gaut,
57 2017 Evolutionary genomics of grape (*Vitis vinifera* ssp. *vinifera*)
58 domestication. Proc. Natl. Acad. Sci. U.S.A. **114**: 11715–11720.

Appendix

Full description of the model

Observed and hidden variables The data consist of genotyped individuals for genes $k \in \{1, 2, \dots, K\}$, each with biallelic sites $t \in \{1, 2, \dots, T_k\}$. For each biallelic site, the observed alleles are randomly named a and b , such that three genotypes are possible: homozygote for allele a (“ aa ”), heterozygote (“ ab ”), and homozygote for allele b (“ bb ”). At each biallelic site t of gene k , and for each individual i , $\text{OG}_{i,h}^{tk}$ is an indicator of the individual having observed genotype $g \in \{1 = aa, 2 = ab, 3 = bb\}$ and sex $h \in \{1, 2\}$. $h = 1$ for females, $h = 2$ males; for convenience, we will write φ and σ . N_{gh}^{kt} is the number of individuals with sex h and observed genotype g ; note that the total number of observations (i.e., genotyped individuals) can vary between sites. The vector \mathbf{OG}^{kt} describes all observations at a site.

We seek under which segregation type S_j , $j \in 1..7$, these observed genotypes are most likely. These segregation types will be listed in fixed order, and a specific number corresponds to each of them:

1. Diploid autosomal segregation
2. Haploid sequences
3. Paralogs
4. X-hemizygous segregation
5. XY gametologous segregation
6. Z-hemizygous segregation
7. ZW gametologous segregation

For some segregation types, as will be detailed below, several sub-types of segregation have to be specified; these are denoted A_l with $l \in \{1..L\}$. The conditional likelihood of observing the genotypes under each segregation type depends on the allele frequencies and a genotyping error rate.

We introduce a hidden variable $\text{TG}_{i,h}^{kt}$, which is an indicator for the true genotype $g' \in \{aa, ab, bb\}$ of an individual i with sex h . The conditional probabilities of observing a true genotype for an individual, given the fully specified segregation type, are

$$\mathbf{TG}_{i,h}^{kt} | S_j, A_l \sim \mathcal{M} \left(1; P_{1hjl}^{kt}, P_{2hjl}^{kt}, P_{3hjl}^{kt} \right).$$

\mathbf{P}_{hjl}^{kt} is the vector of the probabilities $(P_{1hjl}^{kt}, P_{2hjl}^{kt}, P_{3hjl}^{kt})$ for each genotype at a site, given the sex of the individual and the segregation type and subtype. The genotype probabilities P_{ghjl}^{kt} are calculated from the empirical allele frequencies \hat{f}_{jl}^{kt} using the following population genetic expectations.

1. For autosomal segregation, the Hardy-Weinberg equilibrium should hold in both sexes. Thus,

$$\mathbf{P}_{\varphi, j=1}^{kt} = \mathbf{P}_{\sigma, j=1}^{kt} = \begin{pmatrix} (\hat{f}_{j=1}^{kt})^2 \\ 2\hat{f}_{j=1}^{kt}(1 - \hat{f}_{j=1}^{kt}) \\ (1 - \hat{f}_{j=1}^{kt})^2 \end{pmatrix}$$

where

$$\hat{f}_{j=1}^{kt} = \frac{2N_{aa}^{kt} + N_{ab}^{kt}}{2N^{kt}}.$$

2. Haploid segregation is modeled by

$$\mathbf{P}_{\varphi, j=2}^{kt} = \mathbf{P}_{\sigma, j=2}^{kt} = \begin{pmatrix} \hat{f}_{j=2}^{kt} \\ 0 \\ 1 - \hat{f}_{j=2}^{kt} \end{pmatrix}$$

and

$$\hat{f}_{j=2}^{kt} = \frac{N_{aa}^{kt}}{N_{aa}^{kt} + N_{bb}^{kt}}.$$

3. Paralogy is caused by the mapping of the reads of two more or less recently duplicated genes on one locus in the reference. There is no recombination between the copies, that thus evolve independently. For simplicity, we assume that one of the copies is fixed for one of the alleles. The genotype probabilities depend on which allele is considered fixed in one of the copies, and two sub-types have to be modeled.

(a) First, we consider that allele a is fixed in one of the copies. $\hat{f}_{j=3,l=1}^{kt}$ is the frequency of allele b in the other copy. In reality, such sites have four copies; thus, the genotypes are $aaaa$, $aaab$, $aabb$, with frequencies $(1 - \hat{f}_{j=3,l=1}^{kt})^2$, $2\hat{f}_{j=3,l=1}^{kt}(1 - \hat{f}_{j=3,l=1}^{kt})$ and $(\hat{f}_{j=3,l=1}^{kt})^2$. Genotypes $aaab$ and $aabb$ will probably be considered as ab by the genotyper that expects only diploids; thus, the genotype probabilities are:

$$\mathbf{P}_{\varnothing,j=3,l=1}^{kt} = \mathbf{P}_{\sigma,j=3,l=1}^{kt} = \begin{pmatrix} (1 - \hat{f}_{j=3,l=1}^{kt})^2 \\ (\hat{f}_{j=3,l=1}^{kt})^2 + 2\hat{f}_{j=3,l=1}^{kt}(1 - \hat{f}_{j=3,l=1}^{kt}) \\ 0 \end{pmatrix}$$

To estimate the empirical allele frequency, note that the ab genotype counts that are obtained from the genotyper (N_{ab}) will likely be a mixture of $aaab$ and $aabb$. The expected proportions N_{aaab} and N_{aabb} can be calculated depending on the frequency $\hat{f}_{j=3,l=1}^{kt}$ that we concisely denote \hat{f} here: $2\hat{f}(1 - \hat{f}) / (\hat{f}^2 + 2\hat{f}(1 - \hat{f}))$ and $\hat{f}^2 / (\hat{f}^2 + 2\hat{f}(1 - \hat{f}))$. While in reality, $\hat{f} = 0.5(N_{aaab} + 2N_{aabb}) / (N_{aaaa} + N_{aaab} + N_{aabb})$, we instead calculate

$$\hat{f} = \frac{1}{2(N_{aa} + N_{ab})} \left(\frac{2\hat{f}(1 - \hat{f})N_{ab}}{\hat{f}^2 + 2\hat{f}(1 - \hat{f})} + \frac{2\hat{f}^2N_{ab}}{\hat{f}^2 + 2\hat{f}(1 - \hat{f})} \right)$$

This yields

$$\hat{f}_{j=3,l=1}^{kt} = 1 - \sqrt{1 - \frac{N_{ab}^{kt}}{N_{aa}^{kt} + N_{ab}^{kt}}}.$$

(b) Alternatively, allele b could be fixed in one of the copies. $\hat{f}_{j=3,l=2}^{kt}$ is the frequency of allele a in the other copy. The genotype probabilities and empirical allele frequency are

$$\mathbf{P}_{\varnothing,j=3,l=2}^{kt} = \mathbf{P}_{\sigma,j=3,l=2}^{kt} = \begin{pmatrix} 0 \\ (\hat{f}_{j=3,l=2}^{kt})^2 + 2\hat{f}_{j=3,l=2}^{kt}(1 - \hat{f}_{j=3,l=2}^{kt}) \\ (1 - \hat{f}_{j=3,l=2}^{kt})^2 \end{pmatrix}$$

$$\hat{f}_{j=3,l=2}^{kt} = 1 - \sqrt{1 - \frac{N_{ab}^{kt}}{N_{ab}^{kt} + N_{bb}^{kt}}}.$$

4. For X-hemizyously segregating genes, the males are haploid while the females are diploid.

$$\mathbf{P}_{\varnothing,j=4}^{kt} = \begin{pmatrix} (\hat{f}_{j=4}^{kt})^2 \\ 2\hat{f}_{j=4}^{kt}(1 - \hat{f}_{j=4}^{kt}) \\ (1 - \hat{f}_{j=4}^{kt})^2 \end{pmatrix} ; \quad \mathbf{P}_{\sigma,j=4}^{kt} = \begin{pmatrix} \hat{f}_{j=4}^{kt} \\ 0 \\ 1 - \hat{f}_{j=4}^{kt} \end{pmatrix}$$

$$\hat{f}_{j=4}^{kt} = \frac{2N_{aa,\varnothing}^{kt} + N_{ab,\varnothing}^{kt} + N_{aa,\sigma}^{kt}}{2(N_{aa,\varnothing}^{kt} + N_{ab,\varnothing}^{kt} + N_{bb,\varnothing}^{kt}) + N_{aa,\sigma}^{kt} + N_{bb,\sigma}^{kt}}$$

5. XY gametologous segregation is characterized by the presence of two independent copies in males, and two copies of the X gene in females. We assume that an allele is fixed in at least one of the copies.

(a) X-polymorphism, allele 1 fixed on Y. f is the frequency of allele 2 on X.

$$\mathbf{P}_{\varnothing,j=5,l=1}^{kt} = \begin{pmatrix} (1 - \hat{f}_{j=5,l=1}^{kt})^2 \\ 2\hat{f}_{j=5,l=1}^{kt}(1 - \hat{f}_{j=5,l=1}^{kt}) \\ (\hat{f}_{j=5,l=1}^{kt})^2 \end{pmatrix} ; \quad \mathbf{P}_{\sigma,j=5,l=1}^{kt} = \begin{pmatrix} 1 - \hat{f}_{j=5,l=1}^{kt} \\ \hat{f}_{j=5,l=1}^{kt} \\ 0 \end{pmatrix}$$

$$\hat{f}_{j=5,l=1}^{kt} = \frac{2N_{bb,\varnothing}^{kt} + N_{ab,\varnothing}^{kt} + N_{ab,\sigma}^{kt}}{2(N_{aa,\varnothing}^{kt} + N_{ab,\varnothing}^{kt} + N_{bb,\varnothing}^{kt}) + N_{aa,\sigma}^{kt} + N_{ab,\sigma}^{kt}}$$

(b) X-polymorphism; allele 2 fixed on Y. f is the frequency of allele 1 on X :

$$\mathbf{P}_{\varphi,j=5,l=2}^{kt} = \begin{pmatrix} (\hat{f}_{j=5,l=2}^{kt})^2 \\ 2\hat{f}_{j=5,l=2}^{kt}(1-\hat{f}_{j=5,l=2}^{kt}) \\ (1-\hat{f}_{j=5,l=2}^{kt})^2 \end{pmatrix} ; \quad \mathbf{P}_{\sigma,j=5,l=2}^{kt} = \begin{pmatrix} 0 \\ \hat{f}_{j=5,l=2}^{kt} \\ 1-\hat{f}_{j=5,l=2}^{kt} \end{pmatrix}$$

$$\hat{f}_{j=5,l=2}^{kt} = \frac{2N_{aa\varphi}^{kt} + N_{ab\varphi}^{kt} + N_{ab\sigma}^{kt}}{2(N_{aa\varphi}^{kt} + N_{ab\varphi}^{kt} + N_{bb\varphi}^{kt}) + N_{bb\sigma}^{kt} + N_{ab\sigma}^{kt}}$$

(c) Y-polymorphism, allele 1 fixed on X. f is the frequency of allele 2 on Y:

$$\mathbf{P}_{\varphi,j=5,l=3}^{kt} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} ; \quad \mathbf{P}_{\sigma,j=5,l=3}^{kt} = \begin{pmatrix} 1-\hat{f}_{j=5,l=3}^{kt} \\ \hat{f}_{j=5,l=3}^{kt} \\ 0 \end{pmatrix}$$

$$\hat{f}_{j=5,l=3}^{kt} = \frac{N_{ab\sigma}^{kt}}{N_{aa\sigma}^{kt} + N_{ab\sigma}^{kt}}$$

(d) Y-polymorphism, allele 2 fixed on X. f is the frequency of allele 1 on Y:

$$\mathbf{P}_{\varphi,j=5,l=4}^{kt} = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} ; \quad \mathbf{P}_{\sigma,j=5,l=4}^{kt} = \begin{pmatrix} 0 \\ \hat{f}_{j=5,l=4}^{kt} \\ 1-\hat{f}_{j=5,l=4}^{kt} \end{pmatrix}$$

$$\hat{f}_{j=5,l=4}^{kt} = \frac{N_{ab\sigma}^{kt}}{N_{bb\sigma}^{kt} + N_{ab\sigma}^{kt}}$$

6. Z-hemizygous segregation is similar to X-hemizygous segregation:

$$\mathbf{P}_{\varphi,j=6}^{kt} = \begin{pmatrix} \hat{f}_{j=6}^{kt} \\ 0 \\ 1-\hat{f}_{j=6}^{kt} \end{pmatrix} ; \quad \mathbf{P}_{\sigma,j=6}^{kt} = \begin{pmatrix} (\hat{f}_{j=6}^{kt})^2 \\ 2\hat{f}_{j=6}^{kt}(1-\hat{f}_{j=6}^{kt}) \\ (1-\hat{f}_{j=6}^{kt})^2 \end{pmatrix}$$

$$\hat{f}_{j=6}^{kt} = \frac{2N_{aa,\sigma}^{kt} + N_{ab,\sigma}^{kt} + N_{aa,\varphi}^{kt}}{2(N_{aa,\sigma}^{kt} + N_{ab,\sigma}^{kt} + N_{bb,\sigma}^{kt}) + N_{aa,\varphi}^{kt} + N_{bb,\varphi}^{kt}}$$

7. ZW gametologous segregation is modeled similar to XY gametologous segregation, for both Z and W polymorphism, and two asymmetrical cases for each.

(a) Z-polymorphism, allele 1 fixed on W. f is the frequency of allele 2 on Z.

$$\mathbf{P}_{\varphi,j=7,l=1}^{kt} = \begin{pmatrix} 1-\hat{f}_{j=7,l=1}^{kt} \\ \hat{f}_{j=7,l=1}^{kt} \\ 0 \end{pmatrix} ; \quad \mathbf{P}_{\sigma,j=7,l=1}^{kt} = \begin{pmatrix} (1-\hat{f}_{j=7,l=1}^{kt})^2 \\ 2\hat{f}_{j=7,l=1}^{kt}(1-\hat{f}_{j=7,l=1}^{kt}) \\ (\hat{f}_{j=7,l=1}^{kt})^2 \end{pmatrix}$$

$$\hat{f}_{j=7,l=1}^{kt} = \frac{2N_{bb\sigma}^{kt} + N_{ab\sigma}^{kt} + N_{ab\varphi}^{kt}}{2(N_{aa\sigma}^{kt} + N_{ab\sigma}^{kt} + N_{bb\sigma}^{kt}) + N_{aa\varphi}^{kt} + N_{ab\varphi}^{kt}}$$

(b) Z-polymorphism; allele 2 fixed on W. f is the frequency of allele 1 on Z:

$$\mathbf{P}_{\varphi,j=7,l=2}^{kt} = \begin{pmatrix} 0 \\ \hat{f}_{j=7,l=2}^{kt} \\ 1-\hat{f}_{j=7,l=2}^{kt} \end{pmatrix} ; \quad \mathbf{P}_{\sigma,j=7,l=2}^{kt} = \begin{pmatrix} (\hat{f}_{j=7,l=2}^{kt})^2 \\ 2\hat{f}_{j=7,l=2}^{kt}(1-\hat{f}_{j=7,l=2}^{kt}) \\ (1-\hat{f}_{j=7,l=2}^{kt})^2 \end{pmatrix}$$

$$\hat{f}_{j=7,l=2}^{kt} = \frac{2N_{aa\sigma}^{kt} + N_{ab\sigma}^{kt} + N_{ab\varphi}^{kt}}{2(N_{aa\sigma}^{kt} + N_{ab\sigma}^{kt} + N_{bb\sigma}^{kt}) + N_{bb\varphi}^{kt} + N_{ab\varphi}^{kt}}$$

(c) W-polymorphism, allele 1 fixed on Z. f is the frequency of allele 2 on W:

$$\mathbf{P}_{\Phi, j=7, l=3}^{kt} = \begin{pmatrix} 1 - \hat{f}_{j=7, l=3}^{kt} \\ \hat{f}_{j=7, l=3}^{kt} \\ 0 \end{pmatrix} ; \quad \mathbf{P}_{\sigma, j=7, l=3}^{kt} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

$$\hat{f}_{j=7, l=3}^{kt} = \frac{N_{ab\Phi}^{kt}}{N_{aa\Phi}^{kt} + N_{ab\Phi}^{kt}}$$

(d) W-polymorphism, allele 2 fixed on Z. f is the frequency of allele 1 on W:

$$\mathbf{P}_{\Phi, j=7, l=4}^{kt} = \begin{pmatrix} 0 \\ \hat{f}_{j=7, l=4}^{kt} \\ 1 - \hat{f}_{j=7, l=4}^{kt} \end{pmatrix} ; \quad \mathbf{P}_{\sigma, j=7, l=4}^{kt} = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

$$\hat{f}_{j=7, l=4}^{kt} = \frac{N_{ab\Phi}^{kt}}{N_{bb\Phi}^{kt} + N_{ab\Phi}^{kt}}$$

In some cases, calculation of \hat{f}_{jl}^{kt} might lead to division by 0. To avoid this problem, counts that are expected to be 0 under a segregation type are added to the numerator and the denominator.

Genotyping errors (whether they are due to sequencing errors, read mapping errors, or violations of the assumptions of the method for genotyping) cause the observed genotype g to be different from the true genotype g' . We define $q_{gg'} = \mathbb{P}(\text{OG}_{ihg}^{kt} | \text{TG}_{ihg'}^{kt})$, i.e., the probability to observe genotype g when the true genotype is g' , and \mathbf{Q} is the matrix of all $q_{gg'}$, such that

$$\mathbf{Q} = \begin{pmatrix} q_{1,1} & q_{1,2} & q_{1,3} \\ q_{2,1} & q_{2,2} & q_{2,3} \\ q_{3,1} & q_{3,2} & q_{3,3} \end{pmatrix}$$

We can now directly calculate the probabilities of the observed genotypes for each segregation type:

$$\begin{aligned} \mathbb{P}(\text{OG}_{ihg}^{kt} | S_j, A_l) &= \sum_{g'} \mathbb{P}(\text{TG}_{ihg'}^{kt} | S_j, A_l) \mathbb{P}(\text{OG}_{ihg}^{kt} | \text{TG}_{ihg'}^{kt}) \\ &= \sum_{g'} P_{g'hjl}^{kt} q_{gg'} \end{aligned}$$

We rename the quantity $\sum_{g'} P_{g'hjl}^{kt} q_{gg'}$ as \tilde{P}_{ghjl}^{kt} ; it is the expected frequency of the observed genotype given the segregation type and a certain genotyping error rate. For each sex, \mathbf{OG} follows a multinomial distribution $\mathcal{M}(N_h^{kt}, \tilde{P}_{1hjl}^{kt}, \tilde{P}_{2hjl}^{kt}, \tilde{P}_{3hjl}^{kt})$. Thus, the conditional likelihood of the data (given the segregation type) at each site, that we name M_{jl}^{kt} , is

$$M_{jl}^{kt} = \mathbb{P}(\mathbf{OG}^{kt} | S_j, A_l) = \prod_{gh} (\tilde{P}_{ghjl}^{kt})^{N_{gh}^{kt}} \quad (1)$$

Parameters The error rates $q_{gg'}$ depend on one error parameter e . We assume all genotyping errors to occur with the same frequency, so $q_{g,g' \neq g} = e$ and $q_{g,g'=g} = 1 - 2e$, which gives the error matrix

$$\mathbf{Q} = \begin{pmatrix} 1 - 2e & e & e \\ e & 1 - 2e & e \\ e & e & 1 - 2e \end{pmatrix}$$

Two more series of parameters are required to model the data; these indicate the proportion of the genome that segregates under each type. There are a maximum of seven segregation types S_j , each occupying a proportion π_j of the genome, such that $\sum_j \pi_j = 1$. $\boldsymbol{\pi}$ is the vector containing all π_j . The segregation types S_j^{kt} are distributed multinomially, thus

$$\mathbf{S} \sim \mathcal{M}(1, \boldsymbol{\pi})$$

Several biologically relevant segregation types (S) have several “subtypes” (A), depending on the fixation of one of the alleles on either of the copies. Thus, for a segregation type S_j , there are L subtypes, and each subtype A_{jl} applies to a proportion α_{jl} of the proportion π_j of the genome (corresponding to the segregation type S_j). For each segregation type with subtypes, $\sum_{l=1}^L \alpha_{jl} = 1$, and

$$\mathbf{A}_j | S_j \sim \mathcal{M}(1, \boldsymbol{\alpha}_j)$$

For the paralogs, the subtype depends uniquely on the choice of what allele is called a , which is random. Thus, no parameter is needed, and

$$\boldsymbol{\alpha}_3 = \left(\frac{1}{2}, \frac{1}{2} \right)$$

For the XY and ZW types, more sites can be polymorphic on one chromosome than on the other. The proportion of XY or ZW sites that are polymorphic on X or on Z is described by the parameter ρ_j , which takes a single value for each segregation type. The (random) choice what allele is called a affects both X (or Z) and Y (or W) polymorphisms, leading to four subtypes

$$\boldsymbol{\alpha}_5 = \left(\frac{\rho_5}{2}, \frac{\rho_5}{2}, \frac{1-\rho_5}{2}, \frac{1-\rho_5}{2} \right)$$

When aggregating the segregation subtypes (A) to biologically relevant types (S), we get

$$\mathbb{P}(\mathbf{OG}^{kt} | S_j) = B_j^{kt} = \sum_l \alpha_{jl} M_{jl}^{kt} \quad (2)$$

Expectation-Maximization algorithm The full log-likelihood of the model is given by

$$\begin{aligned} \log \mathbb{P}(\mathbf{OG}, \mathbf{TG}, \mathbf{S}, \mathbf{A}) &= \log \mathbb{P}(\mathbf{OG} | \mathbf{TG}) \\ &+ \log \mathbb{P}(\mathbf{TG} | \mathbf{S}, \mathbf{A}) \\ &+ \log \mathbb{P}(\mathbf{A} | \mathbf{S}) \\ &+ \log \mathbb{P}(\mathbf{S}) \end{aligned}$$

This likelihood is maximized through an Expectation-Maximization (EM) algorithm.

E-step The posterior segregation types are given by

$$\mathbb{E}(\log \mathbb{P}(\mathbf{S}) | \mathbf{OG}) = \sum_{jkt} \mathbb{E}(S_j^{kt} | \mathbf{OG}^{kt}) \log \pi_j$$

with

$$\mathbb{E}(S_j^{kt} | \mathbf{OG}^{kt}) = \hat{S}_j^{kt} = \frac{\pi_j B_j^{kt}}{\sum_{j'} \pi_{j'} B_{j'}^{kt}} \quad (3)$$

The posteriors for the subtypes are calculated by

$$\begin{aligned} \mathbb{E}(\log \mathbb{P}(\mathbf{A}^{kt} | \mathbf{S}^{k(t)}) | \mathbf{OG}^{kt}) &= \sum_{j l k t} \mathbb{E}(A_{jl}^{kt} S_j^{k(t)} | \mathbf{OG}^{kt}) \log \alpha_{jl} \\ &= \sum_{j l k t} \mathbb{E}(A_{jl}^{kt} | \mathbf{OG}^{kt}, S_j^{k(t)}) \mathbb{E}(S_j^{k(t)} | \mathbf{OG}^{k(t)}) \log \alpha_{jl} \\ &= \sum_{j l k t} \hat{A}_{lj}^{kt} \hat{S}_j^{k(t)} \log \alpha_{jl} \\ \hat{A}_{lj}^{kt} &= \frac{\alpha_{jl} M_{jl}^{kt}}{\sum_{l'} \alpha_{jl'} M_{jl'}^{kt}} \end{aligned}$$

For the true expected true genotypes, we calculate

$$\begin{aligned} \mathbb{E}(\log \mathbb{P}(\mathbf{TG} | \mathbf{S}, \mathbf{A}) | \mathbf{OG}) &= \sum_{(kt)(jl)(i_h g')} \mathbb{E}(S_j^{kt} A_{jl}^{kt} \mathbf{TG}_{i_h g'}^{kt} | \mathbf{OG}_{i_h}^{kt}) \log P_{g' h j l}^{kt} \\ &= \sum_{(kt)(jl)(i_h g')} \mathbb{E}(\mathbf{TG}_{i_h g'}^{kt} | \mathbf{OG}_{i_h}^{kt}, S_j^{kt}, A_{jl}^{kt}) \mathbb{E}(S_j^{kt}, A_{jl}^{kt} | \mathbf{OG}) \log P_{g' h j l}^{kt} \\ &= \sum_{(kt)(jl)(i_h g')} \mathbb{E}(\mathbf{TG}_{i_h g'}^{kt} | \mathbf{OG}_{i_h}^{kt}, S_j^{kt}, A_{jl}^{kt}) \hat{A}_{lj}^{kt} \hat{S}_j^{kt} \log P_{g' h j l}^{kt} \end{aligned}$$

$$\widehat{\text{TG}}_{i_h g' j l}^{kt} = \frac{P_{g' h j l}^{kt} \prod_g q_{g g'}^{\text{OG}_{i_h g}^{kt}}}{\sum_{g''} P_{g'' h j l}^{kt} \prod_g q_{g g''}^{\text{OG}_{i_h g}^{kt}}}$$

As individuals are defined uniquely by their sex and their observed genotype, $\widehat{\text{TG}}_{i_h g' j l}^{kt}$ is the same for two individuals having the same sex and genotype. Thus, we write $\widehat{\text{TG}}_{h g g' j l}^{kt} = \frac{P_{g' h j l}^{kt} q_{g g'}}{\sum_{g''} P_{g'' h j l}^{kt} q_{g g''}}$.

Finally, the conditional likelihood of the observed genotypes is given by

$$\begin{aligned} \mathbb{E}(\log \mathbb{P}(\mathbf{OG} | \mathbf{TG}) | \mathbf{OG}) &= \sum_{(kt)(jl)(i_h g')g} \text{OG}_{i_h g}^{kt} \mathbb{E}(\text{TG}_{i_h g' j l}^{kt} | \mathbf{OG}^{kt}) \log q_{g g'} \\ &= \sum_{(kt)(jl)(i_h g')g} \text{OG}_{i_h g}^{kt} \widehat{\text{TG}}_{i_h g' j l}^{kt} \widehat{A}_{l j}^{kt} \widehat{S}_j^{kt} \log q_{g g'} \end{aligned}$$

M-step The key quantity to be used in the M-step is the conditional expectation of the complete-data likelihood:

$$\begin{aligned} \mathbb{E}(\log \mathbb{P}(\mathbf{OG}, \mathbf{TG}, \mathbf{S}, \mathbf{A}) | \mathbf{OG}) &= \sum_{(kt)(jl)(i_h g')g} \text{OG}_{i_h g}^{kt} \widehat{\text{TG}}_{i_h g' j l}^{kt} \widehat{A}_{l j}^{kt} \widehat{S}_j^{k(t)} \log q_{g g'} \\ &+ \sum_{(kt)(jl)(i_h g')} \widehat{\text{TG}}_{i_h g' j l}^{kt} \widehat{A}_{l j}^{kt} \widehat{S}_j^{k(t)} \log P_{g' j l}^{kth} \\ &+ \sum_{(kt)(jl)} \widehat{A}_{l j}^{kt} \widehat{S}_j^{k(t)} \log \alpha_{j l} + \sum_{(k(t))j} \widehat{S}_j^{k(t)} \log \pi_j \\ &= \sum_{(kt)(h g)} N_g^{kth} \left(\sum_{(jl)} \widehat{A}_{l j}^{kt} \widehat{S}_j^{k(t)} \left(\sum_{g'} \widehat{\text{TG}}_{h g g' j l}^{kt} (\log q_{g g'} + \log P_{g' j l}^{kth}) \right) \right) \\ &+ \sum_{(kt)(j l)} \widehat{A}_{l j}^{kt} \widehat{S}_j^{k(t)} \log \alpha_{j l} + \sum_{(k(t))j} \widehat{S}_j^{k(t)} \log \pi_j \end{aligned}$$

Parameters to estimate are π , α and error rate e . These parameters only involve

$$\begin{aligned} \mathbb{E}(\log \mathbb{P}(\mathbf{OG} | \mathbf{TG}) | \mathbf{OG}) &= \sum_{(kt)} \sum_{(jl)(i_h)} \text{OG}_{i_h g}^{kt} \widehat{\text{TG}}_{i_h g' j l}^{kt} \widehat{A}_{l j}^{kt} \widehat{S}_j^{k(t)} \log q_{g g'} \\ &= \sum_{(kt)} \sum_{(jl)(g h)} N_g^{kth} \widehat{\text{TG}}_{h g g' j l}^{kt} \widehat{A}_{l j}^{kt} \widehat{S}_j^{k(t)} \log q_{g g'} \end{aligned}$$

To simplify notations, let us denote

$$\begin{aligned} \widehat{U}_{g g'} &= \sum_{(kt)} \sum_{(jl)(i_h)} \text{OG}_{i_h g}^{kt} \widehat{\text{TG}}_{i_h g' j l}^{kt} \widehat{A}_{l j}^{kt} \widehat{S}_j^{k(t)} \\ &= \sum_{(kt)} \sum_{(jl)(h)} N_g^{kth} \widehat{\text{TG}}_{h g g' j l}^{kt} \widehat{A}_{l j}^{kt} \widehat{S}_j^{k(t)} \end{aligned}$$

Thus,

$$\mathbb{E}(\log \mathbb{P}(\mathbf{OG} | \mathbf{TG}) | \mathbf{OG}) = \sum_{(g g')} \widehat{U}_{g g'} \log q_{g g'}$$

We find the new values of e by $\frac{\partial \mathbb{E}(\log \mathbb{P}(\mathbf{OG} | \mathbf{TG}) | \mathbf{OG})}{\partial e} = 0$, which gives:

$$\begin{aligned} \widehat{e} &= \frac{\widehat{U}_{ab} + \widehat{U}_{13} + \widehat{U}_{21} + \widehat{U}_{23} + \widehat{U}_{31} + \widehat{U}_{32}}{2(\widehat{U}_{ab} + \widehat{U}_{bb} + \widehat{U}_{32} + \widehat{U}_{ab} + \widehat{U}_{13} + \widehat{U}_{21} + \widehat{U}_{23} + \widehat{U}_{31} + \widehat{U}_{32})} \\ &= \frac{\sum_{(g \neq g')} \widehat{U}_{g g'}}{2 \sum_{(g g')} \widehat{U}_{g g'}} \end{aligned}$$

Similarly, we calculate

$$\hat{\rho} = \frac{\sum_{kt} \hat{S}_3^{kt} (\hat{A}_{13}^{kt} + \hat{A}_{23}^{kt})}{\sum_{kt} \hat{S}_3^{kt}}$$

$$\hat{\pi}_j = \frac{\sum_{kt} \hat{S}_j^{kt}}{\sum_{kt} 1}$$

Monitoring and convergence The likelihood of the data in the model is

$$\begin{aligned} \log \mathbb{P}(\mathbf{OG}) &= \sum_{kt} \log \left(\sum_{jl} \left(\mathbb{P}(A_l^{kt}) \mathbb{P}(S_j^{kt}) \prod_i \mathbb{P}(\text{OG}_{i,g}^{kt} | S_j^{kt}, A_l^{kt}) \right) \right) \\ &= \sum_{kt} \log \left(\sum_{jl} \pi_j \alpha_l M_{jl}^{kt} \right) \end{aligned}$$

Convergence is evaluated as a function of the relative change in parameter value estimations. Optimization is halted when the largest relative change of all parameters has been less than 10^{-4} for 10 iterations, except for the error rate parameter, which is not considered for convergence.

There are $J - 1$ free parameters for the segregation types, one parameter α for each of the XY and ZW types, and one parameter for the error rate. If the number of parameters is ξ , we calculate the Bayesian Information Criterion (BIC) as follows:

$$\text{BIC} = \log \mathbb{P}(\mathbf{OG}; \hat{\theta}) - \frac{1}{2} \log \left(\sum_{kt} 1 \right) \xi$$

The model with the lowest BIC has the best fit.

Site- and contig-wise probabilities The posterior probabilities per site, as given in Equation 3, are calculated using the priors π_j , which are the estimated proportions of each segregation type in the genome. The smaller π_j , the higher the conditional likelihood B_j^{kt} should be to produce a high posterior probability. For the sex-linked segregation types π_j can easily be very small. If, say, 0.1% of the sites are inferred as gametologous and 99.9% as autosomal, the conditional likelihood for the gametologous segregation types should be 1000 times higher than the one for autosomal segregation to obtain comparable posterior likelihoods with this formula. In order to avoid excessive biases against rare segregation types, for inference purposes at the end of the optimization, we calculate the posterior probabilities without priors, which amount to using a uniform prior. Thus, for the output, we compute

$$\hat{S}_j^{kt} = \frac{B_j^{kt}}{\sum_{j'} B_{j'}^{kt}} \quad (4)$$

At the contig level, the goal is to estimate the posterior probability to be sex-linked, autosomal, or not informative (*i.e.*, haploid or paralogous), given the observed data for each of its sites and the optimal parameter values. This probability is the expectation of each segregation type, \hat{S}_j^k , which we calculate from the site-wise probabilities. As sites are treated as unlinked, which they are obviously not within a contig, especially when they are sex-linked, calculating the product of the site likelihoods would lead to ignoring the dependence induced by linkage and to overestimating the effective number of independent observations. This is thus expected to inflate the posterior probability contrasts between alternative hypotheses (segregation types) for a given contig. Instead, we take the geometric mean, which reduces this effect:

$$\hat{S}_{N_j}^k = \frac{\text{GM}(\hat{S}_j^{kt})}{\sum_{j'} \text{GM}(\hat{S}_{j'}^{kt})} = \frac{\text{GM}(B_j^{kt})}{\sum_{j'} \text{GM}(B_{j'}^{kt})} = \frac{\exp\left(\frac{1}{T_k} \sum_t \log B_j^{kt}\right)}{\sum_{j'} \exp\left(\frac{1}{T_k} \sum_t \log B_{j'}^{kt}\right)} \quad (5)$$

The geometric mean has the further advantage to give more weight to informative sites, for which the probabilities for each segregation type are very different (say, 0.1 and 10^{-5}), than to sites with less information (say, 0.4 and 0.6). Thus, a site with all females heterozygous and all males homozygous, which would produce a much higher likelihood to be sex-linked than to be autosomal, has more weight than a site with one female heterozygous and all other individuals homozygous, a pattern compatible with both sex-linkage and autosomal segregation.

For completeness (e.g. to allow additional calibration by expert users), we provide two other ways to calculate the posterior probabilities per contig. First, we provide the posterior probability as the geometric mean of the site-wise probabilities calculated using the estimated genome proportions π_j as priors (as in Equation 3):

$$\hat{S}_{G_j}^k = \frac{\pi_j \text{GM}(B_j^{kt})}{\sum_{j'} \pi_{j'} \text{GM}(B_{j'}^{kt})} = \frac{\pi_j \exp\left(\frac{1}{T_k} \sum_t \log B_j^{kt}\right)}{\sum_{j'} \pi_{j'} \exp\left(\frac{1}{T_k} \sum_t \log B_{j'}^{kt}\right)} \quad (6)$$

Second, we provide the arithmetic mean of the expectations per site, \widehat{S}_j^{kt} from Equation 3:

$$\widehat{S}_j^k = \frac{1}{T_k} \sum_t \widehat{S}_j^{kt} = \frac{1}{T_k} \sum_t \frac{\pi_j B_j^{kt}}{\sum_j \pi_j B_j^{kt}} \quad (7)$$

We recommend that inferences of segregation types should be based on the posterior probabilities that were calculated without the priors, *i.e.* Equation 4 for sites and Equation 5 for contigs.

Population genetic predictions From the allele frequencies and segregation subtypes, it is possible to calculate the expected diversity and divergence of the gametologous copies. For each site, the frequency of allele a on chromosome $v \in \{W, X, Y, Z\}$ is

$$\begin{aligned} \widehat{f}_v^{kt} &= \widehat{A}_{j,l=1}^{kt} (1 - \widehat{f}_{j,l=1}^{kt}) + \widehat{A}_{j,l=2}^{kt} \widehat{f}_{j,l=2}^{kt} + \widehat{A}_{j,l=3}^{kt} \quad \text{for } v \in \{X, Z\}, \\ \widehat{f}_v^{kt} &= \widehat{A}_{j,l=1}^{kt} + \widehat{A}_{j,l=3}^{kt} (1 - \widehat{f}_{j,l=3}^{kt}) + \widehat{A}_{j,l=4}^{kt} \widehat{f}_{j,l=4}^{kt} \quad \text{for } v \in \{W, Y\}. \end{aligned}$$

A different way of predicting the allele frequency on both sex chromosomes is to assign it to be the frequency corresponding to the most probable subtype.

This information can be used to infer the consensus sequences of the X and Y sequences. For a given contig (that can be chosen on the basis of $\widehat{S}_{N_j}^k$, but not necessarily if we have other reasons to believe the contig is sex-linked), each polymorphic site can be considered fixed for an allele if \widehat{f}_X or \widehat{f}_Y are above a threshold U_f ($0.5 \leq U_f \leq 1$) or below $1 - U_f$. A further threshold can be applied to genotype non-fixed sites: if \widehat{f}_X or \widehat{f}_Y are above a threshold u_f ($0.5 \leq u_f \leq U_f$) or below $1 - u_f$.

Nucleotide diversity can be calculated as

$$\pi_v^k = \frac{1}{\tau_k} \sum_t 2 \widehat{f}_v^{kt} (1 - \widehat{f}_v^{kt})$$

where $\tau_k \geq T_k$ is the total length of the contig k , including monoallelic sites. The divergence is

$$D_{XY}^k = \frac{1}{\tau_k} \sum_t \left(\widehat{f}_X^{kt} (1 - \widehat{f}_Y^{kt}) + \widehat{f}_Y^{kt} (1 - \widehat{f}_X^{kt}) \right)$$

in the XY case; extension to ZW chromosomes is trivial.

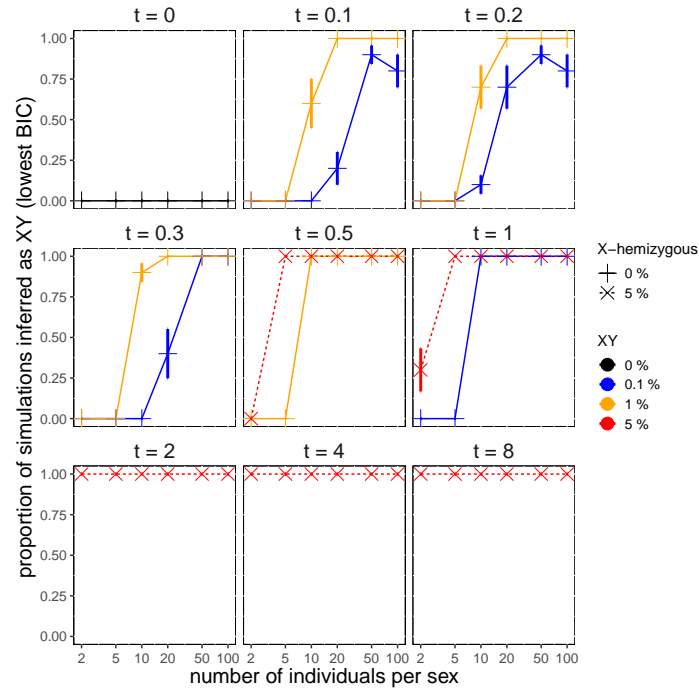


Figure S1 Model choice by SDpop on simulated data. The proportion of simulations for which the XY model had the lowest BIC is indicated; each combination of simulation parameter values was repeated 10 times from a random seed. Vertical bars indicate the expected variance based on the binomial distribution. Different panels represent result for different values of the time since recombination suppression t ; the percentage of simulated X-hemizygous genes is indicated by the line types and symbols (solid lines with "+" symbols indicate no X-hemizygous genes; dashed lines with "x" symbols 5% of X-hemizygous genes); the colors indicate the percentage of XY gametologous genes (black 0%, blue 0.1%, orange 1% and red 5%). These simulations were carried out with error rate $e = 0.001$; results with $e = 0.0001$ are shown in Figure 1.

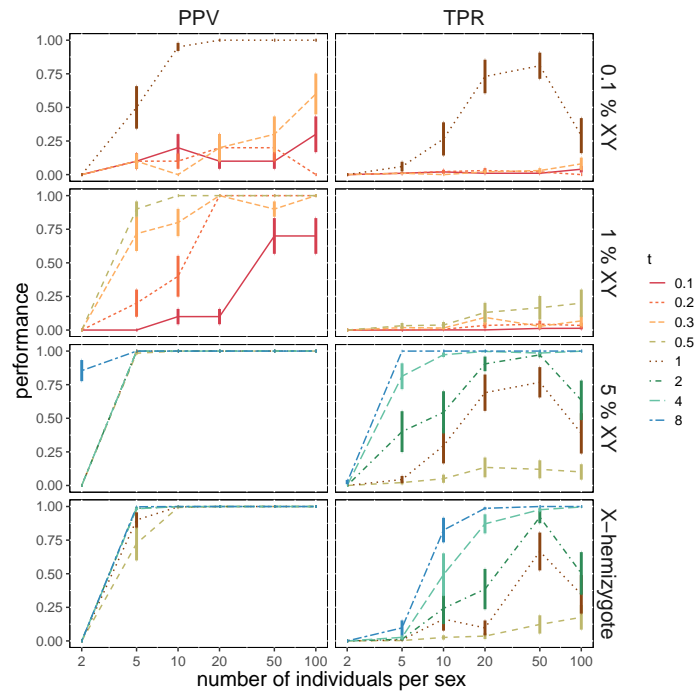


Figure S2 Precision (Positive Predictive Value, left) and power (True Positive Rate, right) of the detection of sex-linked contigs in simulated data, using a threshold for the posterior probability of 0.8. First three rows: XY gametologs, grouped by the proportion of simulated gametologs in the genome (0.1%, 1%, 5%). Bottom graphs: X-hemizygous genes, for which the simulated proportion in the genome was 5%. The color and line scales indicate the simulated time since recombination suppression t . Each point is the average of 100 simulations, with the bars representing the standard error. For all cases shown here, the simulated error rate was 0.001; for $e = 0.0001$, see Figure 2.

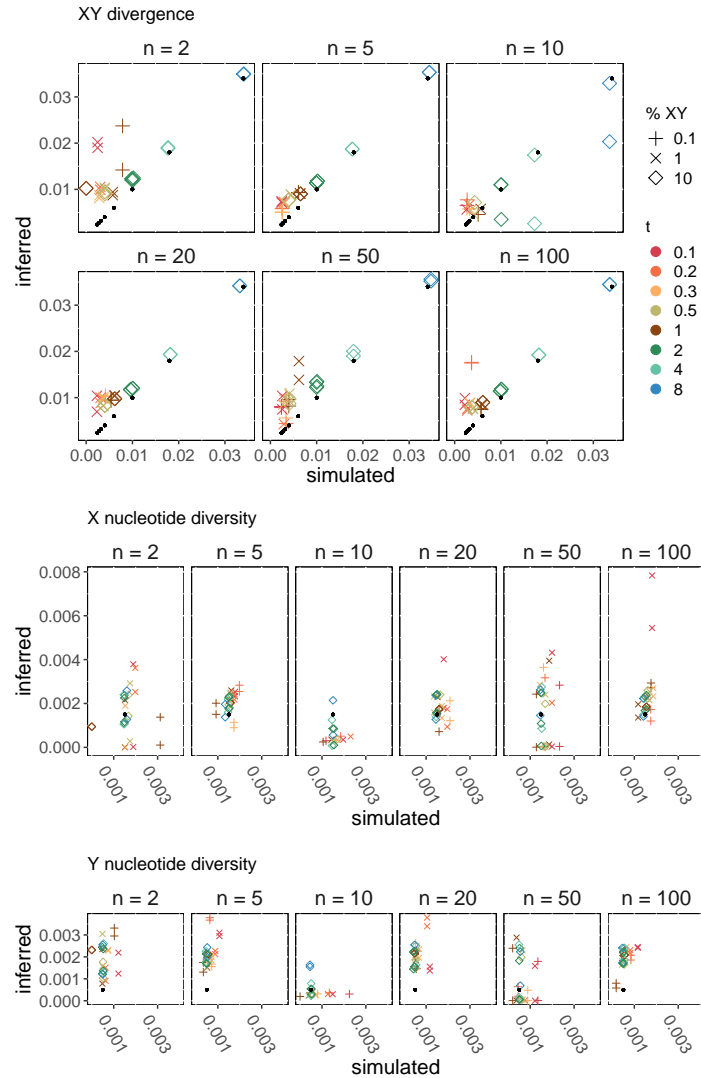


Figure S3 Population genetic inferences of SDpop at higher error rates than Figure 3; here, $e = 0.001$. The values of nucleotide diversity and divergence calculated directly from the simulation results are compared to the values inferred from SDpop's output. Comparisons are based on SDpop's assignment of the genes (i.e. all genes with a posterior probability > 0.8 were used). Top: gametolog divergence (D_{XY}); middle: X nucleotide diversity (π_X); bottom: Y nucleotide diversity (π_Y). Facets are separated by the number of individuals per sex used (n). Color indicates the time since recombination suppression t , and symbols the simulated proportion of gametologs %XY. The black points indicate the theoretical values (D_{XY} : one for each t ; π_X and π_Y : one value for all runs).

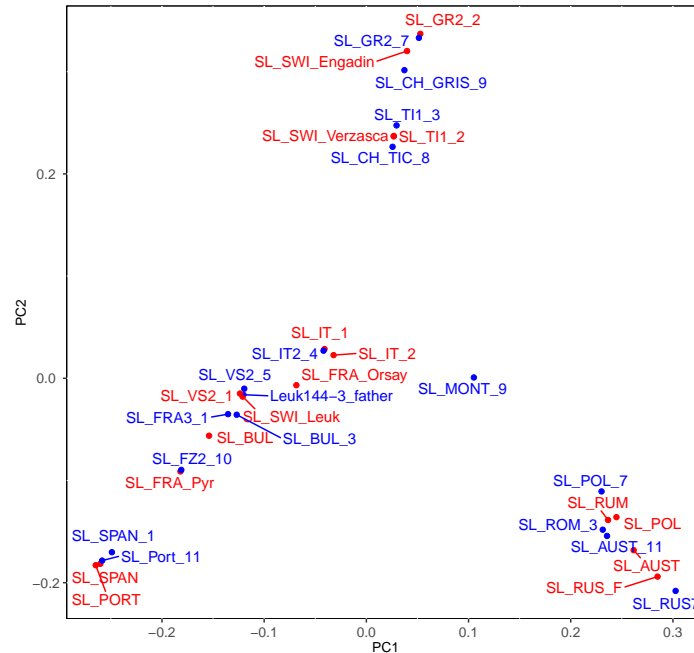


Figure S4 Principal Component Analysis of genetic variation in 34 plants of *Silene latifolia*. Females are colored red, males blue. The exact locations of sampling are given in [Muyle et al. \(2020\)](#). Leuk144-3_father, SL_SWI_Engadin, SL_SWI_Leuk, SL_SWI_Verzasca, SL_CH_GRIS_9, SL_CH_TIC_8, SL_GR2_2, SL_GR2_7, SL_TI1_2, SL_TI1_3, SL_VS2_1 and SL_VS2_5 are from Switzerland; SL_AUST and SL_AUST_11 from Austria; SL_BUL and SL_BUL_3 from Bulgaria; SL_FRA_Orsay, SL_FRA_Pyr, SL_FRA3_1 and SL_FZ2_10 from France; SL_IT_1, SL_IT_2 and SL_IT2_4 from Italy; SL_POL and SL_POL_7 from Poland; SL_PORT and SL_Port_11 from Portugal; SL_RUM and SL_ROM_3 from Romania; SL_RUS_F and SL_RUS7 from Russia; SL_SPAN and SL_SPAN_1 from Spain; and SL_MONT_9 from Montenegro. The central cluster, used as a subsample to test SDpop, extends from SL_FRA_Pyr (lower left) to SL_IT_1 (upper right). The reason for the Bulgarian plants to be located in this cluster remain obscure, but we chose to retain them, as well as to exclude one Italian female (SL_IT_2) to obtain a balanced sample of six females and six males.

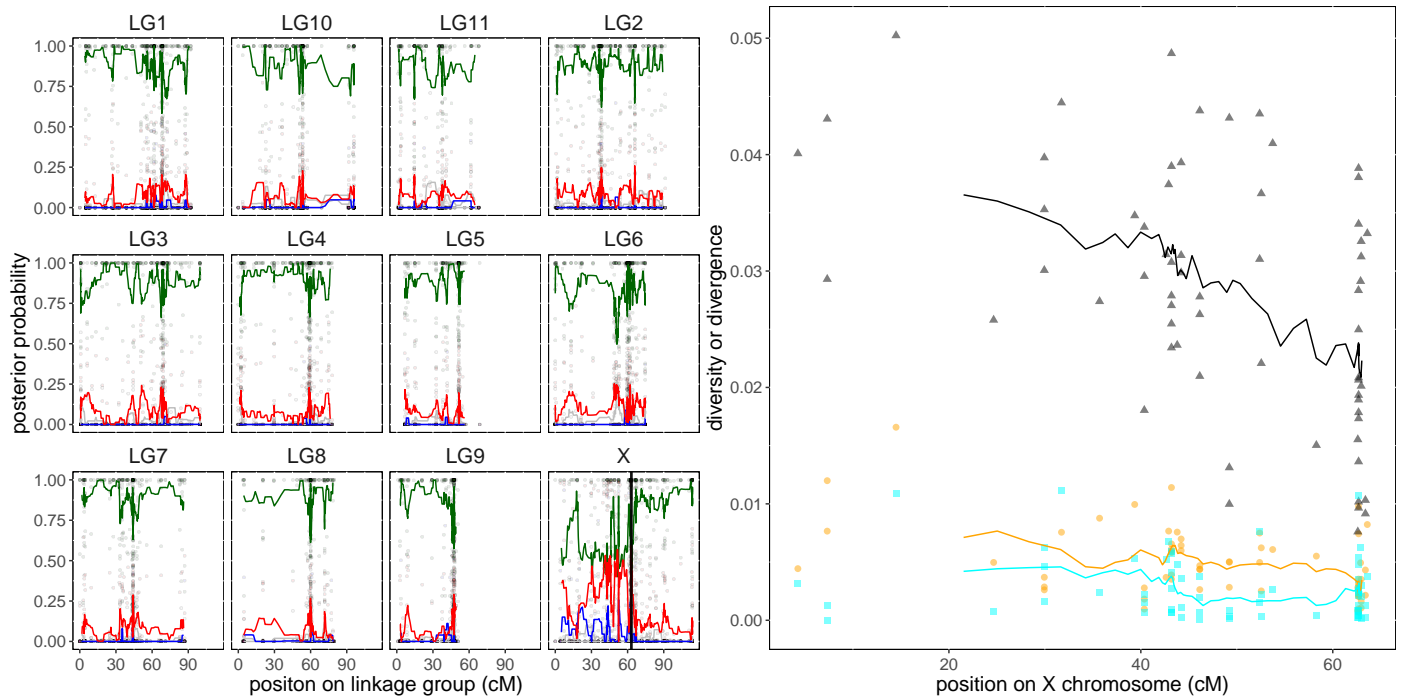


Figure S5 SDpop's inferences of sex-linkage in *Silene latifolia*, using the 12 individuals of the "central cluster" (Figure S4). Contigs were placed on the genetic map of Papadopulos *et al.* (2015). Left panels: posterior probabilities for all placed contigs: autosomal segregation in green, x-hemizyosity in blue, and XY gametology in red; the (uninformative) haploid and paralogous segregation types are indicated in grey. Lines represent running averages, using sliding windows of 10 contigs. The "fuzzy boundary" between the non-recombining region and the pseudoautosomal region on the X chromosome (Krasovec *et al.* 2020) is indicated by the horizontal line. Right panel: predicted divergence (black triangles) and nucleotide diversity of X and Y copies (orange circles and cyan squares) based on SDpop's output. The lines are the running averages over 10 genes.

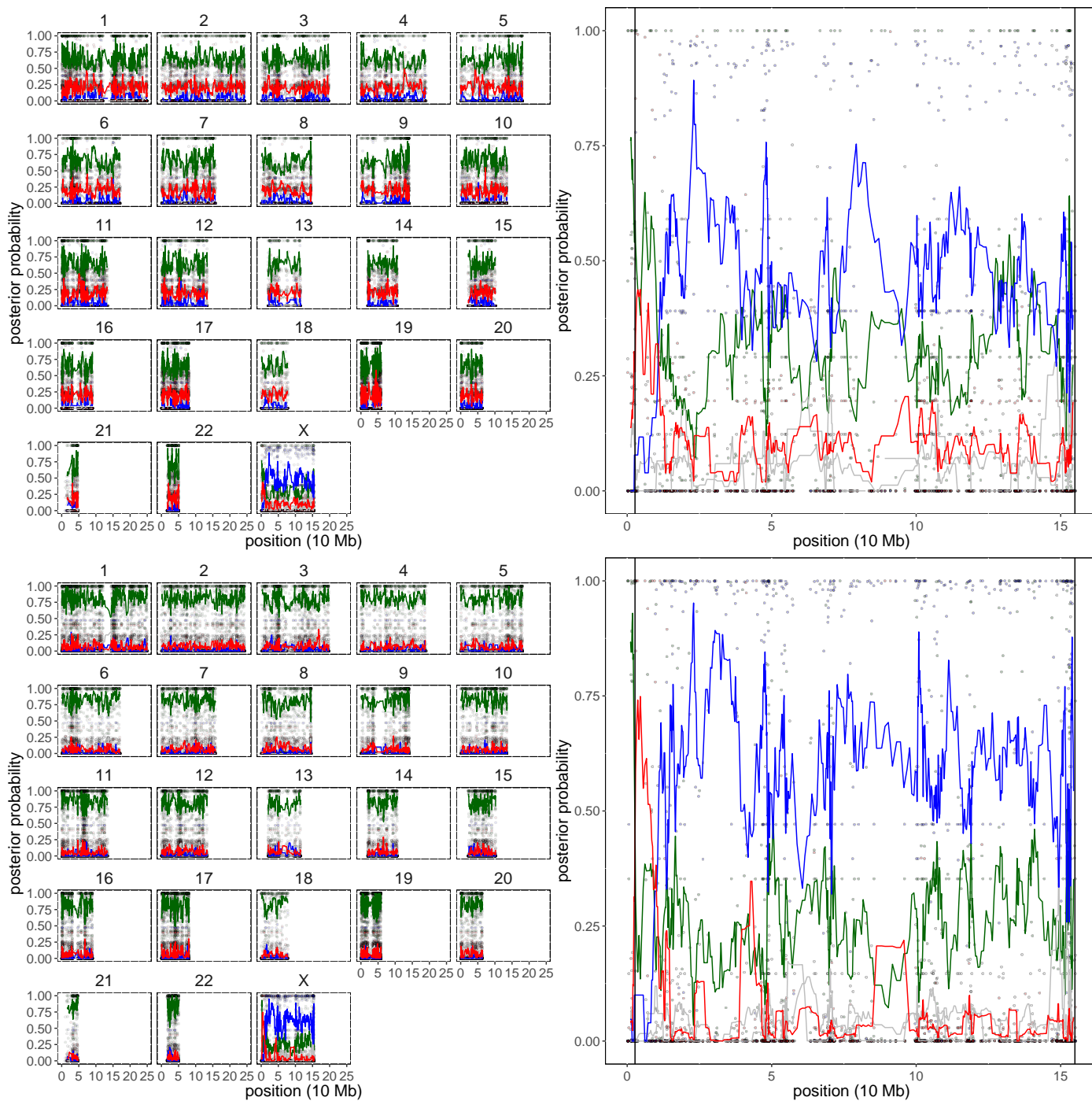


Figure S6 Test of SDpop’s performance on the human exome-targeted sequencing data from the 1000 genome project, using 5 individuals per sex (top graphs) or 20 (bottom graphs). Smoothed gene-level posterior probabilities for autosomal (black), X-hemizygote (blue) and XY (red) segregation are shown; haploid and paralogous posterior probabilities are indicated in gray. The right panels show the results on the X chromosome: the extremities corresponding to the pseudo-autosomal regions are predicted to be autosomal by SDpop, while most XY gametologous genes are found close to the pseudo-autosomal region on the left arm, which represents the youngest stratum where Y copies have not yet been lost. The rest of the chromosome consists of X-hemizygous genes, that lack a Y copy.