

Next-generation inference of past population history by integrating diverse types of genomic markers

Thibaut Sellinger^{1,2}, Frank Johannes³, Aurélien Tellier^{2*}

¹ Department of Environment and Biodiversity,
Paris Lodron University of Salzburg

² Professorship for Population Genetics,
Department of Life Science Systems,
Technical University of Munich

³ Professorship for Plant Epigenomics,
Department of Molecular Life Sciences,
Technical University of Munich

* Corresponding author, aurelien.tellier@tum.de

Abstract

1
2 With the availability of high quality full genome polymorphism
3 (SNPs) data, it becomes feasible to study the past demographic and
4 selective history of populations in exquisite detail. However, such
5 inferences still suffer from a lack of statistical resolution for recent,
6 e.g. bottlenecks, events, and/or for populations with small nucleotide
7 diversity. Additional heritable (epi)genetic markers, such as indels,
8 transposable elements, microsatellites or cytosine methylation, may
9 provide further, yet untapped, information on the recent past popula-
10 tion history. We extend the Sequential Markovian Coalescent (SMC)
11 framework to jointly use SNPs and other hyper-mutable markers. We
12 are able to 1) improve the accuracy of demographic inference in recent
13 times, 2) uncover past demographic events hidden to SNP-based infer-
14 ence methods, and 3) infer the hyper-mutable marker mutation rates
15 under a finite site model. As a proof of principle, we focus on demo-
16 graphic inference in *A. thaliana* using DNA methylation diversity data
17 from 10 European natural accessions. We demonstrate that segregat-
18 ing Single Methylated Polymorphisms (SMPs) satisfy the modelling
19 assumptions of the SMC framework, while Differentially Methylated
20 Regions (DMRs) are not suitable as their length exceeds that of the
21 genomic distance between two recombination events. Combining SNPs
22 and SMPs while accounting for site- and region-level epimutation pro-
23 cesses, we provide new estimates of the glacial age bottleneck and post
24 glacial population expansion of the European *A. thaliana* population.
25 Our SMC framework paves the way for next generation demographic
26 and selection inference by combining information from several herita-
27 ble (epi)genomic markers.

28 **Keywords**— Kingman coalescent, Sequentially Markovian Coalescent, ances-
29 tral recombination graph, epigenetics, hidden markov model

30 Introduction

31 A central goal in population genetics is to reconstruct the evolutionary history
32 of populations from patterns of genetic variation observed in the present. Rele-
33 vant aspects of these histories include past demographic changes as well as sig-
34 natures of selection. Inference methods based on Deep Learning (DL, [37]), Ap-
35 proximate Bayesian Computation (ABC, [9]) or Sequential Markovian Coalescent
36 (SMC, [39, 50]) aim to infer this information directly from full genome sequencing
37 data, which is becoming rapidly available for many (non-model) species due to
38 decreasing costs. The SMC, in particular, offers an elegant theoretical framework
39 that builds on the classical Wright-Fisher and the backward-in-time Kingman coa-
40 lescent stochastic models (*e.g.* [35, 12, 67]). Both models conceptualize Mendelian
41 inheritance as generating the genealogy of a population (or a sample), that is, the
42 unique history of a fragment of DNA passing from parents to offspring. When this
43 genealogy includes the effect of recombination, it is called the Ancestral Recombi-
44 nation Graph (ARG, [26, 71]).

45
46 Under the Kingman coalescent model, the true genealogy of a population (or
47 sample) is defined by its topology and branch length, and contains the information
48 on past demographic changes and life history traits [45, 55, 60, 62] as well as selec-
49 tive events [12, 67]. The genealogical and the mutational processes of any heritable
50 marker can therefore be disentangled, and the frequency of any given marker state
51 is given by the shape of the genealogy in time (see Figure 1A). A central assumption
52 about heritable genomic markers is that they are generated by two homogeneous
53 Poisson mutation processes along the genome as well as through time. This entails
54 that mutations in different genealogies are independent due to the effect of recom-
55 bination [71, 43], and that there are no time periods with a large excess, or a severe
56 lack, of mutations along a genealogy (mutations are independently distributed in
57 time within a DNA fragment). In other words, the frequency of polymorphisms
58 at DNA markers observed across a sample of sequences are constrained by, as well
59 as inform on, the underlying genealogy at this locus (Figure 1A). To clarify these
60 assumptions, we present a schematic representation of a marker 1 (yellow in Figure
61 1) which fulfills both homogeneous Poisson processes in time and along the genome.
62 We also present cases applicable to a second genomic marker 2 that violates the
63 model assumptions, namely by not being heritable (Figure 1B) or not following a
64 non-homogeneous Poisson process in the genome (Figure 1C) or in time (Figure
65 1D).

66

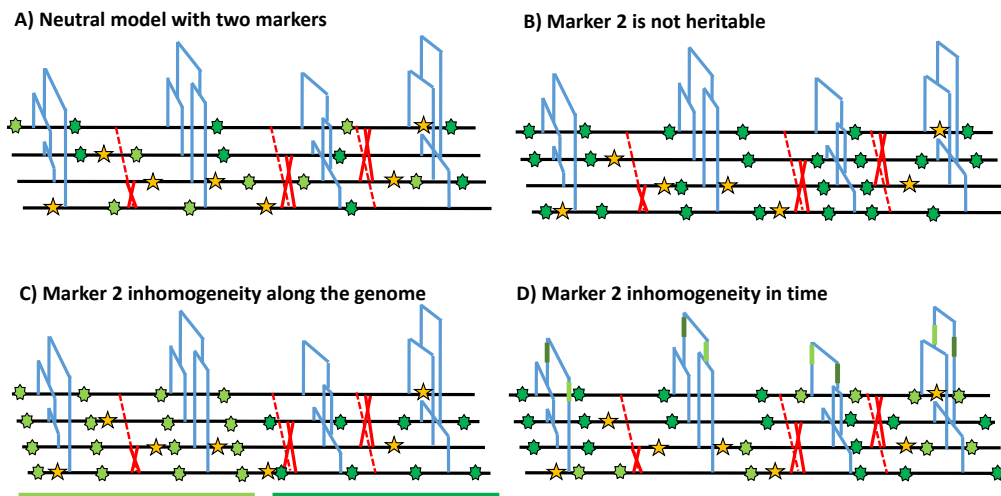


Fig. 1. Schematic distribution of two markers along the genealogy and four genomes. A) Schematic distribution of marker 1 (yellow star) and marker 2 (green star) along the genealogies in a sample of four genomes both following a homogeneous Poisson process. B) The green marker 2 is not heritable, so that its distribution is independent from the genealogy. C) The green marker 2 is spatially structured along the genome, violating the distribution of the Poisson process along the genome and conflicting with the genealogy. D) The green marker 2 does not follows Poisson process through time, *e.g.* burst of mutations at a specific time point represented by given branches of the genealogies in green. The yellow marker 1 has an identical Poisson process along the genome and the genealogy in all four panels, and for readability, marker 2 exhibits light and dark green states.

67 Despite the power of the SMC, well-known model violations such as variation
68 of recombination and mutation rates along the genome [5, 4] or pervasive selection
69 [53, 31, 30] can compromise the accuracy of demographic and selective inference
70 [24, 56]. There are two other important issues that have received less attention in
71 the literature. The first issue occurs when the population recombination rate (ρ)
72 is higher than the population mutation rate (θ). In such cases, inferences can be
73 biased if not erroneous [63, 56, 55], because several recombination events cannot
74 be inferred due to the lack of Single Nucleotide Polymorphisms (SNPs for point
75 mutations). This problem affects many species, though interestingly not humans
76 which have a ratio $\rho/\theta \approx 1$. A second issue occurs when the mutational process
77 along the genealogy is too slow be informative about sudden and strong variation
78 in population size (*i.e.* population bottlenecks), such as during colonization events
79 of novel habitats. The typical low mutation rate of 10^{-9} up to 10^{-8} (per base, per

80 generation) found in most species therefore places strong limitations on SMC anal-
81 ysis of recent bottleneck events (up to ca. 10^{-4} generations ago) when inference is
82 based solely on SNP data. Indeed, bottlenecks are often either not found, or when
83 inferred, their timing and magnitude are not well estimated (inferred smoother
84 than in reality, [31, 56]), even when a large number of samples is used. A typical
85 example is the large uncertainty of the timing and magnitude of the population
86 size bottleneck during the Last Glacial Maximum (LGM) and post-LGM expan-
87 sion in *A. thaliana* European populations based on several studies using different
88 accessions and SMC inference methods [2, 18].

89
90 Nonetheless, current SMC, DL or ABC inference methods making use of full
91 genome sequence data rely almost exclusively on SNPs for inference [50, 63, 55,
92 9, 36]. There are both practical and theoretical reasons for using SNPs: They are
93 easily detectable from short-read re-sequencing data and their mutational process
94 is well approximated by the infinite site model [12, 67], simplifying the inference of
95 the underlying genealogy. However, other heritable genomic markers exists whose
96 mutation rates can be several orders of magnitude higher than that of SNPs, and
97 could thus be more informative about recent demographic events. These include
98 microsatellites, insertions, deletions and transposable elements (TEs). Current
99 technological limitations still impede the easy detection and estimation of allele
100 frequencies for many of these markers [73, 48, 68]. For example, identifying inser-
101 tion/excision variation of transposable elements (TEs) or or copy number variation
102 of microsatellites requires a high quality reference genome and ideally long-read se-
103 quencing approaches [48]. In addition to these genomic markers, DNA cytosine
104 methylation is emerging as a potentially useful epigenetic marker for phylogenetic
105 inference in plants [75, 76]. Stochastic gains and losses of DNA methylation at
106 CG dinucleotides, in particular, arise at a rate of ca. 10^{-4} up to 10^{-5} per site per
107 generation (that is 4 to 5 orders of magnitude faster than DNA point mutations,
108 [65]), and can be inherited across generations [70]. These so-called spontaneous
109 epimutations are likely neutral at the genome-wide scale ([66, 29], but see [44]),
110 and can be easily detected from bisulphite converted short read sequencing data
111 [40, 52]. Recent work suggests that CG methylation data can be used as a molec-
112 ular clock for timing divergence between pairs of lineages over timescales ranging
113 from years to decades [76].

114
115 However, theoretical integration of the above-mentioned (epi)genomic markers
116 into a population genomics and SMC inference framework is not trivial. Because of
117 the high mutation rate, the mutational process at these (hyper-mutable) markers
118 is reversible and more consistent with a finite site, rather than infinite site, model,
119 which can result in extensive homoplasy (as known for microsatellite markers, [16]).
120 Indeed, classic expectations of population genetics diversity statistics, mostly build
121 for SNPs, need to be revised for these hyper-mutable markers [13, 69]. Here we
122 develop the theoretical and methodological inference framework for the inclusion

123 of additional (potentially hyper-mutable) markers into the SMC. We showcase our
124 model using extensive simulations as well as application to published DNA cytosine
125 methylation data from local populations of *A. thaliana* ([52, 66]). We demonstrate
126 that integration of hyper-mutable genomic markers into SMC models significantly
127 improves the inference accuracy of past variation of population size, or can even un-
128 cover demographic events not uncovered using SNPs alone. Our proof-of-principle
129 approach opens up novel avenues for studying population genetic processes over
130 time-scales that have been largely inaccessible using traditional SNP-based ap-
131 proaches. This may prove particularly useful when exploring recent demographic
132 changes of endangered species as a way to assess their potential for extinction in
133 the context of biodiversity loss and global change.

134 Results

135 Theoretical results with two markers underlying the SMC 136 computations

137 We study polymorphic sites across genomes of several sampled individuals which
138 exhibit several possible markers (DNA nucleotides, methylation, TEs, indels, mi-
139 crosatellites,...). We define any marker by 1) its maximum number of possible
140 states (nb_s), for example nucleotide sites have four states (A, T, C and G) while a
141 methylation site has two states (methylated or unmethylated), and 2) its mutation
142 rate μ , *i.e.* the rate at which the state of a marker changes into another state per
143 position and per generation [3]. More specifically, we are interested in two rates:
144 the DNA mutation rate for changes in DNA nucleotides, and epimutation rate for
145 change in methylation state. Furthermore, we assume that at each position on
146 the genome only one type of marker can occur and be observed. We obtain as a
147 first theoretical result the probability for a given site in the genome to be identical
148 ($P(id)$) or segregating ($P(seg)$) (*i.e.* polymorphic) in a sample of size two ($n = 2$,
149 two sampled chromosomes are compared):

$$\begin{aligned} P(id, n = 2) &= \frac{1}{nb_s} + \frac{(nb_s - 1)}{nb_s} e^{-2\mu t_M \frac{(nb_s)}{(nb_s - 1)}} \\ P(seg, n = 2) &= \frac{(nb_s - 1)}{nb_s} - \frac{(nb_s - 1)}{nb_s} e^{-2\mu t_M \frac{(nb_s)}{(nb_s - 1)}} \end{aligned} \quad (1)$$

150 This probability is a function of the time to the most recent common ances-
151 tor (TMRCA in text and t_M in equation 1, details in Supplementary Text). The
152 probability for a mutation to occur for a given marker increases with an increased
153 TMRCA [12, 67], but under high mutation rates the marker may not be polymor-
154 phic in the sample as mutations may be reversed (so-called homoplasy, [16, 13]). In
155 Figure 2 we illustrate these properties by computing the probability 1 for different
156 mutation rates. The inference of recent demographic events and bottlenecks do rely

157 on the presence of polymorphic sites to detect recent coalescent event (TMRCA),
158 and should be improved by using markers with high (or fast) mutation rate (*e.g.*
159 hyper mutable).

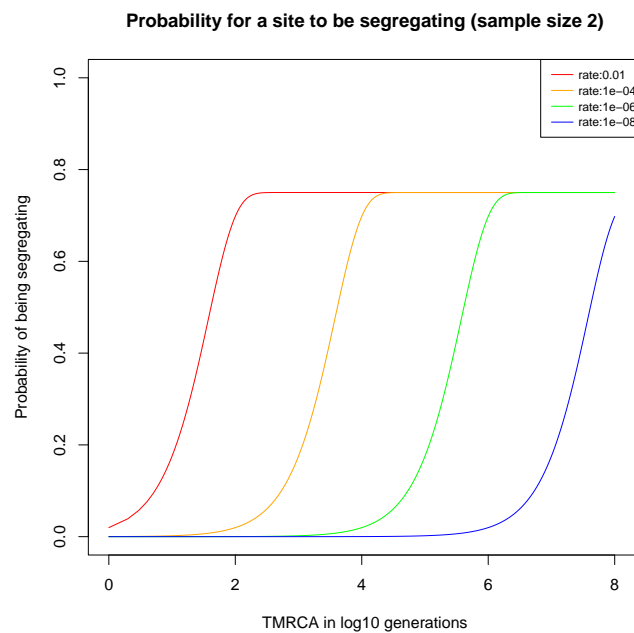


Fig. 2. Probability of a site to be segregating in a sample of size two for different mutation rates. The probability for a site to be segregating in a sample of size two under different mutation rates: 10^{-2} in red, 10^{-4} in orange, 10^{-6} in green and 10^{-8} in blue. The marker is assumed here to have $nb_s = 4$ possible states.

160 In the following, we simulate data under different demographic scenarios using
161 the sequence simulator program *msprime* [6, 33], which generates the ARG of n
162 sampled diploid individuals (set to $n = 5$ throughout this study, leading to 10
163 haploid genomes). This ARG contains the genealogy of a given sample at each
164 position of the simulated chromosomes. We then process the ARG to create DNA
165 sequences according to the model parameters and the type of marker considered.
166 We first assume a set of genomic markers obtained for a sample size n , and mut-
167 ating according an homogeneous Poisson process along the genome and in time
168 (along the genealogy) as in Figure 1A. To simulate the sequence data, we define
169 the number of marker types (any number between 1 and the sequence length) and
170 the proportion of sites of each marker type in the sequence. Each marker is char-
171 acterized by both parameters nb_s and μ . For simplicity, we simulate sequences
172 with two markers, but note that the method can be easily extent to additional
173 markers. Marker 1 represents 98% of the sequence, and has a per site mutation

174 rate $\mu_1 = 10^{-8}$ mimicking nucleotide SNP markers under an infinite site model
175 (thus considered as bi-allelic at a given DNA site, [74]). By contrast, marker 2
176 composes the complementary 2% of the sequence length, with a per site mutation
177 rate of $\mu_2 = 10^{-4}$ per generation between two possible states. Marker 2 is thus
178 hyper-mutable compared to marker 1 and mimics methylation/epimutation sites.
179 Note, that mutation events in Marker 1 and 2 are simulated under a finite site
180 model.

181

182 We use different SMC-based methods throughout this study. These methods
183 include: 1) MSMC2 used as a reference method [19], 2) SMCtheo is an extension
184 of the PSMC' [39, 50] accounting for any number of heritable theoretical mark-
185 ers, and 3) eSMC2 which is equivalent to SMCtheo but accounting only for SNPs
186 markers [56] (to avoid any bias in implementation differences between SMCtheo
187 and MSMC2). All methods are Hidden Markov Models (HMM) derived from the
188 Pairwise Sequentially Markovian Coalescent (PSMC') [50] and assume neutral evo-
189 lution and a panmictic population. The hidden states of these methods are the
190 coalescence time of a sample of size two at a position on the sequence. From the
191 distribution of the hidden states along the genome, all methods can infer population
192 size variation through time as well as the recombination rate [50, 19, 56].

193 **The inclusions of hyper-mutable genomic markers im-** 194 **proves demographic inference**

195 We assume that the mutation rate of marker 1 is $\mu_1 = 10^{-8}$ per generation per
196 bp. We use this information to estimate the mutation rate of marker 2, which
197 we vary from $\mu_2 = 10^{-8}$ to $\mu_2 = 10^{-2}$ per generation per bp. The estimation
198 results based on simulated data under a constant population size of $N = 10,000$
199 are displayed in Table 1. We find that our approach is capable of inferring μ_2 with
200 high accuracy for rates up to $\mu_2 = 10^{-4}$. However, when the mutation rate μ_2 is
201 10^{-2} , our approach underestimates it by a factor three, suggesting the existence of
202 an accuracy limit. To demonstrate that information can be gained by integrating
203 marker 2 (with $\mu_2 = 10^{-4}$), we compared the ability of several inference methods to
204 recover a recent bottleneck (Figure 3A). All methods correctly infer the amplitude
205 of population size variation. When accounting only for marker 1 (with $\mu_1 = 10^{-8}$,
206 MSMC2 and eSMC2 fail to infer accurately the sudden variation of population size.
207 However, with the inclusion of hyper-mutable marker 2, our SMCtheo approach
208 correctly infers the rapid change of population size of the bottleneck (Figure 3A,
209 green). It is encouraging that an accurate estimation of the demography is ob-
210 tained, even when the mutation rate of marker 2 is unknown (Figure 3A, blue).

211

True μ_2 value	Estimated value of μ_2
10^{-8}	9.9×10^{-9} (0.02)
10^{-6}	1.0×10^{-6} (0.008)
10^{-4}	1.4×10^{-4} (0.01)
10^{-2}	3.05×10^{-3} (0.41)

Table 1: Average estimated values of the mutation rate of marker 2 (μ_2), knowing that of marker 1. We use 10 sequences of 100 Mb ($r = \mu_1 = 10^{-8}$ per generation per bp) under a constant population size fixed to $N = 10,000$. The coefficient of variation over 10 repetitions is indicated in brackets.

212 Furthermore, some species or populations might feature small effective popu-
213 lation sizes (ca. $N = 1,000$), potentially resulting in reduced genomic diversity.
214 In such cases the inclusion of hyper-mutable markers should also improve demo-
215 graphic inference. We present the results of such a scenario in Figure 3B, where
216 the population size was divided by a factor 10 compared to the previous scenario in
217 Figure 3A. We find that in the absence of the hyper-mutable marker 2, no approach
218 can correctly infer the variation of population size. From the shape of the inferred
219 demography, methods using only marker 1 do not suggest the existence of a bottle-
220 neck followed by recovery (the "U-shaped" demographic scenario is not apparent
221 with the orange and red lines, Figure 3B). Yet, when integrating both markers,
222 the population size can be recovered, even if the mutation rate of marker 2 is not
223 *a priori* known. In both Figure 3A and B, we assume that the marker 2 occurs
224 at a frequency of 2% in the genome. This percentage may be unrealistically high
225 depending on the marker and the species. To test the impact of reducing marker 2
226 frequency, we repeat the simulations shown in Figure 3A, but set its frequency to as
227 low as 0.1% (a 20-fold reduction). We find that the inclusion of the hyper-mutable
228 marker 2 continues to improve inference accuracy in very recent times, albeit less
229 pronounced than in Figure 3A (see Supplementary Figure 1). This suggests that a
230 very small proportion of hyper-mutable genomic sites is sufficient to significantly
231 improve the accuracy of inferences.

232

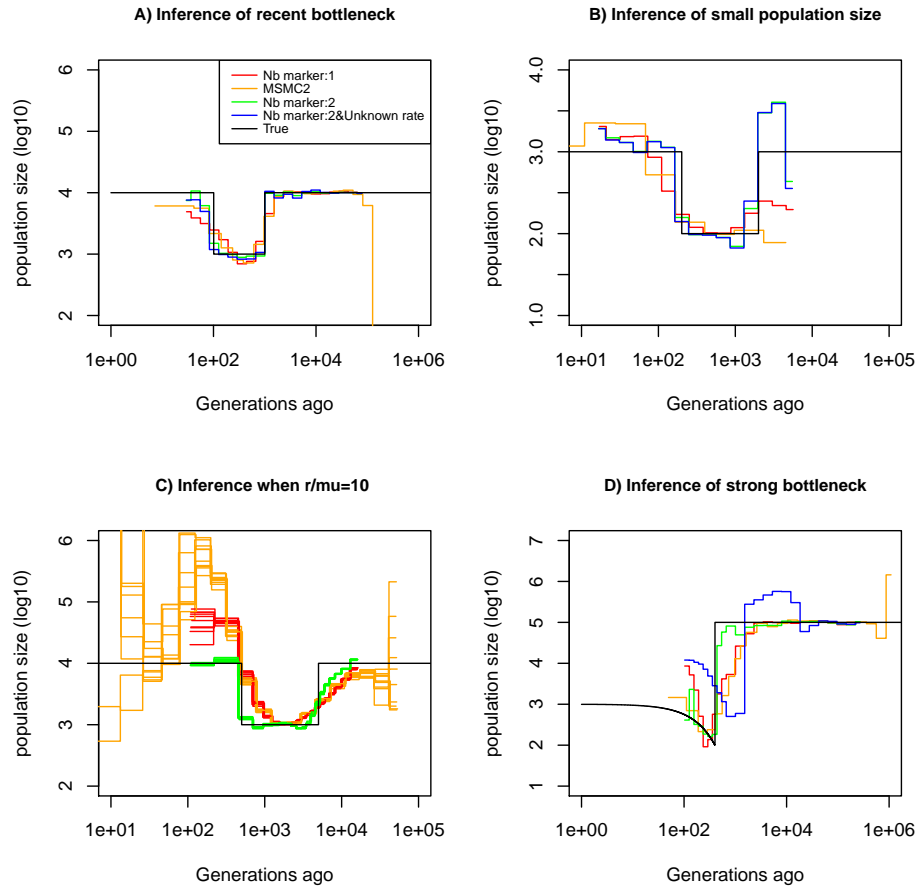


Fig. 3. Performance of SMC approaches using different markers. Estimated demographic history of a bottleneck (black line) by SMC approaches using two genomic markers. In orange and red, are the estimates by MSMC2 and eSMC2 based on only marker 1. Estimates from SMCtheo integrating both markers are in green (with known μ_2), and in blue with unknown μ_2 . The demographic scenarios are A) 10-fold recent bottleneck with an ancestral population size $N = 10,000$, B) 10-fold recent bottleneck with an ancestral population size $N = 1,000$, C) 10-fold bottleneck with an ancestral population size $N = 10,000$, and D) a very severe (1,000 fold) and very recent bottleneck with incomplete size recovery. In A, B and D, we assume $r/\mu_1 = 1$ (with $r = \mu_1 = 10^{-8}$, $\mu_2 = 10^{-4}$ per generation per bp) and in C, $r/\mu_1 = 10$ (with $r = 10^{-7}$, $\mu_1 = 10^{-8}$, and $\mu_2 = 10^{-4}$ per generation per bp).

233 All full genome inference methods, especially SMC approaches, display lower
 234 accuracy when the population recombination rate ($\rho = 4Nr$) is larger than the
 235 population mutation rate of marker 1 ($\theta_1 = 4N\mu_1$). We simulate sequence data

236 under a bottleneck scenario slightly more ancient than in Figure 3A and assume
237 that $\rho/\theta_1 = r/\mu_1 = 10$ and $\rho/\theta_2 = r/\mu_2 = 10^{-3}$. Our results show that by
238 integrating the genomic marker 2 which mutation rate is larger than the recom-
239 bination rate, estimates of the recombination rate as well as past population size
240 variation are substantially improved (Table 2, Figure 3C). Indeed, analyzing only
241 marker 1, eSMC2 and MSMC2 fail to infer the sudden variation of population size,
242 overestimate the population size in recent times (Figure 3D). By integrating the
243 hyper-mutable marker 2, our SMCtheo approach correctly infers the strength and
244 time of the bottleneck when μ_1 and μ_2 are known (Figure 3D, green line), while
245 the timing of the bottleneck is slightly shifted in the past when μ_2 is unknown and
246 estimated by our method (Figure 3D, blue line). Using only marker 1, the esti-
247 mates of the recombination rate are inaccurate (Table 2 under various demographic
248 scenarios in Supplementary Figure S2). We further improve the accuracy of esti-
249 mation by optimizing the likelihood (LH) to estimate the recombination rate and
250 demography compared to the classically used Baum-Welch (BW) algorithm (Table
251 2). Our results demonstrate that SNPs are limiting and insufficient for accurate
252 inferences in recent times and that the inclusion of an additional marker with mu-
253 tation rate higher than the recombination rate generates significant improvements
254 in demographic inference. However, by directly optimizing the likelihood the true
255 recombination rate can be well recovered even with Marker 1 only.

256

Method	True recombination rate	Average estimated recombination rate
MSMC2 (BW)	10^{-7}	0.23×10^{-7} (0.017)
1 Marker : BW	10^{-7}	0.25×10^{-7} (0.012)
2 Marker : BW	10^{-7}	0.90×10^{-7} (0.004)
1 Marker : LH	10^{-7}	0.84×10^{-7} (0.036)
2 Marker : LH	10^{-7}	0.94×10^{-7} (0.01)

Table 2: Estimates of recombination rates with one or both markers. For SMCtheo, BW stands for the use of the Baum-Welch algorithm to infer parameters, and LH to the use of the likelihood. We use 10 sequences of 100 Mb with $r = 10^{-7}$, $\mu_1 = 10^{-8}$ and $\mu_2 = 10^{-4}$ per generation per bp in a population with a past bottleneck event. The coefficient of variation over 10 repetitions is indicated in brackets.

257 Integrating DNA methylation improves the accuracy of 258 inference

259 Definition of the theoretical model for DNA methylation

260 Following the previously encouraging results of demographic inference with SNPs
261 and an hyper-mutable marker under the specific assumptions of Figure 1A, we de-
262 velop a specific SMCm method to jointly analyse SNPs and cytosine methylation as
263 an epigenetic hyper-mutable marker. We focus here on methylation located in CG

264 contexts within genic regions as these are more likely to evolve neutrally [66, 75, 76].
265 The methylation of individual CG dinucleotides presents a biallelic heritable marker
266 with a finite number of (epi)mutable sites (Figure 4). In a sample of several se-
267 quences from a population, variation in the methylation status of individual CGs
268 is known as single methylation polymorphism (SMP, Figure 4A) which could be
269 used for demographic and divergence inference [65, 66]. However, CG methylation
270 sites can also be organized in spatial clusters (of similar state) due to region level
271 epimutation (Figure 4B, [70, 15, 44]. Region level epimutations can have differ-
272 ent epimutation rates than individual CG sites. Population-level variation in the
273 methylation status of these clusters is known as differentially methylated regions
274 (DMRs). Furthermore, when integrating SMP and DMR epimutational processes
275 (*i.e.* what we here call region level epimutation), the methylation status of CG
276 sites is therefore affected by the superposition of both processes. Therefore the
277 simulation and modeling of epimutational processes of SMPs is more complex than
278 in our previous model as we need to account for the effect of region methylation
279 as well as for methylation and demethylation epimutation rates to be different and
280 asymmetrical [65, 15].
281

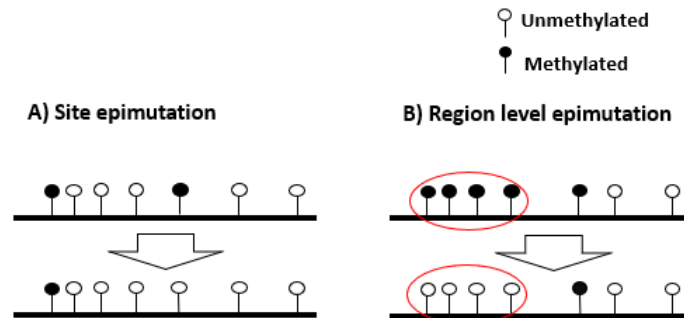


Fig. 4. Schematic representation of site and region epimutations
Schematic representation of a sequence undergoing epimutation at A) the cyto-
sine site level, and B) at the region level. A methylated cytosine in CG context is
indicated in black and an unmethylated cytosine in white.

282 To make our simulations realistic, we use the *A. thaliana* genome sequence as
283 a starting point, and focus on CG dinucleotides within genic regions. To that end,
284 we selected random 1kb regions within genes and choose only those CG sites that
285 are clearly methylated or unmethylated in *A. thaliana* natural populations based
286 on whole genome bisulphite sequencing (WGBS) measurements from the 1001G
287 project (SI text). Our simulator for CG methylation is build in a similar way as
288 the one described above but the epimutation rates are allowed to be asymmetric
289 with the per-site methylation rate (μ_{SM}) and demethylation (μ_{SU}). Region-level
290 epimutations are also implemented, setting the region length to either 1kb [44] or

291 150 bp [15]. The region level methylation and demethylation rates are defined as
292 μ_{RM} and μ_{RU} , respectively. We assume that site-level and region-level epimuta-
293 tion processes are independent. Making this assumption explicit later allow us to
294 test if it is violated in comparisons with actual data. Our simulator also assumes
295 that DNA mutations and epimutations are independent of one another. That is,
296 for simplicity we ignore the fact that methylated cytosines are more likely to tran-
297 sition to thymines as a result of spontaneous deamination [28]. We also ignore
298 the possibility that new DNA mutations could act as CG methylation quantitative
299 trait loci and affect CG methylation patterns in both cis and trans. Such events are
300 extremely rare, and we therefore think that the above assumptions hold reasonably
301 well over short evolutionary time-scales. As the goal is to apply our approach to *A.*
302 *thaliana*, we simulate sequence data for a sample size $n = 10$ (but considering *A.*
303 *thaliana* haploid) from a population displaying 90% selfing [55, 60] under a recent
304 severe population bottleneck demographic scenario. We simulate data assuming
305 previously estimates of the rates of recombination [49], DNA mutation [47], and
306 site- and region-level methylation [65, 15].

307
308 As guidance for future analyses of demographic inference using SNPs and DNA
309 methylation data, the theoretical and empirical analysis of *A. thaliana* methylomes
310 consist of the following five steps: 1) assessing the relevance of region-level methy-
311 lation (DMRs) for inference, 2) inference of site and region epimutation rates, 3)
312 comparing statistics for the SNPs, SMPs and DMRs distributions, 4) demographic
313 inference using SNPs with SMPs or DMRs, and 5) demographic inference using
314 SNPs with SMPs and DMRs.

315

316 **Step 1: assessing the relevance of region-level methylation (DMRs)** 317 **for inference**

318 We determine our ability to detect the existence of spatial correlations between
319 epimutations. That is, we asked if site-specific epimutations can lead to region-
320 level methylation status changes. We assess this across a range of epimutation
321 rates assuming two sequences of 100 Mb ($r = \mu_1 = 10^{-8}$ per generation per bp)
322 under a constant population size fixed to $N = 10,000$ (results in Supplementary
323 Table 1). If site-specific epimutations are independently distributed, the probabil-
324 ity of a given site to be in a certain (methylated or unmethylated) state should
325 be independent from the state of nearby sites (knowing the epimutation rate per
326 site). Conversely, if there is a region effect on epimutation (DMRs), two consecu-
327 tive sites along the genome would exhibit a positive correlation in their methylated
328 states. We therefore calculate from the per-site (de)methylation rates μ_{SM} and
329 μ_{SD} the probability that two successive cytosine positions are identical in their
330 methylation assuming they are independent. This probability can be compared
331 to the one observed (here simulated) methylation data so that we obtain a sta-

332 tistical test for the existence of a positive correlation in the methylation status
333 of nearby sites, interpreted as region-level epimutation process (p-value = 0.05)
334 according to Figure 1A. If the test is non-significant, we validate the existence of a
335 region effect for methylation/demethylation affecting neighbouring cytosines. We
336 find that when region epimutation rates are higher than (or similar to) site-level
337 epimutation rates, namely $\mu_{RM} \gtrsim \mu_{SM}$ and $\mu_{RU} \gtrsim \mu_{SU}$, the existence of regions
338 of consecutive cytosines is detected with high accuracy. However, when site-level
339 epimutation rates are higher ($\mu_{SU} > \mu_{RU}$ and $\mu_{SM} > \mu_{RM}$) than region-level
340 epimutation rates, region-level changes cannot be readily detected (Supplementary
341 Table 1). When methylated regions are detected, we can further determine their
342 length using a specifically developed Hidden Markov Model (HMM) using all pairs
343 of genomes (similarly to [57, 15, 61]). While the length of the methylated region is
344 pre-determined in our simulations (1kb or 150bp) but site-level epimutation occur
345 which can change the distribution of methylation states in that region and across
346 individuals, thus DMR regions can vary in length along the genome and between
347 pairs of chromosomes.

348

349 **Step 2: inference of site- and region-level epimutation rates**

350 As the epimutation rates of most plant species remain unknown, we assess the accu-
351 racy of SMCm to infer epimutation rates at the site- and region-level directly from
352 simulated data. We first assume that either only site- or only region epimutations
353 can occur, and infer their respective rates (see Supplementary Table 2 and 3). Our
354 SMCm approach can accurately recover these rates except when these are higher
355 than 10^{-4} . Next, we assess the accuracy of our approach to simultaneously infer
356 site- and region-level epimutation rates assuming that region and site epimutation
357 rates are equal (Supplementary Table 4). Similar to our previous observation, we
358 find that when the epimutation rates are very high (*e.g.* close to 10^{-2}), accuracy
359 is lost compared to slower epimutation rates. Nonetheless, our average estimated
360 rates are off from the true value by less than an factor 10. Hence, under our model
361 assumptions, we are able to recover the correct order of magnitude for site- and
362 region-level methylation and demethylation rates.

363

364 **Step 3: distribution of statistics for SNPs, SMPs and DMRs**

365 To gain insights on the distribution of epimutations under the assumptions de-
366 scribed in the introduction, we look at key statistics from our simulations: the
367 distribution of distance between two recombination events versus the distribution
368 of the length of estimated DMR regions (Figure 5A), and the LD decay for SMPs
369 (in genic regions) and SNPs (in all contexts) (Figure 5C and D). In our simulations
370 DMRs regions have a maximum fixed size, but their length depends on the inter-
371 action between the region- and site-level epimutation rates. As mentioned in step

372 1, the methylated/demethylated regions are detected using the binomial test and
373 their length estimated by the HMM. Therefore, while variation exist for the length
374 of these regions (Figure 5A), there are shorter than the span of genealogies along
375 the genome, which are defined by the frequency of recombination events along the
376 genome ($r = 3.5 \times 10^{-8}$ as in *A. thaliana*). There is virtually no linkage disequi-
377 librium (LD) between epimutations due to the high epimutation rate (Figure 5C),
378 while the LD between SNPs can range over few kbp (Figure 5D, as observed in
379 *A. thaliana* [27, 52]). Note however, that the region methylation process in itself
380 does not generate LD because this measure can only be computed if SMPs are
381 present in frequency higher than $2/n$ in the sample, *i.e.* there is no LD measure
382 defined for monomorphic methylated/unmethylated regions. In other words, our
383 simulator generates SNPs, SMPs and DMRs which fulfil the three key assumptions
384 of Figure 1A. We note that by using a constant population size $N = 10,000$, the
385 LD decay for SNPs is higher than in the *A. thaliana* data which exhibit an effective
386 population size of ca. $N = 250,000$ [27] and past changes in size.
387

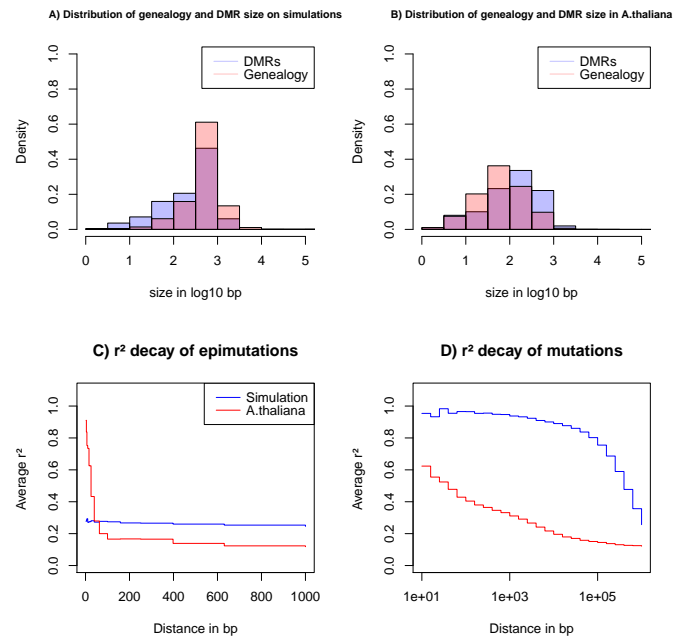


Fig. 5. Key statistics for epimutations and mutations. A) Histogram of the length between two recombination events (genomic span of a genealogy) and DMRs size in bp of the simulated data. B) Histogram of genealogy span and DMRs size in bp from the *A. thaliana* data (10 German accessions). C) Linkage disequilibrium decay of epimutations in our samples of *A. thaliana* (red) and simulated data (blue). D) Linkage disequilibrium decay of mutations in our *A. thaliana* samples (red) and simulated data (blue). The simulations reproduce the outcome of a recent bottleneck with sample size $n = 5$ diploid of 100 Mb, the rates per generation per bp are $r = 3.5 \times 10^{-8}$, $\mu_1 = 7 \times 10^{-9}$, $\mu_{SM} = 3.5 \times 10^{-4}$, $\mu_{SU} = 1.5 \times 10^{-3}$, and per 1kb region $\mu_{RM} = 2 \times 10^{-4}$ and $\mu_{RU} = 1 \times 10^{-3}$.

388 Step 4: demographic inference based on SNPs with SMPs or DMRs

389 We test the usefulness of either SMPs or DMRS for demographic inference. Simulations
 390 under the demographic model from steps 1-3 assume DNA mutations (SNPs)
 391 and only site epimutations (SMPs), *i.e.* no region-level methylation ($\mu_{RM} = \mu_{RU} =$
 392 0). We perform inference of past demographic history under different amount
 393 of potentially methylated sites with and without *a priori* knowing the methyla-
 394 tion/demythlation rates (Figure 6A, B). When the site epimutation rates are *a*
 395 *priori* known, the sharp decrease of population size can be accurately detected.
 396 When epimutation rates are unknown, the shape of the past demographic history
 397 is also well inferred except for a scaling issue (a shift along the x- and y-axes sim-
 398 ilar to that in Figure 6D). When we vary the amount of potentially methylated
 399 sites (2%, 10% and 20%) our inference results remain largely unchanged. This

400 suggests that having methylation measurements for as low as 2% of all CG sites
401 being epimutable in the genome is entirely sufficient to improved SNP-based de-
402 mographic inference (eSMC2 in Figure 6A).

403

404 The amount of sequence data used in Figure 4A and B is fairly large com-
405 pared to real datasets (10 haploid genomes of length 100 Mb). We therefore ran
406 the SMCm and eSMC2 on sequence data simulated under the same scenario but
407 with a reduced sequence length of 10 Mb ($n = 5$ diploid, Figure 6C and D, only
408 3 repetitions are presented for visibility). In this case, we found that inference
409 is significantly affected when using only SNPs (eSMC2 in blue), as we are un-
410 able to correctly recover the demographic scenario. However, incorporating SMPs
411 with known site-level epimutations into the model leads to substantial inference
412 improvements (Figure 6C and D).

413

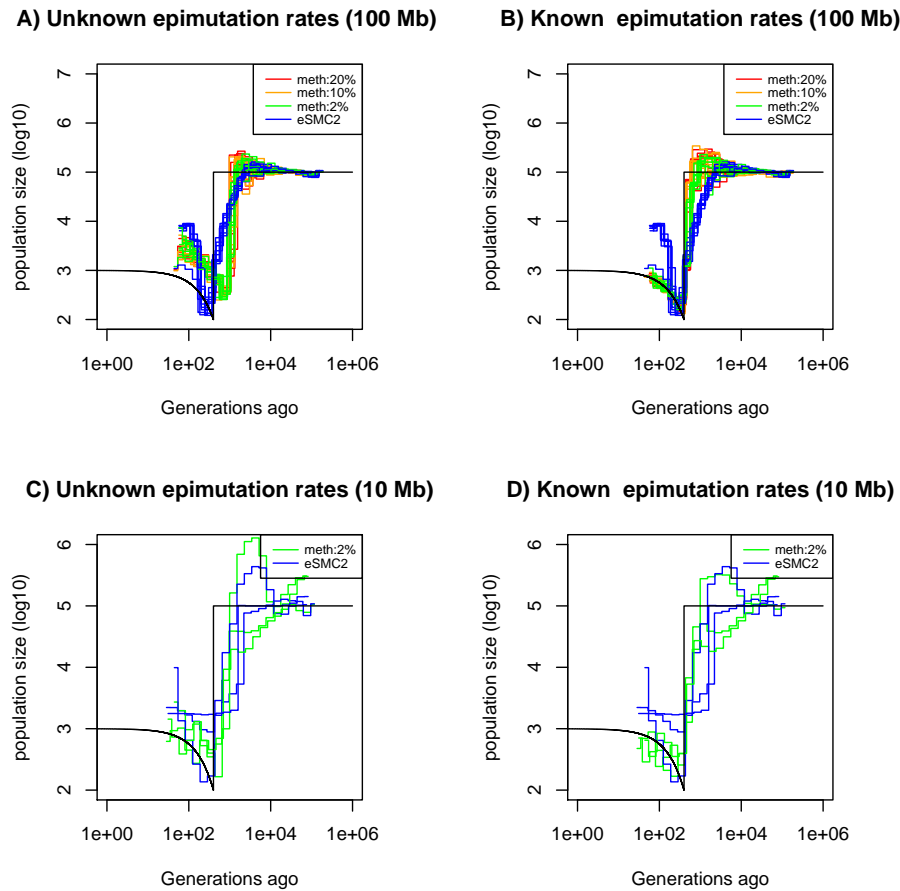


Fig. 6. Performance of SMC approaches using site epimutations (SMPs) and mutations (SNPs) under a bottleneck scenario. Estimated demographic history by eSMC2 (blue) and SMCm assuming the epimutation rate is known (B and D) or not (A and C) where the percentage of CG sites with methylated information varies between 20% (red), 10% (orange) and 2% (green) using 10 sequences of 100 Mb in A and B (with 10 repetitions) and 10 sequences of 10 Mb in C and D (three repetitions displayed) under a recent severe bottleneck (black). The parameters are: $r = 3.5 \times 10^{-8}$ per generation per bp, mutation rate $\mu_1 = 7 \times 10^{-9}$, methylation rate to $\mu_{SM} = 3.5 \times 10^{-4}$ and demethylation rate to $\mu_{SU} = 1.5 \times 10^{-3}$ per generation per bp.

414 We then simulate data under the same demographic scenario, but assume only
 415 region level epimutations (DMRs, $\mu_{SM} = \mu_{SU} = 0$). The results for DMR region
 416 sizes 1kb and 150bp are displayed in Supplementary Figure 4 and 5, respectively.
 417 As in Figure 6, we observed a gain of accuracy in inference when region-level epimuta-
 418 tion rates are known, while the length of the region (1kb or 150bp) does not seem

419 to affect the result. However, no significant gain of information is observed when
420 integrating DMR data with unknown epimutation rates (Supplementary Figure 4
421 and 5). In summary, CG methylation SMPs and to a lesser extend DMRs, can be
422 used jointly with SNPs to improve demographic inference.

423

424 **Step 5: demographic inference based on SNPs with SMPs and** 425 **DMRs**

426 Since site- and region-level methylation processes occur in real data, we run SMCm
427 on simulated data under the same demographic scenario, but now using both site
428 (SMPs) and region (DMRs) epimutations and accounting for both mutation pro-
429 cesses. Inference results are displayed in Supplementary Figure 6. When epimu-
430 tation rates are *a priori* known (in our simulations the rates are fixed and thus
431 known), we find the counter-intuitive result that integrating epimutations decreases
432 the accuracy of inference (compared to SMPs alone, Figure 4). However, when the
433 epimutation rates are set to be inferred by SMCm, integrating SMP and DMR data
434 slightly restores the accuracy of inference (Supplementary Figure 6). Finally, we as-
435 sess the inference accuracy when using SNPs and SMPs but ignoring in SMCm the
436 region methylation effect (DMRs), even though this latter process takes place (Sup-
437 plementary Figure 7). Interestingly, the inference accuracy decreases compared to
438 the previous results (Supplementary Figure 4-6). While the sudden variation of
439 population is somehow recovered, the estimates of the time and magnitude of size
440 change are not well recovered in recent time.

441

442 We demonstrate that our SMCm exhibits an improved statistical power for
443 demographic inference using SNPs and SMPs while accounting for site and region-
444 level methylation processes under the assumptions of Figure 1A. We show that
445 1) using SMPs we can unveil past demographic events hidden by limitations in
446 SNPs, 2) the correct demography can be uncovered irrespective of knowing *a pri-*
447 *ori* the epimutation rates, 3) ignoring site or region-level processes can decrease
448 the accuracy of inference, and 4) knowing the epimutation rates may improve the
449 estimate of demography compared to simultaneously estimating them with SMCm.

450

451 **Joint use of SNPs and SMPs improves the inference of** 452 **recent demographic history in *A. thaliana***

453 **Step 1: assessing the strength of region-level methylation process** 454 **in *A. thaliana***

455 We apply our inference model to genome and methylome data from 10 *A. thaliana*
456 plants from a German local population [27]. We start by assessing the strength of
457 a region effect on the distribution of methylated CG sites along the genome. As

458 expected from [15], for all 10 individual full methylomes we reject the hypothesis of
459 a binomial distribution of methylated and unmethylated sites along the genomes,
460 suggesting the existence of region effect methylation (yielding DMRs) meaning
461 that CG are more likely to be methylated if in a highly methylated region, and
462 conversely for unmethylated CG. This is consistent with the autocorrelations in
463 mCG found in [15, 11]. As a first measure of methylated region length, we test
464 the independence between two annotated CG methylation given a minimum ge-
465 nomic distance between them (within one genome). We observe an average p-value
466 smaller than 0.05 for distances up to 2,000bp but then the p-value rapidly increases
467 (>0.4) (Supplementary Figure 8). As a second measure, our HMM (based on pairs
468 of genomes) yields a DMR average length of 222 bp (distribution in Figure 5B).

469
470 We conclude that the minimum distance for epimutations to be independent
471 along a genome is over 2kb and spans larger distance than the typically proposed
472 DMR size (ca. 150 bp in [15] and 222bp in our analysis) and can therefore cover the
473 size of a gene (see [44]). The simulations and data from *A. thaliana* indicate that
474 the epimutation processes that produces DMRs at the population level in plants
475 cannot simply results from the cumulative action of single-site epimutations. This
476 insights is consistent with recent analyses of epimutational processes in gene bodies,
477 which seems to indicate that the autocorrelation in CG methylation is a function of
478 cooperative methylation maintenance and the distribution of histone modifications
479 [11].

480 **Step 2: site- and region-level epimutation rates**

481 We used the known rates empirically estimated in *A thaliana* and used in simula-
482 tions above ($\mu_{SM} = 3.5 \times 10^{-4}$ and $\mu_{SU} = 1.5 \times 10^{-3}$ per bp per generation and
483 $\mu_{RM} = 2 \times 10^{-4}$ and $\mu_{RU} = 1 \times 10^{-3}$ per region per generation, [65, 15].

485 **Step 3: distribution statistics for SNPs, SMPs and DMRs in *A.*** 486 ***thaliana***

487 Since our SMC model assumes that DNA, SMP and DMR polymorphisms are de-
488 termined by the underlying population/sample genealogy, DMR which span long
489 genomic regions may spread across multiple genealogies and thus violates our as-
490 sumptions. We thus further investigate the potential discrepancies between the
491 data and our model (Figure 5). We infer the DMR sizes from all 10 *A. thaliana*
492 accessions using our ad hoc HMM, and measure the bp distance between a change
493 in the expected hidden state (*i.e.* coalescent time) along the genome, which we
494 interpret as recombination events (called the genomic span of a genealogy). The
495 resulting distributions are found in Figure 5B. We observe that both distributions
496 have a similar shape but DMRs are on average twice as large as the inferred ge-
497 nomic genealogy span: average length of 222 bp (DMR) vs 137 bp (genealogy) and

498 median length of 134 bp (DMR) vs 62 bp (genealogy). This means that on average
499 DMRs are larger than the average distance between two recombination events, thus
500 violating the homogeneous distribution of epimutations along the genome (Figure
501 1C).

502
503 To further unveil potential non-homogeneity of epimutations distribution, we
504 assess the decay of LD of mutations (SNPs) and epimutations (SMPs) (Figure 5C
505 and D) confirming the results in [52]. We find the LD between SMPs in the data
506 to be high (and higher than LD between SNPs) for distance smaller than 100 bp
507 (red line in Figure 5C and D). The LD decay of SMPs is much faster than for
508 SNPs (no linkage between epimutations for distances > 100 bp), likely stemming
509 from 1) epimutation rates being much higher than the DNA mutation rate, and
510 2) the high per site recombination rate in *A. thaliana*. Moreover, the LD between
511 SMPs at distance smaller than 100bp in *A. thaliana* being much higher compared
512 to our simulations (Figure 5C), we suggest that additional local mechanisms of
513 epimutation processes may not be accounted for in our model of the region-level
514 methylation process.

515

516 **Step 4: demographic inference for *A. thaliana* based only on SNPs** 517 **and SMPs**

518 Finally, we apply the SMCm approach to data from the German accessions of *A.*
519 *thaliana*. When using SNP data only, the demographic results are similar to those
520 previously found [55, 60] (Figure 7 purple lines), with no strong evidence for an
521 expansion post-Last Glacial Maximum (LGM) [27]. We then sub-sample and ana-
522 lyze segregating SMPs, which exhibit both methylated and unmethylated states in
523 our sample (as in [65]). Here we ignore DMRs and account only for SMPs. When
524 we use as input the methylation and demethylation rates that have been inferred
525 experimentally [65], a mild bottleneck post-LGM is followed by recent expansion
526 (Figure 7 blue lines). By contrast, letting our SMCm estimate the epimutations
527 rates, we find in recent times a somehow similar but stronger demographic change
528 post-LGM. We find a strong bottleneck event occurring between ca. 5,000 and
529 10,000 generations ago followed by an expansion until today (Figure 7 green lines).
530 The inferred site epimutation rates are 10,000 faster than the DNA mutation rate
531 (Supplementary Table 5) which is close to the expected order of magnitude from
532 experimental measures with and without DMR effects [65, 15]. Both estimates
533 thus yield a post-LGM bottleneck followed by a recent population expansion.

534

535 These results indicate that the inclusion of DNA methylation data can aid in
536 the accurate reconstruction of the evolutionary history of populations, particularly
537 in the recent past where SNPs reach their resolution limit. This is made possible by
538 the fact that the DNA methylation status at CG dinucleotide undergoes stochastic

539 changes at rates that are several orders of magnitude higher than the DNA muta-
540 tion rate, and can be inherited across generations similar to DNA mutations.

541

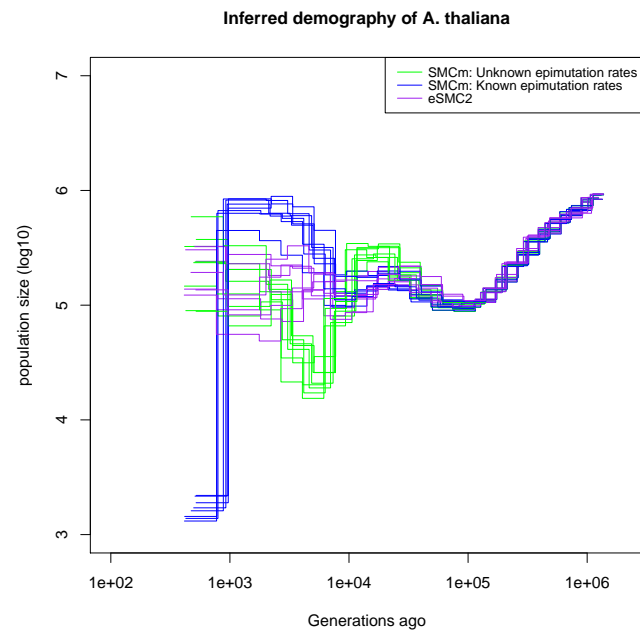


Fig. 7. Integrating epimutations and mutations on German accessions of *A. thaliana*. Estimated demographic history of the German population by eSMC2 (only SNPs, purple) and SMCm when keeping polymorphic methylation sites (SMPs) only: green with epimutation rates estimated by SMCm, blue with epimutation rates fixed to empirical values. The region epimutation effect is ignored. The parameters are $r = 3.6 \times 10^{-8}$, $\mu_1 = 6.95 \times 10^{-9}$, and when assumed known, the site methylation rate is $\mu_{SM} = 3.5 \times 10^{-4}$ and demethylation rate is $\mu_{SU} = 1.5 \times 10^{-3}$.

542 **Step 5: demographic inference accounting for DMRs in *A. thaliana***

543 To assess the robustness of our inference results, we run SMCm using all cytosines
544 (CG) sites with an annotated methylation status (segregating or not) while ac-
545 counting or not for DMRs (Supplementary Figure 9). We fix epimutation rates to
546 the empirically estimated values, and confirm the estimates from Figure 7. When
547 the region-level methylation process is ignored the inferred demography (blue lines
548 in Supplementary Figure 9) is similar to the estimates from SMPs with fixed rates
549 in Figure 7 (blue lines). When the region-level methylation process is taken into
550 account (orange lines in Supplementary Figure 9), the inferred demography is simi-
551 lar to that of the Figure 7 (green lines). In the case where we infer the epimutation

552 rates (sites and region) the demographic history inference is not improved com-
553 pared to that estimated using SNPs only (Supplementary Figure 9, green and red
554 lines) while the inferred epimutation rates are smaller than expected (Supplemen-
555 tary Table 5 and 6), but the ratio of site to region epimutation rates is consistent
556 with empirical estimates [15].

557

558 Discussion

559 Current approaches analyzing whole genome sequences rely on statistics derived
560 from the distribution of ancestral recombination graphs [23, 56, 36, 60, 10, 72,
561 58, 34]. In this study we present a new SMC method that combines SNP data
562 with other types of genomic marker (e.g. TE, microsatellites, DNA methylation).
563 We focus mainly on the inclusion of genomic markers whose mutation rates ex-
564 ceed the DNA point mutation rate, as such (hyper-mutable) markers can provide
565 increased temporal resolution in the recent evolutionary past of populations, and
566 aid in the identification of demographic changes (e.g. population bottlenecks).
567 We demonstrate that by integrating multiple heritable genomic markers, the ARG
568 can be more accurately recovered (outperforming any other methods given the
569 amount of data used in this study [63, 58]). Our simulations demonstrate that
570 if the SNP mutation rate is known, the mutation rate of other markers can be
571 recovered. Moreover, our method accounts for the finite site problem that arises at
572 reversible (hyper-mutable) markers and/or effective population size is high [62, 64].
573 Because model inferences are based on a Baum-Welch algorithm, the accuracy of
574 the method depends on the HMM's capacity to correctly recover the true hidden
575 states [39, 50, 23, 56]. The simulator and SMC methods presented here therefore
576 pave the way for a rigorous statistical framework to test if a common ARG can
577 explain the observed diversity patterns under the model hypotheses laid out in
578 Figure 1. We find that comparisons of LD for different markers along the genome
579 is a useful way to assess violations of our model assumptions.

580 As proof of principle, we apply our approach on data originating from whole
581 genome and methylome data of *A. thaliana* natural accessions (focusing on CG
582 context in genic regions, as in [66, 75, 76]). Our model-based approach provides
583 strong evidence that DMRs cannot simply emerge from site-level epimutations that
584 arise according to a Poisson processes along genome. Instead, stochastic changes
585 in region-level methylation states must be the outcome of spontaneous methyl-
586 ation and demethylation events that operate at both the site- and region-level. Our
587 epimutation model cannot fully describe the observed diversity of epimutations
588 along the genome, meaning that the epimutation processes may indeed be more
589 complex than expected [15, 25]. We observe non-independence between annotated
590 methylation sites spanning genomic regions larger than the span of the underly-
591 ing genealogy (determined by recombination events) which no model can currently

592 describe. Additionally, we find high LD between SMPs over short distances which
593 does not appear in our simulations. Thus, methylation likely violate the assump-
594 tions of a Poisson process distribution along the genome and in time, in line with
595 recent functional studies [25, 41]. We thus further caution against conclusions on
596 the role of natural (purifying) selection [44] or its absence [66] based on population
597 epigenomic data due to the above mentioned assumptions violation. We suggest
598 a possible way forward for modeling epimutations would be to use an Ising model
599 [77] to account for the heterogeneous methylation process along the genome. How-
600 ever, our preliminary work indicates that this model generates non-homogeneous
601 mutation process in space and time which violate strongly our SMC assumptions
602 (Figure 1C and D).

603 Interestingly, the distance of LD decay for SMPs matches quite well the estimated
604 distance between recombination events (Figure 5). In addition to our theoretical
605 results in Table 2, this observation reinforces the usefulness of using SMPs (or any
606 hyper-mutable marker) to improve estimates of the recombination rate along the
607 genome in species where the per site DNA mutation rate (μ) is smaller than the
608 per site recombination rate (r) as in *A. thaliana*. As far as we are aware, our SMC
609 method is the first one to use the forward algorithm output to provide estimates
610 of the position where a change in the expected hidden state (*i.e.* coalescent time)
611 occurs (here interpreted as a recombination event). Future work is needed to im-
612 prove the accuracy of this algorithm based on several markers.

613
614 Nonetheless, we find that a restricted focus on segregating SMPs meets our
615 model assumptions reasonably well, and thus provides a promising way forward.
616 Using these segregating SMPs, we recover a past demographic bottleneck followed
617 by an expansion which could fit the post- Last Glacial Maximum (LGM) coloniza-
618 tion of Europe, a scenario which could not be clearly identified using SNPs only
619 from European (relic and non-relic) accessions [27]. This scenario has been long
620 speculated in *A. thaliana* [21] but strong evidence from inference methods was
621 lacking ([27], Figure 5 in [18]). Furthermore, the absence of highly conflicting de-
622 mography inferred from SNPs and from methylation confirm that, at the time scale
623 of thousands of generations, CG methylation sites are mainly heritable and can be
624 modeled using population genetics theory [13, 66] and used to estimate divergence
625 between lineages [76, 75]. In other words fast ecological local adaptation [51] and
626 response to stresses [59] may likely not be prominent forces reshaping endlessly CG
627 methylation patterns (non-heritability in Figure 1B).

628
629 With the release of new sequencing technology [38], long and accurate reads are
630 becoming accessible, leading to the availability of high quality reference genomes
631 for model and non-model species alike [46, 7]. Additionally, the quality of re-
632 sequencing (population sample) genome data and their annotations is enhanced so
633 that additional markers such as transposable elements, insertion, deletion or mi-
634 crosatellites can be called with increasing confidence. These accurate genomes will

635 provide access to new classes of genomic markers that span the entire mutational
636 spectrum. We therefore suspect in the near future an improvement in our under-
637 standing of the heritability of many markers besides SNPs. Adding other genomic
638 markers besides SNPs will improve full genome approaches, which are currently
639 limited by the observed nucleotide diversity [34, 58, 54]. We predict that our re-
640 sults pave the way to improve the inference of 1) biological traits or recombination
641 rate through time [14, 60], 2) multiple merger events [36], and 3) recombination
642 and mutation rate maps [5, 4]. Our method also should help to dissect the effect of
643 evolutionary forces on genomic diversity [32, 31], and to improve the simultaneous
644 detection, quantification and dating of selection events [1, 8, 30].

645
646 Hence, there is no doubt that extending our work, by simultaneously integrat-
647 ing diverse types of genomic markers into other theoretical framework (*e.g.* ABC
648 approaches), likely represents the future of population genomics. We believe our
649 approach helps to develop more general classes of models capable of leveraging
650 information from any type and amount of diversity observed in sequencing data.
651 Only by doing so can we challenge our current understanding of genomes and unify
652 under a common theory the complex evolution of genomes through generations.

653

654 **Materials and Methods**

655 **Simulating two genomic markers**

656 The sequence is written as a sequence of markers with a given state. Each site is
657 annotated as $MXSY$, where X indicates the marker type and Y the current state
658 of that marker: for example $M1S1$ indicate at this position a marker of type 1 in
659 the state 1.

660 To simulate sequence of theoretical marker we start by simulating an ARG
661 which is then split in a series of genealogies (*i.e.* a sequence of coalescent trees)
662 along the chromosome and create an ancestral sequence (based on equilibrium
663 probability of marker states). Mutation events (nucleotides or epimutations for
664 methylable cytosine) are then added when going along the sequence, *i.e.* along the
665 series of genealogies. The ancestral sequence is thus modified by mutation event
666 assuming a finite site model [74] conditioned to the branch length and topology of
667 the genealogies. Each leaf of the genealogy is one of the n sample. Our model has
668 thus two important features: 1) markers are independent from one another, and
669 2) a given marker has a polymorphism distribution between samples (frequencies
670 of alleles) determined by one given genealogy. The simulator can be found in the
671 latest version of eSMC2 R package (<https://github.com/TPPSellinger/eSMC2>).

672 Simulating methylome data

673 We now focus on methylation data located at cytosine in CG context within genic
674 regions. Only, CG sites in those regions are considered "methylable", and CG
675 sites outside those defined genic regions do not have a methylation status and
676 are considered "unmethylable". We vary the percentage of CG site with methyla-
677 tion state annotated from 2 to 20% of the sequence length. The simulator can in
678 principle simulate epimutations in different methylation context and different rates
679 [40, 17, 78, 20]. We simulate epimutations as described above but with asymmetric
680 rates: the methylation rate per site is $\mu_{SM} = 3.5 \times 10^{-4}$, and the demethyla-
681 tion rate per site is $\mu_{SM} = 1.5 \times 10^{-3}$ [65, 15]. For simplicity and computational
682 tractability, we assume that when an epimutation occurs, it occurs on both DNA
683 strands which then present the same information. In other words, for a haploid
684 individual, a cytosine site can only be methylated or unmethylated (as in [61]).
685 For region level epimutations, the region length is either 1kbp [44] or 150 bp [15].
686 The region level methylation and demethylation rates are set to $\mu_{RM} = 2 \times 10^{-4}$
687 and $\mu_{RU} = 10^{-3}$ respectively (similar to rates measured in *A. thaliana*, [15]). In
688 addition to this, unlike for theoretical marker described above, mutations, site and
689 region epimutations can occur at the same position of the sequence.

690

691 To simulate methylation data, we start with an ancestral sequence of random
692 nucleotide and then randomly select regions in which CG sites have their methy-
693 lation state annotated (representing the genic regions). Cytosine in CG context
694 in those regions are either methylated or unmethylated (noted as M or U). Cy-
695 tosine in other context or regions are considered as unmethylable (and noted as
696 C). The ancestral methylation state is then randomly attributed according to the
697 equilibrium probabilities. Our simulator then introduces DNA mutations, site- and
698 region-epimutations in a similar way as described above.

699 SMC Methods

700 All three methods (eSMC2, SMCtheo and SMCm) are based on the same mathe-
701 matical foundations and implemented in a similar way within the eSMC2 R package
702 (<https://github.com/TPPSellinger>) [60, 36, 56]. This allows to specifically quan-
703 tify the accuracy gained by accounting for multiple genomic markers.

704 SMC optimization function

705 All current SMC approach rely on the Baum-Welch (BW) algorithm for parameter
706 estimation in order to reduce computational load. Yet, the Baum-Welch algorithm
707 is an Expectation-Maximization algorithm, and can hence fall in local extrema
708 when optimizing the likelihood. We alternatively extend SMCtheo to estimate
709 parameters by directly optimizing the likelihood (LH) at the greater cost of com-
710 putation time. We run this approach on a sub-sample of size six haploid genomes

711 to limit the required computational time.

712 **eSMC2 and MSMC2**

713 SMC methods based on the PSMC' [50], such as eSMC2 and MSMC2, focus on the
714 coalescent events between two individuals (*i.e.* two haploid genomes or one diploid
715 genome). The algorithm moves along the sequence and estimates the coalescence
716 time at each position by assessing whether the two sequences are similar or different
717 at each position. If the two sequences are different, this indicates a mutation took
718 place in the genealogy of the sample. The intuition being that the absence of
719 mutations (*i.e.* the two sequences are identical) is likely due to a recent common
720 ancestor between the sequences, and the presence of several mutations likely reflects
721 that the most recent common ancestor of the two sequences is distant in the past.
722 In the event of recombination, there is a break in the current genealogy and the
723 coalescence time consequently takes a new value according to the model parameters
724 [42, 50]. A detailed description of the algorithm can be found in [19, 55].

725 **SMCtheo based on several genomic markers**

726 Our SMCtheo approach is equivalent to PSMC' but take as input a sequence of
727 several genomic markers. The algorithm goes along a pair of haploid genomes and
728 checks at each position which marker is observed and then if both states of the
729 marker are identical or not. The approach is identical to the one described above,
730 except that the probability of both sequences to be identical at one site depends
731 on the mutation rate of the marker at this site (equation 1 and Figure 2). While
732 the mutation rates for many heritable genomic markers are unknown, there is an
733 increasing amount of measures of the DNA (SNP) mutation rate for many species.
734 Our SMCtheo approach is able to leverage the information from the distribution of
735 one theoretical marker (*e.g.* mutations for SNPs) to infer the mutation rate of the
736 other marker 2 (assuming both mutation rates to be symmetrical). If more than 1%
737 of sites are polymorphic in a sequence we use the finite site assumption. If not, then
738 from the diversity observed, the different mutation rates can be recovered by simply
739 comparing Waterson's theta (θ_W) between the reference marker (*i.e.* with known
740 rate) and the marker with the unknown rates. For example, if the diversity (θ_W)
741 at marker 2 is smaller by a factor ten than the reference marker 1 (and no marker
742 violates the infinite site hypothesis), the mutation rate of marker 2 is inferred to
743 be ten times smaller (corrected by the number of possible states). However, if the
744 marker 2 violates the infinite site hypothesis, a Baum-Welch algorithm is run to
745 infer the most likely mutation rates under the SMC to overcome this issue (the
746 Baum-Welch algorithm description can be found in [55]).

747 SMCm

748 When integrating epimutations, the number of possible observations increases com-
749 pare to eSMC2. As in eSMC2, if the two nucleotides (DNA mutation) at one
750 position are identical at a non methylable site, we indicate this as 0. If the two nu-
751 cleotides are different, it is indicated as 1 (*i.e.* a DNA mutation occurred). When
752 assuming site-level epimutation only, three possible observations are possible at a
753 given methylable position: 1) if the two cytosines from the two chromosomes are
754 unmethylated, it is indicated as a 2, 2) if the two cytosines are methylated, it is
755 indicated as a 3, and 3) if at a position a cytosine is methylated and the other
756 one unmethylated, it is indicated as a 4. Depending on the mutation, methylat-
757 ion and, demethylation rates, different frequencies of these states are possible in
758 the sample of sequences, which provide information on the emission rate in the
759 SMC method. When both site- and region-level methylation processes occur, the
760 methylation state is conditioned by the region level methylation state (increasing
761 the number of possible observation to 9)

762 To choose the appropriate settings for SMCm (*i.e.* if there are region level
763 epimutations), we test if the methylation state are distributed independently from
764 one another along one genome. In absence of region methylation effect, the prob-
765 ability at each site (position) to be methylated or unmethylated should be inde-
766 pendent from the previous position (or any other position). Conversely, if there
767 is a region effect on epimutation, two consecutive sites along one genome would
768 exhibit a positive correlation in their methylated states (and across pairs of se-
769 quences). We therefore calculate the probability that two successive positions with
770 an annotated methylation state would be identical under a binomial distribution
771 of methylation along a given genome. We then compare theoretical expectations
772 to the observed data and build the statistical test based on a binomial distribution
773 of probabilities. If existence of region level epimutation is detected, the regions
774 level methylation states are recovered through a hidden markov model (HMM)
775 similarly to [57, 15, 61]. The complete description of the mathematical models and
776 probabilities are in the supplementary material Text S1.

777 We postulate that the epimutation rates remain unknown in most species, while
778 the DNA mutation rate may be known (or approximated based on a closely related
779 species). Hence, we develop an approach based on the SMC capable of leverag-
780 ing information from the distribution of DNA mutations to infer the epimutation
781 rates (similar to what is described above). Our approach first tests if epimutations
782 violates or not the infinite site assumptions. If less than 1% of sites with their
783 methylation state annotated are polymorphic in a sequence we use the infinite site
784 assumption: the site and region level epimutation rates can be recovered straight-
785 forwardly from the observed diversity (θ_W , see above) . Otherwise, a Baum-Welch
786 algorithm is run to infer the most likely epimutation rates (site rate for SMP, and
787 region rates for DMRs) [65, 66, 61].

788 Sequence data of *A. thaliana*

789 We download genome and methylome data of *A. thaliana* from the 1001 genome
790 project [27]. We select 10 individuals from the German accessions respectively
791 corresponding to the accession numbers: 9783, 9794, 9808, 9809, 9810, 9811, 9812,
792 9816, 9813, 9814. We only keep methylome data in CG context and in genic regions
793 [66, 15]. The genic regions are based on the current reference genome TAIR 10.1.
794 The SNPs and epimutations are called according to previously published pipeline
795 [61, 15]. As in previous studies [55, 22, 18], we assume *A. thaliana* data to be
796 haploid due to high homozygosity (caused by high selfing rate). The resulting
797 files are available on GitHub at <https://github.com/TPPSellinger>. To perform
798 analysis we chose $\mu = 6.95 \times 10^{-9}$ per generation per bp as the DNA mutation
799 rate [47] and $r = 3.6 \times 10^{-8}$ as the recombination rate [49] per generation per bp.
800 In order to have the most realistic model, we assume that the methylome of *A.*
801 *thaliana* undergoes both region (RMM) and site (SMM) level epimutations [15].
802 When fixed, we respectively set the site methylation and demethylation rate to
803 $\mu_{SM} = 3.48 \times 10^{-4}$ and $\mu_{SU} = 1.47 \times 10^{-3}$ per generation per bp according to
804 [65]. We additionally set the region level methylation and demethylation rate to
805 $\mu_{RM} = 1.6 \times 10^{-4}$ and $\mu_{RU} = 9.5 \times 10^{-4}$ per generation per bp according to [15].
806 Because we do not account for the effect of variable mutation or recombination rate
807 along the genome, we cut the five chromosome of *A. thaliana* into eight smaller
808 scaffolds [4, 5]. By doing this we remove centromeric regions and limit the effect
809 the variation of mutation and recombination rate along the genome. The selected
810 regions and the SNP density (from the German accessions) are represented in
811 Supplementary Figures 11 to 15.

812 Acknowledgments

813 We thank Zhilin Zhang and Rashmi Hazarika for giving and processing the data
814 of *Arabidopsis thaliana*. TS is supported by the Deutsche Forschungsgemeinschaft,
815 project number 317616126 (TE809/7-1) to AT, and the Austrian Science Fund
816 (project no. TAI 151-B) to Anja Hörger.

817 References

- 818 [1] P. K. Albers and G. McVean. Dating genomic variants and shared ancestry in
819 population-scale sequencing data. *PLOS BIOLOGY*, 18(1), JAN 2020. ISSN
820 1544-9173. doi: 10.1371/journal.pbio.3000586.
- 821 [2] C. Alonso-Blanco, J. Andrade, C. Becker, F. Bemm, J. Bergelson, K. M.
822 Borgwardt, J. Cao, E. Chae, T. M. Dezwaan, W. Ding, et al. 1,135 genomes

- 823 reveal the global pattern of polymorphism in *Arabidopsis thaliana*. *Cell*, 166
824 (2):481–491, 2016.
- 825 [3] T. Anzai, T. Shiina, N. Kimura, K. Yanagiya, S. Kohara, A. Shigenari, T. Ya-
826 magata, J. K. Kulski, T. K. Naruse, Y. Fujimori, et al. Comparative sequenc-
827 ing of human and chimpanzee MHC class I regions unveils insertions/deletions
828 as the major path to genomic divergence. *Proceedings of the National Academy
829 of Sciences*, 100(13):7708–7713, 2003.
- 830 [4] G. V. Barroso and J. Y. Dutheil. Mutation rate variation shapes genome-
831 wide diversity in *Drosophila melanogaster*. Sept. 2021. doi: 10.1101/2021.09.
832 16.460667. URL [http://biorxiv.org/lookup/doi/10.1101/2021.09.16.
833 460667](http://biorxiv.org/lookup/doi/10.1101/2021.09.16.460667).
- 834 [5] G. V. Barroso, N. Puzovic, and J. Y. Dutheil. Inference of recombination maps
835 from a single pair of genomes and its application to ancient samples. *PLoS
836 Genetics*, 15(11), NOV 2019. ISSN 1553-7404. doi: {10.1371/journal.pgen.
837 1008449;10.1371/journal.pgen.1008449.r001;10.1371/journal.pgen.1008449.
838 r002;10.1371/journal.pgen.1008449.r003;10.1371/journal.pgen.1008449.r004}.
- 839 [6] F. Baumdicker, G. Bisschop, D. Goldstein, G. Gower, A. P. Ragsdale,
840 G. Tsambos, S. Zhu, B. Eldon, E. C. Ellerman, J. G. Galloway, A. L.
841 Gladstein, G. Gorjanc, B. Guo, B. Jeffery, W. W. Kretzschmar, K. Lohse,
842 M. Matschiner, D. Nelson, N. S. Pope, C. D. Quinto-Cortes, M. F. Rodrigues,
843 K. Saunack, T. Sellinger, K. Thornton, H. van Kemenade, A. W. Wohns,
844 Y. Wong, S. Gravel, A. D. Kern, J. Koskela, P. L. Ralph, and J. Kelleher.
845 Efficient ancestry and mutation simulation with msprime 1.0. *GENETICS*,
846 220(3), MAR 3 2022. ISSN 0016-6731. doi: 10.1093/genetics/iyab229.
- 847 [7] A. C. Beichman, E. Huerta-Sanchez, and K. E. Lohmueller. Using Genomic
848 Data to Infer Historic Population Dynamics of Nonmodel Organisms. In
849 Futuyma, DJ, editor, *Annual Review of Ecology, Evolution, and Systemat-
850 ics, VOL 49*, volume 49 of *Annual Review of Ecology Evolution and Sys-
851 tematics*, pages 433–456. 2018. ISBN 978-0-8243-1449-1. doi: {10.1146/
852 annurev-ecolsys-110617-062431}.
- 853 [8] G. Bisschop, K. Lohse, and D. Setter. Sweeps in time: leveraging the joint
854 distribution of branch lengths. *GENETICS*, 219(2), OCT 2021. ISSN 0016-
855 6731. doi: 10.1093/genetics/iyab119.
- 856 [9] S. Boitard, W. Rodríguez, F. Jay, S. Mona, and F. Austerlitz. Inferring pop-
857 ulation size history from large samples of genome-wide molecular data - an
858 approximate bayesian computation approach. 12(3):e1005877. ISSN 1553-
859 7404. doi: 10.1371/journal.pgen.1005877. URL [https://dx.plos.org/10.
860 1371/journal.pgen.1005877](https://dx.plos.org/10.1371/journal.pgen.1005877).

- 861 [10] D. Y. C. Brandt, X. Wei, Y. Deng, A. H. Vaughn, and R. Nielsen. Evaluation
862 of methods for estimating coalescence times using ancestral recombination
863 graphs. *GENETICS*, 221(1), MAY 5 2022. ISSN 0016-6731. doi: 10.1093/
864 genetics/iyac044.
- 865 [11] A. Briffa, E. Hollwey, Z. Shahzad, J. D. Moore, D. B. Lyons, M. Howard,
866 and D. Zilberman. Unified establishment and epigenetic inheritance of dna
867 methylation through cooperative met1 activity. *bioRxiv*, pages 2022–09, 2022.
- 868 [12] B. Charlesworth and D. Charlesworth. Elements of evolutionary genetics.
869 2010.
- 870 [13] B. Charlesworth and K. Jain. Purifying Selection, Drift, and Reversible Mu-
871 tation with Arbitrarily High Mutation Rates. *Genetics*, 198(4):1587+, DEC
872 2014. ISSN 0016-6731. doi: {10.1534/genetics.114.167973}.
- 873 [14] Y. Deng, Y. S. Song, and R. Nielsen. The distribution of waiting distances
874 in ancestral recombination graphs. *THEORETICAL POPULATION BIOL-
875 OGY*, 141:34–43, OCT 2021. ISSN 0040-5809. doi: {10.1016/j.tpb.2021.06.
876 003}.
- 877 [15] J. Denkena, F. Johannes, and M. Colome-Tatche. Region-level epimutation
878 rates in arabidopsis thaliana. *HEREDITY*, 127(2):190–202, AUG 2021. ISSN
879 0018-067X. doi: 10.1038/s41437-021-00441-w.
- 880 [16] A. Estoup, P. Jarne, and J.-M. Cornuet. Homoplasmy and mutation model
881 at microsatellite loci and their consequences for population genetics analysis.
882 *Molecular ecology*, 11(9):1591–1604, 2002.
- 883 [17] C. et al. Shotgun bisulphite sequencing of the Arabidopsis genome reveals
884 DNA methylation patterning. *Nature*, 452(7184):215–219, MAR 13 2008. ISSN
885 0028-0836. doi: {10.1038/nature06745}.
- 886 [18] D. et al. African genomes illuminate the early history and transition to selfing
887 in Arabidopsis thaliana. *Proceedings of the National Academy of Sciences of
888 the United States of America*, 114(20):5213–5218, MAY 16 2017. ISSN 0027-
889 8424. doi: {10.1073/pnas.1616736114}.
- 890 [19] M. et al. A genomic history of Aboriginal Australia. *Nature*, 538(7624):207+,
891 OCT 13 2016. ISSN 0028-0836. doi: {10.1038/nature18299}.
- 892 [20] Z. et al. Genome-wide high-resolution mapping and functional analysis of
893 DNA methylation in Arabidopsis. *Cell*, 126(6):1189–1201, SEP 22 2006. ISSN
894 0092-8674. doi: {10.1016/j.cell.2006.08.003}.

- 895 [21] O. François, M. G. B. Blum, M. Jakobsson, and N. A. Rosenberg. De-
896 mographic history of european populations of arabidopsis thaliana. *PLOS*
897 *Genetics*, 4:1–15, 05 2008. URL [https://doi.org/10.1371/journal.pgen.](https://doi.org/10.1371/journal.pgen.1000075)
898 1000075.
- 899 [22] A. Fulgione, M. Koornneef, F. Roux, J. Hermisson, and A. M. Hancock.
900 Madeiran Arabidopsis thaliana Reveals Ancient Long-Range Colonization and
901 Clarifies Demography in Eurasia. *Molecular Biology and Evolution*, 35(3):564–
902 574, MAR 2018. ISSN 0737-4038. doi: {10.1093/molbev/msx300}.
- 903 [23] L. Gattepaille, T. Guenther, and M. Jakobsson. Inferring Past Effective Pop-
904 ulation Size from Distributions of Coalescent Times. *Molecular Biology and*
905 *Evolution*, 204(3):1191+, NOV 2016. ISSN 0016-6731. doi: {10.1534/genetics.
906 115.185058}.
- 907 [24] L. M. Gattepaille, M. Jakobsson, and M. G. B. Blum. Inferring population
908 size changes with sequence and SNP data: lessons from human bottlenecks.
909 *Heredity*, 110(5):409–419, MAY 2013. ISSN 0018-067X. doi: {10.1038/hdy.
910 2012.120}.
- 911 [25] R. R. Hazarika, M. Serra, Z. Zhang, Y. Zhang, R. J. Schmitz, and F. Johannes.
912 Molecular properties of epimutation hotspots. *Nature Plants*, 8(2):146–156,
913 2022.
- 914 [26] R. HUDSON. Properties of a neutral allele model with intragenic recombina-
915 tion. *Theoretical Population Biology*, 23(2):183–201, 1983. ISSN 0040-5809.
916 doi: {10.1016/0040-5809(83)90013-8}.
- 917 [27] e. a. J.E. Cao. Whole-genome sequencing of multiple Arabidopsis thaliana
918 populations. *Nature Genetics*, 43(10):956–U60, OCT 2011. ISSN 1061-4036.
919 doi: {10.1038/ng.911}.
- 920 [28] F. Johannes. DNA methylation makes mutational history. *Nature Plants*, 5
921 (8):772–773, AUG 2019. ISSN 2055-026X. doi: {10.1038/s41477-019-0491-z}.
- 922 [29] F. Johannes and R. J. Schmitz. Spontaneous epimutations in plants. *New*
923 *Phytologist*, 221(3):1253–1259, FEB 2019. ISSN 0028-646X. doi: {10.1111/
924 nph.15434}.
- 925 [30] P. Johri, B. Charlesworth, and J. D. Jensen. Toward an evolutionarily
926 appropriate null model: Jointly inferring demography and purifying selec-
927 tion. *GENETICS*, 215(1):173–192, MAY 2020. ISSN 0016-6731. doi:
928 10.1534/genetics.119.303002.

- 929 [31] P. Johri, K. Riall, H. Becher, L. Excoffier, B. Charlesworth, and J. D.
930 Jensen. The impact of purifying and background selection on the infer-
931 ence of population history: Problems and prospects. *MOLECULAR BIOL-
932 OGY AND EVOLUTION*, 38(7):2986–3003, JUL 2021. ISSN 0737-4038. doi:
933 10.1093/molbev/msab050.
- 934 [32] P. Johri, C. F. Aquadro, M. Beaumont, B. Charlesworth, L. Excoffier, A. Eyre-
935 Walker, P. D. Keightley, M. Lynch, G. McVean, B. A. Payseur, S. P. Pfeifer,
936 W. Stephan, and J. D. Jensen. Recommendations for improving statistical
937 inference in population genomics. *PLOS Biology*, 20(5):e3001669, May 2022.
938 ISSN 1545-7885. doi: 10.1371/journal.pbio.3001669. URL [https://dx.plos.
939 org/10.1371/journal.pbio.3001669](https://dx.plos.org/10.1371/journal.pbio.3001669).
- 940 [33] J. Kelleher, A. M. Etheridge, and G. McVean. Efficient Coalescent Simula-
941 tion and Genealogical Analysis for Large Sample Sizes. *PLOS Computational
942 Biology*, 12(5), MAY 2016. doi: {10.1371/journal.pcbi.1004842}.
- 943 [34] J. Kelleher, Y. Wong, A. W. Wohns, C. Fadil, P. K. Albers, and G. McVean.
944 Inferring whole-genome histories in large population datasets (vol 51, pg 1330,
945 2019). *Nature Genetics*, 51(11):1660, NOV 2019. ISSN 1061-4036. doi: {10.
946 1038/s41588-019-0523-7}.
- 947 [35] J. Kingman. The Coalescent . *Stochastic Processes and their Applications*, 13,
948 1982.
- 949 [36] K. Korfmann, T. P. P. Sellinger, F. Freund, M. Fumagalli, and A. Tellier.
950 Simultaneous inference of past demography and selection from the ancestral
951 recombination graph under the beta coalescent. *bioRxiv*, 2022.
- 952 [37] K. Korfmann, O. E. Gaggiotti, and M. Fumagalli. Deep Learning in Population
953 Genetics. *Genome Biology and Evolution*, 15(2), 01 2023. ISSN 1759-6653.
954 doi: 10.1093/gbe/evad008. URL <https://doi.org/10.1093/gbe/evad008>.
955 evad008.
- 956 [38] D. Lang, S. Zhang, P. Ren, F. Liang, Z. Sun, G. Meng, Y. Tan, X. Li,
957 Q. Lai, L. Han, D. Wang, F. Hu, W. Wang, and S. Liu. Comparison
958 of the two up-to-date sequencing technologies for genome assembly: Hifi
959 reads of pacific biosciences sequel ii system and ultralong reads of oxford
960 nanopore. *GIGASCIENCE*, 9(12), DEC 2020. ISSN 2047-217X. doi:
961 10.1093/gigascience/giaa123.
- 962 [39] H. Li and R. Durbin. Inference of human population history from individual
963 whole-genome sequences. *Nature*, 475(7357):493–U84, JUL 28 2011. ISSN
964 0028-0836. doi: {10.1038/nature10231}.

- 965 [40] R. Lister, R. C. O'Malley, J. Tonti-Filippini, B. D. Gregory, C. C. Berry,
966 A. H. Millar, and J. R. Ecker. Highly integrated single-base resolution maps
967 of the epigenome in *Arabidopsis*. *Cell*, 133(3):523–536, MAY 2 2008. ISSN
968 0092-8674. doi: {10.1016/j.cell.2008.03.029}.
- 969 [41] D. B. Lyons, A. Briffa, S. He, J. Choi, E. Hollwey, J. Colicchio, I. Ander-
970 son, X. Feng, M. Howard, and D. Zilberman. Extensive de novo activity
971 stabilizes epigenetic inheritance of cg methylation in *arabidopsis* transposons.
972 *bioRxiv*, 2022. doi: 10.1101/2022.04.19.488736. URL [https://www.biorxiv.](https://www.biorxiv.org/content/early/2022/04/19/2022.04.19.488736)
973 [org/content/early/2022/04/19/2022.04.19.488736](https://www.biorxiv.org/content/early/2022/04/19/2022.04.19.488736).
- 974 [42] P. Marjoram and J. Wall. Fast “coalescent” simulation. *BMC Genetics*, 7,
975 MAR 15 2006. ISSN 1471-2156. doi: {10.1186/1471-2156-7-16}.
- 976 [43] G. McVean and N. Cardin. Approximating the coalescent with recombina-
977 tion. *Philosophical Transactions of the Royal Society B-Biological Sciences*,
978 360(1459):1387–1393, JUL 29 2005. ISSN 0962-8436. doi: {10.1098/rstb.
979 20053.1673}.
- 980 [44] A. Muyle, J. Ross-Ibarra, D. K. Seymour, and B. S. Gaut. Investigation
981 Gene body methylation is under selection in *Arabidopsis thaliana*. Sept. 2020.
982 doi: 10.1101/2020.09.04.283333. URL [http://biorxiv.org/lookup/doi/10.](http://biorxiv.org/lookup/doi/10.1101/2020.09.04.283333)
983 [1101/2020.09.04.283333](http://biorxiv.org/lookup/doi/10.1101/2020.09.04.283333).
- 984 [45] M. Nordborg. Linkage disequilibrium, gene trees and selfing: An ancestral re-
985 combination graph with partial self-fertilization. *Molecular Biology and Evo-*
986 *lution*, 154(2):923–929, FEB 2000. ISSN 0016-6731.
- 987 [46] S. Nurk, B. P. Walenz, A. Rhie, M. R. Vollger, G. A. Logsdon, R. Grothe,
988 K. H. Miga, E. E. Eichler, A. M. Phillippy, and S. Koren. Hicnu: accurate
989 assembly of segmental duplications, satellites, and allelic variants from high-
990 fidelity long reads. *GENOME RESEARCH*, 30(9):1291–1305, SEP 2020. ISSN
991 1088-9051. doi: 10.1101/gr.263566.120.
- 992 [47] S. Ossowski, K. Schneeberger, J. I. Lucas-Lledo, N. Warthmann, R. M. Clark,
993 R. G. Shaw, D. Weigel, and M. Lynch. The Rate and Molecular Spectrum
994 of Spontaneous Mutations in *Arabidopsis thaliana*. *Science*, 327(5961):92–94,
995 JAN 1 2010. ISSN 0036-8075. doi: {10.1126/science.1180677}.
- 996 [48] S. Ou, W. Su, Y. Liao, K. Chougule, J. R. A. Agda, A. J. Hellinga, C. S. B.
997 Lugo, T. A. Elliott, D. Ware, T. Peterson, N. Jiang, C. N. Hirsch, and M. B.
998 Hufford. Benchmarking transposable element annotation methods for creation
999 of a streamlined, comprehensive pipeline. *GENOME BIOLOGY*, 20(1), DEC
1000 16 2019. ISSN 1474-760X. doi: 10.1186/s13059-019-1905-y.

- 1001 [49] P. A. Salome, K. Bomblies, J. Fitz, R. A. E. Laitinen, N. Warthmann, L. Yant,
1002 and D. Weigel. The recombination landscape in *Arabidopsis thaliana* F-2
1003 populations. *Heredity*, 108(4):447–455, APR 2012. ISSN 0018-067X. doi:
1004 {10.1038/hdy.2011.95}.
- 1005 [50] S. Schiffels and R. Durbin. Inferring human population size and separation his-
1006 tory from multiple genome sequences. *Nature Genetics*, 46(8):919–925, AUG
1007 2014. ISSN 1061-4036. doi: {10.1038/ng.3015}.
- 1008 [51] M. W. Schmid, C. Heichinger, D. Coman Schmid, D. Guthörl, V. Gagliardini,
1009 R. Bruggmann, S. Aluri, C. Aquino, B. Schmid, L. A. Turnbull, et al. Contri-
1010 bution of epigenetic variation to adaptation in *arabidopsis*. *Nature Commu-
1011 nications*, 9(1):1–12, 2018.
- 1012 [52] R. J. Schmitz, M. D. Schultz, M. A. Urich, J. R. Nery, M. Pelizzola, O. Libiger,
1013 A. Alix, R. B. McCosh, H. Chen, N. J. Schork, et al. Patterns of population
1014 epigenomic diversity. *Nature*, 495(7440):193–198, 2013.
- 1015 [53] J. G. Schraiber and J. M. Akey. Methods and models for unravelling human
1016 evolutionary history. *Nature Reviews Genetics*, 16(12):727–740, DEC 2015.
1017 ISSN 1471-0056. doi: {10.1038/nrg4005}.
- 1018 [54] R. Schweiger and R. Durbin. Ultra-fast genome-wide inference of pairwise
1019 coalescence times. *bioRxiv*, 2023.
- 1020 [55] T. P. P. Sellinger, D. Abu Awad, M. Moest, and A. Tellier. Inference of
1021 past demography, dormancy and self-fertilization rates from whole genome
1022 sequence data. *PLOS Genetics*, 16(4), APR 2020. ISSN 1553-7404.
1023 doi: {10.1371/journal.pgen.1008698;10.1371/journal.pgen.1008698.r001;
1024 10.1371/journal.pgen.1008698.r002;10.1371/journal.pgen.1008698.r003;10.
1025 1371/journal.pgen.1008698.r004;10.1371/journal.pgen.1008698.r005;10.1371/
1026 journal.pgen.1008698.r006}.
- 1027 [56] T. P. P. Sellinger, D. Abu-Awad, and A. Tellier. Limits and convergence
1028 properties of the sequentially markovian coalescent. *MOLECULAR ECOL-
1029 OGY RESOURCES*, 21(7):2231–2248, OCT 2021. ISSN 1755-098X. doi:
1030 {10.1111/1755-0998.13416}.
- 1031 [57] Y. Shahryary, A. Symeonidi, R. R. Hazarika, J. Denkena, T. Mubeen,
1032 B. Hofmeister, T. van Gurp, M. Colome-Tatch, K. J. F. Verhoeven, G. Tuskan,
1033 R. J. Schmitz, and F. Johannes. Alphabeta: computational inference of
1034 epimutation rates and spectra from high-throughput dna methylation data
1035 in plants. *GENOME BIOLOGY*, 21(1), OCT 6 2020. ISSN 1474-760X. doi:
1036 10.1186/s13059-020-02161-6.

- 1037 [58] L. Speidel, M. Forest, S. Shi, and S. R. Myers. A method for genome-wide
1038 genealogy estimation for thousands of samples. *Nature Genetics*, 51(9):1321+,
1039 SEP 2019. ISSN 1061-4036. doi: {10.1038/s41588-019-0484-x}.
- 1040 [59] T. Srikant and H.-G. Drost. How stress facilitates phenotypic innovation
1041 through epigenetic diversity. *Frontiers in Plant Science*, 11:606800, 2021.
- 1042 [60] S. Struett, T. Sellinger, S. Glémin, A. Tellier, and S. Laurent. Inference
1043 of evolutionary transitions to self-fertilization using whole-genome sequences.
1044 *bioRxiv*, 2022. doi: 10.1101/2022.07.29.502030. URL <https://www.biorxiv.org/content/early/2022/08/01/2022.07.29.502030>.
1045
- 1046 [61] A. Taudt, D. Roquis, A. Vidalis, R. Wardenaar, F. Johannes, and M. Colome-
1047 Tatche. Methimpute: imputation-guided construction of complete methylomes
1048 from wgbs data. *BMC GENOMICS*, 19, JUN 7 2018. ISSN 1471-2164. doi:
1049 10.1186/s12864-018-4641-x.
- 1050 [62] A. Tellier, S. J. Y. Laurent, H. Lainer, P. Pavlidis, and W. Stephan. Infer-
1051 ence of seed bank parameters in two wild tomato species using ecological and
1052 genetic data. *Proceedings of the National Academy of Sciences of the United*
1053 *States of America*, 108(41):17052–17057, OCT 11 2011. ISSN 0027-8424. doi:
1054 {10.1073/pnas.1111266108}.
- 1055 [63] J. Terhorst, J. A. Kamm, and Y. S. Song. Robust and scalable inference
1056 of population history froth hundreds of unphased whole genomes. *Nature*
1057 *Genetics*, 49(2):303–309, FEB 2017. ISSN 1061-4036. doi: {10.1038/ng.3748}.
- 1058 [64] G. Upadhyia and M. Steinrücken. Robust Inference of Population Size Histories
1059 from Genomic Sequencing Data. May 2021. doi: 10.1101/2021.05.22.445274.
1060 URL <http://biorxiv.org/lookup/doi/10.1101/2021.05.22.445274>.
- 1061 [65] van der Graaf et al. Rate, spectrum, and evolutionary dynamics of sponta-
1062 neous epimutations. *Proceedings of the National Academy of Sciences of the*
1063 *United States of America*, 112(21):6676–6681, MAY 26 2015. ISSN 0027-8424.
1064 doi: {10.1073/pnas.1424254112}.
- 1065 [66] A. Vidalis, D. Zivkovic, R. Wardenaar, D. Roquis, A. Tellier, and F. Johannes.
1066 Methylome evolution in plants. *Genome Biology*, 17, DEC 20 2016. ISSN
1067 1474-760X. doi: {10.1186/s13059-016-1127-5}.
- 1068 [67] J. Wakeley. Coalescent theory: an introduction. roberts and company. *Green-*
1069 *wood Village Wayne AF, Maxwell MA, Ward CG, Vellios CV, Wilson I, Wayne*
1070 *JC, Williams MR (2015) Sudden and rapid decline of the abundant marsupial*
1071 *Bettongia penicillata in Australia. Oryx*, 49:175185Webb, 2008.

- 1072 [68] C. Wang and C. Liang. Msipred: a python package for tumor microsatellite
1073 instability classification from tumor mutation annotation data using a support
1074 vector machine. *SCIENTIFIC REPORTS*, 8, DEC 3 2018. ISSN 2045-2322.
1075 doi: 10.1038/s41598-018-35682-z.
- 1076 [69] J. Wang and C. Fan. A neutrality test for detecting selection on dna methyla-
1077 tion using single methylation polymorphism frequency spectrum. *GENOME*
1078 *BIOLOGY AND EVOLUTION*, 7(1):154–171, JAN 2015. ISSN 1759-6653.
1079 doi: 10.1093/gbe/evu271.
- 1080 [70] D. Weigel and V. Colot. Epialleles in plant evolution. *Genome biology*, 13
1081 (10):1–6, 2012.
- 1082 [71] C. Wiuf and J. Hein. Recombination as a point process along sequences.
1083 *Theoretical Population Biology*, 55(3):248–259, JUN 1999. ISSN 0040-5809.
1084 doi: {10.1006/tpbi.1998.1403}.
- 1085 [72] A. W. Wohns, Y. Wong, B. Jeffery, A. Akbari, S. Mallick, R. Pinhasi, N. Pat-
1086 terson, D. Reich, J. Kelleher, and G. McVean. A unified genealogy of modern
1087 and ancient genomes. *SCIENCE*, 375(6583):836+, FEB 25 2022. ISSN 0036-
1088 8075. doi: 10.1126/science.abi8264.
- 1089 [73] R. Yang, J. L. Van Etten, and S. M. Dehm. Indel detection from dna and rna
1090 sequencing data with transindel. *BMC GENOMICS*, 19, APR 19 2018. ISSN
1091 1471-2164. doi: 10.1186/s12864-018-4671-4.
- 1092 [74] Z. Yang. Statistical properties of a DNA sample under the finite-sites model.
1093 *Genetics*, 144(4):1941–1950, DEC 1996. ISSN 0016-6731.
- 1094 [75] N. Yao, R. J. Schmitz, and F. Johannes. Epimutations define a fast-ticking
1095 molecular clock in plants. *Trends in Genetics*, 37(8):699–710, 2021.
- 1096 [76] N. Yao, Z. Zhang, L. Yu, R. Hazarika, C. Yu, H. Jang, L. M. Smith, J. Ton,
1097 L. Liu, J. Stachowicz, et al. An evolutionary epigenetic clock in plants. *bioRxiv*,
1098 pages 2023–03, 2023.
- 1099 [77] Y. Zhang, S. Wang, and X. Wang. Data-driven-based approach to identifying
1100 differentially methylated regions using modified 1d ising model. *BIOMED*
1101 *RESEARCH INTERNATIONAL*, 2018, 2018. ISSN 2314-6133. doi: 10.1155/
1102 2018/1070645.
- 1103 [78] D. Zilberman, M. Gehring, R. K. Tran, T. Ballinger, and S. Henikoff. Genome-
1104 wide analysis of Arabidopsis thaliana DNA methylation uncovers an interde-
1105 pendence between methylation and transcription. *Nature Genetics*, 39(1):
1106 61–69, JAN 2007. ISSN 1061-4036. doi: {10.1038/ng1929}.