

1 **Inferring genomic landscapes with the**
2 **integrative sequentially Markov coalescent**
3 **(iSMC)**

4 Gustavo V. Barroso^[0000-0002-1943-9297] and
5 Julien Y. Dutheil^[0000-0001-7753-4121]

6 **Abstract** The integrative Sequentially Markovian Coalescent (iSMC) is an extension
7 of the sequentially Markovian Coalescent (SMC) model allowing for parameter het-
8 erogeneity along the genome, such as recombination and mutation rates. Heteroge-
9 neous parameters follow an autocorrelation process that modulates the genealogical
10 process, extending the hidden state space and adding as few as two extra parameters
11 per heterogeneous rate. Classical hidden Markov chain methodology is used to in-
12 fer the posterior estimate of the rate landscape. In this chapter, we demonstrate the
13 use of iSMC to infer both recombination and mutation landscapes, using data from
14 *Homo sapiens* and *Homo neanderthalensis* genomes. We further indicate how to use
15 simulations to assess statistical power and investigate possible sources of inference
16 noise.

17 **Keywords:** sequentially Markovian coalescent, coalescent hidden Markov mod-
18 els, recombination landscape, mutation landscape

19 **Running head:** Inferring genomic landscapes with iSMC

Gustavo V. Barroso
University of Wisconsin-Madison, 447 Birge Hall, Madison, Wisconsin, USA e-mail: gvbarroso@
gmail.com

Julien Y. Dutheil
Max Planck Institute for Evolutionary Biology, August-Thienemann-Str. 2, 24306 Plön, Germany
e-mail: dutheil@evolbio.mpg.de

20 **1 Introduction**

21 The ancestral recombination graph (ARG) is the collection of coalescence and re-
22 combination events in the history of a sample of haplotypes. As such, it is a complete
23 description of the genealogical history in a sample of genomes. With knowledge of
24 the true ARG, it would be straightforward to estimate key evolutionary parameters
25 such as ancestral population sizes and recombination rates. However, because many
26 events do not leave a signature in the sequence data, a large number of configurations
27 are compatible with a given dataset. Therefore, inference methods typically rely on
28 population genetic models to integrate over possible ARGs to account for the uncer-
29 tainty in its reconstruction. In this chapter, we focus on describing one such class
30 of models and their application to inferring molecular processes that vary along the
31 genome, such as mutation and recombination rates, while accounting for temporal
32 variation in population sizes.

33 The sequential coalescent frames the genealogical process as unfolding spatially
34 along the genome [1]. Starting at the leftmost position and moving site by site, the co-
35 alescent tree subtending a sample of haplotypes remains the same until an ancestral
36 recombination event separates lines of descent. The ensuing ‘floating lineage’ then
37 follows its own coalescent trajectory, resulting in a new genealogy that spans until
38 the next recombination event, and this process is iterated down to the last position
39 in the genome, allowing the construction of the ARG. The sequentially Markovian
40 coalescent (SMC) is an approximation of this process in which some parts of the
41 ARG (which do not have any direct descendants in the sample) are discarded [2].
42 Therefore, unlike the classic (and unwieldy) backward-in-time model, the SMC en-
43 forces that recombination events always happen rightward, restricting the uncertainty
44 to the temporal dynamics of recombination and coalescence events, for which there
45 are accurate population genetic predictions. This allows the SMC to be formulated

46 as a hidden Markov model (HMM) with genealogies and nucleotide configurations
47 as hidden and observed states, respectively [3]. As the coalescence trajectory of re-
48 combining lineages is a function of ancestral population sizes, this model forms the
49 basis of many demographic inference tools that have been developed in the past 15
50 years.

51 Besides demography, the SMC framework has been used to model evolutionary
52 processes that vary in magnitude along the genome. In 2018, Palamara et al. [4]
53 used hundreds of genomes to identify regions with consistently shallow coalescence
54 events as candidates for natural selection. In the following year, Barroso et al. [5]
55 added further complexity to the SMC by casting it as a Markov-modulated HMM,
56 where a prior distribution of recombination rates modulates the frequency of ge-
57 nealogical transitions along the genome. These local changes in the auto-correlation
58 of the ARG are recorded in the sequence data because patterns of diversity are a func-
59 tion of the underlying genealogies. This flexible integrative SMC framework (*ismc*)
60 was later used to incorporate spatial variation in mutation rates into the model [6].

61 The approach of *ismc* is to first estimate the distributions of recombination and/or
62 mutation rates, then use well-established HMM techniques to reconstruct their ge-
63 nomic landscapes. In practice, the statistical power to perform fine-scale inferences
64 is a function (among other things) of the level of diversity maintained in the popula-
65 tion. In humans, where pairwise diversity is on the order of 10^{-4} , the resolution of
66 inferred maps is low, with estimates generally noisy on scales $< 100kb$. However,
67 with larger window sizes, the number of recombination events recorded is enough
68 that the landscapes are well recovered, even from a single diploid genome. In this
69 chapter, we guide the user through a typical application of the iSMC to infer re-
70 combination maps in the Altai Neanderthal. We then infer the human mutation map
71 for chromosome 1 and compare our results with mutation rates that were recently
72 predicted using a collection of genomic features.

73 **2 Getting and running iSMC**

74 The `ismc` source code can be obtained from GitHub [https://github.com/gvbarroso/](https://github.com/gvbarroso/iSMC/)
75 `iSMC/`. The package is written in C++ and depends on the BOOST and Bio++ li-
76 braries [7]. Standalone executable programs can be downloaded for Linux 64-bit
77 systems.

78 `ismc` uses the Bio++ Option (BppO) syntax to specify models and options, such
79 as input and output file paths. The options can be conveniently saved in one or more
80 option files. `ismc` can then be run using the command line

81 **Terminal command**

81

82

```
83 ismc param=options.bpp
```

1

84

85 where `options.bpp` contains a list of program options. One of these options is
86 `dataset_label`, which defines a prefix tag that will be used for all output files. In
87 the following, we will refer to this suffix as [PREFIX] when mentioning the output
88 files.

89 **3 Input data**

90 The core information used by `ismc` is a list of homozygous (coded '0') and heterozy-
91 gous positions (coded '1') in at least one pair of haplotypes, possibly with missing
92 data (coded '2'). Following the approach introduced by MSMC2 [8], `ismc` imple-
93 ments a composite likelihood approach where the likelihood of a data set is the like-
94 lihood product of multiple pairs of haplotypes. These pairs may correspond to several

95 (possibly unphased) diploid individuals or multiple combinations of individual hap-
 96 lotypes (up to $\frac{2n \times (2n-1)}{2}$ possible pairs for n phased diploid individuals). Sequence
 97 data can be provided in different formats. Individual haplotypes can be provided as
 98 complete sequences in a file in FASTA format, with missing data coded 'N'. Alter-
 99 natively, variable positions for each sampled individual can be provided as a file in
 100 Variant Call Format (VCF). Genome VCF (gVCF) files provide both variable and
 101 non-variable positions, other positions being considered missing data. Conversely,
 102 simple VCF files only contain variable positions, and an additional mask file should
 103 be provided to distinguish non-variable positions from missing data.

104 To run ismc from a Fasta file, the following options are required:

105 **Option file**

106
 107 `input_file_type = FASTA` 1
 108 `sequence_file_path = sequences.fasta.gz` 2
 109 `seq_compression_type = gzip` 3

111 Input files can be provided as compressed, either by gzip or bgzip. To use the VCF
 112 format instead, these options are required, including the specification of a mask file:

113 **Option file**

114
 115 `input_file_type = VCF //or gVCF` 1
 116 `sequence_file_path = variants.vcf.gz` 2
 117 `mask_file_type = BED` 3
 118 `mask_file_path = mask.bed` 4
 119 `seq_compression_type = bgzip` 5
 120 `mask_compression_type = none` 6

122 The mask file can be in BED or Fasta format and must match (in content and
123 order) the chromosomes in the sequence file. In the BED case, `ismc` will mask out
124 (that is, convert to missing data) regions that are present in this file (i.e., a negative
125 mask). In the Fasta case, `ismc` will mask out sites that are not represented by either
126 '1' or 'P' characters. When a genome VCF file is used (`input_file_type = gVCF`),
127 no mask file is required.

128 Only diploid individuals are allowed in the VCF file, whose phasing, if any, will
129 be ignored. Each diploid will be considered as a pair of haplotypes to be analyzed.
130 Conversely, when the data are input as a Fasta file, each sequence corresponds to a
131 single haplotype. By default, `ismc` will then create all possible pairs of haplotypes
132 and fit the model to all of them. The option `diploid_indices` allows the user to restrict
133 the haplotype pairs that are formed. For example, `diploid_indices = (0,1,1,2)` will
134 create only two pairs: 0 with 1 and 1 with 2.

135 **4 Specifying the demographic model**

136 Similarly to other SMC methods like MSMC2 and eSMC (see Chapters 6 and 7),
137 `ismc` discretizes time in intervals. The number of categories to use is specified by
138 the option `number_intervals`, with a value typically between 30 and 40. More inter-
139 vals lead to more precise inferences, but longer execution times. Setting a too large
140 number may lead to some intervals never being observed in the data if the number
141 of variable positions in the genome is too small, leading to optimization issues.

142 Unlike MSMC2, where time intervals represent free population sizes to be esti-
143 mated, `ismc` adopts cubic spline interpolation to reduce the number of model param-
144 eters. The idea is to have a smooth, continuous curve that spans the entire coales-
145 cent history, then map TMRCAs at the prescribed time boundaries, converting to a
146 piecewise-constant history. This ensures that we have a discrete-space Markov chain.

147 To allow flexible size changes, ismc can incorporate multiple spline curves that meet
148 at each end. This is set by the number of ‘knots’ k in the model (default = 2), where
149 $k + 1$ curves are fit. The height of the splines at each knot (in units of scaled coalescent
150 rates) are the demographic parameters to be estimated by ismc (labeled y_0 to y_{k+1}).
151 ismc can estimate and compare models with different numbers of knots and select the
152 best-fitting model according to Akaike’s information criterion (AIC). For instance,
153 setting the options `init_number_knots = 2` and `max_number_knots = 5` will
154 compare four models with two, three, four, and five knots, respectively.

155 **5 Estimating parameters**

156 ismc uses maximum likelihood inference to estimate parameters. Optimization is
157 started using the `optimize = true` option. The optimizer is set with the option
158 `numerical_optimizer`, with two available options: Powell or NewtonRaphson.
159 The Powell optimizer will fit most situations. The Newton-Raphson optimizer may
160 converge faster, but requires computationally expensive numerical derivatives. In
161 both cases, optimization will stop when the (relative) increase in logarithmic likeli-
162 hood between two iterations is less than a given threshold, specified by the option
163 `function_tolerance`. A threshold of 0.1 or below is recommended. The maximum
164 likelihood framework also enables us to compute confidence intervals for the esti-
165 mated parameter values using the Fisher information matrix. This calculation can be
166 enabled using the option `compute_confidence_interval = true`. In this case, ismc
167 outputs `ConfidenceIntervals.txt` which contains lower and upper 95% intervals for
168 each model parameter.

169 ismc produces two log files that allow the user to assess the proper convergence of
170 the parameter estimation procedure: a `[PREFIX]_profile.txt` file containing all pa-
171 rameter and log-likelihood values at each iteration, and a `[PREFIX]_messages.txt`

172 file with potential warnings generated during the optimization process. In case sev-
 173 eral demographic models are tested, one pair of files will be generated per model.
 174 Furthermore, ismc produces a [PREFIX]_backup_params.txt file that contains
 175 the current parameter value at each iteration. This file is updated as the optimization
 176 progresses. Since some ismc models can take a rather long time to fit, this temporary
 177 file allows us to restart the parameter estimation in case of hardware failure or mis-
 178 configuration of the job scheduling system when using a computer grid. This can be
 179 done by setting the option `resume_optim = true`. The option can also be used to
 180 provide initial parameter values, for instance, obtained from a simpler, faster model.

181 The final parameters are outputted as a [PREFIX]_estimates.txt file. The coa-
 182 llescence rates for each time interval are in a separate file [PREFIX]_demography.txt,
 183 which follows the same format as MSMC2 (see Chapter 6).

184 **6 Example 1: homogeneous model**

185 In order to illustrate the options that we have seen so far, we fit a simple homoge-
 186 neous model to the Altai Neanderthal genome, restricting the analysis to chromo-
 187 some 1. The complete data preparation pipeline can be found on the companion
 188 GitHub repository.

189 **Option file**

```

190
191 DATA=chr1 1
192 NB_KNOTS=2 2
193 dataset_label = $(DATA)_homogeneous_$(NB_KNOTS) knots 3
194 4
195 sequence_file_path = ../VCFnorm/$(DATA).norm.vcf.gz 5
196 mask_file_path = ../Bed/$(DATA)_negmask.bed.gz 6
197 mask_file_type = BED 7
198 input_file_type = VCF 8

```

199	<code>seq_compression_type = gzip</code>	9
200	<code>mask_compression_type = gzip</code>	10
201	<code>tab_file_path = \$(DATA).tab</code>	11
202	<code>diploid_indices = (0,1)</code>	12
203		13
204	<code>splines_type = Sigmoidal</code>	14
205	<code>init_number_knots = \$(NB_KNOTS)</code>	15
206	<code>max_number_knots = \$(NB_KNOTS)</code>	16
207		17
208	<code>number_theta_categories = 1</code>	18
209	<code>number_rho_categories = 1</code>	19
210	<code>number_ne_categories = 1</code>	20
211	<code>number_intervals = 30</code>	21
212		22
213	<code>numerical_optimizer = Powell</code>	23
214	<code>function_tolerance = 1e-1</code>	24
215	<code>confidence_interval = false</code>	25
216		26
217	<code>optimize = true</code>	27

218

219 The script makes use of the BppO variable syntax, defining two global variables:
 220 DATA, which specifies chromosome 1 (this allows to easily run the same model on
 221 other chromosomes with no further modification), and NB_KNOTS, which spec-
 222 ifies the number of knots in the spline model. (For convenience, we run the model
 223 4 times with four different numbers of knots instead of letting it run all models at
 224 once.) The options `tab_file_path` will be described in section 13.

225 The inferred demography consistently reveals a decrease in population size (Fig-
 226 ure 1). As `ismc` and `MSMC2` use an identical population genetic model (albeit with
 227 different implementations and estimation procedures), they produce similar results.
 228 `MSMC2` displays a typical off-the-charts high population size in the recent past,
 229 while `ismc` provides a smoother curve due to the use of splines. However, the de-
 230 mography inferred by `ismc` has a lower resolution and does not capture intermediate
 231 peaks. Using a higher number of knots provides more nuance, but a large number

232 of knots ($n = 5$) yields a most likely artifact peak in the distant past. The lowest
233 AIC value is obtained with the four-knot model (see Table 1). Using more data (all
234 chromosomes) would allow us to fit more precise models.

235 In the next sections, we will see how to run `ismc` to estimate two types of land-
236 scape: recombination and mutation.

237 **7 Specifying a heterogeneous model with variable recombination** 238 **rate**

239 The strength of `ismc` lies in its ability to model heterogeneous parameters along
240 the genome. Fitting such models requires specifying (1) the a priori discrete distri-
241 bution of ρ along the genome and (2) the transition probabilities between ρ cate-
242 gories. Property (1) is achieved with option `rho_var_model = Gamma`, specify-
243 ing a (discretized) gamma distribution with mean 1. The shape of this distribution
244 is estimated together with the demography parameters. The number of ρ categories
245 is specified with option `number_rho_categories`. Just like for the time intervals,
246 there is a trade-off between accuracy (more categories) and computational efficiency
247 (less categories). We recommend using a number of five categories. `ismc` currently
248 assumes a simple auto-correlation model with a single rate of transitions between
249 categories (Property 2).

250 **8 Example 2: variable recombination rate**

251 Here we extend the example from section 6 to fit a model with variable recombination
252 rate:

Option file

```

253
254
255 DATA=chr1 1
256 NB_KNOTS=2 2
257 NB_RHO=5 3
258 dataset_label = $(DATA)_gamma$(NB_RHO)_$(NB_KNOTS) knots 4
259 5
260 sequence_file_path = ../VCFnorm/$(DATA).norm.vcf.gz 6
261 mask_file_path = ../Bed/$(DATA)_negmask.bed.gz 7
262 mask_file_type = BED 8
263 input_file_type = VCF 9
264 seq_compression_type = gzip 10
265 mask_compression_type = gzip 11
266 tab_file_path = $(DATA).tab 12
267 diploid_indices = (0,1) 13
268 14
269 splines_type = Sigmoidal 15
270 init_number_knots = $(NB_KNOTS) 16
271 max_number_knots = $(NB_KNOTS) 17
272 18
273 number_theta_categories = 1 19
274 number_rho_categories = $(NB_RHO) 20
275 number_ne_categories = 1 21
276 number_intervals = 30 22
277 23
278 rho_var_model = Gamma 24
279 25
280 numerical_optimizer = Powell 26
281 function_tolerance = 1e-1 27
282 confidence_interval = false 28
283 29
284 optimize = true 30
285 resum_optim = false 31

```

286

287 The option file is mostly identical; we only added the options for the heteroge-
288 nous recombination rate. The inferred demographic scenarios are very similar to
289 those inferred with the homogeneous model (Figure 2). When comparing all mod-

Table 1 Models fitted to the Altai neanderthal dataset.

Model	<i>a</i>	<i>b</i>	Log likelihood	AIC
Homogeneous	2	5	-369098.22838652721839	738206.5
Homogeneous	3	6	-369109.84457955235848	738231.7
Homogeneous	4	7	-369029.35426597035257	738072.7
Homogeneous	5	8	-369105.81665211619111	738227.6
Gamma-5	2	7	-368181.47791081061587	736377.0
Gamma-5	3	8	-368165.42031291732565	736346.8
Gamma-5	4	9	-368130.59860143717378	736279.2
Gamma-5	5	10	-368119.60966349433875	736259.2

^a Number of knots in the demography model

^b Total number of parameters

290 els, Akaike's information criterion favors the most complex one, with five knots and
 291 a heterogeneous landscape (Table 1).

292 **9 Inference of recombination maps**

293 Once a heterogeneous model has been fitted, it is possible to use posterior decoding
 294 to obtain site-specific estimates of recombination rates, which can be averaged in
 295 windows along the genome.

296 The option `decode = true` will trigger the posterior decoding and output the
 297 results in multiple files for each pair of genomes. As posterior decoding may re-
 298 quire a large amount of memory, slicing the chromosome in blocks is important for
 299 computational efficiency. The maximum length of such blocks is specified by the
 300 `fragment_size` option, with a value set to 10,000,000 that generally enables a rea-
 301 sonable compromise (see section 13). In case the model was fitted but not decoded,
 302 it is possible to rerun `ismc` by enabling the decoding options and setting the option
 303 `optimize = false`; this will use the previously estimated parameters to generate the
 304 posterior decoding, without further optimization.

305 The `ismc` package contains an auxiliary program named `ismc_mapper` that
 306 takes the single-nucleotide landscapes output by `ismc` and uses the information
 307 present in the tab file to output rate maps (averaged within larger genomic windows).
 308 `ismc_mapper` requires its own option file, with a syntax similar to the main `ismc`
 309 program:

310 **Option file**

```

311
312 DATA=chr1 1
313 NB_KNOTS=2 2
314 dataset_label = $(DATA)_gamma5_$(NB_KNOTS) knots 3
315 tab_file_path = $(DATA).tab 4
316 bin_sizes = (10000,100000,1000000) 5
317 bin_rate = rho 6
    
```

318 _____

319 The two arguments `dataset_label` and `tab_file_path` mirror `ismc`'s options file
 320 and must be provided identically. The remaining options are specific to `ismc_mapper`
 321 and tell the program how to aggregate the single-nucleotide landscapes. The
 322 `bin_sizes` argument specifies multiple window sizes, here 10 kb, 100 kb and 1 Mb.
 323 The argument `bin_rate` specifies the variable of interest, in our case ρ (but see Note
 324 1). The resulting files will contain averaged estimates of ρ per window, over all in-
 325 dividuals and all positions in each window.

326 Figure 3A shows the resulting window estimates plotted along the chromosome,
 327 showing the typical ‘‘smiling’’ pattern with higher recombination rates close to the
 328 telomeres. Recombination estimates are robust to the demography model and are
 329 highly correlated when two or five knots are used (Figure 3B).

330 **10 Power analysis using simulations**

331 While in this example the inferred demography and recombination map appear plau-
332 sible and in agreement with current knowledge about the neanderthal history and
333 primates biology in general, it might be difficult to make sense out of parameters
334 inferred from non-model organisms where little is known. To assess whether the es-
335 timates of ismc (or any other method) are reliable, it is recommended to conduct a
336 power analysis. The idea is to simulate datasets with identical characteristics to the
337 real one (alignment length, number of contigs, genetic diversity) under the inferred
338 scenario (population size variation + recombination landscape) to assess how likely
339 we are to recover the scenario if it were the true one.

340 We use msprime [9] (see Chapter 15) to simulate data, with the demography
341 inferred with the selected model and the 10 kb-scale recombination map. The cor-
342 responding code is available in the companion repository at [https://github.com/](https://github.com/StatisticalPopulationGenomics-2ndEd/iSMC/)
343 [StatisticalPopulationGenomics-2ndEd/iSMC/](https://github.com/StatisticalPopulationGenomics-2ndEd/iSMC/). The resulting VCF files are in-
344 dexed on the original neanderthal VCF by creating an artificial SNP at the first and
345 last positions of the original Neanderthal chromosome. The original mask file was
346 then applied to the simulated data during the inference.

347 Demography could be well recovered from the simulation, demonstrating suf-
348 ficient statistical power for demography inference (Figure 4). In contrast, despite
349 showing significant correlations (Table 2), recombination rates were largely under-
350 estimated (Figure 5A). Recombination rates are bounded by the values of the min-
351 imum and maximum categories of the discretized prior distribution. The effect is
352 notable here, as ismc infers a shape parameter higher than 0.5, ten times higher than
353 the simulated value. This is an artifact of binning the recombination map in windows,
354 resulting in a normalization of the landscape.

355 Increasing the number of categories did not improve the estimation (Table 2, sec-
356 ond row). Using a manual discretization scheme that allows for lower and higher
357 categories, yet with low probabilities, improves the estimate (Figure 5B) but not the
358 correlation with the true map (Table 2, third row). Note that the reported correlations
359 are rank-based, which indicates that the underestimation is only a scaling issue, rela-
360 tive rates are largely unaffected. Custom discretization can be achieved by specifying
361 the following argument to `ismc`:

Option file

```
362  
363  
364 rho_boundaries=(0.0001, 0.001, 0.01, 0.1, 0.5, 1, 10) 1  
365
```

366 This argument supersedes the `number_rho_categories` argument as the number of
367 categories is set to the number of bounds minus one (see also Note 2). `ismc` recovers
368 the 1-Mb recombination landscape with a good accuracy, but this accuracy decreases
369 at smaller scales, becoming lower than previously reported values [5]. We can use
370 simulations to investigate possible causes.

371 Ancient genomes are typically of lower quality compared to genomes of extant
372 species. This leads to an increased proportion of missing data (*i.e.*, more masked
373 positions). To investigate how much this reduces the accuracy of the inference, we
374 re-analyzed the simulated data without masking positions. The correlations increase
375 at all scales, but only marginally (Table 2, fourth row), indicating that missing data
376 are not the main cause of lower precision in this case.

377 The inferred ρ value was relatively high, with a ratio $\rho/\theta > 1$, which can lead to
378 a reduction in statistical power [5]. Surprisingly, simulating under the same scenario

Table 2 Power analysis for the Altai neanderthal dataset. Spearman correlation coefficients between the inferred and true map at three different window sizes. All correlations are significant with P values lower than $1e-200$.

Method	10 kb	100 kb	1 Mb
Gamma (5 categories)	0.5345	0.5865	0.6890
Gamma (10 categories)	0.5323	0.5857	0.6860
Gamma (custom discretization)	0.5201	0.5721	0.6587
No mask	0.5353	0.5870	0.6883
Low recombination	0.4742	0.5391	0.6447
Constant population size	0.6946	0.7498	0.7772

379 but with a lower recombination rate led to even lower correlations (Table 2, fifth
380 row), indicating that the high recombination rate is not the source of lower accuracy.

381 Finally, we investigated the singular “collapsing” demography of Neanderthals
382 as a source of power reduction. Simulating under the same recombination map but
383 with a flat demography of $N_e = 30,000$ resulted in much higher correlations (Table
384 2, bottom row), identifying demography as a main cause of power reduction.

385 **11 Inference of mutation maps**

386 Several, if not most, evolutionary models are sensitive to variation in mutation rate
387 along the genome. Despite their importance, the development of tools for inferring
388 mutation maps is lagging behind other statistical models. On the one hand, identifi-
389 cation of *de novo* mutations (e.g., in family trios or biobank-scale datasets) is partic-
390 ularly sensitive to sequencing artifacts, besides being impractical in most species. On
391 the other hand, population genetic methods struggle because disentangling the effect
392 of mutation rate variation from genetic drift and linked selection is challenging.

393 ismc can jointly model mutation and recombination landscapes, alongside the
394 TMRCA distribution. Under the neutral SMC model, variation in the TMRCA repre-
395 sents genetic drift. Informally, the hope is that it also absorbs linked selection effects,

396 such that ismc learns the mutation landscape from the remaining variation in diver-
 397 sity along the genome. However, due to the lack of empirically validated mutation
 398 maps, this has been difficult to benchmark. Although simulations suggest that ismc
 399 recovers the mutation landscape under neutrality and that it can be robust to linked
 400 selection, questions remain about performance in real data. Fortunately, due to its
 401 importance in evolutionary and medical genetics, the human mutation map has re-
 402 cently been predicted from genomic features such as sequence context and epigenetic
 403 markers [10].

404 **12 Example 3: variable mutation rate**

405 To compare our model with this independent estimate of mutation rate variation, we
 406 fit ismc to high-quality sequence data from four individuals of the Mbuti population
 407 (chromosome 1).

408 **Option file**

```

409
410 dataset_label=Mbuti_chr1                1
411 input_file_type = VCF                    2
412 sequence_file_path = mbuti_chr1.vcf.gz   3
413 mask_file_path = sgdp_chr1_mask.fa.gz    4
414 mask_file_type = FASTA                   5
415 tab_file_path = Mbuti_chr1.tab           6
416                                           7
417 seq_compression_type = gzip              8
418 mask_compression_type = gzip             9
419                                           10
420 diploid_indices = (0,1,2,3,4,5,6,7)     11
421                                           12
422 optimize = $(OPTIM)                      13
423 numerical_optimizer = Powell              14
424 confidence_interval = false              15
425 function_tolerance = 1e-3                16
426                                           17

```

427	<code>number_threads = 4</code>	18
428		19
429	<code>number_rho_categories = 5</code>	20
430	<code>number_theta_categories = 5</code>	21
431	<code>number_intervals = 30</code>	22
432		23
433	<code>decode = \$(DECODE)</code>	24
434	<code>decode_diploids_parallel = true</code>	25
435	<code>decode_breakpoints_parallel = false</code>	26
436	<code>fragment_size = 800000</code>	27
437		28
438	<code>init_number_knots = 2</code>	29
439	<code>max_number_knots = 2</code>	30
440	<code>splines_type = Sigmoidal</code>	31
441	<code>rho_var_model = Gamma</code>	32
442	<code>theta_var_model = Gamma</code>	33

443

444 This option file specifies a model with heterogeneous mutation and recombina-
 445 tion landscapes in a $30 \times 5 \times 5$ configuration. Since fitting such a complex model
 446 requires hefty computational resources, we run optimization and posterior decoding
 447 separately (see section 13), parallelizing each step over individuals. (Due to the large
 448 run-times, in this example we refrain from comparing models.)

449 **Command line**

450

451 `ismc param=opt.bpp OPTIM=true DECODE=false` 1

452

453 Although optimization does not require large amounts of memory, it takes about
 454 two weeks to fit a model with 750 hidden states on the largest human chromosome.
 455 (A workaround is to use the tab file to specify breakpoints around the centromere,
 456 effectively treating chromosome arms with independent genealogical processes and

457 accordingly doubling the number of computing threads.) The estimated parameters
 458 shaping the mutation landscape suggest an approximately bell-shaped distribution
 459 (Gamma $\alpha = \beta \approx 4$) with auto-correlation $t_{ii} \approx 0.99998$, i.e., changing mutation
 460 rates on average every 50 kb. In contrast, the inferred recombination distribution is
 461 exponential-like (Gamma $\alpha = \beta \approx 0.1$) and less auto-correlated ($r_{ii} \approx 0.99995$).

462 We then use the estimated parameters to reconstruct the posterior landscapes:

463 **Command line**

```
465 ismc param=opt.bpp OPTIM=false DECODE=true 1
```

466

467 To parallelize the decoding over diploids, we specified a short `fragment_size`
 468 (800 kb) in the options file. This option tunes peak memory use, since `ismc` resets
 469 the matrix of posterior probabilities after extending it over the specified length. After
 470 another two weeks and peaking at around 128 Gb of memory, `ismc` output single-
 471 nucleotide landscapes that were binned by `ismc_mapper` with the following op-
 472 tions:

473 **Options file**

```
474
475 dataset_label = Mbuti_chr1 1
476 bin_rate = theta 2
477 bin_sizes = (10000,100000,1000000) 3
478 tab_file_path = Mbuti_chr1.tab 4
```

480 The correlations between the mutation maps predicted by `ismc` and Roulette in-
 481 crease dramatically with window size (Figure 6). This has at least two potential ex-

482 planations. First, noisy estimates at smaller scales are smoothed out when averaging
483 over larger genomic bins. This effect should be observed in both methods, but is espe-
484 cially strong in *ismc*, which extracts the signal from a small sample of genomes with
485 average pairwise diversity $< 10^{-4}$ (see fluctuations in mutation maps as a function
486 of scale, Figure 7). Second, different biological factors shape mutation rate variation
487 at different scales, with triplet context in one extreme and others (e.g. replication
488 timing) likely making a difference in the order of dozens of kilobases. Because it
489 ultimately relies on sparse polymorphisms as its source of information, *ismc* is more
490 sensitive to capture larger-scale effects.

491 Overall, the correlations between *ismc* and Roulette mutation maps are encour-
492 aging. They indicate that although fitting these models is computationally expensive,
493 *ismc* is capable of inferring mutation maps from a small sample with fairly low di-
494 versity, especially in large genomic windows. To this end, high-quality sequence data
495 is critical. The model accumulates information as it traverses the genome and, there-
496 fore, relies on long blocks of informative sites (which are interrupted if missing data
497 are widespread). However, further benchmarking is required if strong model viola-
498 tions are suspected. We suggest caution in the inference and analyses of mutation
499 maps, but we are optimistic that *ismc* opens new avenues in the population genetics
500 of non-model organisms.

501 **13 Practical considerations**

502 **13.1 Computer resource management**

503 Depending on genome length, sample size, and model complexity (i.e., the total num-
504 ber of hidden states), *ismc* can use a large amount of computer resources. To mitigate
505 execution time, *ismc* uses multi-threading to compute the likelihood, spreading the

506 calculation over several cores (set by `number__threads` option). This parallelization
507 is performed on the individual and their chromosomes. That is, if data consist of,
508 e.g. two diploids with three chromosomes each, `ismc` will benefit from using up to
509 six threads. While likelihood calculation does not require extensive live memory (in-
510 termediate results are not stored), posterior decoding requires recording conditional
511 likelihoods for every hidden state and position in the genome. As this can be quite
512 large, `ismc` offers the ability to decode the genome in windows. This saves mem-
513 ory, but implies resetting the Markov chain at every window boundary. The size of
514 the window can be set using the `fragment__size` argument. The larger the value,
515 the more accurate the decoding is, but the more memory is required. The minimal
516 value depends on the recombination rate and genetic diversity; as a rule of thumb,
517 we recommend not getting below 1 Mb. Similarly to parameter estimation by maxi-
518 mum likelihood, posterior decoding can also be sped up with multi-threading. Since
519 this comes at the cost of a linear increase in memory, `ismc` provides two boolean
520 options to specify axes of parallelization and fine-tune resource management: `de-`
521 `code__diploids__parallel` and `decode__breakpoints__parallel` (see Example 3: vari-
522 able mutation rate), both of which default to false.

523 **13.2 Two-pass estimation**

524 `ismc` employs a two-step approach to reconstruct the rate maps. The first step is the
525 maximum likelihood estimation of the parameters governing the prior distributions
526 of mutation, recombination, and TMRCA along the genome. (This includes fitting
527 a demographic model.) The second step is to use the inferred distributions to recon-
528 struct the posterior average μ and r for every site in the genome.

529 In practice, we can exploit this two-step procedure to better allocate computa-
530 tional resources. The optimization step employs several evaluations of the likelihood

531 function (the compressed “forward” algorithm of the HMM) and therefore con-
 532 sumes substantial amounts of time but not memory. On the other hand, computing
 533 the posterior average requires a single execution of the full “forward-backward”
 534 HMM algorithm, which not only can be time-consuming but also requires consider-
 535 able memory. It is therefore desirable to run these steps separately, first by specifying

536 Option file

```
537
538 optimize = true           1
539 decode = false           2
```

540

541 in the option file, subsequently inverting both booleans before running iSMC again.
 542 The reconstructed rate landscapes are scaled by the effective population size, $\rho =$
 543 $4 \times N_e \times r$ and $\theta = 4 \times N_e \times \mu$ (but see Note 3).

544 13.3 Tab files

545 The options file specifies the various inputs for ismc. Among these is a tab-separated
 546 text file that summarizes the boundaries of contiguous “blocks” (e.g., chromosomes
 547 or scaffolds) that are labeled in the sequence data. ismc treats each block indepen-
 548 dently, i.e., assumes free recombination among them. They are listed one per line,
 549 with seven columns in total (coordinates are 1-indexed):

- 550 1. Block ID (e.g. 'chr1')
- 551 2. The start coordinate of the block, mapped to the reference genome
- 552 3. The end coordinate of the block, mapped to the reference genome
- 553 4. 0 for all blocks

- 554 5. The difference between the third and second columns
- 555 6. Bottom cut-off position from where `ismc_mapper` should bin single-nucleotide
- 556 landscapes into larger genomic windows
- 557 7. Top cut-off position until where `ismc_mapper` should bin single-nucleotide land-
- 558 scapes into larger genomic windows.

559 The last two columns are used by `ismc_mapper` and are convenient for trimming

560 the final maps into a range of sites that match that of another file. For example, if

561 the input sequence data for chromosome 1 spans positions 5,000–10,000,000, and

562 we want to align it with an experimental genetic map that ranges from 10,000 to

563 8,000,000, columns 6-7 should be 10,000 and 8,000,000. Similarly, if one wishes to

564 align rate maps from different VCFs files, column 6 should have the largest of the

565 starting positions, whereas column 7 should have the smallest of the end positions

566 among the VCFs. If such trimming is not desirable, columns 6 and 7 should be the

567 same as columns 2 and 3. The Python script `generate_tab.py`, provided with the

568 `ismc` package, generates a default tab file from a VCF file.

Competing Interests The authors have no conflicts of interest to declare that are relevant to the content of this chapter.

569 **14 Notes**

570 Note 1: multiple decoding

571 If multiple variables (e.g. ρ and θ) have been decoded in the same model, `ismc_mapper`
572 must be run separately for each of them.

573 Note 2: discretization of prior distribution

574 The default behavior is to consider categories with equal mass, which is usually the
575 best option for estimating the hyper-parameters (parameters of the prior distribution).
576 The number of categories is specified by the argument `number_rho_categories`. A
577 drawback of this approach is that some category values might become very close, if
578 not identical, in the case of an extremely skewed distribution (low alpha). It is also
579 possible to use an equally distributed grid by specifying `max_rho_value`, which
580 will create n equally spaced categories between 0 and the specified value.

581 Note 3: definition of ρ

582 Although the population recombination rate is typically defined as $\rho = 4 \cdot Ne \cdot r$,
583 simulation revealed that ρ in `ismc` is off by a factor of 2, which means that it is
584 defined as $\rho = 2 \cdot Ne \cdot r$. This implementation was kept for the sake of backward
585 compatibility. Estimates should be corrected accordingly when converting ρ to r :

$$\frac{\rho}{\theta} = \frac{2 \cdot Ne \cdot r}{4 \cdot Ne \cdot u} \quad (1)$$

$$= \frac{r}{2u} \quad (2)$$

$$r = 2u \frac{\rho}{\theta} \quad (3)$$

586 **References**

- 587 1. Carsten Wiuf and Jotun Hein. Recombination as a point process along sequences. *Theoretical*
588 *population biology*, 55(3):248–259, 1999.
- 589 2. Gilean AT McVean and Niall J Cardin. Approximating the coalescent with recombination.
590 *Philosophical Transactions of the Royal Society B: Biological Sciences*, 360(1459):1387–1393,
591 2005.
- 592 3. Jeffrey P Spence, Matthias Steinrücken, Jonathan Terhorst, and Yun S Song. Inference of
593 population history using coalescent hmms: review and outlook. *Current opinion in genetics &*
594 *development*, 53:70–76, 2018.
- 595 4. Pier Francesco Palamara, Jonathan Terhorst, Yun S. Song, and Alkes L. Price. High-throughput
596 inference of pairwise coalescence times identifies signals of selection and enriched disease
597 heritability. *Nat. Genet.*, 50(9):1311–1317, September 2018.
- 598 5. Gustavo V. Barroso, Nataša Puzović, and Julien Y. Dutheil. Inference of recombination
599 maps from a single pair of genomes and its application to ancient samples. *PLOS Genetics*,
600 15(11):e1008449, November 2019.
- 601 6. Gustavo V. Barroso and Julien Y. Dutheil. The landscape of nucleotide diversity in *Drosophila*
602 *melanogaster* is shaped by mutation rate variation. *Peer Community Journal*, 3, 2023.
- 603 7. Laurent Guéguen, Sylvain Gaillard, Bastien Boussau, Manolo Gouy, Mathieu Groussin, Nico-
604 las C Rochette, Thomas Bigot, David Fournier, Fanny Pouyet, Vincent Cahais, Aurélien
605 Bernard, Céline Scornavacca, Benoît Nabholz, Annabelle Haudry, Loïc Dachary, Nicolas
606 Galtier, Khalid Belkhir, and Julien Y Dutheil. Bio++: Efficient Extensible Libraries and Tools
607 for Computational Molecular Evolution. *Mol. Biol. Evol.*, June 2013.
- 608 8. Anna-Sapfo Malaspinas, Michael C Westaway, Craig Muller, Vitor C Sousa, Oscar Lao, Isabel
609 Alves, Anders Bergström, Georgios Athanasiadis, Jade Y Cheng, Jacob E Crawford, et al. A
610 genomic history of aboriginal australia. *Nature*, 538(7624):207–214, 2016.
- 611 9. Franz Baumdicker, Gertjan Bisschop, Daniel Goldstein, Graham Gower, Aaron P. Rags-
612 dale, Georgia Tsambos, Sha Zhu, Bjarki Eldon, E. Castedo Ellerman, Jared G. Galloway,
613 Ariella L. Gladstein, Gregor Gorjanc, Bing Guo, Ben Jeffery, Warren W. Kretzschmar, Konrad
614 Lohse, Michael Matschiner, Dominic Nelson, Nathaniel S. Pope, Consuelo D. Quinto-Cortés,
615 Murillo F. Rodrigues, Kumar Saunack, Thibaut Sellinger, Kevin Thornton, Hugo van Keme-
616 nade, Anthony W. Wohns, Yan Wong, Simon Gravel, Andrew D. Kern, Jere Koskela, Peter L.

- 617 Ralph, and Jerome Kelleher. Efficient ancestry and mutation simulation with msprime 1.0.
618 *Genetics*, 220(3):iyab229, March 2022.
- 619 10. Vladimir Seplyarskiy, Evan M Koch, Daniel J Lee, Joshua S Lichtman, Harding H Luan, and
620 Shamil R Sunyaev. A mutation rate model at the basepair resolution identifies the mutagenic
621 effect of polymerase iii transcription. *Nature genetics*, 55(12):2235–2242, 2023.

622 **Figures**

Fig. 1 Inferred demography from Altai Neanderthal's chromosome 1. A mutation rate of $1.25e-8 \text{ bp}^{-1}$ was used to scale the axes.

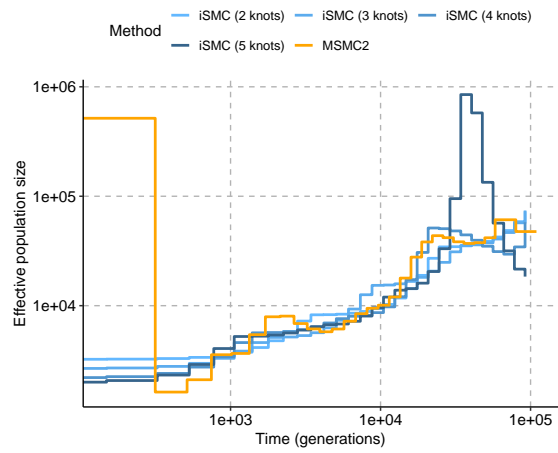
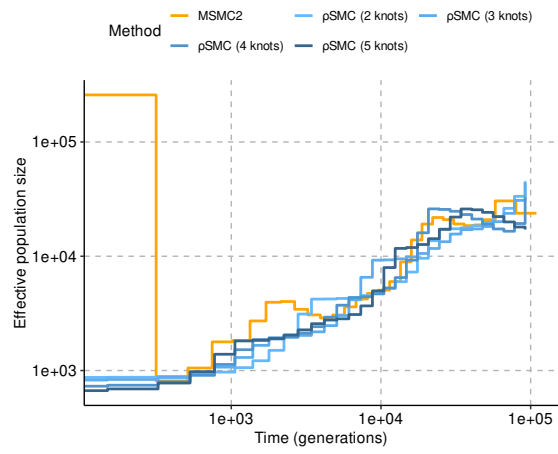


Fig. 2 Inferred demography from Altai Neanderthal chromosome 1, with a heterogeneous recombination rate. A mutation rate of $1.25e-8 \text{ bp}^{-1}$ was used to scale the axes.



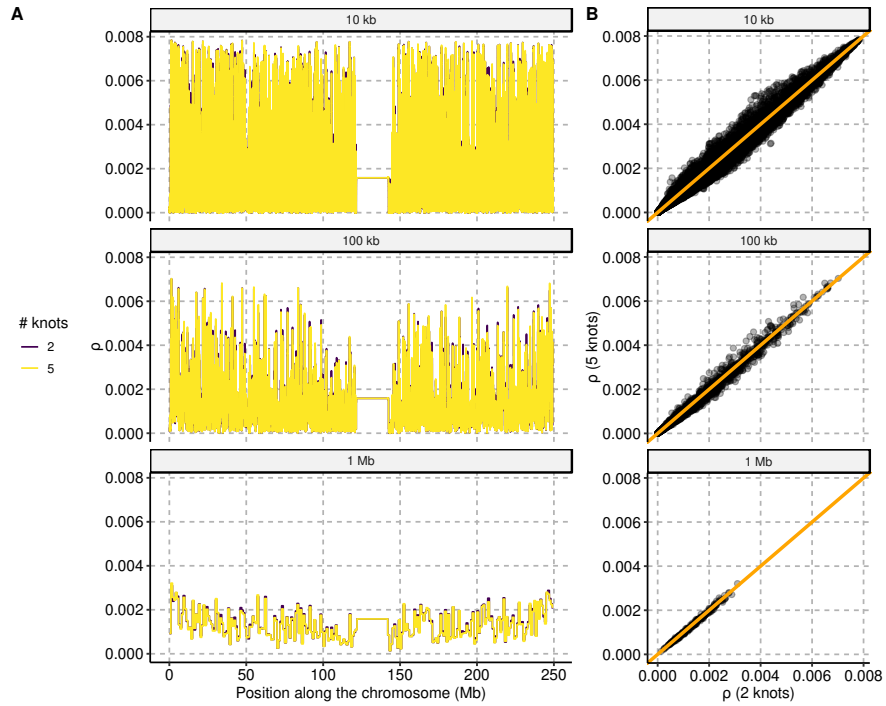
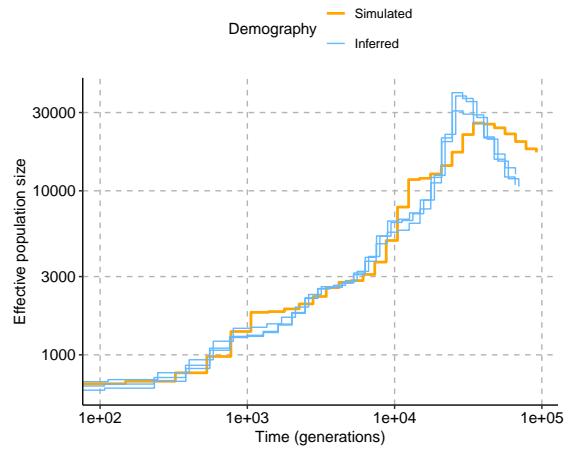


Fig. 3 Inferred recombination landscape of Altai Neanderthal chromosome 1, with distinct demography models and averaged over three different window sizes. A) Landscapes plotted along the chromosome, for two-node and five-node spline models. B) Comparison of window estimates under the two spline models. The orange line indicates the 1:1 ratio.

Fig. 4 Re-inferred demography from three simulated data sets using the model (demography + recombination map averaged in 10 kb windows) inferred from Altai Neanderthal's chromosome 1, with a heterogeneous recombination rate. A mutation rate of $1.25e-8 \text{ bp}^{-1}$ was used in the simulations and to scale the axes.



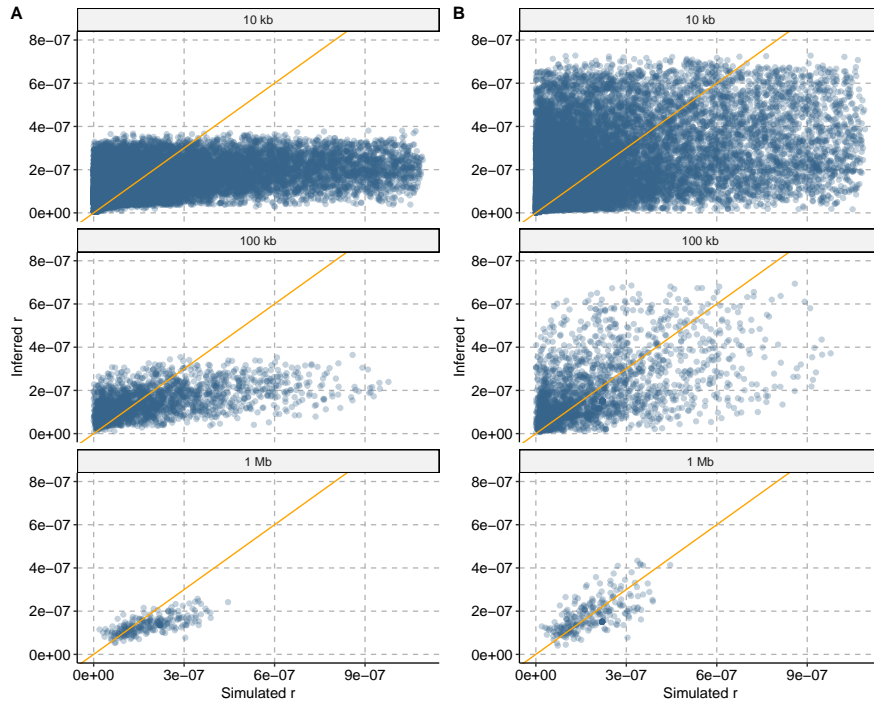


Fig. 5 Inferred recombination landscape from one simulation output (see Figure 4), plotted against the simulated landscape. A) Posterior decoding using the default discretization scheme. B) Posterior decoding using a custom discretization scheme with additional point mass at the tails of the distribution. The orange lines indicate the 1:1 ratio.

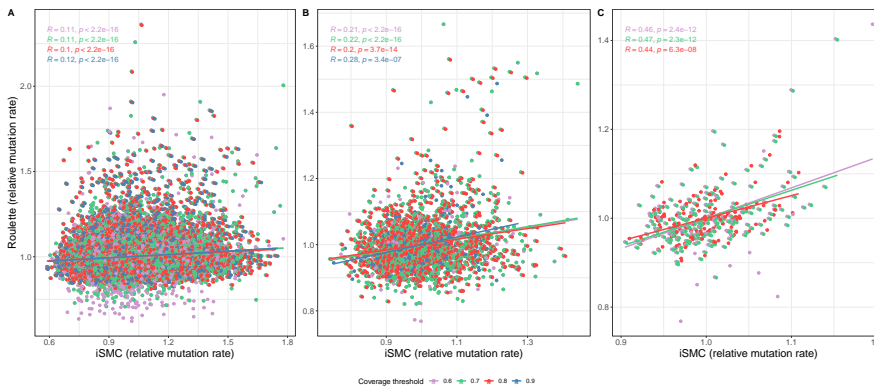


Fig. 6 Pearson correlations between iSMC-inferred and Roulette-predicted mutation rates at 10 kb (A), 100 kb (B) and 1 Mb (C) scales. Colors denote filtering thresholds, i.e., the minimum coverage required for keeping each window in the analysis.

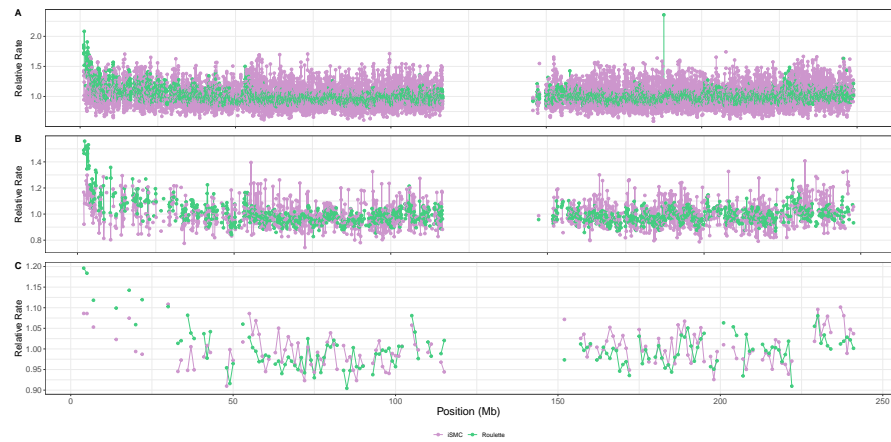


Fig. 7 Mutation maps for human chromosome 1, inferred with iSMC (purple) or predicted by Roulette (green) at 10 kb (A), 100 kb (B) and 1 Mb (C) scales. A coverage threshold of 80% is applied in each. Shown are relative rates for each method. iSMC displays larger variance than Roulette in A and B, likely due to estimation noise in these scales.