



CoMuS: Simulating coalescent histories and polymorphic data from multiple species

Journal:	<i>Molecular Biology and Evolution</i>
Manuscript ID:	Draft
Manuscript Type:	Article
Date Submitted by the Author:	n/a
Complete List of Authors:	Papadantonakis, Stefanos; University of Crete, Biology Poirazi, Panayiota; Foundation for Research and Technology - Hellas (FORTH), Institute for Molecular Biology and Biotechnology (IMBB) Pavlidis, Pavlos; Foundation of Research and Technology - Hellas (FORTH), Institute of Molecular Biology and Biotechnology (IMBB)
Key Words:	Coalescent, Multi-Species Coalescent, Simulations, Approximate Bayesian Computation, Parameter Inference

SCHOLARONE™
Manuscripts

CoMuS: Simulating coalescent histories and polymorphic data from multiple species

Stefanos Papadantonakis¹, Panayiota Poirazi², Pavlos Pavlidis^{2*}

¹Department of Biology, University of Crete, PO Box 2208, 71409 Heraklio, Greece

²Institute of Molecular Biology and Biotechnology (IMBB), Foundation for Research and Technology-Hellas (FORTH), 70013 Heraklio, Greece

*Corresponding author:

Pavlos Pavlidis

pavlidisp@gmail.com; pavlidis@imbb.forth.gr

Institute of Molecular Biology and Biotechnology (IMBB)

Foundation of Research and Technology - Hellas (FORTH)

Nikolaou Plastira 100, Vassilika Vouton

GR - 711 10, Heraklion, Crete, Greece

Keywords: Coalescent, Multi-Species Coalescent, simulations, Approximate Bayesian Computation, parameter inference

Running Title: CoMuS

Nonstandard abbreviations: MSC, ABC

31 ***Abstract***

32 The advent of modern DNA sequencing technology is the driving force in obtaining complete
33 intra- and inter-specific genomes. The simultaneous analysis of intra- and inter-species
34 variation is difficult mainly because our knowledge about patterns of polymorphisms in
35 samples where both intra- and inter-species sequences co-exist is limited. In the present study,
36 we implement CoMuS (Coalescent of Multiple Species), a multi-species coalescent software
37 that can simulate intra- and inter-species polymorphisms. CoMuS supports a variety of birth-
38 death models for speciation events and demographic scenarios related to the history of each
39 species. In CoMuS, speciation can be characterized either by instant isolation or by partial
40 isolation between sister species. Our software supports both the infinite and the finite site
41 model. As an application, we estimate parameters from simulated data, as well as the
42 speciation time and mutation rate using human-chimpanzee alignments. Furthermore, we
43 study the relations of summary statistics in multi-species scenarios. We expect that CoMuS
44 will be particularly useful for studies where species have been separated recently from their
45 common ancestor and phenomena such as incomplete lineage sorting or introgression still
46 occur.

47 Availability: CoMuS is implemented in C programming language and source code is
48 available on www.pop-gen.eu. Most updated and experimental versions of the code are
49 available from an online repository at [bitbucket.org \(git@bitbucket.org:idaios/comus.git\)](https://bitbucket.org/idaios/comus.git).

50

51 ***Introduction***

52 Coalescent simulation represents a Monte Carlo process to generate samples drawn from a
53 Wright-Fisher model of evolution. The most famous implementation of single species
54 coalescent is Hudson's *ms* software (Hudson 2002). *ms* can simulate samples from neutrally
55 evolving populations, allowing for migration, recombination, gene conversion and ancestral
56 changes of the population size. Due to its efficiency and flexibility to simulate a multitude of
57 evolutionary scenarios, *ms* has been used in Approximate Bayesian Computation (ABC)
58 inference methods to construct distributions of summary statistics under various evolutionary
59 models (e.g. Pavlidis et al. 2010; Saminadin-Peter et al. 2012).

60 Here, we extend Hudson's *ms* to simulate multiple species, thus implementing the multi-
61 species coalescent (MSC). Multi-species coalescent (Hobolth et al. 2007; Heled and
62 Drummond, 2010) represents a retrospective model of evolution when speciation events have
63 occurred in the ancestry of the sample. Thus, similar to single-species coalescent it effectively
64 models random genetic drift within species, but it also allows for speciation events at a
65 specific point in time, after which the two sub-species evolve in total isolation. Recently,
66 Heled et al. (2013) proposed a modification of MSC to model speciation more realistically.
67 Thus, Heled et al. (2013) allowed lineages from different species to remain in partial contact
68 for a period after the speciation event. The study by Heled et al. (2013) focuses on
69 phylogenetic applications of MSC and specifically on the effect of gene flow on phylogenetic
70 inference and migration rate estimation. MSC has been exploited in several different
71 applications of evolutionary biology. In phylogenetics, it has been used to study the effect of
72 incomplete lineage sorting (Degnan and Rosenberg 2006; Hobolth et al. 2007; Heled and
73 Drummond 2010; Mossel and Roch (unpublished work at *arxiv*) 2007), which is the failure of
74 gene copies to coalesce within the species time boundaries and, thus, they coalescent in

1
2
3 75 ancestral species. In population genetics, Hobolth et al. (2007) analyzed a sample of human,
4
5 76 chimpanzee, and gorilla sequences by considering the MSC of these three species and
6
7 77 applying a hidden Markov model to infer ancestral effective population sizes and speciation
8
9 78 times. Furthermore, Zhang et al. (2013) applied MSC simulations to assess the performance of
10
11 79 their species delimitation method under various speciation rates. In this article, we propose
12
13 80 CoMuS, an open-source MSC simulation software written in C programming language.
14
15 81 CoMuS combines Hudson's *ms* with Rambaut's Seq-Gen (Rambaut and Grassly, 1997) and
16
17 82 various birth-death speciation models presented by T. Stadler (Stadler 2009; Hartmann et al.
18
19 83 2010). Since it is based on Hudson's *ms*, it can efficiently implement recombination,
20
21 84 population size changes, multiple-populations per species, gene flow (or species
22
23 85 introgression). By implementing CoMuS, our aim is to provide a tool that can generate
24
25 86 simultaneously intra- and inter-species samples. Thus, it can be used to infer parameters when
26
27 87 both intra- and inter-species variation is available (for example, in an ABC framework), to
28
29 88 study the effect of incomplete lineage sorting in phylogenetics tree reconstruction or to test
30
31 89 species delimitation methods (e.g. Zhang et al. 2013).

32 33 34 35 36 37 38 39 40 *New Approaches*

41
42 91 CoMuS extends coalescent simulations to multiple species. It allows simulations under either
43
44 92 the infinite site or the finite site model. Furthermore, most of the flexibility of Hudson's *ms*
45
46 93 has been preserved. Thus, complicated evolutionary scenarios can be simulated efficiently.
47
48 94 Simulations are guided by a phylogenetic tree, which specifies the species boundaries. The
49
50 95 guide phylogenetic tree can be either simulated, under a variety of speciation models, or it is
51
52 96 provided by the user. By using simulated guide trees, uncertainty about the speciation process
53
54 97 is incorporated in the results. One important addition in CoMuS is the implementation of the
55
56 98 gradual isolation model after a speciation event. Thus, two species can be generated at a
57
58
59
60

1
2
3 99 certain time point, but decreasing gene flow can continue until complete isolation. Finally,
4
5
6 100 we have developed CoMuStats, a C++ software that can calculate summary statistics from the
7
8 101 output of CoMuS, or from user's single- or multi-FASTA alignment files. For more details
9
10 102 please consult the Materials and Methods section, as well as the manual and the examples that
11
12 103 are provided with the software.

16 104 **Results**

17
18 105 CoMuS can be used in studies where random intra- and inter-species samples are required. We
19
20 106 focus on two related applications: first, using an Approximate Bayesian Computation (ABC;
21
22 107 Beaumont et al. 2002) framework we estimate the speciation time either on real data, between
23
24 108 human and chimpanzee or on simulated data. Second, we study the pairwise relations of
25
26 109 summary statistics that are commonly used in population genetics studies. We show that even
27
28 110 though we simulate inter-species variation, summary statistics that are commonly used in
29
30 111 population genetics studies are informative for the inference of the simulation parameters.
31
32 112 ABC is used widely in population genetics studies (e.g. Excoffier et al. 2005; Duchon et al.
33
34 113 2013; Gray et al. 2014). The second application demonstrates the usability of CoMuS on
35
36 114 studying the relations of summary statistics using simulations. All scripts and command lines
37
38 115 used for the generation of simulated data and analysis can be downloaded from [## 46 117 **Estimation of parameters from simulated data**](http://pop-
39
40 116 <u>gen.eu/wordpress/software/comus-coalescent-of-multiple-species.</u></p></div><div data-bbox=)

47
48 118 For this application, we have simulated two scenarios: (i) inference of the birth rate of the
49
50 119 speciation process using data from two species (10 sequences per species); and (ii) inference
51
52 120 of the birth rate and first speciation time (i.e. time of the most recent common ancestor,
53
54 121 TMRCA) using data from 10 species (10 sequences per species).
55
56
57
58
59
60

1
2
3 122 **Scenario (i):** For the first scenario, mutation and recombination rates were fixed and assumed
4
5 123 to be known. The prior distribution of the birth rate b follows a uniform distribution, with
6
7 124 parameters 0 and 500 (i.e. $b \sim U(0,500)$). To assess the accuracy of the method we produced
8
9 125 1,000 pseudo-observed datasets with $b = 5$ for each of them. In the ABC framework, a
10
11 126 parameter can be inferred by using the mean, the median or the mode of the posterior
12
13 127 distribution. Figure 1 shows the density of the inferred values when the mean, the median or
14
15 128 the mode of the posterior distribution is used. Note that b is overestimated by all approaches.
16
17 129 This overestimation reflects the prior $U(0,500)$, which puts most of the density in large values
18
19 130 of b .

20
21
22 131 **Scenario (ii):** For the second scenario, mutation and recombination rates were fixed and
23
24 132 assumed known. Unknown parameters comprise the time of the root of the phylogenetic tree
25
26 133 (TMRCA) and the birth rate of the speciation process. The prior for the TMRCA is log-
27
28 134 uniform within the range $[0.0001, 1]$. Log-uniform priors are used in ABC for distributions
29
30 135 that span several orders of magnitude and we wish to assign the same density on each order of
31
32 136 magnitude. The prior for the birth rate b follows a uniform distribution $U(0, 100)$. As in
33
34 137 scenario (i), we generated 1,000 pseudo-observed datasets; b and TMRCA are 80 and 0.033,
35
36 138 respectively, for each dataset. Figure 2 shows the distribution of the median, mean and mode
37
38 139 values for each of the parameters.

47 48 140 **Estimation of the speciation time between human and chimpanzee**

49
50 141 We used 50 homologous gene regions between human and chimpanzee (kindly provided by
51
52 142 Qihui Zhu, personal communication; Supplementary Table 1). Using all 50 fragments we
53
54 143 calculated common population genetics summary statistics (Supplementary Table 2).
55
56 144 Simulations were performed with CoMuS assuming a finite site model (HKY; Hasegawa et al.
57
58 145 1985). The length of *simulated* fragments equals the average length of real fragments (1,763
59
60

1
2
3 146 bp). Simulation parameters are given in Supplementary Text 1.
4
5

6 147 In this demo application, we have been focused on estimating the speciation time between two
7
8 148 sister species as well as a total mutation parameter θ ; therefore, we neglect the demographic
9
10 149 history of each species. In a real application, however, demographic parameters for each
11
12 150 species should also be inferred since they may affect the inference of the speciation time,
13
14 151 especially in recent speciation events. Using CoMuS and ABC it is possible to infer
15
16 152 simultaneously both the speciation time and the demographic history of each species.
17
18 153 However, inferring such a complicated evolutionary history is out of the scope of the current
19
20 154 study.
21
22
23
24

25 155 Simulations were performed by drawing random variables (log-uniform) for the speciation
26
27 156 time (TMRCA) and the total mutation rate (θ). Note that the units of the speciation time is the
28
29 157 usual phylogenetic unit (i.e. expected substitutions per site). The prior distribution for both
30
31 158 parameters was uniform on the log-scale. However, we have conditioned on the existence of
32
33 159 SNPs in the simulations to be able to compute summary statistics. Thus, instances of very
34
35 160 recent speciation times that produced no SNPs were not included in the analysis (i.e. the prior
36
37 161 density of recent speciation times is lower). The median, mean and mode of the TMRCA is
38
39 162 0.0051, 0.0053 and 0.0026 expected substitutions per site, respectively. The total mutation
40
41 163 rate parameter θ corresponds to the parameter $4N_e\mu$, where N_e is the effective population size
42
43 164 and μ the mutation rate for the whole genomic region. In our context, though, N_e reflects a
44
45 165 measurement of the 'total' population size. The median, mean and mode of θ is 7.6×10^{-5} , 2.1
46
47 166 $\times 10^{-4}$, and 3.5×10^{-5} per base pair, respectively (Figure 3). In other words, θ reflects the rate at
48
49 167 which mutations occur on the ancestral lineages of both human and chimpanzee from the
50
51 168 present until their most recent common ancestor.
52
53
54
55
56
57
58
59
60

169 **Generating scatterplots of summary statistics under various evolutionary** 170 **models**

171 CoMus, together with summary statistics calculation software, can be used to produce
172 distributions and scatterplots of summary statistics under various evolutionary scenarios. We
173 have implemented a C++ software, CoMuStats, based on Libsequence (Thornton 2003) that is
174 able to calculate commonly used population genetics summary statistics. Thus, it is possible to
175 study marginal distributions or the relations between summary statistics under various inter-
176 and intra-species models. CoMuStats can read the output of CoMuS, i.e. a file with multiple
177 FASTA alignments separated by '//'. It produces a table of summary statistics where each line
178 is associated with one dataset and each column with one summary statistic. The output of
179 CoMuStats can be readily processed by R and produce either marginal distributions or
180 pairwise scatterplots. R scripts are available from [http://pop-
181 gen.eu/wordpress/software/comus-coalescent-of-multiple-species](http://pop-gen.eu/wordpress/software/comus-coalescent-of-multiple-species). CoMuStats can be also
182 used in the ABC framework for the calculation of summary statistics. Supplementary Figure 1
183 demonstrates both CoMuS and CoMuStats on producing summary statistics scatterplots. The
184 simulated scenario refers to 5 species with 15, 10, 10, 2, and 20 sampled sequences. The
185 length of the simulated genomic region is 10,000 bp, recombination rate is 100 and mutation
186 rate 20. The birth rate for the speciation process is 20.

187 ***Limitations and Discussion***

188 The current version of CoMuS has some limitations that we will be improved in future
189 versions of the software. CoMuS can run in a single CPU and is not able yet to exploit
190 multiple cores. Furthermore, the current version cannot simulate insertions and deletions
191 (indels). Indels represent an important class of mutations. However, incorporating them in the
192 current version of the software was challenging since their implementation requires a total
193 redesign of the software to be able to properly handle overlapping indel events. Future

1
2
3 194 versions of CoMuS will be able to fully utilize multi-core computers and simulate indel
4
5 195 events.

6
7
8 196 Even though it is possible to generate inter- and intra-species datasets using *ms* and Seq-Gen
9
10 197 in a pipeline (i.e. using the coalescent trees of *ms* as an input for Seq-Gen), this process has
11
12 198 some limitations and cannot fully substitute CoMuS. CoMuS has several advantages over the
13
14 199 simple pipeline approach: (i) it is simpler since the user does not need to implement his own
15
16 200 scripts to combine *ms* and Seq-Gen; (ii) it can simulate the phylogenetic tree under speciation
17
18 201 models with various values of birth rate, death rate, and species sampling proportion. It is also
19
20 202 possible to read in a species tree from the command line; (iii) it automatically performs the
21
22 203 time scale conversion; (iv) it implements the partial isolation model after a speciation event,
23
24 204 thus allowing two sister species to exchange genetic material for some time after the
25
26 205 speciation event. The rate of gene flow is maximum at speciation time and decreases linearly
27
28 206 until the complete isolation. This model is more realistic than an instantaneous isolation after
29
30 207 a speciation event (Heled et al. 2013).

31
32
33
34
35
36
37 208 Heled and Drummond (2010) have implemented *BEAST (StarBEAST) to infer species trees
38
39 209 from multilocus data. *BEAST implements the multi-species coalescent but their study
40
41 210 focuses on the inference of the species tree. In CoMuS we have focused on generating
42
43 211 polymorphic data from multiple species, allowing the full flexibility of Hudson's *ms* software.
44
45 212 Thus, CoMuS allows recombination, exponential size changes, introgression between
46
47 213 different species (i.e. horizontal gene transfer), gene flow between different populations.

48
49
50
51
52 214 CoMus can simulate efficiently intra- and inter-species genomic variability under different
53
54 215 evolutionary scenarios and mutation models. It is written in C and it is freely available under
55
56 216 the GNU GPLv3 license. Its core machinery is based on Hudson's *ms* (Hudson 2002) and Seq-
57
58 217 Gen (Rambaut and Grassly 1997). Together with CoMuS we have developed CoMuStats, a
59
60

1
2
3 218 C++ Libsequence software that can calculate common population genetics summary statistics
4
5 219 from multi-FASTA-alignment files. We demonstrate the usage of CoMuS and CoMuStats in
6
7 220 ABC and pairwise scatterplots of summary statistics. Both softwares as well as scripts used in
8
9
10 221 the analysis are available from [http://pop-gen.eu/wordpress/software/comus-coalescent-of-](http://pop-gen.eu/wordpress/software/comus-coalescent-of-multiple-species)
11
12 222 [multiple-species](http://pop-gen.eu/wordpress/software/comus-coalescent-of-multiple-species).

16 223 ***Material and Methods***

18 224 **Outline of the workflow**

19 225 The workflow during a typical simulation process performed by CoMuS is described in
20
21 226 Figure 4.

27 227 **The guide phylogenetic tree**

28 228 CoMuS implements the multispecies coalescent using a guide phylogenetic tree as input. The
29
30 229 phylogenetic tree delimits the species boundaries; specifically, it determines the speciation
31
32 230 time points. After a speciation event the user defines whether isolation happened
33
34 231 instantaneously or gradually as in Heled et al. (2013). In case of gradual isolation, the gene
35
36 232 flow between sister species is reduced linearly until species are completely isolated. The guide
37
38 233 phylogenetic tree is provided either by the user as input or it can be simulated. To simulate a
39
40 234 phylogenetic tree we assume a Yule process with birth and death events. As in INDELible
41
42 235 (Fletcher and Yang 2009), the following parameters are needed to simulate a phylogenetic
43
44 236 tree: birth rate b , death rate d , the number of species n , and the proportion r of the number of
45
46 237 species sampled.

52
53 238 **Generating the guide phylogenetic tree:** To simulate the guide phylogenetic tree, the
54
55 239 following assumptions are made (as in TreeSim; Stadler 2011): (i) all sequences are sampled
56
57 240 simultaneously at the present time point; (ii) both the birth and death rates are constant; (iii)
58
59
60

1
2
3 241 since the birth-death process does not require a definite start or an end, one of the following
4
5 242 conditions must hold to be able to simulate it according to our needs: (a) the number of
6
7 243 sampled species is fixed, (b) the process starts at a specified time in the past, (c) the age of the
8
9 244 most recent common ancestor of the sampled species is known, (d) the time that the process
10
11 245 starts is fixed at T_{origin} and we condition on the sampled number of species (e) the number of
12
13 246 sampled species is fixed, and we condition on the fact that the process is not older than a
14
15 247 certain time T_{oldest} . Considering the case (d) we have also implemented a special case where
16
17 248 $T_{\text{origin}} = 1.0$. This condition has been developed firstly by Yang and Rannala (1997) and is
18
19 249 implemented in INDELible software by Fletcher and Yang (2009). The guide phylogenetic
20
21 250 tree can also be provided by the user in newick tree format, and it is required to be rooted.
22
23
24
25
26
27

28 251 **Mutation model**

29
30 252 Hudson's *ms* assumes the infinite site model to simulate mutation events on ancestral lineages.
31
32 253 Consequently, each polymorphic site hosts two states. The infinite site model is justified in *ms*
33
34 254 because its goal is to model events that happen within a species boundaries, thus events that
35
36 255 are relatively young in the evolutionary time scale. Indeed, the probability of more than one
37
38 256 mutation at a given site is negligible for realistic mutation rate values. Contrary to single-
39
40 257 species coalescent, multiple mutation events may occur frequently in a multi-species setup.
41
42 258 Thus, CoMuS is able to simulate DNA sequences (and not just polymorphisms as in *ms*) given
43
44 259 an evolutionary mutation model (JC, Jukes and Cantor 1969; F84, Felsenstein 1984; HKY,
45
46 260 Hasegawa et al. 1985; GTR, Rodríguez et al. 1990).
47
48
49
50
51
52

53 261 **Time unit conversion**

54
55 262 The usual time unit in coalescent theory is the number of generations divided by the effective
56
57 263 population size N_e (or by $4N_e$ in *ms*). In other words, if the effective population size is N_e , then
58
59
60

1
2
3 264 a time period $t = 1$ corresponds to N_e generations. Consequently, the generation of a
4
5
6 265 coalescent tree becomes independent of the effective population size. On the other hand, in
7
8 266 phylogenetics, time is measured in expected numbers of substitutions per site. For example, a
9
10 267 branch of length 0.1 corresponds to a time period in which 0.1 substitutions per site occur on
11
12 268 average. To model coalescent processes on a phylogenetic tree, it is necessary to use identical
13
14 269 units for both the phylogenetic and the coalescent tree. Assume a branch of length l expected
15
16 270 substitutions per site. If the mutation rate per site and per generation is μ , and the number of
17
18 271 generations that correspond to the branch are γ , then:

22
23 272 $l = \mu \gamma \rightarrow \frac{\gamma}{4N_e} = \frac{l}{4N_e \mu} = \frac{l}{\theta}$. Since, $\frac{\gamma}{4N_e}$ represents time in $4N_e$ generations, we can divide

24
25
26 273 the branch length (in phylogenetic units) by θ to obtain the branch length in coalescent units.
27
28
29

30 274 **Implementation of speciation events as population merge events**

31
32 275 CoMuS is using the *ms* machinery to build genealogies. Since in *ms* the concept of species is
33
34 276 absent, we treat species as distinct populations of a common origin. Thus, a speciation event
35
36 277 involving species A and B is translated to a merge event between populations. When each
37
38 278 species A and B contain a single population, the process of merging them is simple and
39
40 279 equivalent to a single population merge event in accordance to *ms* 'j' events (implemented by
41
42 280 the -ej flag; see Figure 5A). On the other hand, the situation is more complicated when a
43
44 281 species contains more than one populations. In this case, by convention, we merge
45
46 282 simultaneously each population to the population with the largest index (Figure 5B).
47
48
49
50

51 52 53 283 **Gradual Isolation after a speciation event**

54 284 Recently, Heled et al. (2013) presented a model where speciation may occur over an extended
55
56 285 time period. During this period, sister species are able to exchange genetic material, thus
57
58 286 creating a gradual isolation model after a speciation event. Gradual isolation may be a more
59
60

1
2
3 287 realistic model for speciation for several species including the speciation between human and
4
5 288 chimpanzee (Patterson et al. 2006). Backward in time, there is a time period (specified by the
6
7
8 289 user) that ends at the speciation time point, in which gene flow between different species is
9
10 290 allowed. The rate of gene flow is increasing linearly (backward in time). To sample times
11
12 291 until the next lineage migration event we use a non-homogeneous Poisson process with linear
13
14 292 intensity (gene flow rate changes as a linear function of time from 0 to λ_{\max}). In brief, we first
15
16 293 simulate a Poisson process with a maximum arbitrary intensity λ_{\max} (provided by the user; by
17
18 294 default the maximum rate value equals to 1.0). Then at time t , we accept values with a
19
20 295 probability $p(t) = \lambda(t)/\lambda_{\max}$, where $\lambda(t) = \lambda_{\max} (t - t_{\text{isolation}}) / (t_{\text{speciation}} - t_{\text{isolation}})$, where $t_{\text{isolation}}$ is the
21
22 296 time where the two sister species become completely isolated; $t_{\text{speciation}}$ is the speciation time,
23
24 297 i.e. the age of the node. It has been shown (Ross, 2006) that the process of counted events
25
26 298 corresponds to a non-homogeneous Poisson process with intensity function $\lambda(t) = \lambda_{\max}p(t)$. In
27
28 299 our case, we are interested only in the first event and the time needed to occur. This event
29
30 300 represents gene-flow during the post-speciation period that the sister species have not been
31
32 301 completely isolated.
33
34
35
36
37
38
39

40 **Availability**

41
42 303 CoMus as well as CoMuStats are freely available under the GPLv3.0 license and they are
43
44 304 available from <http://pop-gen.eu/wordpress/software/comus-coalescent-of-multiple-species>.
45
46 305 CoMuS is available also from bitbucket.org (git clone [git@bitbucket.org:idaio/comus.git](https://bitbucket.org/idaio/comus.git))
47
48 306 with public access level. Also, we have constructed a public group forum on Google
49
50 307 (<https://groups.google.com/forum/#!forum/popgen-comus>) where users can ask questions or
51
52 308 report problems.
53
54
55
56
57
58
59
60

1
2
3 309 *Acknowledgements*
4

5 310 We would like to thank Dr. Qihui Zhu (Harvard Medical School) for providing 50 gene
6
7 311 alignments between human and chimpanzee. This study was supported by the FP7 REGPOT-
8
9 312 InnovCrete (No. 316223) grant, as well as by the FP7-PEOPLE-2013-IEF EVOGREN
10
11 313 (625057) to P. Pavlidis.
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

314 **Bibliography**

- 315 1. Beaumont, MA., Zhang, W., & Balding, DJ. (2002). Approximate Bayesian
316 computation in population genetics. *Genetics*, 162(4), 2025–35.
- 317 2. Degnan, J.H., Rosenberg, N. (2006) Discordance of species trees with their most likely
318 gene trees. *PLoS Genet* 2:e68.
- 319 3. Duchen, P., Zivkovic, D., Hutter, S., Stephan, W., Laurent, S. (2013) Demographic
320 inference reveals African and European admixture in the North American *Drosophila*
321 *melanogaster* population. *Genetics* 193:291–301.
- 322 4. Excoffier, L., Estoup, A., Cornuet, J-M. (2005) Bayesian analysis of an admixture
323 model with mutations and arbitrarily linked markers. *Genetics* 169:1727–38.
- 324 5. Felsenstein, J. (1984) Distance methods for inferring phylogenies - a justification.
325 *Evolution* 38:16-24.
- 326 6. Fletcher, W., & Yang, Z. (2009). INDELible: a flexible simulator of biological
327 sequence evolution. *Mol Biol Evol*, 26:1879–88.
- 328 7. Gray, MM., Wegmann, D., Haas R.J., White, M., Gabriel, S.I., Searle, J.B., Cuthbert,
329 R.J., Ryan, P.G., Payseur, B.A. (2014) Demographic History of a Recent Invasion of
330 House Mice on the Isolated Island of Gough. *Mol Ecol* 23:1923–39.
- 331 8. Hartmann, K., Wong, D., Stadler, T. (2010) Sampling trees from evolutionary models.
332 *Syst Biol* 59:465–76.
- 333 9. Hasegawa, M., Kishino, H., Yano, T. (1985) Dating of the human-ape splitting by a
334 molecular clock of mitochondrial DNA. *J Mol Evol* 22:160–74.
- 335 10. Heled, J., Bryant, D., Drummond, AJ. (2013) Simulating gene trees under the
336 multispecies coalescent and time-dependent migration. *BMC Evol Biol* 13:44.
- 337 11. Heled, J., Drummond, AJ. (2010) Bayesian inference of species trees from multilocus
338 data. *Mol Biol Evol* 27:570–80.
- 339 12. Hobolth, A., Christensen, OF., Mailund, T., Schierup, MH. (2007) Genomic
340 relationships and speciation times of human, chimpanzee, and gorilla inferred from a
341 coalescent hidden Markov model. *PLoS Genet* 3:e7.
- 342 13. Hudson, R. R. (2002). Generating samples under a Wright-Fisher neutral model of
343 genetic variation. *Bioinformatics*, 18:337–8.
- 344 14. Jukes. H., and Cantor, CR. (2006), *Evolution of Protein Molecules*, New York:
345 Academic Press
346
- 347 15. Mossel, E., Roch, S. (2013) Incomplete Lineage Sorting : Consistent Phylogeny
348 Estimation From Multiple Loci. *arxiv.org*.
- 349 16. Patterson, N., Richter, DJ., Gnerre, S., Lander, ES., Reich, D. (2006) Genetic evidence
350 for complex speciation of humans and chimpanzees. *Nature* 441:1103–8.
- 351 17. Rambaut A, Grassly NC (1997) Seq-Gen: an application for the Monte Carlo
352 simulation of DNA sequence evolution along phylogenetic trees. *Comput Appl Biosci*
353 13:235–8.

- 1
2
3 354 18. Rodríguez, F., Oliver, J. L., Marín, A., & Medina, J. R. (1990). The general stochastic
4 355 model of nucleotide substitution. *Journal of Theoretical Biology*, 142:485–501.
5
6 356 19. Ross, SM. (2006) *Simulation*, New York, Academic Press.
7
8 357 20. Saminadin-Peter, SS., Kemkemer, C., Pavlidis, P., Parsch, J. (2012) Selective sweep of
9 358 a cis-regulatory sequence in a non-African population of *Drosophila melanogaster*.
10 359 *Mol Biol Evol* 29:1167–74.
11
12 360 21. Stadler, T. (2009) On incomplete sampling under birth-death models and connections
13 361 to the sampling-based coalescent. *J Theor Biol* 261:58–66.
14
15 362 22. Stadler, T. (2011) Simulating trees with a fixed number of extant species. *Syst Biol*
16 363 60:676–84.
17
18 364 23. Thornton, K. (2003) Libsequence: a C++ class library for evolutionary genetic
19 365 analysis. *Bioinformatics* 19:2325–7.
20
21 366 24. Yang, Z., & Rannala, B. (1997). Bayesian phylogenetic inference using DNA
22 367 sequences: a Markov Chain Monte Carlo Method. *Mol Biol Evol*, 14(7), 717–24.
23
24 368 25. Zhang, J., Kapli, P., Pavlidis, P., Stamatakis, A. (2013) A general species delimitation
25 369 method with applications to phylogenetic placements. *Bioinformatics* 29:2869–76.
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4 370 **Figure Legends**
5

6 371 **Figure 1:** Densities of the inferred mean, median and mode of the birth rate. We generated
7
8 372 1,000 pseudo-observed datasets with birth rate $b = 5$. To infer b , we constructed 100,000
9
10 373 simulated data in which $b \sim U(0, 500)$. For each of the 1,000 pseudo-observed datasets the
11
12 374 posterior distribution of b was computed using the ABC framework, and the mean, median
13
14 375 and mode values were extracted and used to generate the plot. All three measurements
15
16 376 overestimate the value of b . This is because the prior distribution of b heavily favors large
17
18 377 values of b ($b \sim U(0, 500)$).
19
20
21
22

23 378
24
25

26 379 **Figure 2:** Densities of the inferred mean, median and mode of TMRCA and birth rate. We
27
28 380 generated 1,000 pseudo-observed datasets with birth rate $b = 80$ and TMRCA = 0.033. To
29
30 381 infer the parameters b and TMRCA, we constructed 100,000 simulated data in which $b \sim U(0,$
31
32 382 $100)$ and the $\log_{10}(\text{TMRCA}) \sim U(0.0001, 1)$, i.e. log-uniform in $[0.0001, 1]$. For each of the
33
34 383 1,000 pseudo-observed datasets the posterior distributions of b and TMRCA were computed
35
36 384 using the ABC framework, and the mean, median and mode values were extracted and used to
37
38 385 generate the plots. Both of the phylogenetic parameters are precisely estimated using the ABC
39
40 386 framework.
41
42
43
44

45 387
46
47

48 388 **Figure 3:** Estimation of the speciation time (TMRCA) and the total mutation rate parameter
49
50 389 (Theta). Parameters were inferred using the ABC framework as it is implemented in the 'abc'
51
52 390 package of the R statistical language. Simulations were performed by drawing random
53
54 391 variables (log-uniform) for the speciation time (TMRCA) and the total mutation rate (θ or
55
56 392 Theta). The x-axis in both plots are in log-scale. A) Inference of the speciation time. Dotted-
57
58
59
60

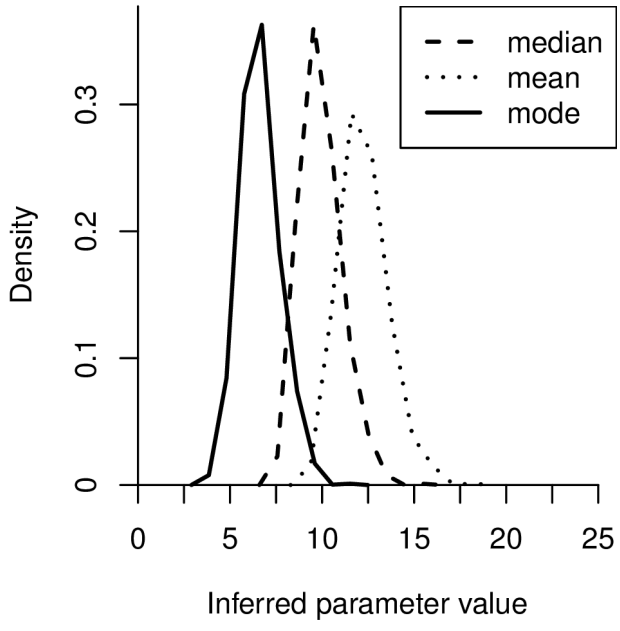
1
2
3 393 line is the prior distribution of the speciation time, dashed-line the distribution of the
4
5 394 speciation time after the rejection step of ABC and the solid line represents the posterior
6
7 395 distribution. Prior was uniform on the log-scale but we have conditioned on the existence of
8
9 396 SNPs in the simulations. Thus, instances of very recent speciation times that produced no
10
11 397 SNPs were not included in the analysis (i.e. the prior density of recent speciation times is
12
13 398 lower). The median, mean and mode of the TMRCA is 0.0051, 0.0053 and 0.0026 expected
14
15 399 substitutions, respectively. B) Estimation of the total mutation rate θ . θ corresponds to the
16
17 400 parameter $4N_e\mu$, where N_e is the effective population size and μ the mutation rate for the
18
19 401 whole genomic region. In our context, though, N_e reflects a measurement of the 'total'
20
21 402 population size. Since the species are isolated for a long period of time, N_e is very large (since
22
23 403 coalescent is not allowed). The median, mean and mode of θ is 7.6×10^{-5} , 2.1×10^{-4} , and $3.5 \times$
24
25 404 10^{-5} per base pair, respectively.
26
27
28
29
30
31
32
33
34

35 406 **Figure 4:** The workflow during a typical multi-species simulation with CoMuS. Dashed-line
36
37 407 boxes denote optional steps. Splits in the workflow denote alternative paths that a user may
38
39 408 follow.
40
41
42
43
44

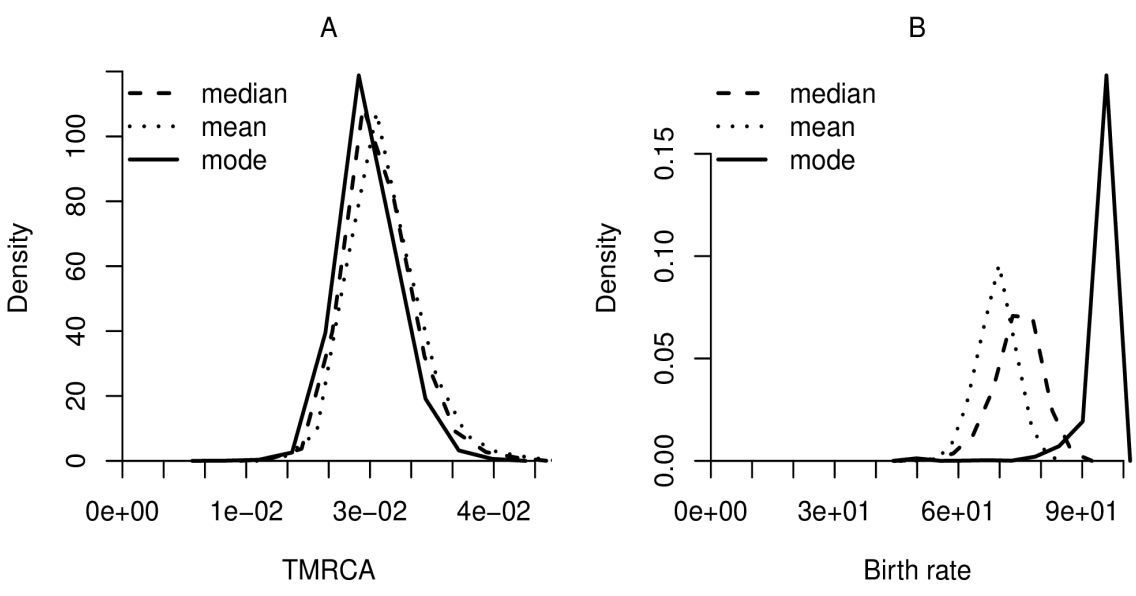
45 410 **Figure 2:** Speciation model as implemented in CoMuS. In (A), each species comprises one
46
47 411 population. Therefore, all lineages from population 1 (pop1) 'migrate' to population 2 (pop2)
48
49 412 at the speciation time. In (B), species one consists of two populations, pop1 and pop2. Thus,
50
51 413 two events take place at speciation time: The lineages of pop1 migrate to pop3 and at the same
52
53 414 time lineages from pop2 migrate to pop3. By convention, at each node of the guide tree
54
55 415 (nodes define speciation events) lineages migrate to the population with the largest index
56
57 416 (here, pop3).
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

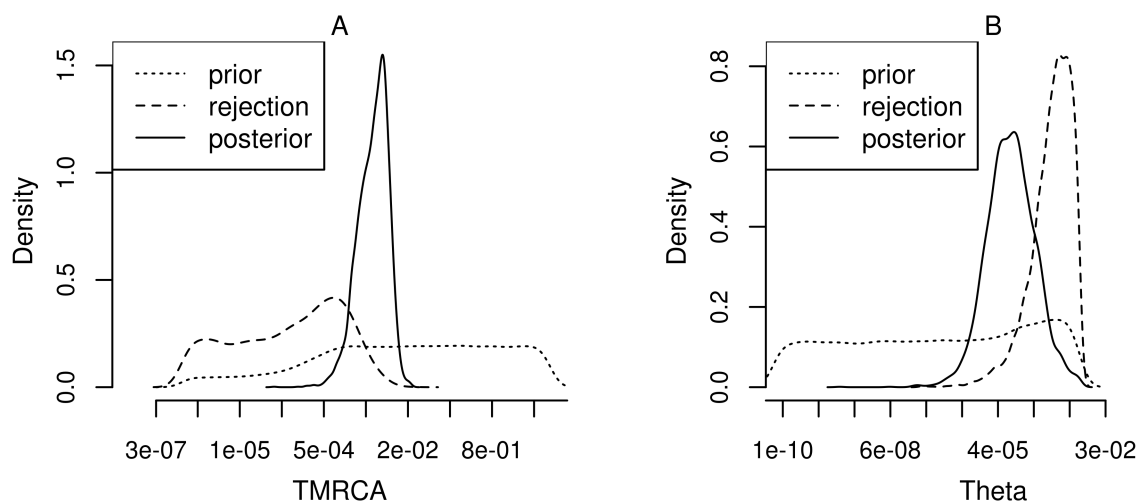
417 **Figures**
418 **Figure 1**



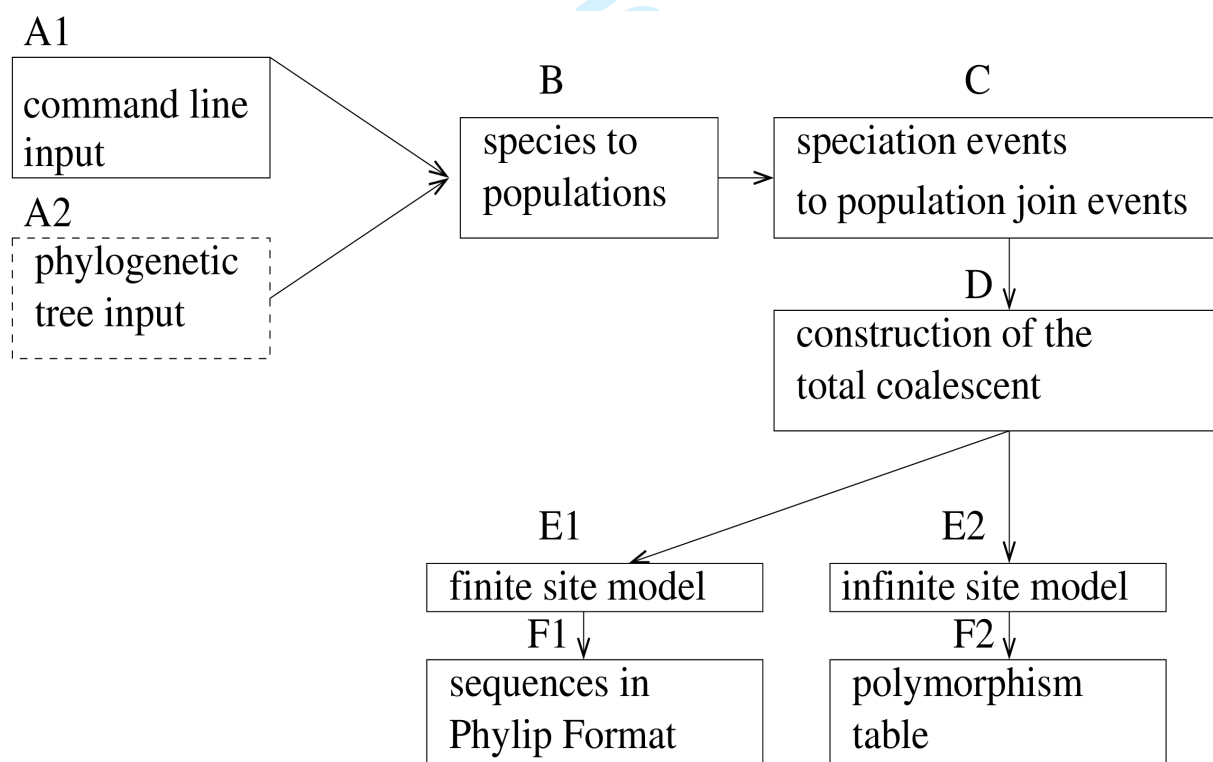
420 **Figure 2**



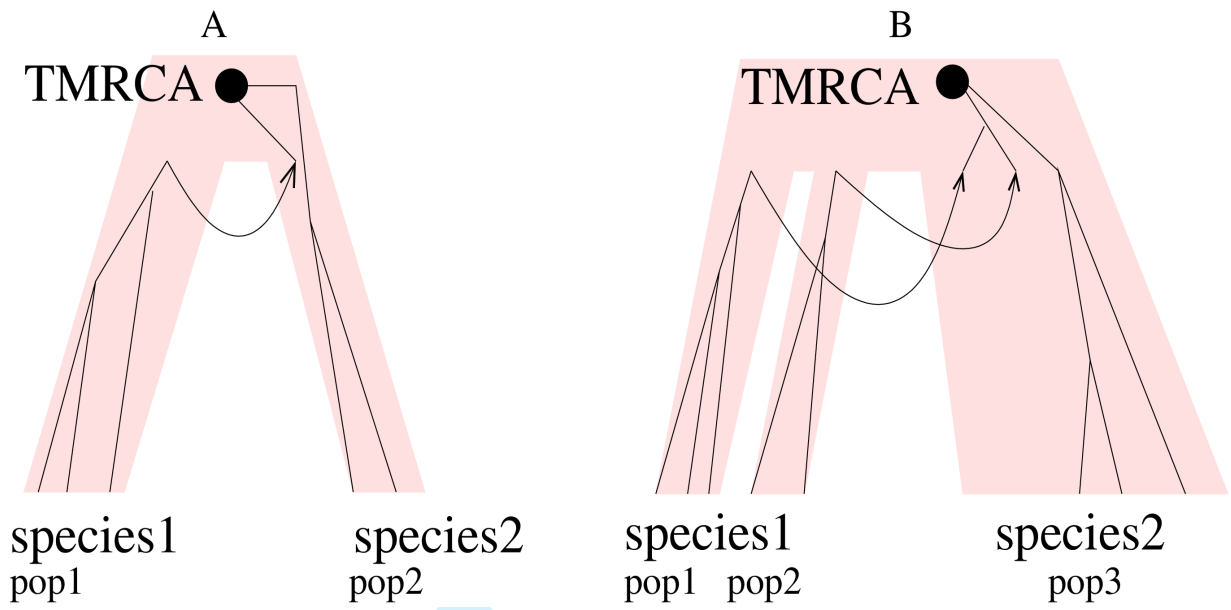
421 **Figure 3**



423 **Figure 4**



425 **Figure 5**



Of: Mol. Biol. Evol.