

1 Cophylogeny reconstruction allowing for 2 multiple associations through approximate 3 Bayesian computation – Supplementary 4 Material

5 BLERINA SINAIMERI^{1,2}, LAURA URBINI^{2*}, MARIE-FRANCE SAGOT² AND
6 CATHERINE MATIAS³

7 ¹ *LUISS University, Rome, Italy*

8 ² *Inria Lyon, 56 Bd Niels Bohr, 69100 Villeurbanne, France, and Université de*
9 *Lyon, F-69000, Lyon; Université Lyon 1; CNRS, UMR5558; 43 Boulevard du 11*
10 *Novembre 1918, 69622 Villeurbanne cedex, France*

11 ³ *Sorbonne Université, Université de Paris Cité, Centre National de la Recherche*
12 *Scientifique, Laboratoire de Probabilités, Statistique et Modélisation, Paris,*
13 *France*

14 **Corresponding author:** Blerina Sinaimeri, LUISS University, Rome, Italy;

15 E-mail: bsinaimeri@luiss.it.

16 Contents

18	A The event-based model	2
19	A.1 Tree-related basic definitions	3

*First co-authors.

20	A.2 Reconciliation model from Tofigh et al.	4
21	A.3 Reconciliation model allowing for spreads	6
22	A.4 Pre-estimating probabilities for the spread events	9
23	B AMOCOALA algorithm	14
24	B.1 Simulation algorithm in AMOCOALA	14
25	B.2 ABC-SMC inference method in AMOCOALA	17
26	B.3 Distance measure in AMOCOALA	19
27	B.4 A proof that dMASST is a distance	21
28	B.5 Polynomial time algorithm for computing the dMASST distance	24
29	C Additional results for the self-test	25
30	D Biological datasets	25
31	D.1 Results on biological datasets	26
32	D.2 Running times	27
33	D.3 Robustness analysis wrt the pre-estimated spread probabilities	27

34 A The event-based model

35 AMOCOALA relies on the event-based model presented in Charleston (2002);
 36 Tofigh *et al.* (2011). For the sake of completeness, we detail the model here.
 37 We first start with some basic definitions related to phylogenetic trees.

38 A.1 Tree-related basic definitions

39 A rooted phylogenetic tree is a leaf-labelled tree that models the evolution
 40 of a set of taxa from their most recent common ancestor (placed at the
 41 root). The internal vertices of the tree correspond to the speciation events.
 42 In a rooted phylogenetic tree, a direction is assumed from the root to the
 43 leaves that corresponds to the direction of evolutionary time. Specifically,
 44 a phylogenetic tree is a rooted tree with labelled leaves where the root has
 45 in-degree 0 and out-degree 2, the leaves have in-degree 1 and out-degree 0
 46 and every internal vertex has in-degree 1 and out-degree 2. For such a tree
 47 T , the set of vertices is denoted by $V(T)$, the set of arcs by $A(T)$, and the set
 48 of leaves by $L(T)$. The cardinality of set A is denoted by $|A|$. The root of T
 49 is denoted by $r(T)$. For a vertex v in a tree T , we denote by T_v the subtree
 50 of T rooted in v (often referred to as a *clade*), and we write $L(v)$ for the set
 51 $L(T_v)$. For a vertex $v \in V(T)$, we denote by $Des(v)$ the set of *descendants* of
 52 v , *i.e.* the set of vertices in the subtree of T_v . Similarly, we denote by $Anc(v)$
 53 the set of *ancestors* of v , that is the set of vertices in the unique path from
 54 $r(T)$ to v (including the end points). For a vertex $v \in V(T)$ different from
 55 the root, we call its *parent*, denoted by $par(v)$, the vertex x for which there
 56 is the arc $(x, v) \in A(T)$. We denote by $mrca(v, w)$ the most recent common
 57 ancestor of v and w in T . Finally, we denote by \leq the partial order induced
 58 by the ancestry relation in the tree. Formally, for $x, y \in V(T)$, we say that
 59 $x \leq y$ if $x \in Anc(y)$. If neither $x \in Anc(y)$ nor $y \in Anc(x)$, the vertices x

60 and y are said to be *incomparable*.

61 For any tree T and any set of leaves t_1, \dots, t_n , we denote by $T_{|\{t_1, \dots, t_n\}}$
 62 the phylogenetic subtree of T induced by the leaves t_1, \dots, t_n and eventually
 63 suppressing the vertices of out-degree 1. When a vertex u with parent vertex
 64 v and child vertex w is suppressed, both vertex u and arcs $(v, u), (u, w)$ are
 65 removed and the arc (v, w) is added to the tree.

66 **A.2 Reconciliation model from Tofigh et al.**

67 In this section, we describe the classical reconciliation model, where 4 co-
 68 evolutionary events are allowed, producing no multiple associations. Let H
 69 and S be respectively the rooted phylogenetic trees of the host and sym-
 70 biont species, both binary and full (*i.e.* each internal vertex has exactly
 71 two children). Let ϕ be a function from $L(S)$ to $L(H)$, representing the
 72 symbiont/host associations between extant species. A reconciliation is a
 73 function λ that assigns, for each symbiont vertex $s \in V(S)$, a host vertex
 74 $\lambda(s) \in V(H)$, and satisfies the conditions stated in Definition A.1.

75 In its classical form, a reconciliation associates to each vertex s in $V(S)$
 76 an event $E(\lambda(s))$ among cospeciation (\mathbb{C}), duplication (\mathbb{D}) and host switch
 77 (\mathbb{S}).

78 **Definition A.1.** *Given two phylogenetic trees S and H , and a function $\phi :$
 79 $L(S) \rightarrow L(H)$, a reconciliation of (S, H, ϕ) is a function $\lambda : V(S) \rightarrow V(H)$
 80 satisfying the following:*

- 81 1. For every leaf vertex $s \in L(S)$, we have $\lambda(s) = \phi(s)$.
- 82 2. For every internal vertex $s \in V(S) \setminus L(S)$ with children s_1, s_2 , exactly
- 83 one of the following applies:
 - 84 (a) $E(\lambda(s)) = \mathbb{S}$, that is, either $\lambda(s_1)$ and $\lambda(s)$ are incomparable and
 - 85 $\lambda(s_2)$ is a descendant of $\lambda(s)$, or $\lambda(s_2)$ and $\lambda(s)$ are incomparable
 - 86 and $\lambda(s_1)$ is a descendant of $\lambda(s)$,
 - 87 (b) $E(\lambda(s)) = \mathbb{C}$, that is, $\text{mrca}(\lambda(s_1), \lambda(s_2)) = \lambda(s)$, and $\lambda(s_1)$ and
 - 88 $\lambda(s_2)$ are incomparable,
 - 89 (c) $E(\lambda(s)) = \mathbb{D}$, that is, $\lambda(s_1)$ and $\lambda(s_2)$ are both descendants of
 - 90 $\lambda(s)$, and the previous two cases do not apply.

91 The loss event is denoted by \mathbb{L} and is identified by a multiset (general-
 92 isation of a set where the elements are allowed to appear more than once)
 93 whose elements are in $V(H)$ containing all the vertices $h \in V(H)$ that are
 94 in the path between the image of a vertex $s \in V(S)$ and the image of one of
 95 its children. The images themselves are not included in the count, except for
 96 the duplication event, where one of the images is included.

97 The function λ partitions the set of internal symbiont tree vertices into
 98 three disjoint subsets according to the coevolutionary event occurring at that
 99 vertex. The number of occurrences of each of the three events and the number
 100 of losses make up the *event vector* of the reconciliation. The *event vector* of
 101 a reconciliation is a vector of integers consisting of the total number of each
 102 type of events \mathbb{C} , \mathbb{D} , \mathbb{S} , \mathbb{L} .

103 We say that a reconciliation is *time-feasible* if it does not violate the time-
 104 feasibility constraints. The exact criterion we use to assess time-feasibility
 105 is the one defined in Stolzer *et al.* (2012) and that was already in force in
 106 COALA.

107 **A.3 Reconciliation model allowing for spreads**

108 The introduction of spread events modifies the previous setting in the fol-
 109 lowing way. Let again H and S be respectively the rooted phylogenetic trees
 110 of the host and symbiont species, both binary and full (*i.e.* every internal
 111 vertex has exactly two children). Now, let ϕ be a relation between $L(S)$ and
 112 $L(H)$, representing the symbiont/host associations between extant species.
 113 More precisely, let us denote $\mathcal{P}(L(H))$ the set of all subsets of $L(H)$. Then
 114 ϕ is now a function from $L(S)$ to $\mathcal{P}(L(H))$. For any extant symbiont species
 115 $s \in L(S)$, whenever the cardinality $|\phi(s)| \geq 2$ (*i.e.* whenever the symbiont
 116 is associated to more than one host), we say that this symbiont has multiple
 117 associations and we count the total number of multiple associations in the
 118 dataset as:

$$\text{Nb of multiple associations} = \sum_{s \in L(S)} (|\phi(s)| - 1).$$

119 A reconciliation is now a function λ from $V(S)$ to $\mathcal{P}(V(H))$ that assigns,
 120 for each symbiont vertex $s \in V(S)$, a set of host vertices $\lambda(s) \subset V(H)$,
 121 and satisfies the conditions stated in Definition A.2. A reconciliation now

122 associates to each vertex s in $V(S)$ an event $E(\lambda(s))$ among cospeciation (\mathbb{C}),
 123 duplication (\mathbb{D}), host switch (\mathbb{S}), vertical spread (\mathbb{VS}) and horizontal spread
 124 (\mathbb{HS}).

125 **Definition A.2.** *Given two phylogenetic trees S and H , and a function*
 126 *$\phi : L(S) \rightarrow \mathcal{P}(L(H))$, a reconciliation of (S, H, ϕ) is a function $\lambda : V(S) \rightarrow$*
 127 *$\mathcal{P}(V(H))$ satisfying the following:*

- 128 1. *For every leaf vertex $s \in L(S)$, we have $\lambda(s) = \phi(s)$.*
- 129 2. *For every internal vertex $s \in V(S) \setminus L(S)$ with children s_1, s_2 , such*
 130 *that $\lambda(s)$ is a singleton, exactly one of the following applies:*

131 (a) *$E(\lambda(s)) = \mathbb{S}$, that is, either $\lambda(s)$ and one element of $\lambda(s_1)$ are*
 132 *incomparable and $\lambda(s_2)$ contains a descendant of $\lambda(s)$, or $\lambda(s)$*
 133 *and one element of $\lambda(s_2)$ are incomparable and $\lambda(s_2)$ contains a*
 134 *descendant of $\lambda(s)$,*

135 (b) *$E(\lambda(s)) = \mathbb{C}$, that is, there is some $h_1 \in \lambda(s_1)$ (resp. $h_2 \in \lambda(s_2)$)*
 136 *such that $\text{mrca}(h_1, h_2) = \lambda(s)$, and h_1 and h_2 are incomparable,*

137 (c) *$E(\lambda(s)) = \mathbb{D}$, that is, there is some $h_1 \in \lambda(s_1)$ (resp. $h_2 \in \lambda(s_2)$)*
 138 *such that both h_1, h_2 are descendants of $\lambda(s)$, and the previous two*
 139 *cases do not apply.*

- 140 3. *For every internal vertex $s \in V(S) \setminus L(S)$ such that $\lambda(s)$ is not a*
 141 *singleton, exactly one of the following applies:*

- 142 (a) $E(\lambda(s)) = \mathbb{VS}$, that is $\lambda(s)$ is a clade in H , and all the descendants
 143 s' of s are also associated to the same clade, i.e. $\lambda(s') = \lambda(s)$.
- 144 (b) $E(\lambda(s)) = \mathbb{HS}$, that is $\lambda(s)$ is the union of two clades in H whose
 145 respective roots are incomparable. Moreover, all the descendants
 146 s' of s are also associated to the same clades, i.e. $\lambda(s') = \lambda(s)$.
- 147 (c) s is the descendant of a node s' where a spread (either vertical or
 148 horizontal) occurred (cases (3a) and (3b)). Then $\lambda(s) = \lambda(s')$. In
 149 that case, no additional coevolutionary event is recorded at that
 150 vertex.

151 The loss event denoted by \mathbb{L} is identified by a multiset (generalisation
 152 of a set where the elements are allowed to appear more than once) whose
 153 elements are in $V(H)$ containing all the vertices $h \in V(H)$ that are in the
 154 path between the image of a vertex $s \in V(S)$ which is a singleton and the
 155 image of one of its children. Note that no other event and thus no losses can
 156 happen below spread events.

157 Now, the function λ partitions the set of internal symbiont tree vertices
 158 into five disjoint subsets according to the coevolutionary event occurring at
 159 that vertex, plus an additional subset of all internal symbiont vertices that
 160 descend from a vertex where a spread occurred. The number of occurrences
 161 of each of the five events and the number of losses make up the *event vector* of
 162 the reconciliation. The *event vector* of a reconciliation is a vector of integers
 163 consisting of the total number of each type of events \mathbb{C} , \mathbb{D} , \mathbb{S} , \mathbb{L} , \mathbb{VS} , \mathbb{HS} .

164 Note that in the case of spread events (either vertical or horizontal) occurring
 165 at internal vertex $s \in V(S) \setminus L(S)$, the event is counted only once and the
 166 internal vertices s' descendants of s have no coevolutionary event associated
 167 to them.

168 The time feasibility condition is unchanged when adding spreads in the
 169 list of coevolutionary events.

170 A.4 Pre-estimating probabilities for the spread events

171 Given an input dataset (H, S, ϕ) , we rely on frequency estimators for the
 172 spread probabilities that will be used in our algorithm. Note that the “clas-
 173 sical events” (cospeciation, duplication, host switch and loss) have the same
 174 probability to occur everywhere in the tree, while the probability of a vertical
 175 or horizontal spread is specific to each vertex of the host tree. These prob-
 176 abilities are pre-estimated based on the input (H, S, ϕ) as described below
 177 rather than in the full ABC procedure. They are estimated through heuris-
 178 tic frequencies observed in the associations of the two trees. In Section D.3,
 179 we explore the robustness of our results with respect to these pre-computed
 180 estimators.

181 **Probability that a vertical spread occurs at host h .** A probability
 182 $p_{\text{vs}}(h)$ is associated to a vertical spread event at host h as follows. If $h \in$
 183 $L(H)$, then $p_{\text{vs}}(h)$ is estimated to 1. Otherwise, for any internal vertex h of

184 the host tree H , the probability $p_{\text{vs}}(h)$ is estimated to

$$p_{\text{vs}}(h) = \left(\frac{1}{|S^{L(h)}|} \right) \frac{\sum_{s \in S^{L(h)}} |\phi(s) \cap L(h)| - 1}{|L(h)| - 1} \quad (\text{A.1})$$

185 where $L(h)$ is the set of leaves in H_h (the subtree of H rooted in h), $S^{L(h)}$ is
 186 the set of leaves in the symbiont tree S that are associated with at least one
 187 leaf of H_h (formally $S^{L(h)} = \{s \in L(S) : \phi(s) \cap L(h) \neq \emptyset\}$), and $|\phi(s) \cap L(h)|$
 188 is the number of host leaves in H_h associated with a symbiont s .

189 Intuitively, the probability $p_{\text{vs}}(h)$ is large whenever a large proportion
 190 of the symbionts in $S^{L(h)}$ are associated to a large proportion of the hosts
 191 $L(h)$ (*i.e.* most of the symbionts are generalists) and is low when most of
 192 those symbionts are associated only with a few hosts of $L(h)$ (*i.e.* most of
 193 the symbionts are specialists). Notice that for a host h that is high in the
 194 tree, *i.e.* that is near to the root of H , the set $L(h)$ is large. Thus, a vertical
 195 spread to occur at h with high probability requires that some symbiont leaves
 196 are associated to an unrealistically large set of hosts $L(h)$. Hence usually the
 197 probability of a vertical spread is lower in hosts that are high in the tree. As
 198 explained in the next paragraph, the same holds for the horizontal spread
 199 event.

200 **Probability that a symbiont present in h invades an incomparable**
 201 **host h' .** For two incomparable vertices h and h' , a probability $p_{\text{jump}}(h \rightarrow h')$

202 is estimated as follows

$$p_{\text{jump}}(h \rightarrow h') = \frac{|S^{L(h)} \cap S^{L(h')}|}{|S^{L(h)} \cup S^{L(h')}|}. \quad (\text{A.2})$$

203 The notion of “jump” does not refer to a coevolutionary event and should
 204 not be confused with a host switch. The jump probability is specific to each
 205 pair of vertices of the host tree. It is a symmetric quantity, *i.e.* $p_{\text{jump}}(h \rightarrow$
 206 $h') = p_{\text{jump}}(h' \rightarrow h)$. It is high whenever the leaves of the subtrees H_h and
 207 $H_{h'}$ share a large proportion of associated symbionts. In particular, it is
 208 zero when they do not share any associated symbiont, and 1 when they have
 209 exactly the same set of associated symbionts.

210 **Probability that a horizontal spread occurs at host h .** From the
 211 probabilities $p_{\text{jump}}(h \rightarrow h')$, we estimate a probability of horizontal spread
 212 at each vertex h . The associated probability depends on all the vertices h'
 213 that are incomparable with h . Indeed, such vertices are all those that may
 214 be reached from h through a horizontal spread event. In fact, a horizontal
 215 spread corresponds to a jump combined with two vertical spreads. We thus
 216 associate a probability of horizontal spread $p_{\text{hs}}(h)$ to each vertex h of the
 217 host tree that takes into account both a jump and two vertical spreads and
 218 is set as

$$p_{\text{hs}}(h) = \min\{1, p^*(h)\}, \quad (\text{A.3})$$

219 where

$$p^*(h) = p_{\text{vs}}(h) \sum_{\substack{h' \in V(H) \\ h, h' \text{ incomparable}}} p_{\text{vs}}(h') p_{\text{jump}}(h \rightarrow h').$$

220 The probability of a horizontal spread $p_{\text{hs}}(h)$ is high whenever $p_{\text{vs}}(h)$ is high
 221 and there exist vertices h' incomparable to h with large $p_{\text{vs}}(h')$ and large
 222 value $p_{\text{jump}}(h \rightarrow h')$ (so that the leaves below h and h' share many symbionts).
 223 Observe that $p^*(h)$ is not a probability but a positive value, that in particular
 224 may be larger than 1.

Probability for sampling a horizontal spread to some specific host h' . In the simulation process, once a horizontal spread is sampled for symbiont s at vertex h , we need to choose an incomparable vertex h' where the symbiont s has to jump to. In this case, we need to guarantee that the jump satisfies the time-feasibility constraints as given in Stolzer *et al.* (2012) and Baudet *et al.* (2015). This constraint depends on the symbionts mapped so far (see Section *Simulation algorithm in AMOCOALA* below). For a current partial mapping λ from the vertices of S to the subsets of vertices of H , the probability $p_{\text{invasion}}(h \rightarrow h', \lambda)$ of a vertex h' to be invaded by a symbiont s mapped in h is estimated as

$$\begin{aligned} p_{\text{invasion}}(h \rightarrow h', \lambda) &= \frac{p_{\text{jump}}(h \rightarrow h') 1\{E_{h, h', \lambda}\} p_{\text{vs}}(h) p_{\text{vs}}(h')}{p_{\text{vs}}(h) \sum_{h''} p_{\text{vs}}(h'') p_{\text{jump}}(h \rightarrow h'') 1\{E_{h, h'', \lambda}\}}, \\ &= \frac{p_{\text{jump}}(h \rightarrow h') 1\{E_{h, h', \lambda}\} p_{\text{vs}}(h')}{\sum_{h''} p_{\text{vs}}(h'') p_{\text{jump}}(h \rightarrow h'') 1\{E_{h, h'', \lambda}\}}, \end{aligned} \quad (\text{A.4})$$

where $1\{E_{h,h',\lambda}\} = 1$ whenever the horizontal spread of the symbiont mapped in h to the new host h' induces a time feasible reconciliation, and the sum in the denominator is restricted to the vertices h'' that are incomparable to h . If no vertex induces a time feasible reconciliation (namely $p_{\text{invasion}}(h \rightarrow h', \lambda) = 0$ for any h' incomparable to h), the horizontal spread is not applied and another event is sampled. Otherwise, as the probabilities $p_{\text{invasion}}(h \rightarrow h', \lambda)$ sum up to one, a vertex h' is necessarily chosen.

Computing the pre-estimated spread probabilities. The estimated spread probabilities are calculated at the beginning of the algorithm. These values depend only on the host tree H , the symbiont tree S and the associations between the leaves ϕ . In a first step, we start by setting to 1 the probabilities p_{vs} for the leaves. Then, for the internal vertices h , these probabilities are computed as in Equation (A.1). In a second step, the probabilities of a jump are calculated for each pair of incomparable vertices h and h' as in Equation (A.2). In the last step, the probabilities of a horizontal spread for vertex h are computed as in Equation (A.3). Observe that the probabilities of invasion (Equation (A.4)) depend on the current simulation. Indeed, one has to take into account the time-feasibility in order to choose the target h' of a horizontal spread. Therefore, it may happen that the invasion $p_{\text{invasion}}(h \rightarrow h', \lambda) > 0$ for the current partial mapping λ but after some steps $p_{\text{invasion}}(h \rightarrow h', \lambda') = 0$ for the new mapping λ' . These probabilities are then updated, during the simulation algorithm, each time a horizontal

247 spread is selected.

248 **B AMOCOALA algorithm**

249 **B.1 Simulation algorithm in AMOCOALA**

250 The simulation of a symbiont tree \tilde{S} together with its reconciliation $\tilde{\lambda}$ starts
 251 with the creation of its root vertex \tilde{s}_{root} . This vertex is positioned before
 252 the root of H on the arc $a = (\rho, H_{root})$. We add the arc (ρ, H_{root}) to allow
 253 the simulation of events that happened in the symbiont tree before the most
 254 recent common ancestor of all host species in H . Figure ?? in main text
 255 depicts this starting configuration.

256 For any vertex \tilde{s} of \tilde{S} that is not yet mapped and whose position is
 257 $\langle \tilde{s} : a \rangle$ (see Figure ?? in main text), AMOCOALA successively considers the
 258 six allowed operations, and chooses one depending on the probability of each
 259 event (once an event is picked, the others are not considered). In what
 260 follows, we denote by a_1, a_2 the arcs outgoing from the head $h(a)$ of the arc
 261 a .

- 262 I. If $h(a)$ is a leaf, we *STOP* the evolution of \tilde{s} .
- 263 II. We first sample a horizontal spread according to the probability $p_{hs}(h(a))$.
 264 When a horizontal spread occurs (Figure ?? in main text), we apply the
 265 mapping $\tilde{\lambda}(\tilde{s}) = H_{h(a)} \cup H_{h(a')}$. The choice of the incomparable vertex
 266 $h(a')$ varies in order to preserve time feasibility (Stolzer *et al.*, 2012;

267 Baudet *et al.*, 2015), thus the probabilities described in Equation (A.4)
 268 are updated according to the new set of incomparable vertices. If there
 269 is no incomparable vertex, it is not possible for a horizontal spread to
 270 occur and we go to Step III. To select the ghost subtree rooted in \tilde{s} ,
 271 we mimic the real symbiont tree as shown in Figure ?? in main text.

272 III. If a horizontal spread did not occur, we sample a vertical spread ac-
 273 cording to the probability $p_{vs}(h(a))$. When a vertical spread occurs
 274 (Figure ?? in main text), we apply the mapping $\tilde{\lambda}(\tilde{s}) = H_{h(a)}$. To se-
 275 lect the ghost subtree rooted in \tilde{s} , we mimic the real symbiont tree as
 276 shown in Figure ?? from main text.

277 In both cases of vertical and horizontal spreads, the evolution of \tilde{s}
 278 stops after the creation of the ghost subtree and its descendants are
 279 not processed anymore.

280 IV. If a spread was not sampled, then we sample with a multinomial distri-
 281 bution a classical event according to the probabilities $\theta = \langle p_c, p_d, p_s, p_l \rangle$.
 282 Notice that $p_c + p_d + p_s + p_l = 1$ so that one of the four events is se-
 283 lected. This case is handled identically as in COALA and the symbiont
 284 is associated to a single host. We briefly recall the procedure below.

285 – Cospeciation (Figure ??(b) in main text): We apply the mapping
 286 $\tilde{\lambda}(\tilde{s}) = \{h(a)\}$ and we create the vertices \tilde{s}_1 and \tilde{s}_2 as children
 287 of \tilde{s} . We position them as follows: $\langle \tilde{s}_1 : a_1 \rangle$ and $\langle \tilde{s}_2 : a_2 \rangle$. This
 288 operation is executed with probability p_c .

- 289 – Duplication (Figure ??(c) in main text): We apply the mapping
 290 $\tilde{\lambda}(\tilde{s}) = \{h(a)\}$ and we create the vertices \tilde{s}_1 and \tilde{s}_2 as children of
 291 \tilde{s} . Both \tilde{s}_1 and \tilde{s}_2 are positioned on a . This operation is executed
 292 with probability p_d .
- 293 – Host switch (Figure ??(e) in main text): We apply the mapping
 294 $\tilde{\lambda}(\tilde{s}) = \{h(a)\}$ and we create the vertices \tilde{s}_1 and \tilde{s}_2 as children of \tilde{s} .
 295 We then randomly choose one of the two children and position it
 296 on a . Finally, we randomly choose an arc a' that does not violate
 297 the time feasibility of the reconstruction so far (Stolzer *et al.*, 2012;
 298 Baudet *et al.*, 2015). If such an arc does not exist, it is not possible
 299 for a host switch to take place. In this case, we choose between
 300 the three remaining events with probability $p_i/(p_c + p_d + p_l)$ with
 301 $i \in \{c, d, l\}$. Otherwise, we position \tilde{s}_2 on a' . This operation is
 302 executed with probability p_s .
- 303 – Loss (Figure ??(e) in main text): This operation consists of ran-
 304 domly choosing an arc outgoing from the head $h(a)$ of a and po-
 305 sitioning \tilde{s} on it. This operation is executed with probability p_l .

306 In any of these four cases, the simulation process recursively continues
 307 with the new vertices created (back to Step I).

308 Note that in our modelling, losses never occur after a spread event. In-
 309 deed, in the case of a vertical spread, a symbiont and its entire clade are
 310 associated to one host clade, while in the case of a horizontal spread, they

are then associated to two host clades. This might appear unrealistic. However, this choice is made for computational reasons. Indeed, as mentioned in the Main Manuscript, there is no simple way of simulating the symbiont tree below a symbiont where a spread occurs.

B.2 ABC-SMC inference method in AMOCOALA

AMOCOALA is based on the same ABC-SMC method as the one developed in COALA (Baudet *et al.*, 2015). For the sake of completeness, we now recall the procedure.

The ABC-SMC procedure is composed of a sequence of $R > 1$ rounds. At each round, parameter vectors θ are sampled in a specific way, symbiont trees \tilde{S}_θ are generated under the reconciliation model allowing for spreads with parameter values given by θ (and relying on the simulation algorithm described in the previous section). Then, these symbiont trees are compared to the original dataset through a summary distance d whose details are given in the next section. The parameters with the smallest discrepancies are selected.

For each of these rounds, we define a tolerance value τ_r ($1 \leq r \leq R$) which determines the percentage of parameter vectors to be accepted. Associated with a tolerance value τ_r , we have a threshold ϵ_r which is the largest value of the summary distance associated with the accepted parameter vectors.

- Initial round ($r = 1$):

-
- 332 – Draw an initial set of N parameter vectors $\{\theta_1^i\}_{(1 \leq i \leq N)}$ from the
 333 prior π .
- 334 – Then, for each θ_1^i , simulate M trees $\{\tilde{S}_j(\theta_1^i)\}_{(1 \leq j \leq M)}$. Compute the
 335 corresponding discrepancies $\{d_j(\theta_1^i)\}_{(1 \leq j \leq M)}$ and summarise them
 336 into the summary discrepancy $d_{\theta_1^i}$ through the mean value.
- 337 – Select $Q_1 = \tau_1 \times N$ parameter vectors θ_1 that have the smallest
 338 value d_{θ_1} , thus defining the threshold ϵ_1 and the set A_1 of accepted
 339 parameter vectors.
- 340 • Following rounds ($2 \leq r \leq R$):
- 341 1. Sample a parameter vector θ^* from the set $A_{(r-1)}$.
- 342 2. Create a parameter vector θ^{**} by perturbing θ^* (through a kernel
 343 proposal).
- 344 3. Simulate M trees relying on the parameter value θ^{**} and compute
 345 $d_{\theta^{**}}$. If $d_{\theta^{**}} \leq \epsilon_{(r-1)}$, add θ^{**} into the quantile set \mathcal{Q}_r . If $|\mathcal{Q}_r| <$
 346 Q_{r-1} , return to Step 1.
- 347 4. Based on the set \mathcal{Q}_r , select $Q_r = \tau_r \times Q$ parameter vectors θ_r that
 348 have the smallest d_{θ_r} , thus defining the threshold ϵ_r and the set
 349 A_r of accepted parameters.

350 **Prior distribution.** We sample from a uniform distribution on the simplex
 351 $\mathcal{S}_3 = \{(p_1, p_2, p_3, p_4); p_i \geq 0 \text{ and } \sum_i p_i = 1\}$ (we recall that $p_c + p_d + p_s + p_l =$
 352 1).

353 **Kernel proposal.** We add to each coordinate of θ a randomly chosen value
 354 in $[-0.01, +0.01]$ and normalise the result. The final set of accepted parame-
 355 ter vectors is the result of the ABC-SMC procedure and characterises the list
 356 of vectors that may explain the evolution of the pair of host and symbiont
 357 trees given as input. Observe that, since in all experiments a uniform prior
 358 distribution is assumed and also the perturbations are performed in a uni-
 359 form way, the weights induced by the proposals will also appear to be uniform
 360 (Beaumont *et al.*, 2009). However, in the case of a different prior, weights
 361 should be used in the process in order to correct the posterior distribution
 362 according to the perturbation made.

363 **Clustering of the vectors.** The final list of accepted vectors are clustered
 364 using a hierarchical clustering procedure implemented in COALA (Baudet
 365 *et al.*, 2015). As final result, we therefore obtain a list of clusters to each one
 366 of which a representative vector is associated.

367 B.3 Distance measure in AMOCOALA

368 The discrepancy between the simulated and the original datasets is measured
 369 through a distance between set-labelled phylogenetic trees which can be cal-
 370 culated in polynomial time. Similarly as in COALA, this distance contains
 371 two components: (i) d_1 , that describes how much the simulated tree \tilde{S}_θ is
 372 representative of the vector θ , and (ii) d_2 that measures how much is \tilde{S}_θ (and
 373 its labels) topologically similar to S (and its labels).

Let us recall the definition of this first component. For a given vector $\theta = \langle p_c, p_d, p_s, p_l \rangle$ and for each simulated tree \tilde{S}_θ that was simulated according to this vector, we keep track of the vector of the number of classical cophylogeny events $\langle o_c, o_d, o_s, o_l \rangle$ associated to this simulation. We compute the corresponding expected vector $\langle e_c, e_d, e_s, e_l \rangle$ as follows

$$\forall \text{event} \in \{c, d, s, l\}, \quad e_{\text{event}} = |S| \times \theta_{\text{event}} = |S| \times p_{\text{event}},$$

where $|S|$ is the size of the symbiont tree, *i.e.* its number of internal leaves. Then by comparing the observed and expected vectors, we define a measure $d_1(S, \tilde{S}_\theta)$ as follows:

$$d_1(S, \tilde{S}_\theta) = \frac{1}{4} \times \sum_{\text{event} \in \{c, d, s, l\}} \frac{|e_{\text{event}} - o_{\text{event}}|}{\max\{e_{\text{event}}, o_{\text{event}}\}}.$$

Note that we did not consider the number of observed spread events, which does not depend on the choice of θ as the corresponding probabilities are pre-estimated before applying the ABC-SMC approach.

As concerns point (ii), we extend the well-known *maximum agreement subtree* (MAST) distance (Finden and Gordon, 1985; Farach-Colton *et al.*, 1995) to handle set-labelled trees. This part is the novelty with respect to the proposal in COALA and details were given in the Main Manuscript. We establish in the next sections that d_{MAST} is a distance and that it can be computed in polynomial time.

391 We use a normalised version of d_{MASSST} and define the distance d_2 (see
 392 Main Manuscript). The two components are then combined to form the
 393 following distance

$$d_\theta = \alpha_1 d_1(S, \tilde{S}_\theta) + \alpha_2 d_2(S, \tilde{S}_\theta).$$

394 According to our experiments and also the ones presented in COALA, the
 395 most appropriate values are $\alpha_1 = 0.7$ and $\alpha_2 = 0.3$.

396 B.4 A proof that dMASST is a distance

397 We show that the distance d_{MASSST} is a metric. For this, we check that
 398 d_{MASSST} satisfies the following properties:

- 399 1. $d_{MASSST}(T_1, T_2) \geq 0$ for all T_1, T_2 : this is trivial.
- 400 2. $d_{MASSST}(T_1, T_2) = 0$ if and only if $T_1 = T_2$. Clearly if $T_1 = T_2$
 401 then $d_{MASSST}(T_1, T_2) = 0$. Otherwise, let $d_{MASSST}(T_1, T_2) = 0$. Then
 402 $\max\{w(T_1), w(T_2)\} = MASSST(T_1, T_2)$. The proof follows by observ-
 403 ing that if T^* is a subtree of T such that $w(T^*) = w(T)$ then $T^* = T$.
- 404 3. $d_{MASSST}(T_1, T_2) = d_{MASSST}(T_2, T_1)$: this is trivial.
4. For any triplet of trees T_1, T_2, T_3 , it holds that $d_{MASSST}(T_1, T_2) +$
 $d_{MASSST}(T_2, T_3) \geq d_{MASSST}(T_1, T_3)$. For simplicity, we set $w_i = w(T_i)$
 and $w_{i,j} = w(MASSST(T_i, T_j))$. Hence $d_{MASSST}(T_i, T_j) = \max\{w_i, w_j\} -$
 $w_{i,j}$. Furthermore, we denote by $w_{1,2,3}$ the weight of the maximum

agreement subtree that is common to the three trees T_1, T_2, T_3 . We then have:

$$\begin{aligned}
& d_{MASST}(T_1, T_2) + d_{MASST}(T_2, T_3) \\
&= \max\{w_1, w_2\} - w_{1,2} + \max\{w_2, w_3\} - w_{2,3} \\
&= \max\{w_1, w_2\} + \max\{w_2, w_3\} - (w_{1,2} + w_{2,3} - w_{1,2,3} + w_{1,2,3}) \\
&\geq \max\{w_1, w_2, w_3\} + w_2 - (w_2 + w_{1,2,3}) \\
&\geq \max\{w_1, w_3\} - w_{1,3},
\end{aligned}$$

405 where for the first inequality, we use the fact that $\max\{w_1, w_2\} +$
406 $\max\{w_2, w_3\} \geq \max\{w_1, w_2, w_3\} + w_2$ and we show in the next Lemma
407 that $w_{1,2} + w_{2,3} - w_{1,2,3}$ is at most w_2 . The last inequality uses $w_{1,2,3} \leq$
408 $w_{1,3}$.

409 This concludes the proof.

410 **Lemma.** *For any three set-labelled trees T_1, T_2, T_3 (using the notation from*
411 *the above proof) it holds that $w_{1,2} + w_{2,3} - w_{1,2,3} \leq w_2$.*

Proof. Let $T_{1,2}$ and $T_{2,3}$ be maximum agreement set-labelled subtrees (MASST) of T_1, T_2 and T_2, T_3 , respectively. Consider any pair of leaf, label that belongs to T_2 , i.e. $(l, lab) \in T_2$. There are only four possibilities: (i) $(l, lab) \in T_{1,2}$ and $(l, lab) \notin T_{2,3}$ (we call these leaves of type A), (ii) $(l, lab) \notin T_{1,2}$ and $(l, lab) \in T_{2,3}$ (we call these leaves of type B), (iii) $(l, lab) \in T_{1,2}$ and $(l, lab) \in T_{2,3}$ (we call these leaves of type C), (iv) $(l, lab) \notin T_{1,2}$ and

$(l, lab) \notin T_{2,3}$ (we call these leaves of type D). Then we have

$$\begin{aligned} w_2 &= |A| + |B| + |C| + |D| \\ &= w_{12} - |C| + w_{23} - |C| + |C| + |D| \\ &= w_{12} + w_{23} - |C| + |D|. \end{aligned}$$

412 Or equivalently

$$w_{12} + w_{23} = w_2 + |C| - |D|. \quad (\text{B.1})$$

Moreover, we define the tree \tilde{T} as the subtree obtained from T_2 by taking all the pairs of leaf, label that belong to T_{12} and T_{23} . Notice that \tilde{T} is also a subtree of T_1 and of T_3 . Thus, \tilde{T} is included in T_{123} . This implies that $|C| \leq w_{123}$. Going back to (B.1), we thus obtain

$$\begin{aligned} w_{12} + w_{23} &= w_2 + |C| - |D| \\ &\leq w_2 + |C| \\ &\leq w_2 + w_{123}. \end{aligned}$$

413 This concludes the proof of the lemma. □

414 **Remark.** *The previous proof and comments show that the MASST distance*
 415 *d_{MASST} is very similar to the MAAC one (Ganapathy et al., 2005) for multi-*
 416 *labelled trees. Thus, it is natural to ask whether comparing two set-labelled*
 417 *trees can be reduced to comparing two multi-labelled trees. One idea is to*

transform a set-labelled tree into a multi-labelled tree. However, the straight-forward transformation seems not to work well for our purpose. For instance, we can transform each set-labelled tree into a multi-labelled tree by substituting each set-labelled leaf by a subtree with a fixed topology (say a complete binary tree, or a multifurcating vertex) as in Figure A. However, in these cases the two trees in Figure A would be considered equivalent, but in our context they are different. In fact, the set-labelled tree in Figure A(a) indicates that there is a symbiont that infects 4 different hosts h_1, h_2, h_3, h_4 , while in Figure A(b), we will have 4 different symbionts infecting each a different host.

B.5 Polynomial time algorithm for computing the dMASST distance

We show that it is possible to calculate the distance $d_{MASST}(T_1, T_2)$ in polynomial time with respect to the size of the trees. This boils down to computing the weight of the maximum agreement subtree $w(MASST(T_1, T_2))$ in polynomial time. The algorithm is based on dynamic programming and extends quite straightforwardly the algorithm for calculating the MAAC distance (Ganapathy *et al.*, 2005). We abbreviate to $w(v_1, v_2)$ the weight of the maximum agreement subtree between the two trees T_1 and T_2 rooted in v_1 and v_2 , respectively. For a leaf v , we denote by $l(v)$ the set of labels associated with it. Finally, for an internal vertex v , we denote by $ch_1(v)$ and

439 $ch_2(v)$ the two children of v .

440 The dynamic programming algorithm starts from the leaves and ends in
 441 the roots of T_1 and T_2 following a recursion. We have that $w(v_1, v_2)$ is given
 442 by:

- 443 • If v_1 and v_2 are both leaves then $w(v_1, v_2) = |l(v_1) \cup l(v_2)|$
- 444 • If v_1 or v_2 (could be both) are internal vertices, $w(v_1, v_2)$ is the maxi-
 445 mum value among the following three quantities
 - 446 1. $\max\{w(ch_1(v_1), v_2), w(ch_2(v_1), v_2)\}$;
 - 447 2. $\max\{w(v_1, ch_1(v_2)), w(v_1, ch_2(v_2))\}$;
 - 448 3. $\max\{w(ch_1(v_1), ch_1(v_2)) + w(ch_2(v_1), ch_2(v_2)), w(ch_1(v_1), ch_2(v_2))$
 449 $+ w(ch_2(v_1), ch_1(v_2))\}$.

450 C Additional results for the self-test

451 The results for parameter values θ_2^* to θ_8^* are presented in Figures Ba to Cd.

452 D Biological datasets

453 We provide here a description of the 4 datasets used. The corresponding
 454 phylogenetic trees are shown in Figures D - G.

455 *Dataset 1: AP - Acacia & Pseudomyrmex.* This dataset was extracted
 456 from Gómez-Acevedo *et al.* (2010) and displays the interaction between *Aca-*

457 *cia* plants and *Pseudomyrmex* species of ants. The host and symbiont trees
 458 include 9 and 7 leaves, respectively. The dataset has 22 multiple-associations.

459 *Dataset 2: MP - Myrmica & Phengaris.* This dataset was extracted from
 460 Jansen *et al.* (2011) and is composed of a pair of host and symbiont trees
 461 which have each 8 leaves. The dataset has 8 multiple-associations.

462 *Dataset 3: SBL - Seabirds & Lice.* This dataset was extracted from
 463 Paterson *et al.* (1997). The host and symbiont trees include 15 and 8 leaves,
 464 respectively. The dataset has 15 multiple-associations.

465 *Dataset 4: SFC - Smut Fungi & Caryophyllaceus plants.* This dataset was
 466 extracted from Refrégier *et al.* (2008). The host and symbiont trees include
 467 15 and 16 leaves, respectively. The dataset has 4 multiple-associations.

468 **D.1 Results on biological datasets**

469 We ran AMOCOALA on all the real datasets and plotted in Figures H to S the
 470 histograms of the summary discrepancies and event probabilities (except for
 471 the spread probabilities which are not inferred) obtained at the end of each
 472 one of the 3 rounds, for each of the 4 datasets. We see on the histograms
 473 that the summary discrepancies for the accepted parameter vectors decrease
 474 after each round. We recall that the summary discrepancy measures the
 475 similarity between the simulated trees and the original symbiont tree, and
 476 hence is related to the quality of the vectors. Thus, our result shows that
 477 the set of accepted vectors is refined at each round, leading to vectors which
 478 can generate trees that are increasingly more similar to the original symbiont

479 tree (and its host associations).

480 **D.2 Running times**

481 Table D.1 shows the running times obtained on the 4 biological datasets,
 482 together with their sizes (as expressed by the number of leaves in the host and
 483 symbiont trees) and the number of multiple associations. The results have
 484 been obtained on a computer with a AMD EPYC 7542 32-Core processor and
 485 128 CPU (2 sockets of 32 double threads cores) and 675Gb RAM. We used
 486 just one core ('nthreads 1', though AMOCOALA has a parallelized version)
 487 and AMOCOALA was run with default values on these datasets.

Dataset	(Host,Symbiont) leaves	Multiple associations	Running time
AP	(9,7)	22	23m20.859s
MP	(8,8)	8	21m25.631s
SBL	(15,8)	15	28m53.597s
SFC	(15,16)	4	117m45.919s

Table D.1: For each of the 4 biological datasets, we indicate the pairs of numbers of host and symbiont trees leaves (2nd column), the number of multiple associations (3rd column) and the running time of AMOCOALA on this dataset (4th column).

488 **D.3 Robustness analysis wrt the pre-estimated spread** 489 **probabilities**

490 In this section, we explore the robustness of our results with respect to the
 491 pre-estimated values of the spread events probabilities. On each of the 4 bio-

492 logical datasets, we ran AMOCOALA with perturbed values of $p_{\text{hs}}(h)$, $p_{\text{vs}}(h)$.
 493 More precisely, to each non zero probability $p_{\text{hs}}(h)$ or $p_{\text{vs}}(h)$, we added a noise
 494 value uniformly drawn in $[-0.1; 0.1]$ (and then took the infimum with 1 and
 495 the supremum with 0, in order to ensure the modified probabilities remain in
 496 $[0, 1]$). With these perturbed values, we ran AMOCOALA and output (after 3
 497 rounds) 50 accepted vectors $\theta = \langle p_c, p_d, p_s, p_l \rangle$. The results are presented in
 498 Figures T to W. Let us recall that AMOCOALA is a stochastic algorithm and
 499 any two runs will give similar but not identical results. The results obtained
 500 adding these perturbations are qualitatively the same for the first 3 datasets
 501 (namely AP, MP and SBL) as the ones without perturbations (see Figures J
 502 to P). The results for dataset SFC show more variability wrt those of the
 503 unperturbed version (Figure S). Thus we also looked at the clusters output
 504 by AMOCOALA in this case in Table D.2. We recall that in Refrégier *et al.*
 505 (2008), the different analyses performed indicated that the most plausible
 506 reconciliations presented for the SFC dataset have from 0 to 3 cospeciations,
 507 no duplication, 12 to 15 host switches and 0 to 2 losses. Here we find that
 508 the first main cluster (31 vectors out of 50) has a representative vector with
 509 around 50% of cospeciations (about 7 or 8 events), almost no duplication
 510 (about 0 or 1 event), 31% of host switches (about 4 or 5 events) and 18% of
 511 losses (about 2 or 3 losses). The second main cluster has a higher probability
 512 of cospeciation and less switches. Only the third cluster could correspond to
 513 Refrégier *et al.* (2008)'s scenario, with 1 or 2 cospeciations, no duplication, 8
 514 or 9 host switches and 4 to 5 losses; though it is supported by only 3 selected

515 vectors out of 50. Thus for the SFC dataset, the detection of the biological
516 scenario presented in Refrégier *et al.* (2008) is more difficult to detect with
517 perturbed values of the spread probabilities. To conclude, our results are
518 overall robust with respect to potential errors in the estimation of the spread
519 events probabilities.

520 References

- 521 Baudet, C., Donati, B., Sinaimeri, B., Crescenzi, P., Gautier, C., Matias, C.,
522 and Sagot, M.-F. 2015. Cophylogeny reconstruction via an Approximate
523 Bayesian Computation. *Systematic Biology*, 64(3): 416–31.
- 524 Beaumont, M. A., Cornuet, J.-M., Marin, J.-M., and Robert, C. P. 2009.
525 Adaptive approximate Bayesian computation. *Biometrika*, 96: 983–990.
- 526 Charleston, M. A. 2002. *Biological Evolution and Statistical Physics*, volume
527 585 of *Lecture Notes in Physics*, chapter Principles of cophylogenetic maps,
528 pages 122–147. Springer Berlin Heidelberg.
- 529 Farach-Colton, M., Przytycka, T. M., and Thorup, M. 1995. On the agree-
530 ment of many trees. *Inform. Process. Lett.*, 55: 297–301.
- 531 Finden, C. R. and Gordon, A. D. 1985. Obtaining common pruned trees. *J.*
532 *Classif.*, 2: 255–276.
- 533 Ganapathy, G., Goodson, B., Jansen, R., Ramachandran, V., and Warnow,

- 534 T. 2005. Pattern Identification in Biogeography. In R. Casadio and G. Myers,
535 editors, *Algorithms in Bioinformatics*, volume 3692 of *Lecture Notes*
536 *in Computer Science*, pages 116–127. Springer Berlin Heidelberg.
- 537 Gómez-Acevedo, S., Rico-Arce, L., Delgado-Salinas, A., Magallón, S., and
538 Eguiarte, L. E. 2010. Neotropical mutualism between Acacia and Pseu-
539 domyrmex: Phylogeny and divergence times. *Molecular Phylogenetics and*
540 *Evolution*, 56(1): 393–408.
- 541 Jansen, G., Vepsäläinen, K., and Savolainen, R. 2011. A phylogenetic test of
542 the parasite-host associations between *Maculinea* butterflies (Lepidoptera:
543 Lycaenidae) and *Myrmica* ants (Hymenoptera: Formicidae). *European*
544 *Journal of Entomology*, 108(1): 53–62.
- 545 Paterson, A., Gray, R. D., Clayton, D. H., and Moore, J. 1997. Host-parasite
546 co-speciation, host switching, and missing the boat. In D. H. Clayton and
547 J. Moore, editors, *Host-parasite evolution: General principles and avian*
548 *models*, pages 236–250. Oxford University Press.
- 549 Refrégier, G., Le Gac, M., Jabbour, F., Widmer, A., Shykoff, J. A., Yock-
550 teng, R., Hood, M. E., and Giraud, T. 2008. Cophylogeny of the anther
551 smut fungi and their Caryophyllaceae hosts: Prevalence of host shifts and
552 importance of delimiting parasite species for inferring cospeciation. *BMC*
553 *Evolutionary Biology*, 8(1): 100.
- 554 Stolzer, M. L., Lai, H., Xu, M., Sathaye, D., Vernot, B., and Durand, D. 2012.

- 555 Inferring duplications, losses, transfers and incomplete lineage sorting with
556 nonbinary species trees. *Bioinformatics*, 28(18): i409–i415.
- 557 Tofigh, A., Hallett, M. T., and Lagergren, J. 2011. Simultaneous identifica-
558 tion of duplications and lateral gene transfers. *IEEE/ACM Trans. Comput.*
559 *Biology Bioinform.*, 8(2): 517–535.

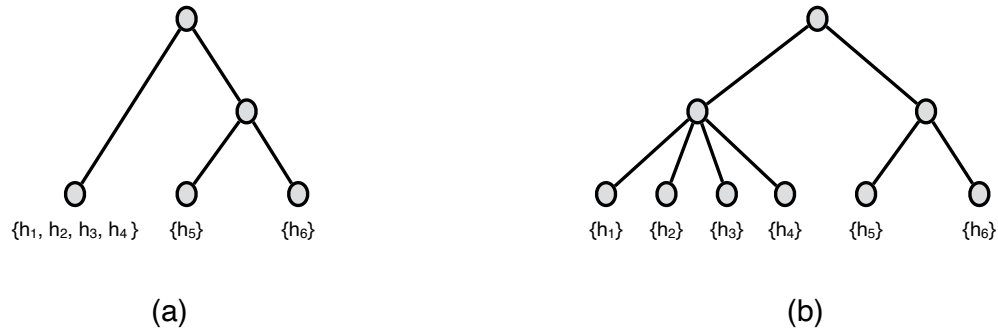


Figure A: The two phylogenetic trees will be considered at distance 0 if we substitute the vertex labelled by the set h_1, h_2, h_3, h_4 by a multifurcated vertex.

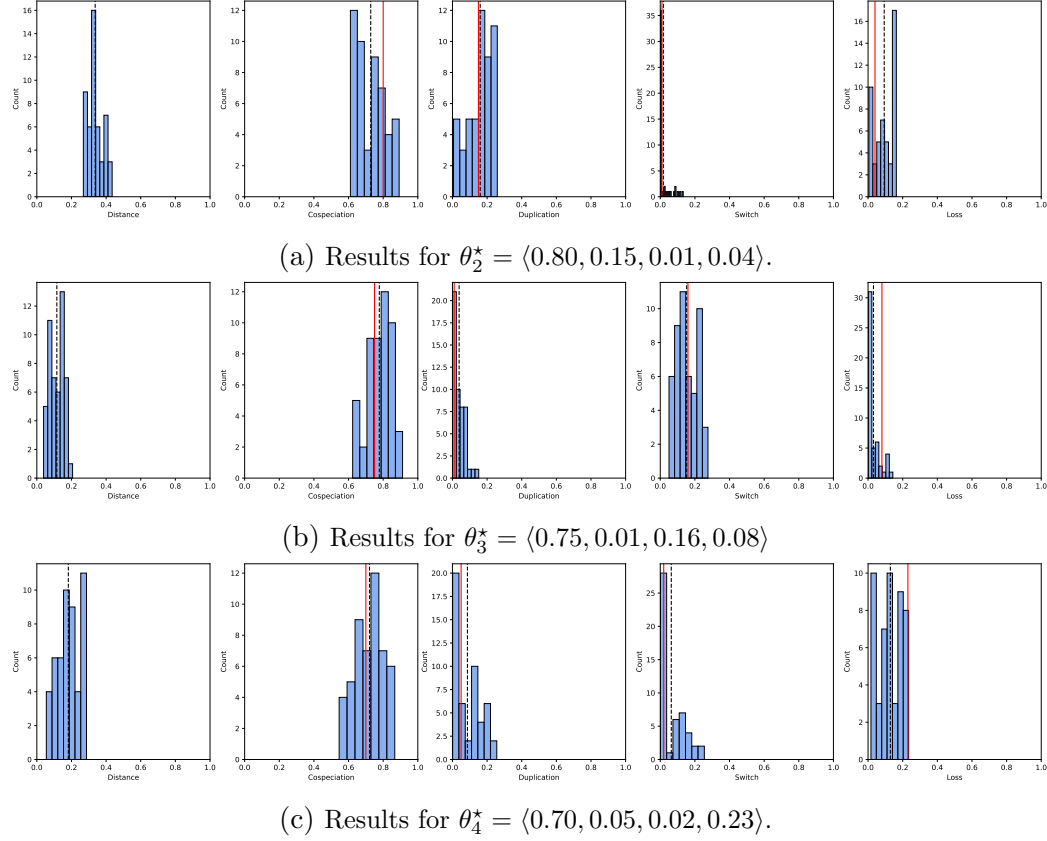


Figure B: For each simulated dataset with true parameter value θ_i^* and $2 \leq i \leq 8$, we ran AMOCOALA 50 times and, at the end of the third round, we took note of the cluster whose representative parameter vector had the smallest euclidean distance (histograms shown in the first column) to θ_i^* . Columns 2 to 5 show the histograms of the distributions of the event probabilities in these "best" clusters. The dashed vertical black line indicates the mean value. The solid vertical red line indicates the true parameter value.

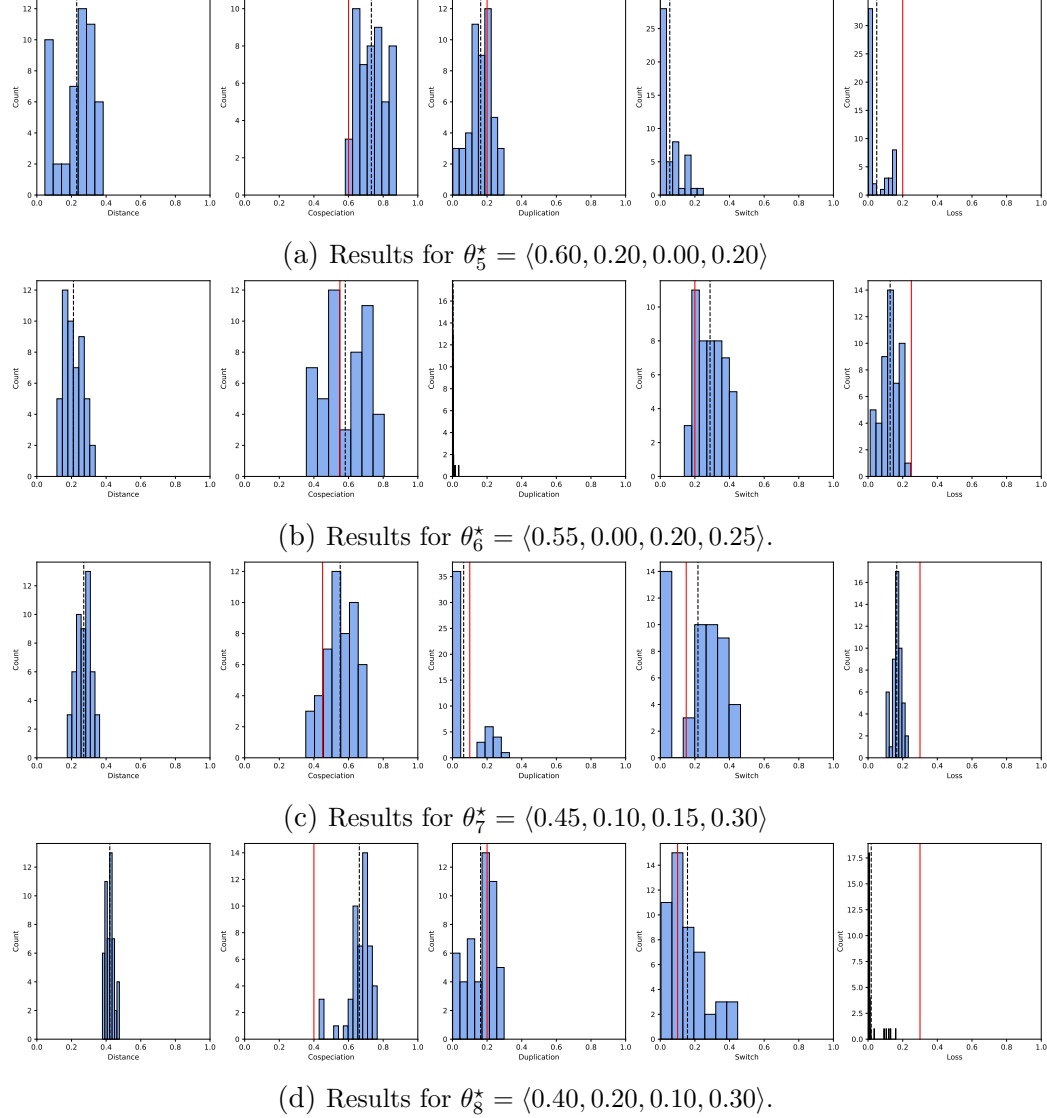


Figure C: For each simulated dataset with true parameter value θ_i^* , we ran AMOCOALA 50 times and, at the end of the third round, we took note of the cluster whose representative parameter vector had the smallest euclidean distance (histograms shown in the first column) to θ_i^* . Columns 2 to 5 show the histograms of the distributions of the event probabilities in these "best" clusters. The dashed vertical black line indicates the mean value. The solid vertical red line indicates the true parameter value.

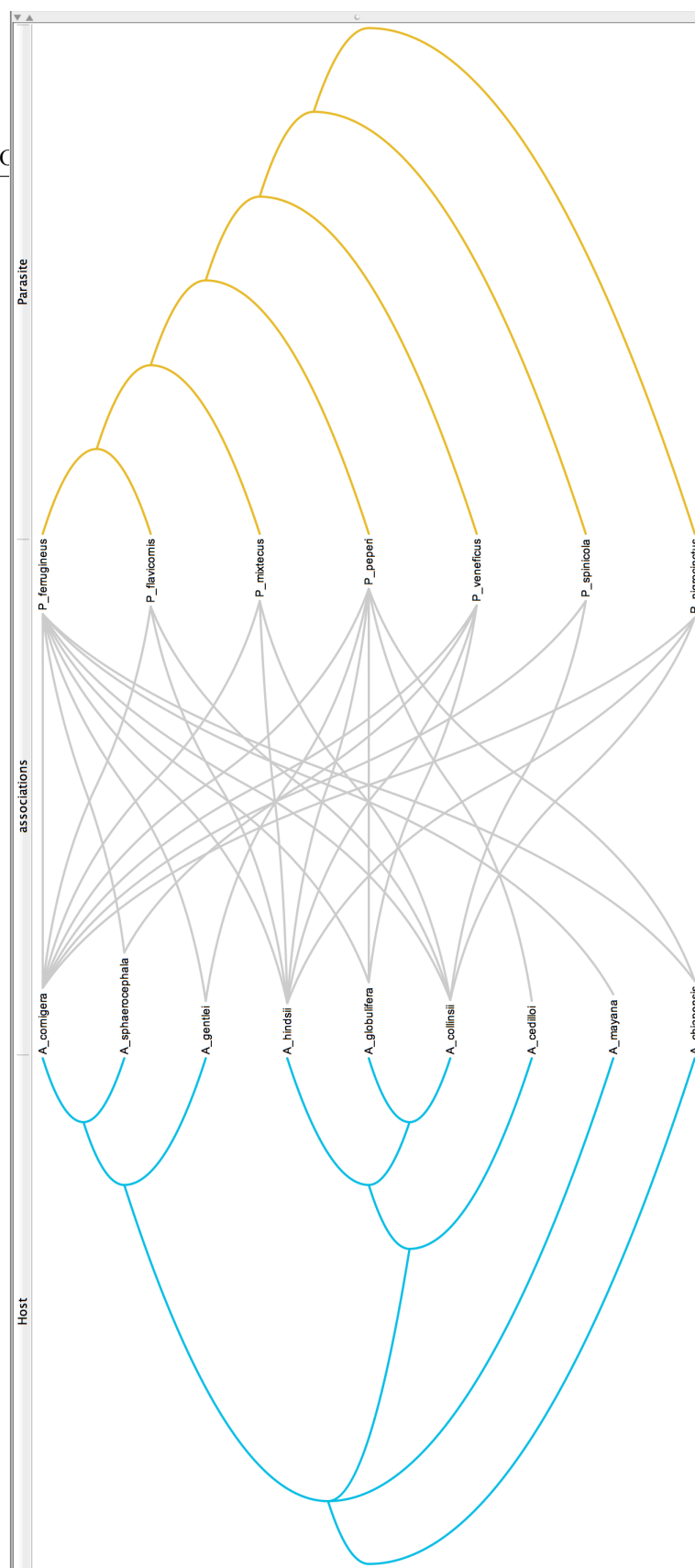


Figure D: AP dataset.

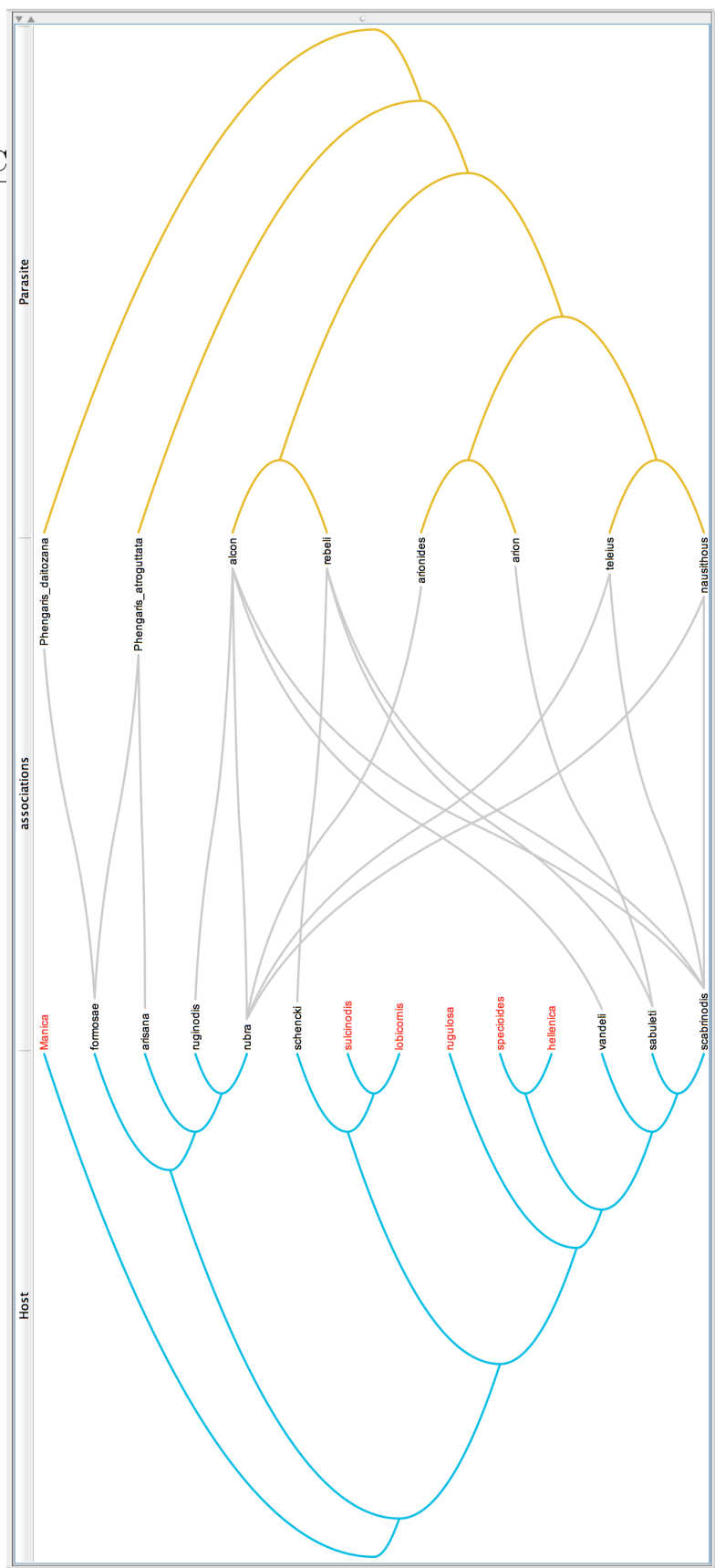


Figure E: MP dataset.

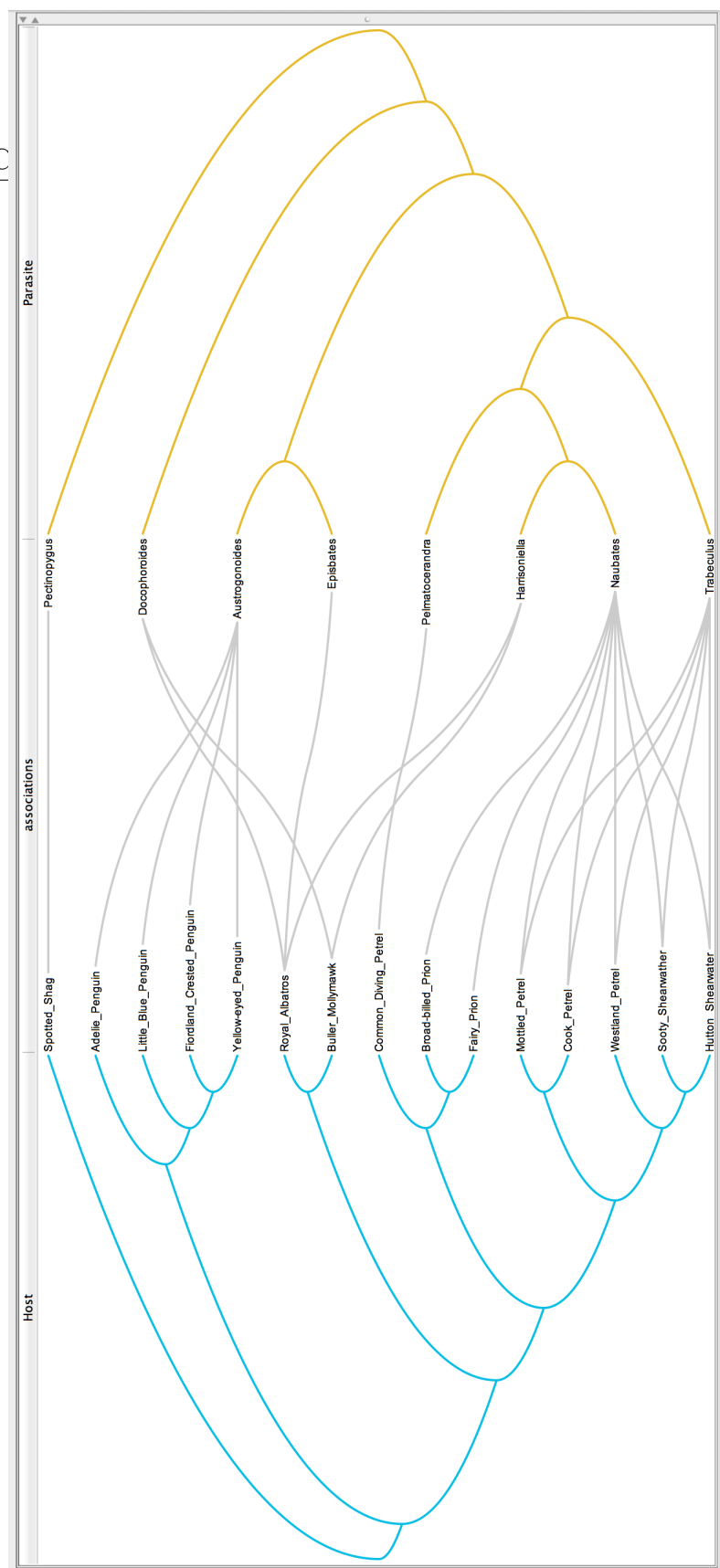


Figure F: SBL dataset.

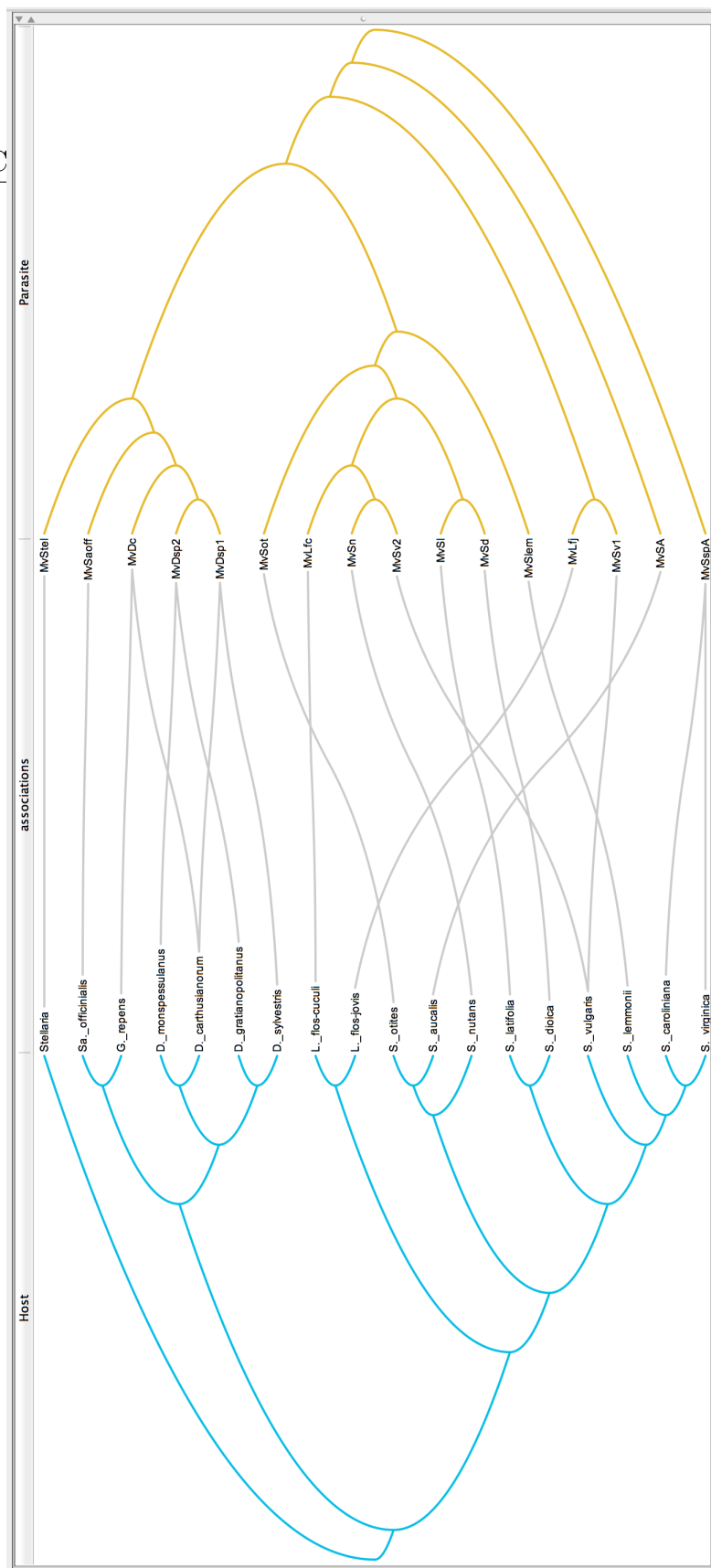


Figure G: SFC dataset.

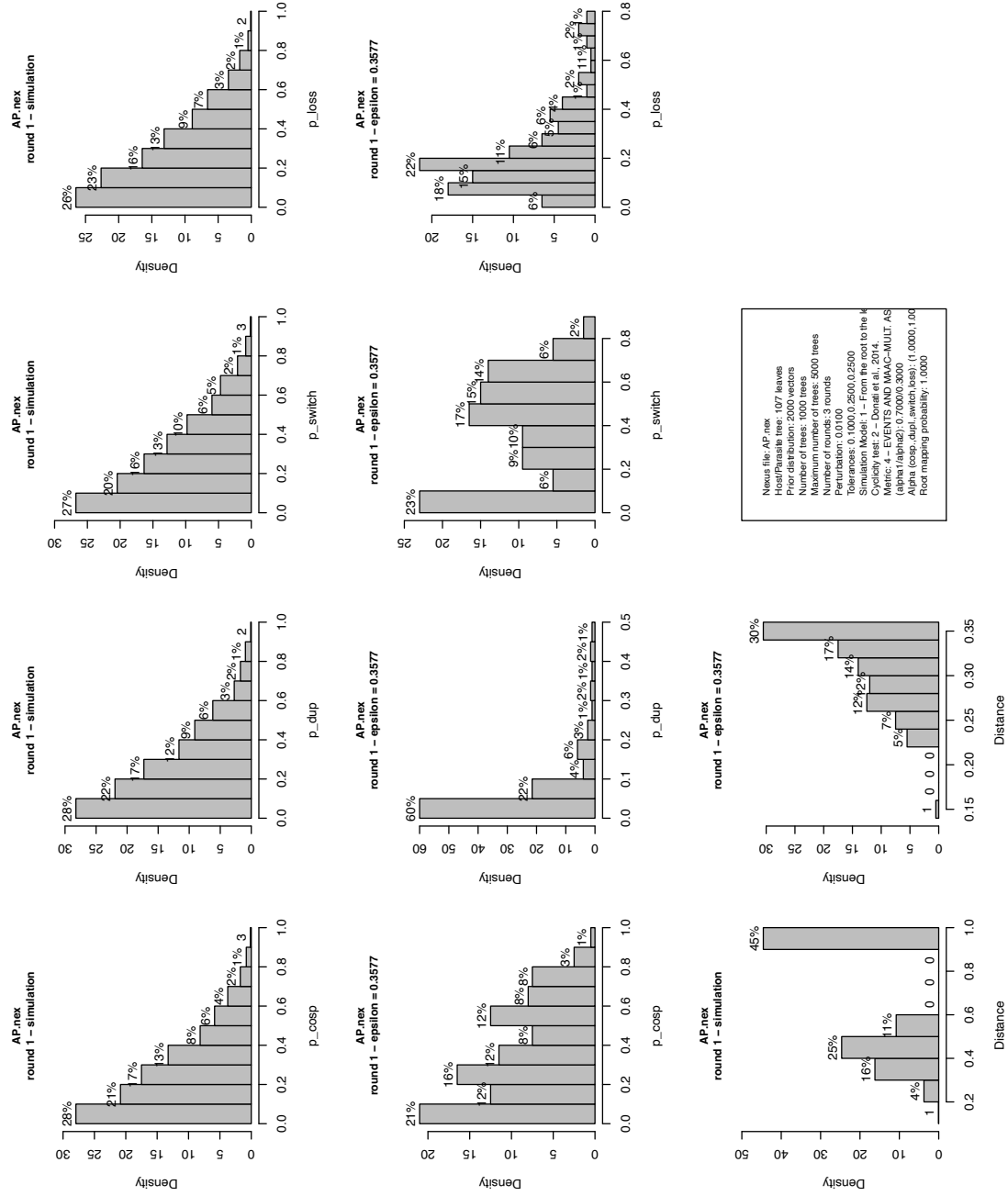


Figure H: AP dataset. First row: histograms of the input parameters. Second row: histograms of the parameters after round 1. Third row: summary discrepancies of the input parameters and of the parameters after round 1.

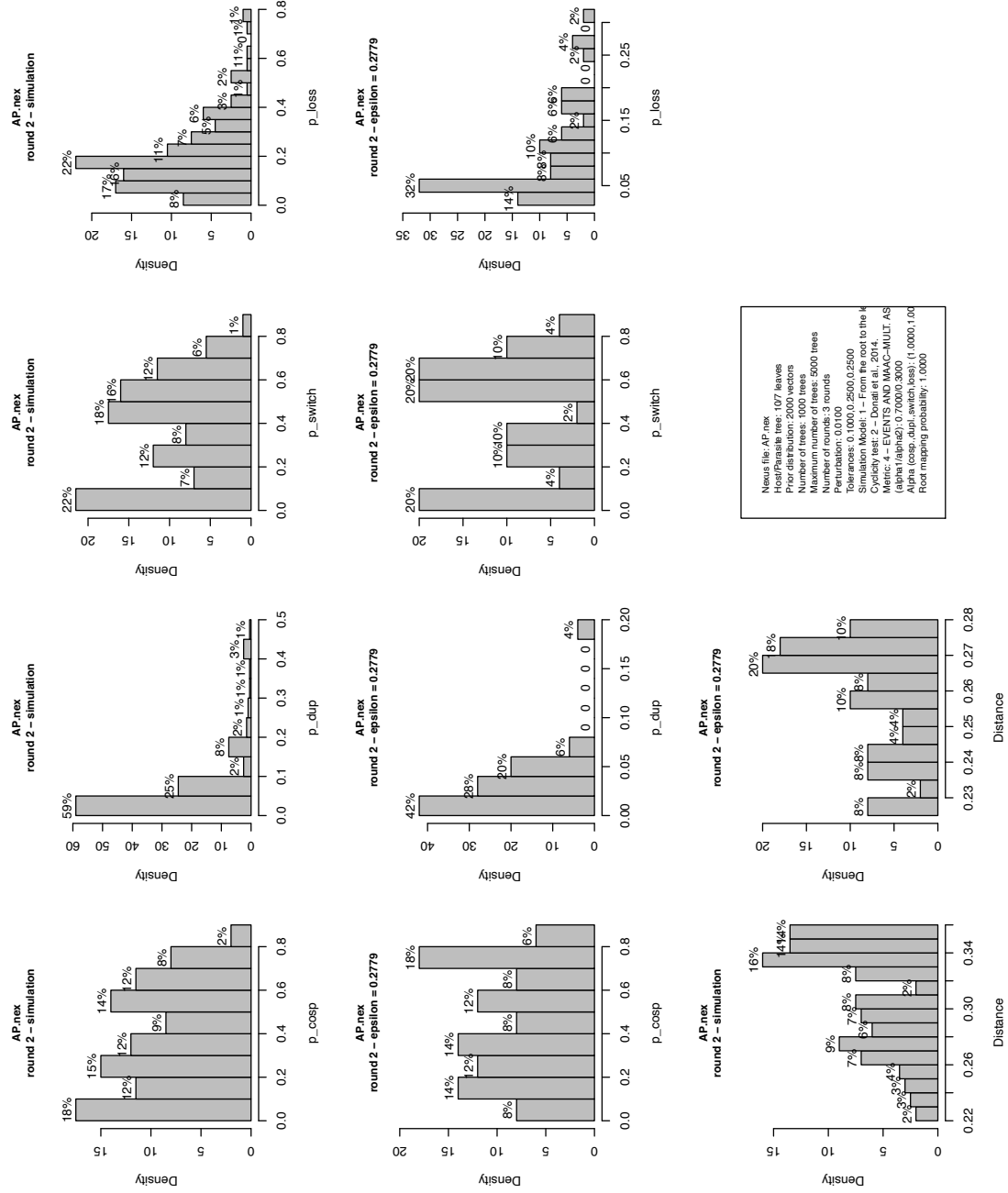


Figure I: AP dataset. First row: histograms of the input parameters. Second row: histograms of the parameters after round 2. Third row: summary discrepancies of the input parameters and of the parameters after round 2.

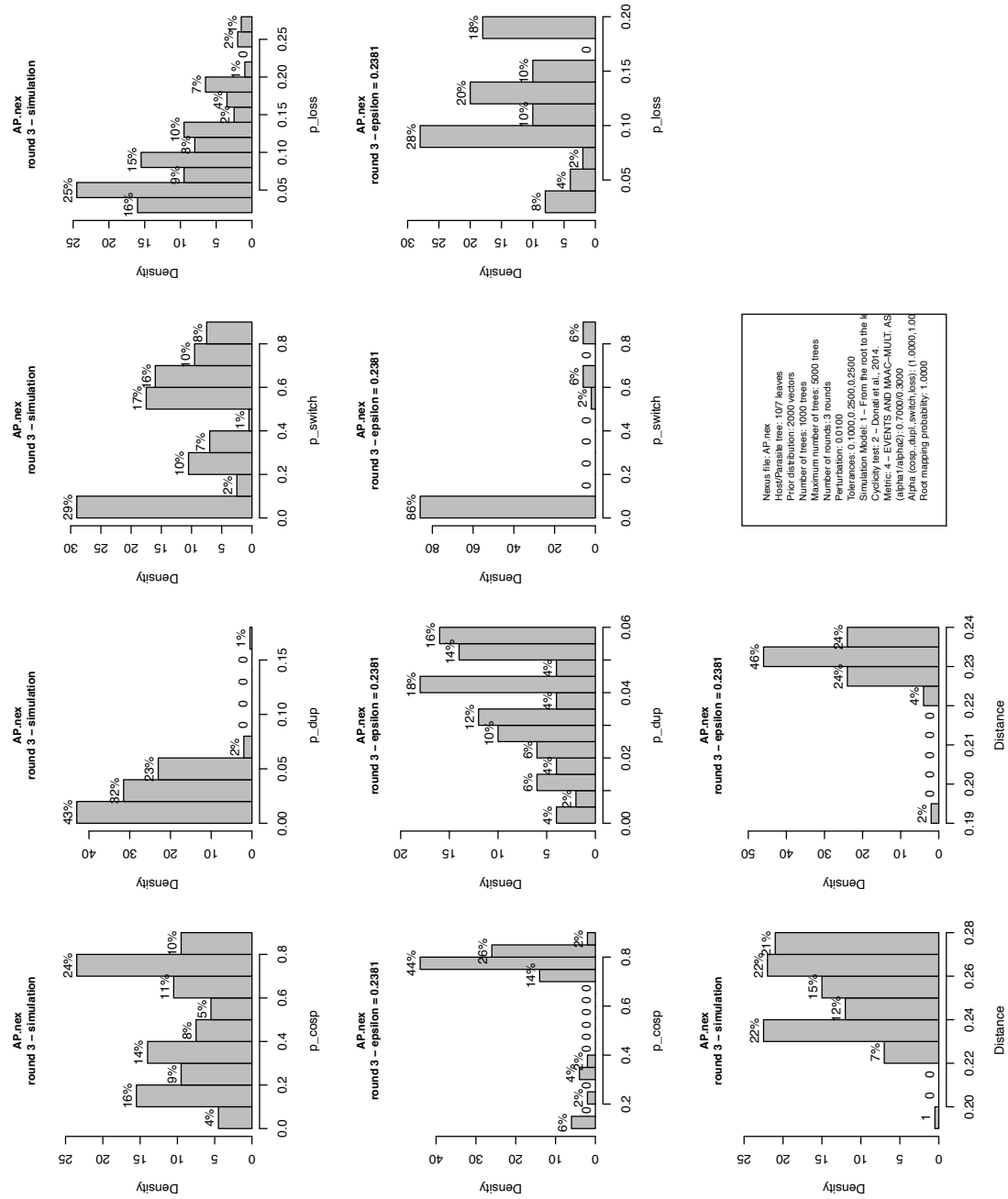


Figure J: AP dataset. First row: histograms of the input parameters. Second row: histograms of the parameters after round 3. Third row: summary discrepancies of the input parameters and of the parameters after round 3.

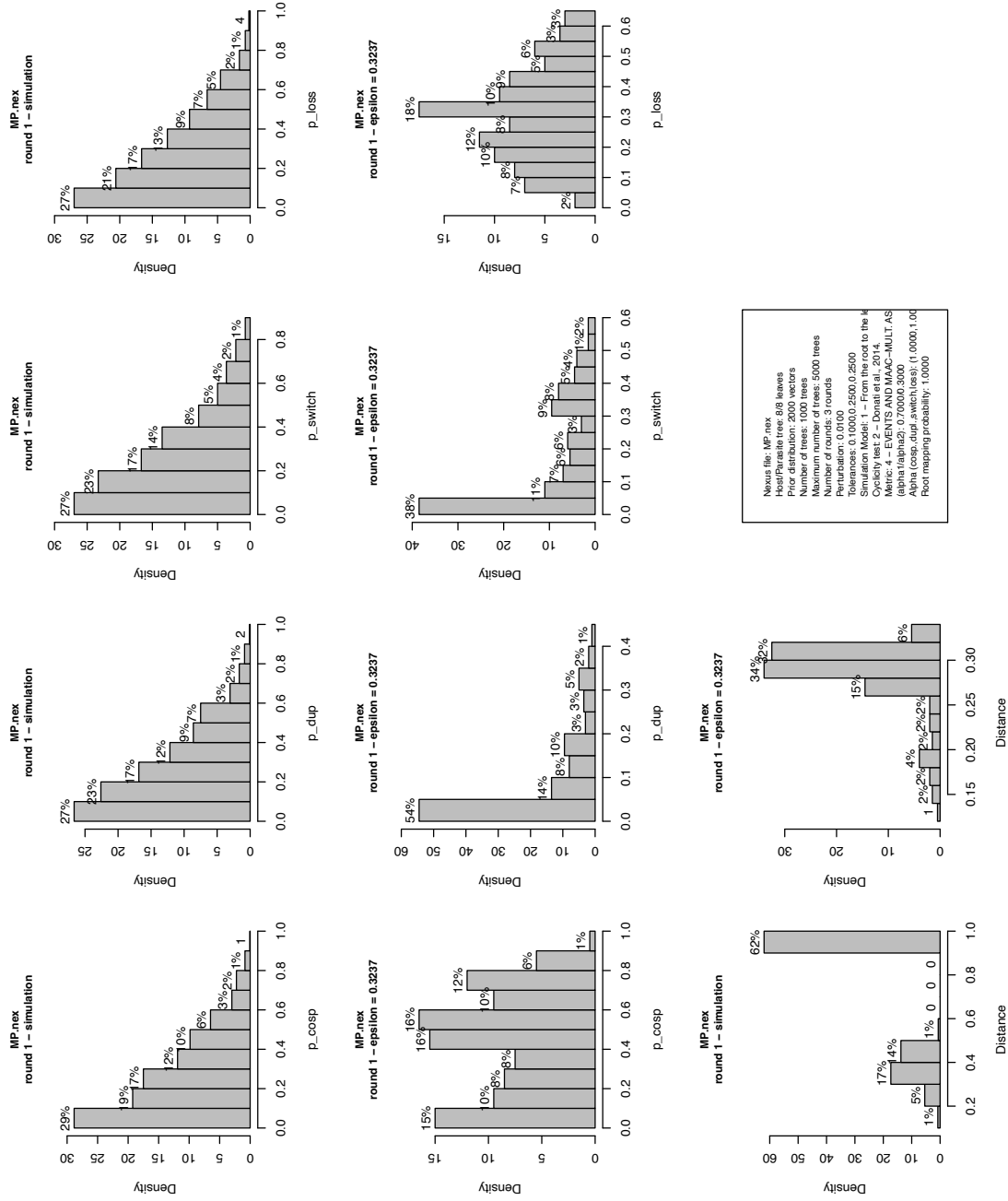


Figure K: MP dataset. First row: histograms of the input parameters. Second row: histograms of the parameters after round 1. Third row: summary discrepancies of the input parameters and of the parameters after round 1.

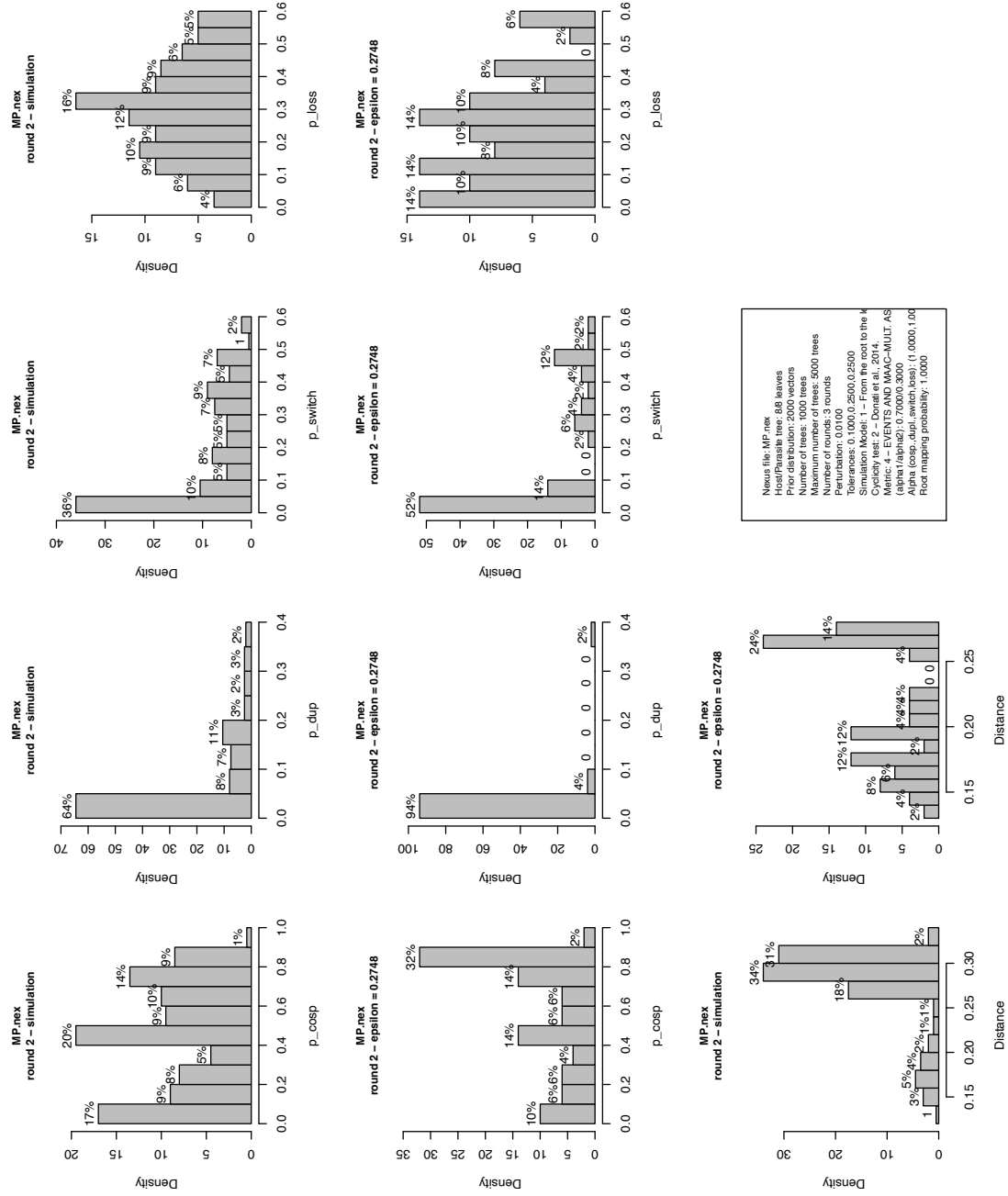


Figure L: MP dataset. First row: histograms of the input parameters. Second row: histograms of the parameters after round 2. Third row: summary discrepancies of the input parameters and of the parameters after round 2.

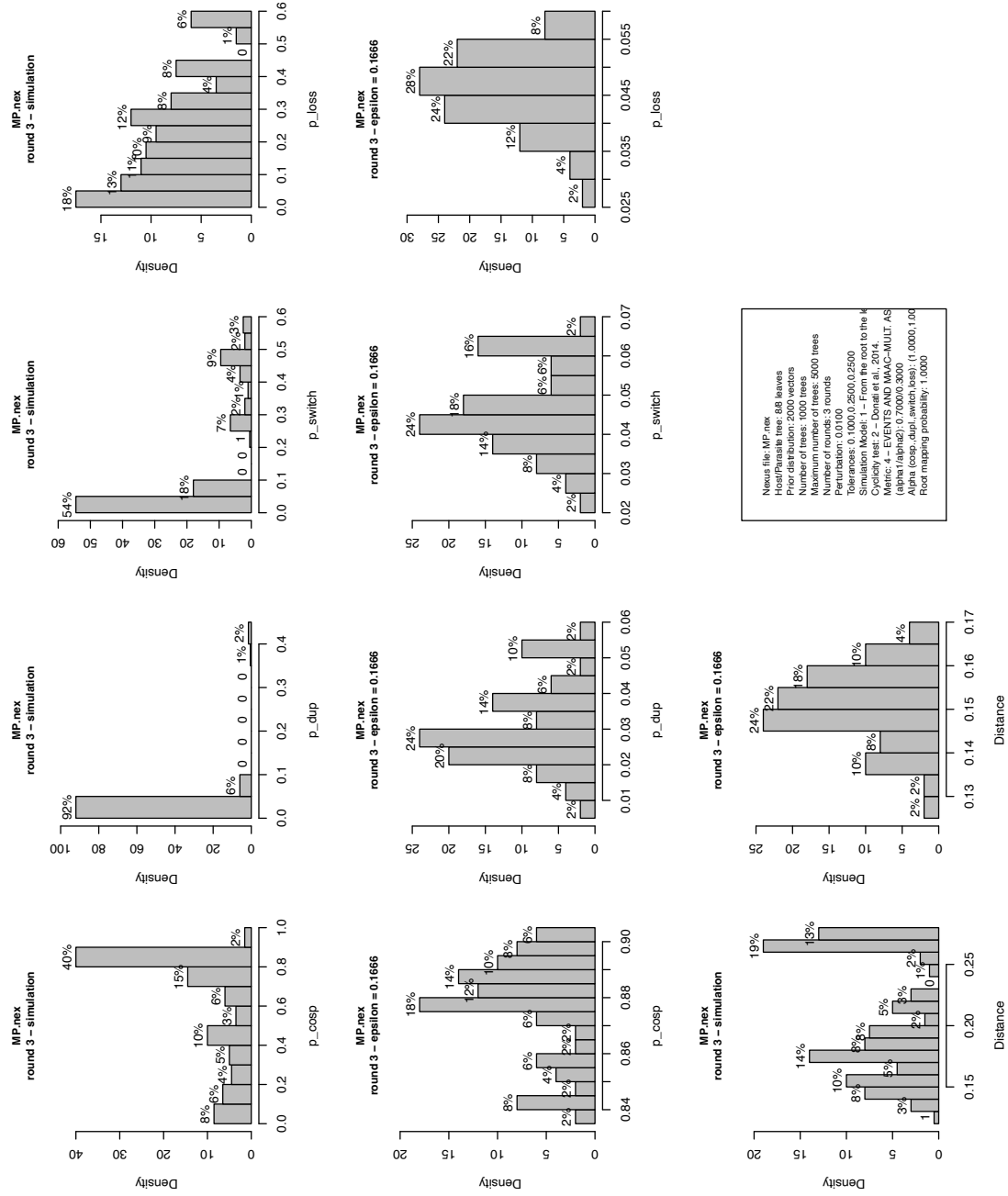


Figure M: MP dataset. First row: histograms of the input parameters. Second row: histograms of the parameters after round 3. Third row: summary discrepancies of the input parameters and of the parameters after round 3.

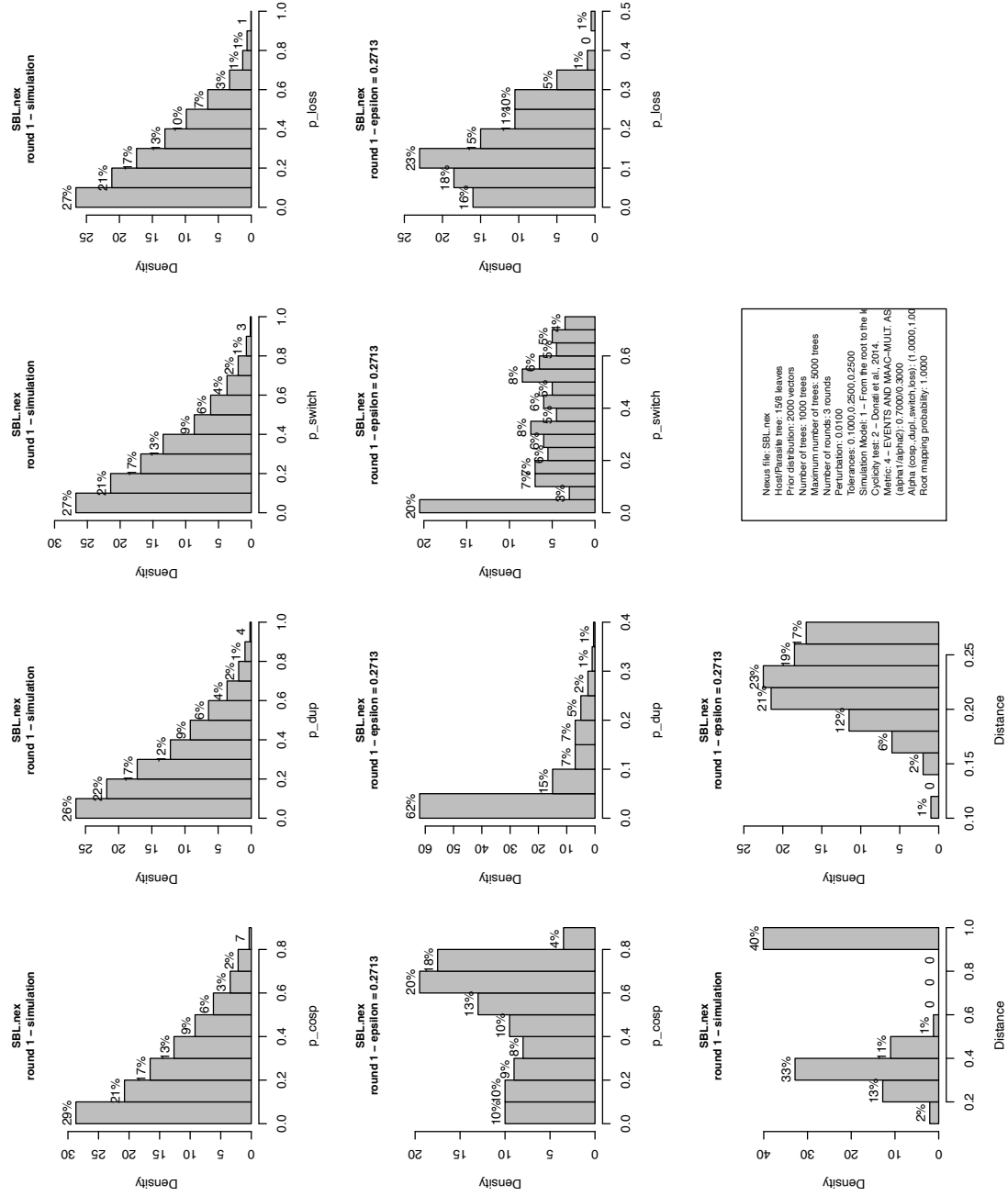


Figure N: SBL dataset. First row: histograms of the input parameters. Second row: histograms of the parameters after round 1. Third row: summary discrepancies of the parameters and of the parameters after round 1.

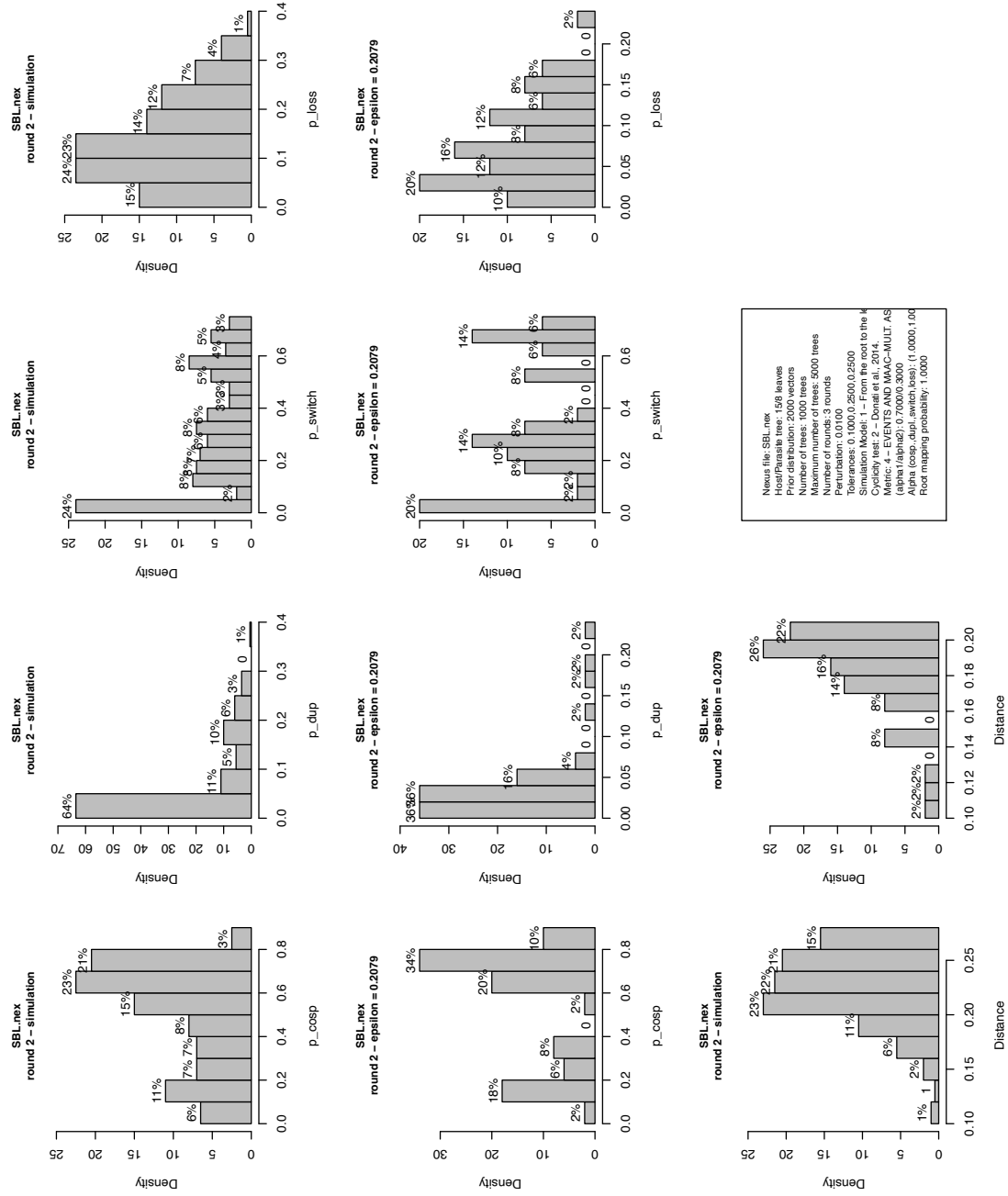


Figure O: SBL dataset. First row: histograms of the input parameters. Second row: histograms of the parameters after round 2. Third row: summary discrepancies of the parameters and of the parameters after round 2.

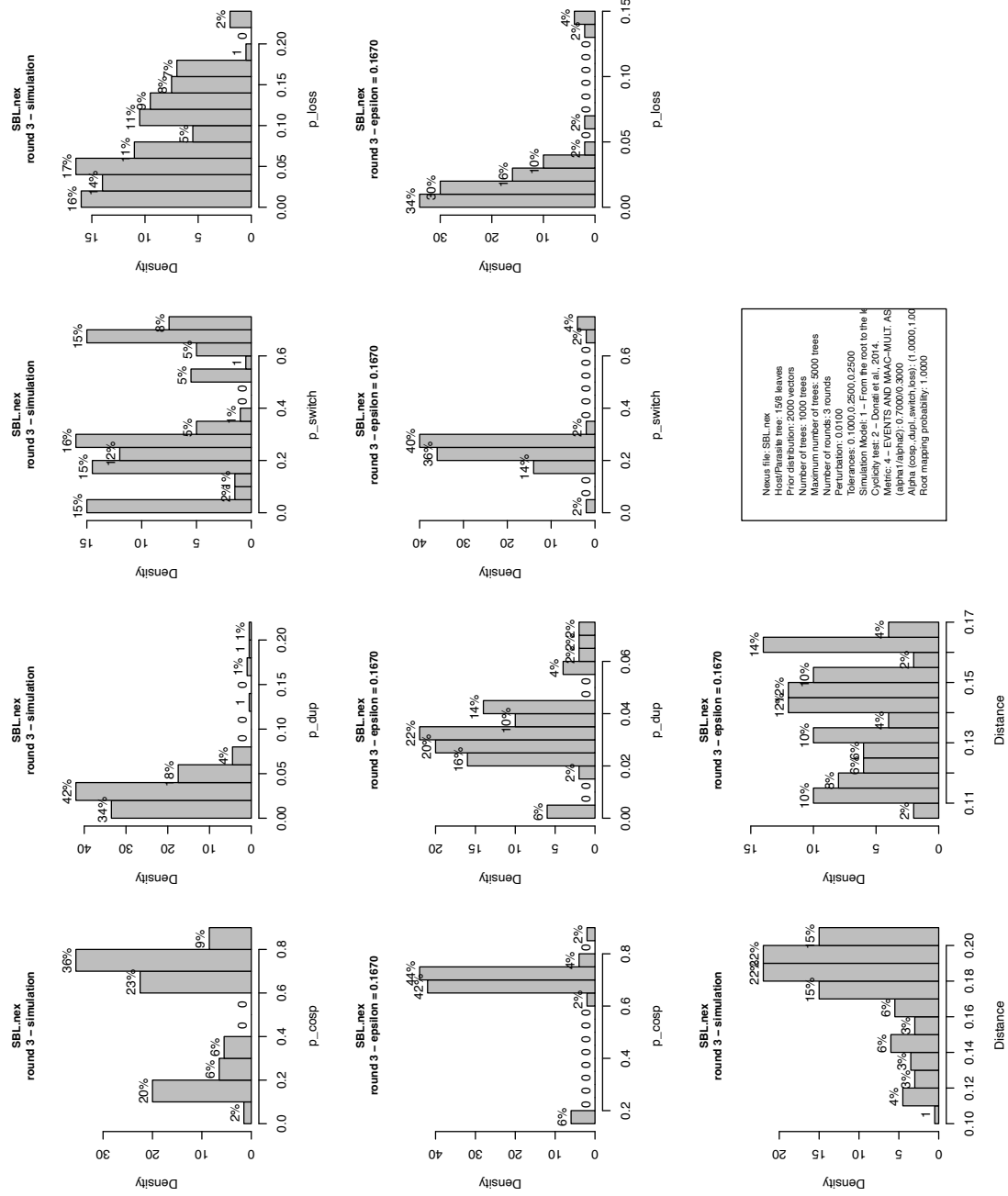


Figure P: SBL dataset. First row: histograms of the input parameters. Second row: histograms of the parameters after round 3. Third row: summary discrepancies of the parameters and of the parameters after round 3.

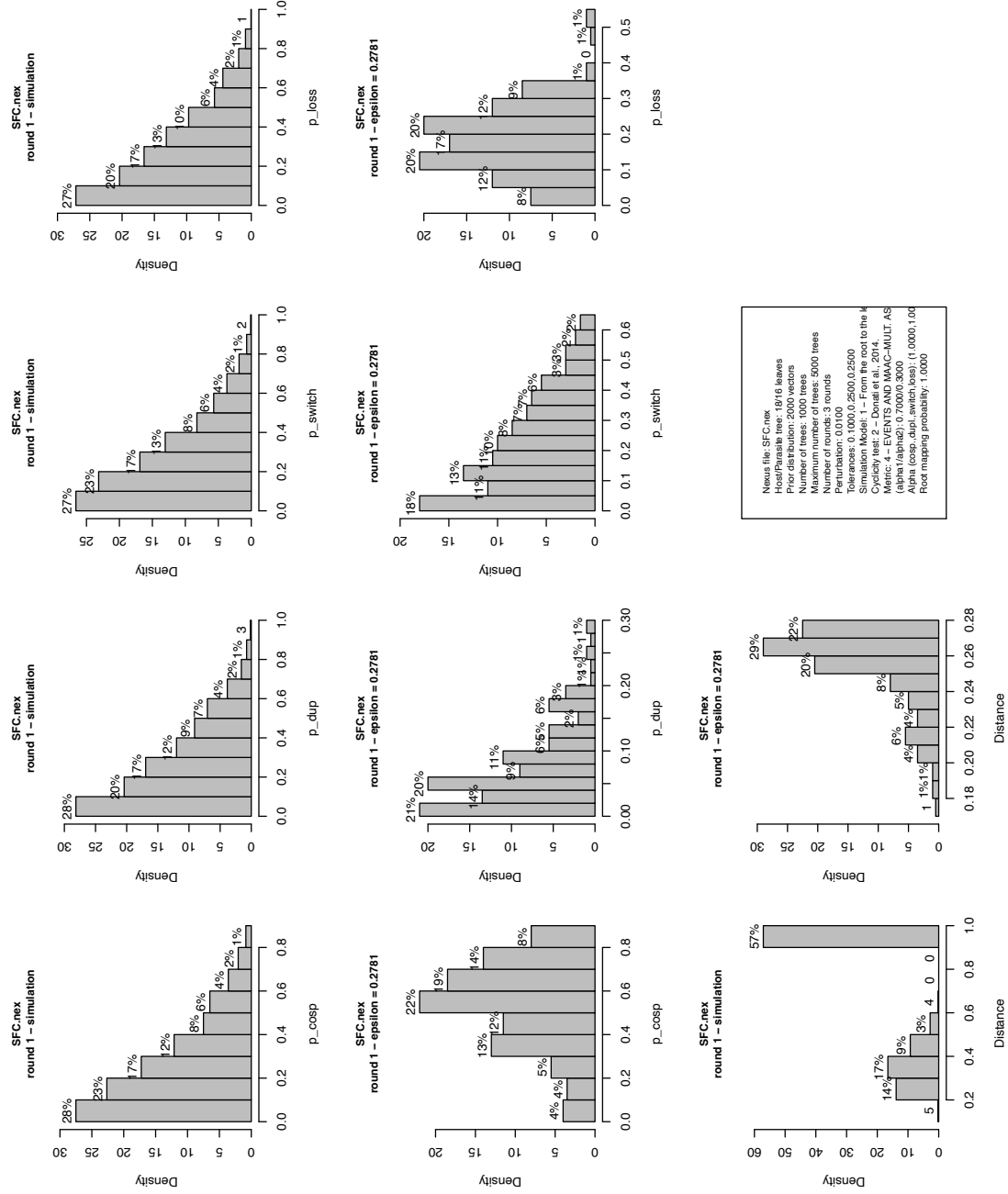


Figure Q: SFC dataset. First row: histograms of the input parameters. Second row: histograms of the parameters after round 1. Third row: summary discrepancies of the input parameters and of the parameters after round 1.

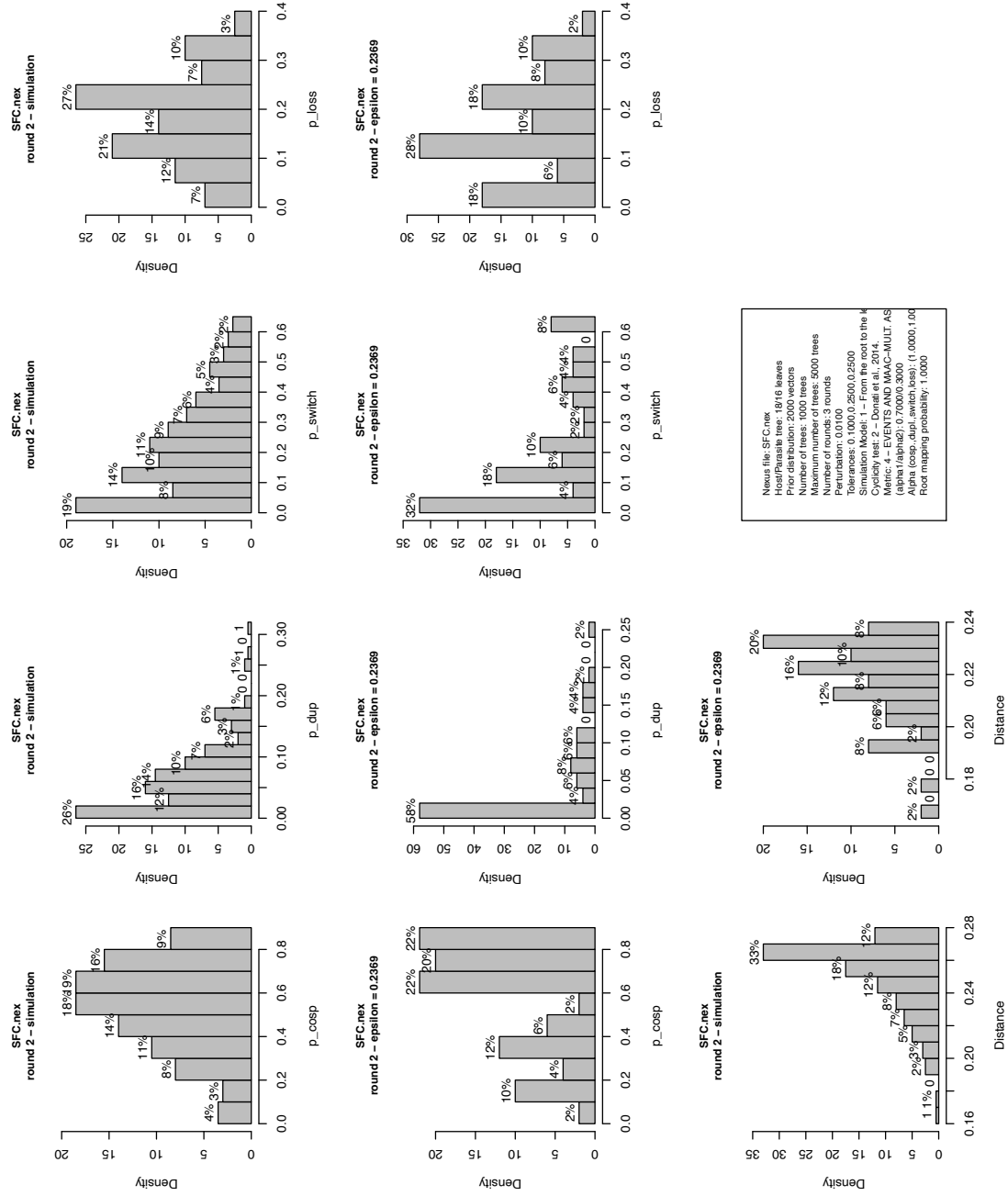


Figure R: SFC dataset. First row: histograms of the input parameters. Second row: histograms of the parameters after round 2. Third row: summary discrepancies of the input parameters and of the parameters after round 2.

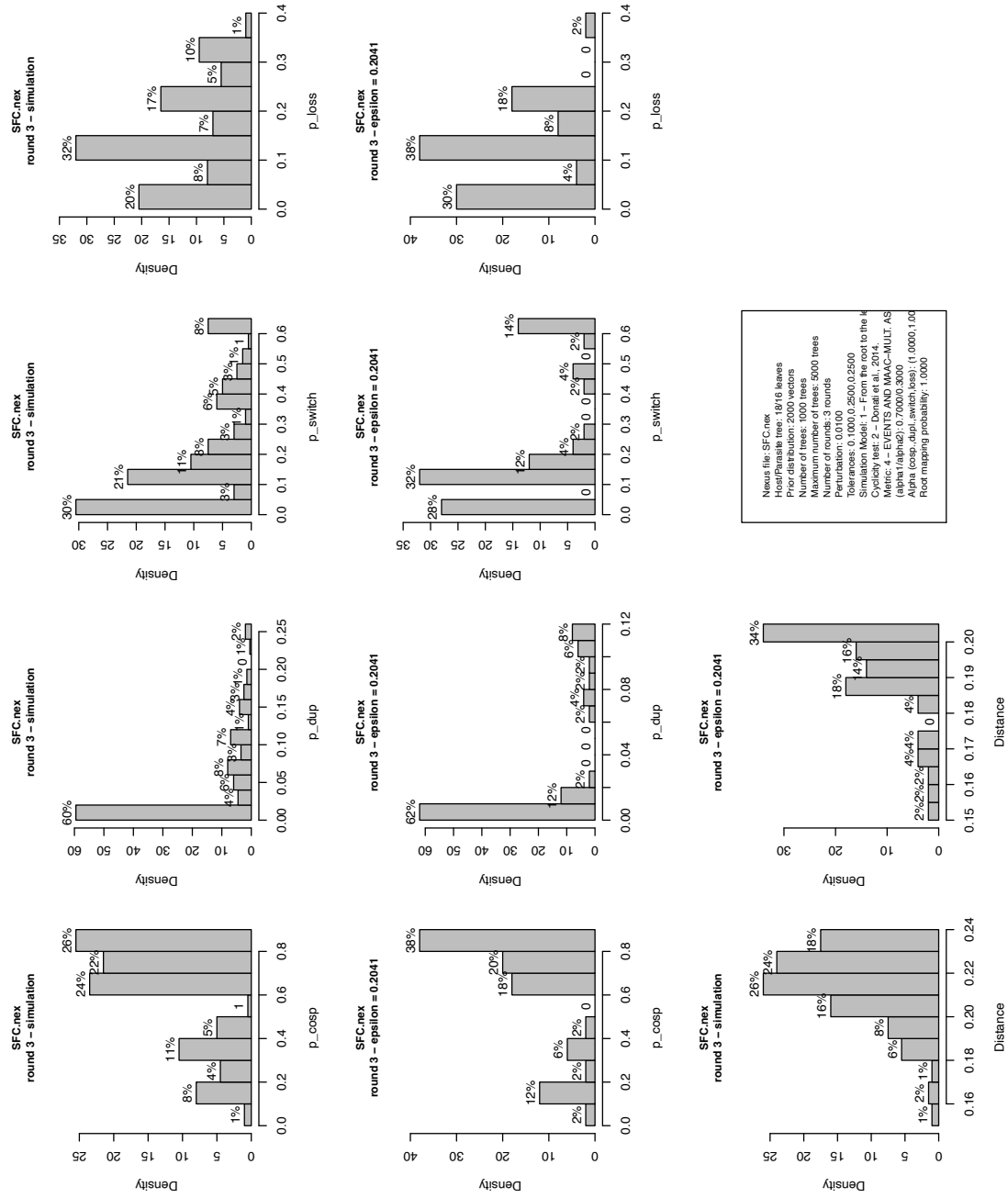


Figure 5: SFC dataset. First row: histograms of the input parameters. Second row: histograms of the parameters after round 3. Third row: summary discrepancies of the input parameters and of the parameters after round 3.

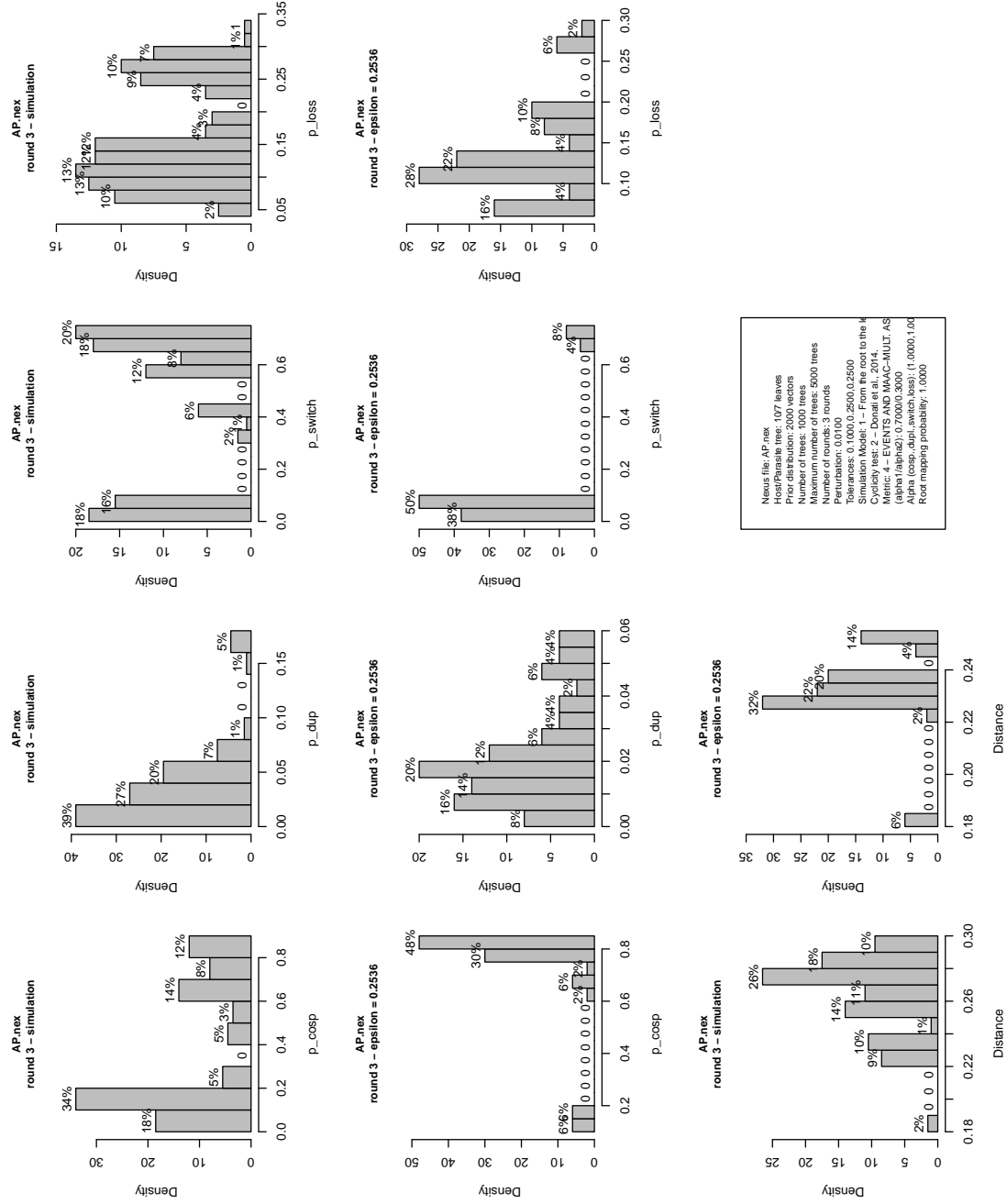


Figure T: AP dataset with perturbed spread probabilities. First row: histograms of the input parameters. Second row: histograms of the parameters after round 1. Third row: summary discrepancies of the input parameters and of the parameters after round 1.

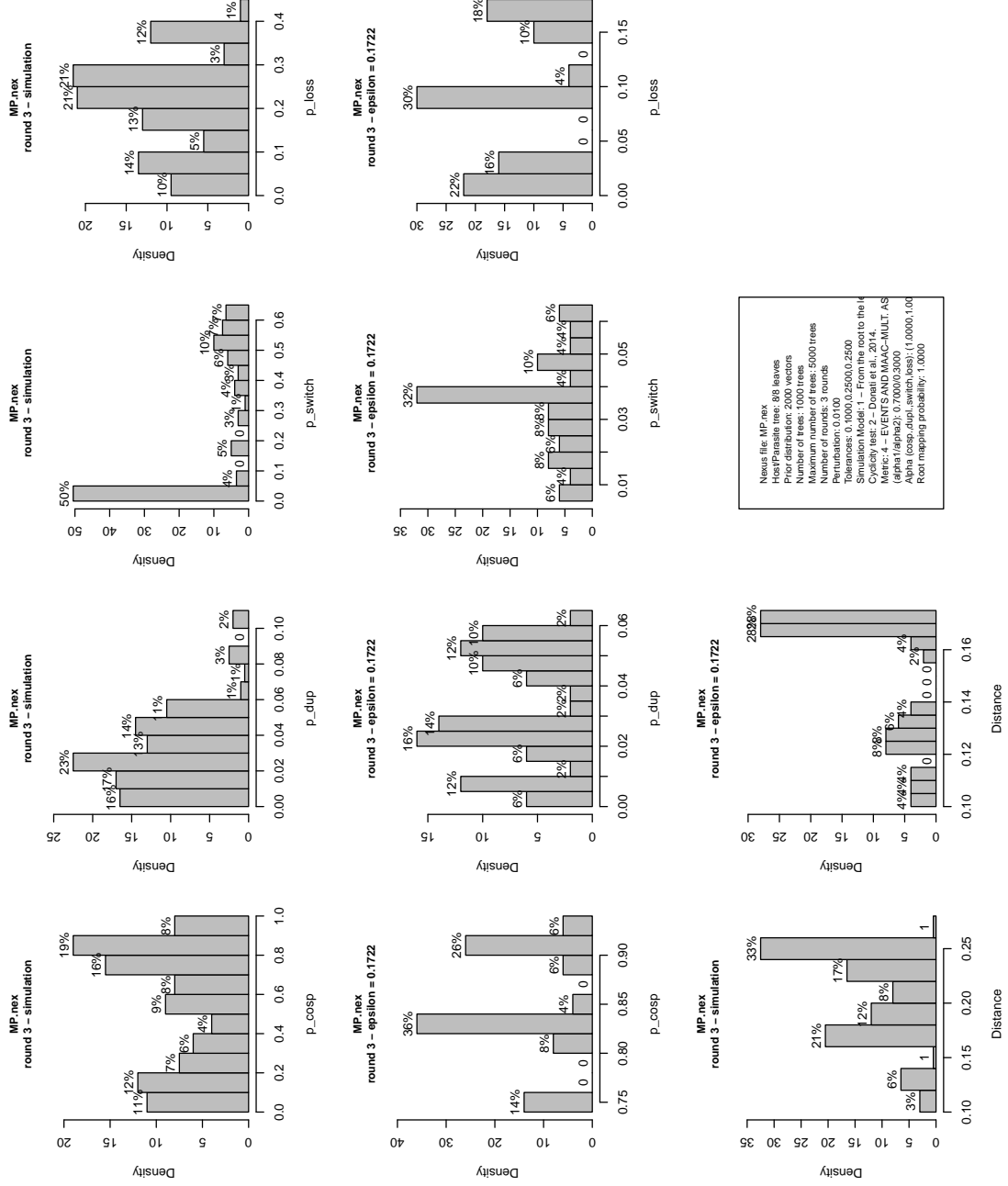


Figure U: MP dataset with perturbed spread probabilities. First row: histograms of the input parameters. Second row: histograms of the parameters after round 1. Third row: summary discrepancies of the input parameters and of the parameters after round 1.

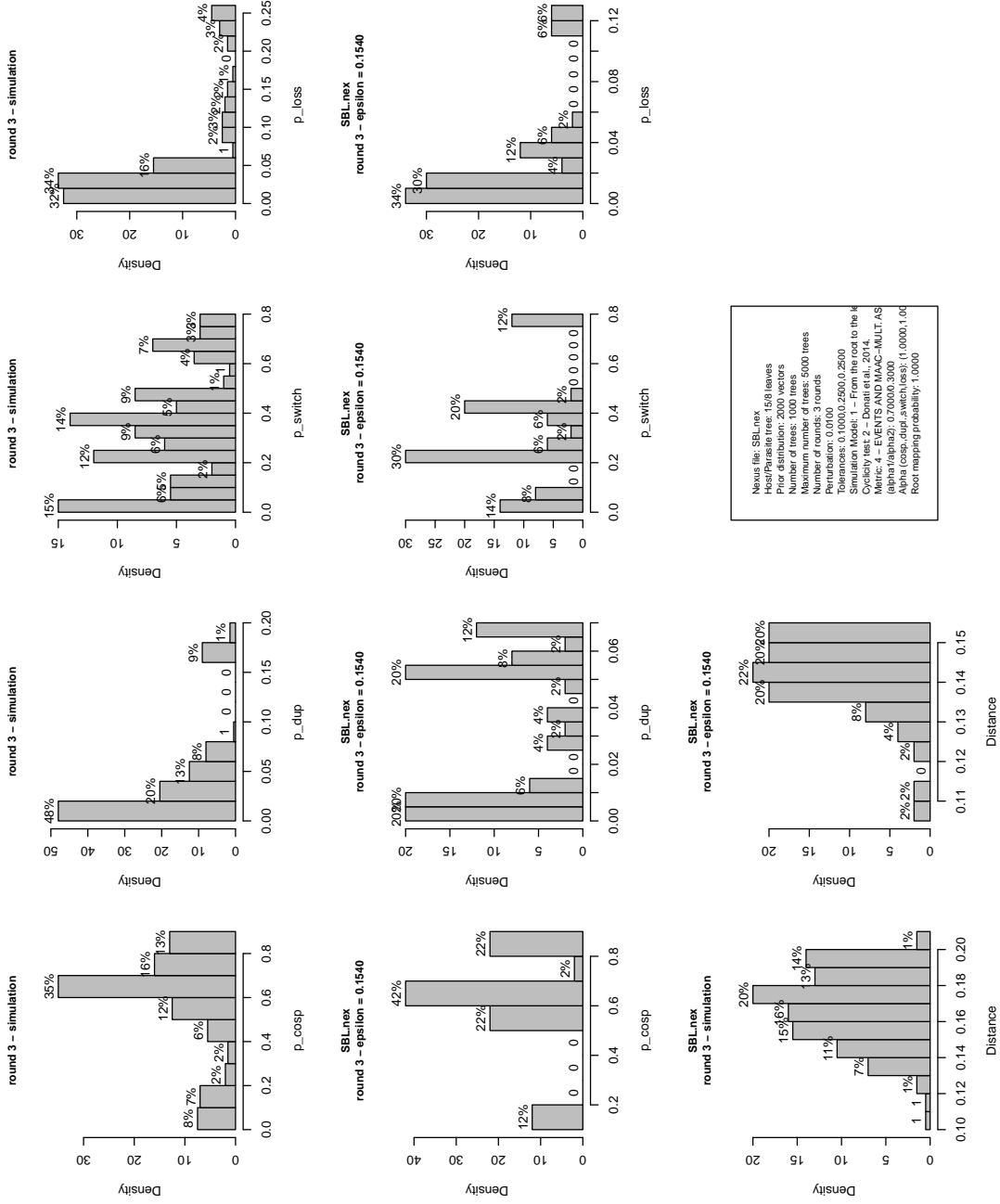


Figure V: SBL dataset with perturbed spread probabilities. First row: histograms of the input parameters. Second row: histograms of the parameters after round 1. Third row: summary discrepancies of the input parameters and of the parameters after round 1.

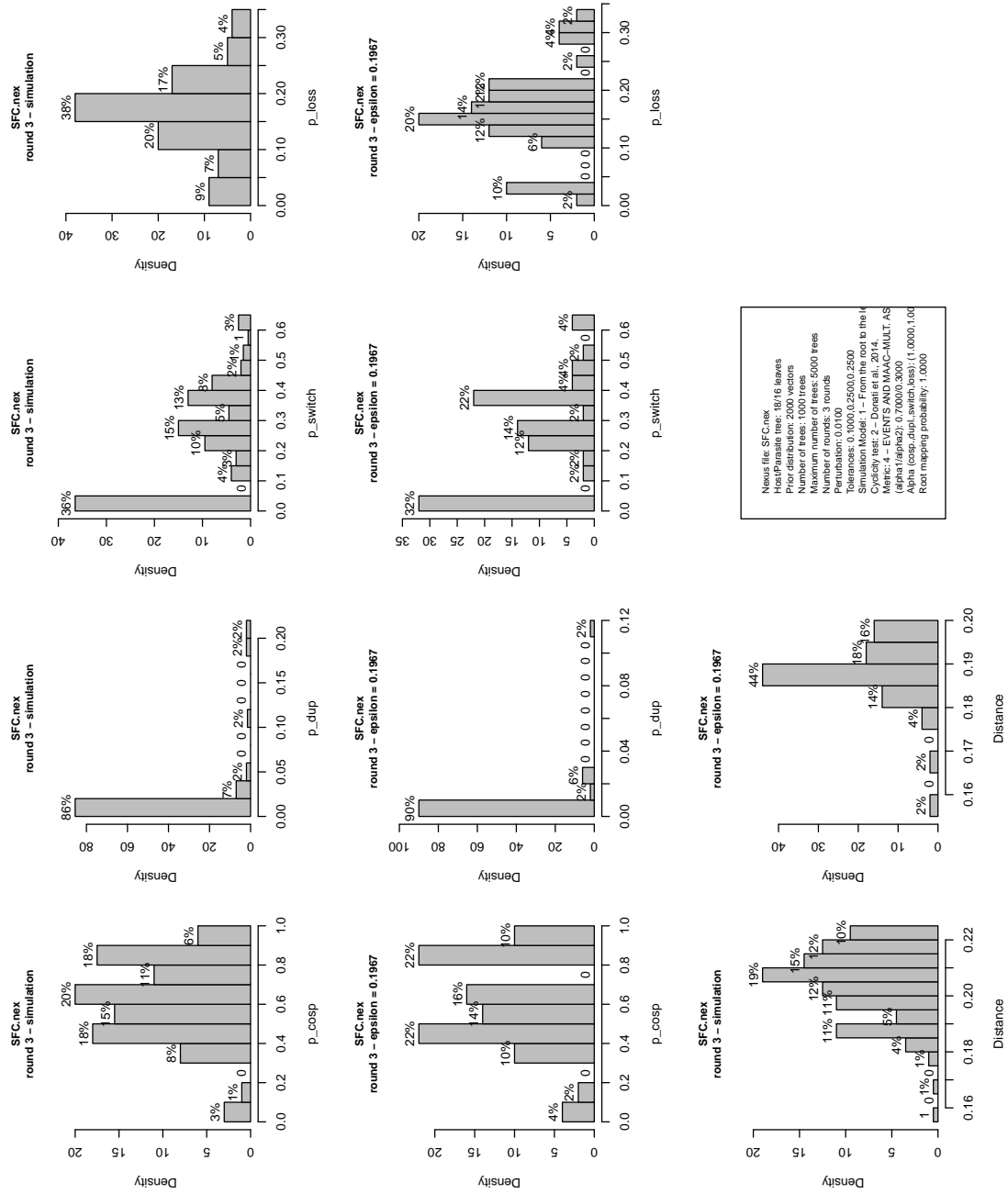


Figure W: SFC dataset with perturbed spread probabilities. First row: histograms of the input parameters. Second row: histograms of the parameters after round 1. Third row: summary discrepancies of the input parameters and of the parameters after round 1.

Table D.2: Representative vectors of the clusters produced by AMOCOALA with perturbations for the SFC dataset. The column *#vectors* indicates the number of vectors in the cluster.

<i>Dataset</i>	<i>Cluster</i>	p_c	p_d	p_s	p_l	<i>#vectors</i>
SFC	1	0.4985	0.0024	0.3162	0.1829	31
	2	0.8738	0.0147	0.0180	0.0935	16
	3	0.1087	0.0012	0.5770	0.3131	3