# A new (old) approach to genotype-based phylogenomic inference within species, with an example from the saguaro cactus (*Carnegiea gigantea*)

Michael J. Sanderson[1,*], Alberto Búrquez[2], Dario Copetti[3], Michelle M. McMahon[4], Yichao Zeng[1], and Martin F. Wojciechowski[5]

[1] *Department of Ecology and Evolutionary Biology, University of Arizona, Tucson, AZ 85721, USA*
[2] *Instituto de Ecología, Unidad Hermosillo, Universidad Nacional Autónoma de México, Hermosillo, Sonora, Mexico*
[3] *Molecular Plant Breeding, Institute of Agricultural Sciences, ETH Zurich, 8092, Switzerland*
[4] *School of Plant Sciences, University of Arizona, Tucson, AZ 85721, USA*
[5] *School of Life Sciences, Arizona State University, Tempe, AZ 85287, USA*

*\*Michael J. Sanderson, Dept. of Ecology and Evolutionary Biology, University of Arizona; Tucson, AZ*

*85721, USA; 520-626-6848; sanderm@email.arizona.edu*

## Abstract

1  Genome sequence data are routinely being used to infer phylogenetic history within and

2  between closely related diploid species, but few tree inference methods are specifically

3  tailored to diploid genotype data. Here we re-examine the method of "polymorphism

4  parsimony" (Inger 1967; Farris 1978; Felsenstein 1979), originally introduced to study

5  morphological characters and chromosome inversion polymorphisms, to evaluate its utility

6  for unphased diploid genotype data in large scale phylogenomic data sets. We show that it

7  is equivalent to inferring species trees by minimizing deep coalescences—assuming an

8  infinite sites model. Two potential advantages of this approach are scalability and

9  estimation of a rooted tree. As with some other single nucleotide polymorphism (SNP)

10  based methods, it requires thinning of data sets to statistically independent sites, and we

11  describe a genotype-based test for phylogenetic independence. To evaluate this approach in

12  genome scale data, we construct intraspecific phylogenies for 10 populations of the saguaro

cactus using 200 Gbp of resequencing data, and then use these methods to test whether the population with highest genetic diversity corresponds to the root of the genotype trees. Results were highly congruent with the (unrooted) trees obtained using SVDquartets, a scalable alternative method of phylogenomic inference.

*Key words*: Polymorphism parsimony, deep coalescence, infinite sites model, phylogenomics

Most published phylogenetic trees of eukaryotes have diploid taxa at their leaves, but most phylogenetic methods implicitly assume input sequence data are haploid—a single row of individual character states, such as nucleotide base calls, for each leaf label. The significance of this disconnect is growing, especially in studies that include multiple individuals within species *and* have high enough coverage to call diploid genotypes reliably. Judging by the number of recent reconstructions of phylogenies of individual diploid accessions within species, there is a gowing demand to characterize the phylogenetic structure in these genotype data (e.g. Durvasula et al. 2017; Wang et al. 2018; Zhao et al. 2019). Yet, even recent comprehensive reviews on phylogenomics (Bravo et al. 2019; Liu et al. 2019) make little mention of "genotypes", despite the fact that diploidy, and polyploidy in general, allows intra-individual polymorphism, which is a well known complicating factor in phylogenetic inference long known in gene families, for example (Goodman et al. 1979; Potts et al. 2014).

Various "workarounds" have been used to recode genotype data into a form that can be used by existing phylogenetic methods. In maximum parsimony, for example, the three genotypes for a biallelic SNP, $\{AA, Aa, aa\}$, can be recoded (i) as the integer states $\{0, 1, 2\}$ and treat them as either unordered or ordered multistate characters (Rheindt et al. 2014), possibly with weighting (Buckler and Holtsford 1996); (ii) as the states $\{0, ?, 1\}$, treating the heterozygote as missing; (iii) as the states $\{0, \{0, 1\}, 1\}$, where the

39  heterozygote is treated as "polymorphic" (in the sense of PAUP* or MacClade, which does

40  not allow an ancestor to be polymorphic, Maddison and Maddison 2000); or (iv) as a pair

41  of binary presence-absence characters, one per allele (Schmidt-Lebuhn et al. 2017). Some

42  of these recoding schemes can be used in likelihood inference (e.g., method (i) in

43  VCFtoTree: Xu et al. 2017) or distance methods. The widely used identity-by-state (IBS)

44  distance (Chang et al. 2015; Subramanian et al. 2019) is actually the same as that implied

45  by the ordered integer coding in (i) above when applied to a single site. It is also the same

46  as the Manhattan distance in allele frequency space if the "population" is just the two

47  alleles in a diploid individual.The latter was used for individual site branch lengths in the

48  FreqPars method of Swofford and Berlocher (1987).

49       Three other frameworks are more explicit about the nature of evolving genotypes.

50  The method of *polymorphism parsimony* (PP) (Inger 1967; Farris 1978; Felsenstein 1979,

51  2004, 2005) infers phylogenies from discrete characters that can be ancestrally polymorphic

52  by minimizing the extent of polymorphism on the tree. Though originally applied to

53  morphological traits and chromosome inversions, it seems clearly applicable to diploid

54  biallelic genotypes derived from genome sequence data. The complexity of morphological

55  traits and chromosome structure led these authors to invoke specific assumptions about

56  character evolution: a unique derivation with no losses of the derived state. This allowed

57  polymorphism to be gained once but lost many times. These assumptions turn out to

58  mirror the infinite sites model of sequence evolution (Kimura 1969), a similarity we return

59  to below. PP is implemented in Phylip (Felsenstein 2005) but has has been used

60  infrequently (Baum 1983; Suh et al. 2015).

61       Two other frameworks, gene tree reconciliation in phylogenetics (Goodman et al.

62  1979; Page 1994), and coalescent theory in population genetics (Kingman 1982; Hein et al.

63  2005), also provide toolkits for modeling genotypes more explicitly by analyzing the

64  genealogy of alleles. Yet, existing implementations are rather agnostic about using

65  genotype vs. haplotype data. For example, considerable headway has been made in tying

phylogenetic inference between species/populations to allele phylogenies using either the gene tree reconciliation framework, which aims to infer species/population trees by minimizing the deep coalescence (DC) score in allele trees (Maddison 1997; Than and Nakhleh 2009); or a more explicit probabilistic model provided by the multi-species coalescent (MSC) (Degnan and Rosenberg 2009). In the MSC framework, several methods use biallelic SNP data to build trees, such as SNAPP (Bryant et al. 2012); SVDquartets (Chifman and Kubatko 2014; Vachaspati and Warnow 2018); Phylonet (Zhu and Nakhleh 2018). Yet, none of these methods embrace genotype data per se. Instead, they assign single sequences—haplotypes—to predefined populations or species (possibly one sequence per species) and infer the phylogenies of those taxa.

These approaches (and others) can be "spoofed" into taking diploid genotype inputs by converting them to pairs of pseudo-haplotypes and then aggregating them into populations or species. At the lowest levels these would be "populations" of two pseudo-haplotypes. We say "pseudo-haplotypes" unless they are explicitly phased either computationally (e.g. Chifman and Kubatko 2014) or experimentally. For example, VCFtoTree (Xu et al. 2017) can use phased haplotypes, but if the data are unphased it simply ignores heterozygous sites. Alternatively, if all sites are treated as coalescent independent sites (Tian and Kubatko 2017), then the two haplotypes in each diploid can be constructed by segregating the alternative alleles in all heterozygotes randomly.

However, this raises another issue. At the species level and below, even a phylogenetic inference method properly focused on genotypes still runs into the problem of gene flow and non-tree-like history. In a Wright-Fisher population of genotypes, the paired alleles within a diploid individual most likely do not coalesce with each other more recently than they do with alleles from other individuals in that population. In other words, individuals (genotypes) do not act like populations. However, this does not invalidate the concept of genotype tree or genotype-specific tree inference methods, per se; it means they run the same risks as many other phylogenetic tools applied to the infraspecific level.

93    Existing approaches to building trees with genotype data, whether workarounds, or

94    implicit genotype-based tools, are all ad hoc to various degrees, and it may prove useful to

95    examine all of them in a common allele tree framework to expose assumptions and

96    limitations. In this paper we take a step toward this goal by showing conditions under

97    which polymorphism parsimony, and allele tree reconciliation by minimizing deep

98    coalescences, are equivalent. We then exploit the existing implementation of PP in Phylip

99    (Felsenstein 2005) to study phylogenetic structure among genotypes of saguaro cactus

100    (*Carnegiea gigantea*), across its range, illustrating its scalability to millions of SNPs and

101    utility in estimating a rooted genotype tree. By adopting the allele tree framework shared

102    by reconciliation methods and coalescent theory, we hope to illuminate the strengths of

103    genotype data while exploring the limits to infraspecific phylogenetic inference imposed by

104    interbreeding.

## Materials and Methods I. Theory

105

### *Genotypes, characters and trees*

106

107    Assuming each site is diploid and biallelic with alternate alleles labelled 0 and 1,

108    and has therefore three possible diploid genotype states plus "missing", we will use one of

109    two notation schemes for the genotypes of individual $i$ and site $j$:

$$G_{ij} \in \{0/0, 0/1, 1/1, ./.\} \text{ or } G_{ij} \in \{(+,-), (+,+), (-,+), (-,-)\} \tag{1}$$

110    The first is the notation typical for biallelic SNPs in modern genotyping pipelines (e.g.,

111    VCF format). The second (useful in proofs in the Appendix) is a vector in which the first

112    and second components denote the presence ($+$) or absence ($-$) of the 0 and 1 alleles,

113    respectively, in that genotype. Define a *genotype character* for site $j$ in $n$ individuals as the

114    vector, $(G_{1j}, ..., G_{nj})$. We will often use the shorthand, $G$, for a genotype character.

115    We assume that genotype characters are statistically independent—that is, sites are

116    *coalescent independent* (Tian and Kubatko 2017)—but we test for dependence to obtain

117 such data. Genotypes at different sites in an individual are assumed to be *unphased*

118 (genotypes are *phased* if the alternate alleles in heterozygous genotypes at different sites

119 are ordered with respect to each other as *haplotypes*).

120       A *genotype tree*, $\Psi$ for $n$ individuals is a rooted tree with $n$ leaves. A genotype

121 character, $G$, for these individuals has an *allele tree*, $t(G)$, with at most two leaves per

122 individual (fewer if genotype data are missing), which is imbedded within $\Psi$, in keeping

123 with the gene tree reconciliation framework (Goodman et al. 1979; Page 1994; Page and

124 Charleston 1997). Both trees are assumed to be binary.

125 *Equivalence of Polymorphism Parsimony and Minimizing Deep Coalescence in Allele Trees*

126       Our main theoretical result is that for any $\Psi$ and $G$, the PP score is the same as the

127 deep coalescence score *for an appropriately specified allele tree, $t(G)$*. This underlying allele

128 tree is always implicitly present, but its coalescent history (Degnan and Salter 2005) and

129 therefore deep coalescence score, are not identifiable without additional constraints, which

130 we discuss below.

131       *Polymorphism parsimony for genotypes.—* Polymorphism parsimony assumes a

132 particular model underling the evolution of genotypes at a site: (i) the 0/0 genotype is

133 ancestral; (ii) the 0/1 genotype can originate only once, and (iii) both 0/0 and 1/1 can

134 evolve from 0/1 any number of times (Farris 1978; Felsenstein 1979). The original rationale

135 was that a novel state, such as a new chromosome inversion, arises once in an ancestral

136 edge of the tree, at which point the edge is polymorphic along with the ancestral

137 chromosome type, and then later the *polymorphism* may be repeatedly lost, by fixation of

138 one or the other type, but, importantly, the novel trait itself is never reversed.

139       Let $y_\Psi(G) \in Y_\Psi^{01}(G)$ be a set of ancestral states that jointly satisfy these

140 assumptions for genotype character, $G$, where $Y_\Psi^{01}(G)$ is the set of all such joint

141 reconstructions. The number of polymorphic edges, $m(y_\Psi(G))$, which is the number of

¹⁴² edges having 0/1 genotypes at both ends, provides an optimality criterion.

¹⁴³ The *polymorphism parsimony (PP) score* is then found by solving the following

¹⁴⁴ problem:

¹⁴⁵ **Problem (PP Score).** Given $\Psi$ and $G$, find the ancestral state set, $y_{opt}$ over all

¹⁴⁶ $y \in Y_\Psi^{01}(G)$ that minimizes $m(y)$. Then the PP score is

$$c_{PP}(\Psi, G) = m(y_{opt}). \tag{2}$$

¹⁴⁷ Felsenstein (1979) outlined a two-pass algorithm to compute these ancestral states

¹⁴⁸ and the PP score. Details are provided in the Appendix, but the intuition is as follows: in

¹⁴⁹ the downpass each node is assigned a pair of presence/absence flags that indicate whether

¹⁵⁰ allele 0 and 1, respectively, are found among any descendant leaves. Because 0 is ancestral,

¹⁵¹ this immediately means that the 0 allele is present at this node. However, the 1 allele

¹⁵² might not yet have evolved. The second phase is an uppass from the root, delaying as long

¹⁵³ as possible the final assignment of the first appearance of a 1 allele. This ensures a minimal

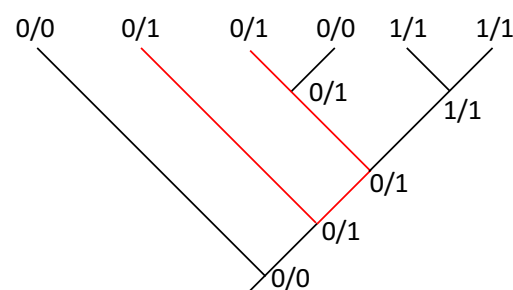¹⁵⁴ number of edges of the tree having both 0 and 1. Figure 1 illustrates an example.



Fig. 1. Ancestral state reconstruction under polymorphism parsimony. The ancestral states are jointly valid under the assumptions of the PP model. Optimal reconstruction indicates four polymorphic edges (red lines) and therefore a PP score of four.

¹⁵⁵ PP can return different ancestral state reconstructions depending on whether 0/0 or

¹⁵⁶ 1/1 is assumed to be the ancestral genotype. The Phylip implementation computes the PP

157  score for both choices and selects the smaller of the two. Although neither Farris nor

158  Felsenstein explicitly consider missing data in their papers, the modification to the

159  Felsenstein algorithm is slight (see Appendix). The Felsenstein algorithm runs in $O(n)$

160  time, where $n$ is the number of leaves of $\Psi$.

161       *Allele Tree Framework.—*  Minimizing the number of polymorphic edges in a tree

162  of genotypes sounds suspiciously like "minimizing deep coalescences", but the latter

163  usually refers to gene trees imbedded within species or population trees. The connection

164  between these seemingly different frameworks can be made quite strong for genotype data

165  under suitable assumptions. First, we assume the existence of a genotype tree, $\Psi$,

166  containing an allele tree, $t(G)$, giving rise to the observed genotypes at the leaves of $\Psi$.

167  Second, because any allele tree can be concordant with any genotype tree trivially by

168  allowing some or all coalescences to occur below the root (Degnan and Salter 2005, p. 25),

169  we will constrain the allele tree to evolve according to the infinite sites model in population

170  genetics. This assumption is consistent with the genotype evolution model described above

171  for polymorphism parsimony. Finally, we will invoke an optimality criterion to select

172  among all allele trees that fit this model.

173       The first assumption is operational. Even in a randomly mating population of

174  individuals, we can postulate the existence of $\Psi$, but it would bear little relationship to the

175  allele tree, because coalescence events would occur at random, and could therefore be

176  imbedded in $\Psi$ only by forcing most coalescences to occur below the root of $\Psi$. However,

177  with increasing amounts of historical isolation between these individuals, there will be a

178  bias in the distribution of allele trees imbedded within the genotype tree away from

179  expectations of a random coalescent.

180       The second assumption requires allele trees, $t$, that obey the PP model. Because the

181  three model assumptions stated above for PP now involve an interplay between $t$ and $\Psi$ in

182  this allele tree framework, it is helpful to leverage the gene tree reconciliation framework

183  from here on. The key model assumption is that allele states 0 and 1 on $t$ are constrained

184 such that there can only be one $0 \to 1$ transition and no $1 \to 0$ transitions. A second

185 model assumption is that this unique origin of the 1 allele arises just "after" a coalescence

186 event in the allele tree (i.e., more recently). This corresponds to the PP assumption that

187 when the novel allele arises it is in a polymorphic condition in $\Psi$—that is, there is

188 co-occurring allele tree edge within $\Psi$ having the 0 allele. In the reconciliation framework

189 this can be thought of as a "duplication". Finally, from a polymorphic edge of $\Psi$, in which

190 allele tree edges having both the 0 and 1 alleles are imbedded, one or the other allele tree

191 edge can go extinct, or be "lost", leaving the remaining allele tree edge present. Formally,

192 these duplication and loss events are the same as those described in gene tree

193 reconciliation (Zhang 2011). They are also equivalent to the infinite sites model in

194 population genetics (Kimura 1969). Let $\mathcal{T}_G^{01}$ be the set of all rooted binary allele trees that

195 (i) have leaf states with the correct mapping from $G$ to binary allele states; and (ii) satisfy

196 the infinite sites assumption.

197      Finally, as an optimality criterion to select among allele trees $t \in \mathcal{T}_G^{01}$, we use the

198 *deep coalescence* (DC) score. An allele tree's DC score, $c_{\mathrm{DC}}(\Psi, t)$, is the total number of

199 "extra" allele tree edges per edge of $\Psi$ (Maddison 1997). See Appendix for further details.

200 Now we are in a position to generalize the DC score to genotypes by solving the following:

201 **Problem (DC-G Score).** Given $\Psi$ and $G$, find the allele tree, $t_{min}$ over all $t \in \mathcal{T}_G^{01}$ that

202 minimizes $c_{\mathrm{DC}}(\Psi, t)$. The corresponding DC-G (deep coalescence–genotype) score is

$$c_{\mathrm{DC\text{-}G}}(\Psi, G) = c_{\mathrm{DC}}(\Psi, t_{min}). \tag{3}$$

203      The solution, $t_{min}$, can be found efficiently because of the constraints of the

204 genotype evolution model (Fig. 2). Intuitively, because the 1 allele evolves exactly once and

205 is unreversed, $t_{min}$ has to have two subtrees, $t_{min}^0$ and $t_{min}^1$, each having only leaves labelled

206 with one of the alleles. In addition, $t_{min}^1$ is a clade. The subtree $t_{min}^0$ must be present at the

207 root of $\Psi$, because 0 is assumed ancestral and must be present there. Moreover, to

208 minimize the number of edges of $\Psi$ in which both allele subtrees are present (and hence

209 have an "extra edge" in the sense counted by the DC score) all nodes in $t_{min}$ are placed as

close to the leaves as possible. Finally, a homozygous genotype state at a leaf is

represented by a coalescence in the allele tree at the leaf node, which ensures that no deep

coalescence occurs on a terminal edge having a homozygous leaf (see Appendix for details).

Given $t_{min}$, the number of deep coalescences can be easily computed based on the

subtree of $\Psi$ in which the two subtrees of $t_{min}$ overlap, $\psi$ (Fig. 2):

$$c_{\text{DC-G}}(\Psi, G) = |E(\psi)|, \tag{4}$$

where $|E(\psi)|$ is the number of edges of $\psi$.

The form of the optimal allele tree is of less interest than the fact that it leads to

the following result:

**Theorem 1.** *Given a rooted genotype tree, $\Psi$, and genotype character, $G$, and assuming that the pair of alleles underlying the genotypes are encoded by a binary character evolving on a rooted allele tree according to the infinite sites model (i.e., at most one $0 \to 1$ transition permitted, with no reversals), then*

$$c_{DC\text{-}G}(\Psi, G) = c_{PP}(\Psi, G).$$

In other words, the PP score is the same as the DC score of an optimal imbedded

allele tree for those genotypes. See Appendix for Proof. This also means that the DC-G

problem can be solved in $O(n)$ time, because it is equivalent to solving the PP problem.

*MDC-G: Inferring a genotype tree from genotype characters*

The "Minimize Deep Coalescences" (MDC) problem, which infers a species tree

from a collection of gene trees by minimizing the sum of DC scores across the gene trees

(Maddison 1997; Ma et al. 2001; Maddison and Knowles 2006; Than and Nakhleh 2009;

Bansal et al. 2010; Zhang 2011), can be modified to use genotypes alone.

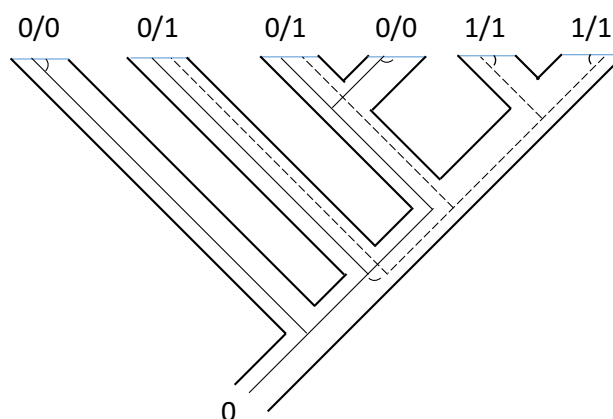**Problem (MDC-G Tree).** Given a collection of genotype characters, $\{G_i\}$, find the

Fig. 2. Genotype tree, $\Psi$, containing the optimal imbedded allele tree, $t_{min}$. Tree $t_{min}$ is constructed from the two minimal subtrees for the two alleles, $t^1_{min}$ (dashed line) and $t^0_{min}$ (fine solid line). These two trees join at a coalescence (duplication) event indicated by a small curved edge. Note the coalescence of homozygotes is shown by very short curved edges near the leaves of $\Psi$. A deep coalescence on $\Psi$ is an edge containing "extra edges" of the imbedded allele tree for its full length. There are four such edges on $\Psi$, each having one extra imbedded edge. This tree has the fewest possible deep coalescences, among all allele trees with these genotypes that obey the infinite sites model. Its DC score, 4, is the same as its PP score (see Fig. 1), which is also equal to the number of "overlapping" edges of $\Psi$—those containing both allele subtrees. Note also that the edges showing "polymorphism" or "deep coalescence" under the two approaches are the same. See Appendix.

genotype tree, $\hat{\Psi}$, that minimizes

$$\sum_i c_{\text{DC-G}}(\hat{\Psi}, G_i) \tag{5}$$

By Theorem 1, this is the same tree inferred by optimizing the PP criterion across these genotypes. Since an implementation of PP is available in Phylip (Felsenstein 2005), it is not necessary to implement the more involved computation implied by the definition of the DC-G score. The decision problem version of finding a tree using the PP score is NP-complete (Day and Sankoff 1987), and thus MDC-G is at least as computationally hard. However, heuristic search implementations benefit from the $O(n)$ runtime of the inner loop score calculation in the PP algorithm (Felsenstein 1979).

## Materials and Methods II. Data Analysis

### Genome Sequence Data

We surveyed population genomic variation in the saguaro cactus (*Carnegiea gigantea* (Engelman.) Britton & Rose) using short read genome sequence data for 20 diploid individuals in 10 populations across its range. We used FASTQC v.0.11.8 (Andrews 2018) to examine all read sets and Trimmomatic 0.38 (Bolger et al. 2014) to trim low quality ends. Options for the latter were either "ILLUMINACLIP:TruSeq3-PE:2:30:10 HEADCROP:3 LEADING:20 TRAILING:20 SLIDINGWINDOW:4:20 MINLEN:36" or "ILLUMINACLIP:TruSeq3-PE:2:30:10 CROP:147 HEADCROP:3 LEADING:20 TRAILING:20 SLIDINGWINDOW:4:20 MINLEN:36" (for the Orégano and San Marcial samples), depending on quality control reports. Trimmed reads were mapped to the saguaro SGP5 v.1.3 assembly ("SGP5" henceforth) (Copetti et al. 2017) using Bowtie2 v. 2.3.4.3 (Langmead and Salzberg 2012) in paired end mode, discarding unaligned, mixed and discordant pairs. Deduplication was implemented with Samtools v. 1.9. (Li et al. 2009) with the command sequence (sort, fixmate, sort, markup -r) described in their documentation, and resulting BAM files were indexed with samtools. Statistics for resulting BAM files were obtained using the samtools stats command, and coverage was computed based on a SGP5 genome size of 1.403 Gbp from the C-value database (Bennett and Leitch 2012).

Subsets of sites in three genomic regions were extracted from BAM files based on the SGP5 assembly. The first consisted of all 1893 scaffolds larger than 100 kb, which totalled 281 MB, or about 1/4 of the assembled genome ("Large100k" data set). The second and third subsets comprised zero-fold and four-fold degenerate sites in codons of all protein coding genes ("0-Fold" and "4-Fold" data sets), which were extracted using a custom PERL script (adapted from "identify_4D_sites.pl", see https://github.com/tsackton/linked-selection/blob/master/misc_scripts/Identify_4D_Sites.pl). The SGP5 annotation v.1.3 was

262 used as input.

263 Genotypes were called from the 20 BAM files with bcftools (options: mpileup -q 20

264 -Q 20 -a 'FORMAT/DP'), followed by bcftools (Li et al. 2009) (options: call -m) to create

265 files with all reads. The FORMAT option allows downstream filtering by individual and

266 site depth. Each of these files were then hard filtered by bcftools (options: -filter –SnpGap

267 3 -i '%QUAL>20 && TYPE != "indel"'). Then VCF files were normalized with bcftools

268 (options: norm -d all) and checked again for duplicates with a custom PERL script. Further

269 filtering by coverage depth and missing genotype calls was effected with a custom PERL

270 script, VCFParse.pl. In general, a genotype call was required to have coverage of between

271 2-100x to be considered non-missing. Analyses in the pooled sample of 20 were required to

272 have 13 of 20 genotype calls present at a site. The individual population samples were each

273 required to have two of two genotypes present to keep a site. These were obtained by using

274 bcftools (option: view -S) to subsample the VCF file for the entire data set.

275 *Genetic diversity.*— Genetic diversity estimators were obtained using two

276 protocols. First, we used the genotype calls derived from the bcftools pipeline described

277 above and calculated Tajima's ($\pi_T$) and Watterson's nucleotide diversity estimator ($\pi_W$)

278 with our VCFParse.pl script. Watterson's estimator was corrected for missing genotypes

279 using the formula in Ferretti et al. (2012). Missing individuals at a site were omitted in the

280 $\pi_T$ calculation.

281 Second, we used the Angsd package (Korneliussen et al. 2014) on only the

282 Large100k data to estimate these parameters directly from the BAM files. This approach

283 sidesteps genotype calling and works directly at the level of genotype likelihoods. Angsd

284 was run with settings chosen to match those used the bcftools pipeline as closely as

285 possible, including hard filters to require "minQ 20 minMapQ 1 minIndDepth 2" and

286 "minInd" set appropriately. Because ancestral states of SNPs were unavailable, folded site

287 frequency spectra were used. BAM files for the species-wide 10 population sample

288 consisted of the original 20 BAM files restricted to the appropriate genomic regions, but

the 2-individual analyses for each population separately were done using just the two BAM

files for that population. This differed from the genotype calling strategy using bcftools, in

which the smaller samples were obtained by subsetting the genotype calls in the large

sample.

The spatial distribution of genomic diversity was interpolated using point kriging

with no drifts with a linear variogram and no nugget effect as recommended for small,

noisy samples. The grid was generated using Surfer (Golden Software, Golden, Colorado,

USA) overlaid on a shaded relief base map derived from GMTED2010

(https://topotools.cr.usgs.gov/gmted_viewer/viewer.htm).

### *Phylogenetic Analysis*

*Complete data sets.—* "Complete" phylogenetic data matrices for all 20

individuals were prepared for each of the three genomic subsamples, and then these were

thinned to account for statistical nonindependence between sites (see below). All sites in

which genotype calls from the bcftools pipeline displayed exactly two alleles among the

individuals in the sample were included in phylogenetic data sets. Data sets were

formatted for use in Phylip (Felsenstein 2005) and PAUP* (Swofford 2002).

*Phylogenetic non-independence across the genome.—* The equivalence of the PP

and MDC-G score summed across multiple sites holds when sites are independent, so that

each site can theoretically evolve on its own allele tree. If a set of nearby sites are linked as

haplotypes evolving on the same allele tree, then the two scores could differ. Complete

data sets should therefore be "thinned" to sites that are approximately independent of one

another (Lee et al. 2014). We estimated the rate of decay of statistical dependence between

sites using two procedures. In the first, we used PLINK v. 1.9 (Chang et al. 2015) to

extract sites in approximate linkage equilibrium using a stringent requirement of 1.0 for

the VIF parameter (option: –indep 50 5 1).

314    In the second procedure, we implemented a direct phylogenetic assay based on

315 character compatibility (Felsenstein 2004) and a generalization of Hudson and Kaplan's

316 (1985) "four-gamete test" for a pair of unphased diploid genotypes. A pair of sites for $n$

317 individuals is *pairwise genotype compatible* if and only if there exist $2n$ haplotypes evolving

318 on some "perfect phylogeny" (one with no homoplasy). This is what is expected under the

319 infinite sites model. Wang (2013) used this concept to identify blocks of compatible sites,

320 but did not explicitly describe it; for completeness, we show the computation in the

321 Appendix. For a sample of 10000 regularly spaced sites in each of the three data sets, we

322 computed the fraction of sites at lag distance of $\lambda$ sites downstream ($1 \leqslant \lambda \leqslant 100$) that are

323 pairwise genotype compatible. This provides an estimate of autocorrelation as a function of

324 coordinate distance. We then assigned a *thinning distance*, $k_{\text{thin}}$, by inspection from these

325 plots, defined as the minimum distance at which the pairwise compatibility fraction

326 decreases to its genome-wide level.


*Phylogenetic data set and tree construction.—*    For a complete data set, $D$, of

length $m$ sites, we built $k_{\text{thin}}$ thinned data sets, $D_i$, of sites spaced evenly at $k_{\text{thin}}$ intervals:

$$D_i = \{s_j : j = i, i + k_{\text{thin}}, i + 2k_{\text{thin}}, ...\}, j \leqslant m, 1 \leqslant i \leqslant k_{\text{thin}}$$

327 For each $D_i$, we estimated a rooted genotype tree, $\hat{\Psi}_i$, using the PP tree search

328 implemented in dollop in Phylip 3.695 (2013) with settings of 'Polymorphism parsimony',

329 'Ancestral' states set to '?', and 100 replicated random addition sequences using Phylip's

330 'Jumble' option. Prior to search, the taxon order in each matrix was also randomized by a

331 PERL script. A bug in the program's output of the total PP score was fixed (Felsenstein,

332 pers. comm.). Because dollop is not optimized for searching for the best *rooting* per se

333 (Felsenstein, pers. comm.), the optimal tree reported by dollop was rerooted in all possible

334 ways using a PERL wrapper script around the Newick Utilities program nw_reroot (Junier

335 and Zdobnov 2010), and each rerooted tree was scored with dollop as described above. The

336 best rooted tree was retained.

For a single thinned data set, such as those generated by running PLINK on each of our complete data sets, we estimated standard bootstrap support using Phylip's seqboot in conjunction with dollop as described. However, for the data sets constructed from the compatibility-based thinning, which produced *multiple* data sets subsampled from each complete data set, we constructed an overall estimate of the tree by computing a (rooted) majority rule consensus of these:

$$\hat{\Psi}_{MR}(\hat{\Psi}_1, \hat{\Psi}_2, ..., \hat{\Psi}_{k_{\text{thin}}})$$

However, trees in $\{\hat{\Psi}_1, \hat{\Psi}_2, ..., \hat{\Psi}_{k_{\text{thin}}}\}$ may not be independent of each other, because the sites used to construct, say, $\hat{\Psi}_j$ are not independent of the sites used to construct $\hat{\Psi}_{j+1}$. In fact, the sites at position $i$ in these two data sets are only one character apart in the original complete matrix. We computed a measure of autocorrelation of the estimated trees as a function of distance between sites in $D$ ("lag") using pairwise tree-to-tree rooted Robinson-Foulds distances. This was only sensible for the Large100k data set, which has a sizable maximum lag distance of $k_{\text{thin}} = 100$. For each lag distance $\lambda : 1 \leqslant \lambda \leqslant k_{\text{thin}}/2$, we computed $L(\lambda) = \sum_i^{k_{\text{thin}}/2} d_{RF}(\hat{\Psi}_i, \hat{\Psi}_{i+\lambda})/(k_{\text{thin}}/2)$, where $d_{RF}(\Psi_i, \Psi_j)$ is the rooted Robinson-Foulds distance between trees $\Psi_i$ and $\Psi_j$.

Because PP trees consistently estimated the root to be near population samples that had higher missing data than average, we performed a numerical experiment to check for the influence of low coverage individuals on the method. For each of the three complete data sets, we downsampled sites by preferentially keeping sites that were present in the three individuals with the most missing data, until the overall percentage of missing data in those three samples was decreased below 5-8%, which was equivalent to levels found in several other samples. Resulting data sets were about 1/3 the size of the originals. We then reran the thinning and PP searches as described above, but changed the $k_{\text{thin}}$ values to maintain independence but correct for the smaller sizes of these matrices.

Finally, we compared our results with trees obtained with SVDquartets (Chifman and Kubatko 2014) in PAUP*, a scalable phylogenomic inference method. To force

357 SVDquartets to treat individuals as genotypes, each individual was recoded as two

358 sequences, randomly phasing all heterozygous positions, and then each pair was forced to

359 be treated as a "species" under the multi-species coalesent assumption of the algorithm.

360 Optimal trees from replicate thinned data sets were combined as described above for PP

361 analyses, with random phasing done separately for each thinned data set.

362 RESULTS

363 *Sequence Data and Genetic Diversity*

364 We obtained a total of 198 Gbp of short read sequence for 20 individuals in 10

365 populations (Table 1). Average coverage ranged from 2.2x to 17.5x (mean = 7.6x). The

366 San Marcial coverage was the lowest by population.

367 The bcftools pipeline inferred from just under 50,000 to nearly 3 million variants in

368 the three data sets (Table 2). Almost all were biallelic: the fraction of variants having 3 or

369 4 alleles ranged from only 0.0035 - 0.0045.

370 Table 3 summarizes all nucleotide diversity estimates obtained with the bcftools

371 pipeline, subsets of the genome and different populations. An overall estimate of neutral

372 genetic diversity based on 4-fold degenerate sites in protein coding genes is 0.0025, which is

373 quite low compared to many plant species. The ratio of 0-fold to 4-fold diversity is quite

374 high, around 50%, indicating a high proportion of nonsynonymous diversity. Tajima's and

375 Watterson's estimators of nucleotide diversity are in broad agreement, as are the estimates

376 obtained from the bcftools pipeline vs. the Angsd pipeline, which uses genotype likelihoods

377 directly (Table 4).

378 A geographic perspective on genetic diversity can be seen in Figure 3. The highest

379 diversity is around Kino Bay in Sonora (Table 3), but it drops off relatively quickly to the

380 east and south, and more slowly to the north. Masiaca and Wickenburg, with generally the

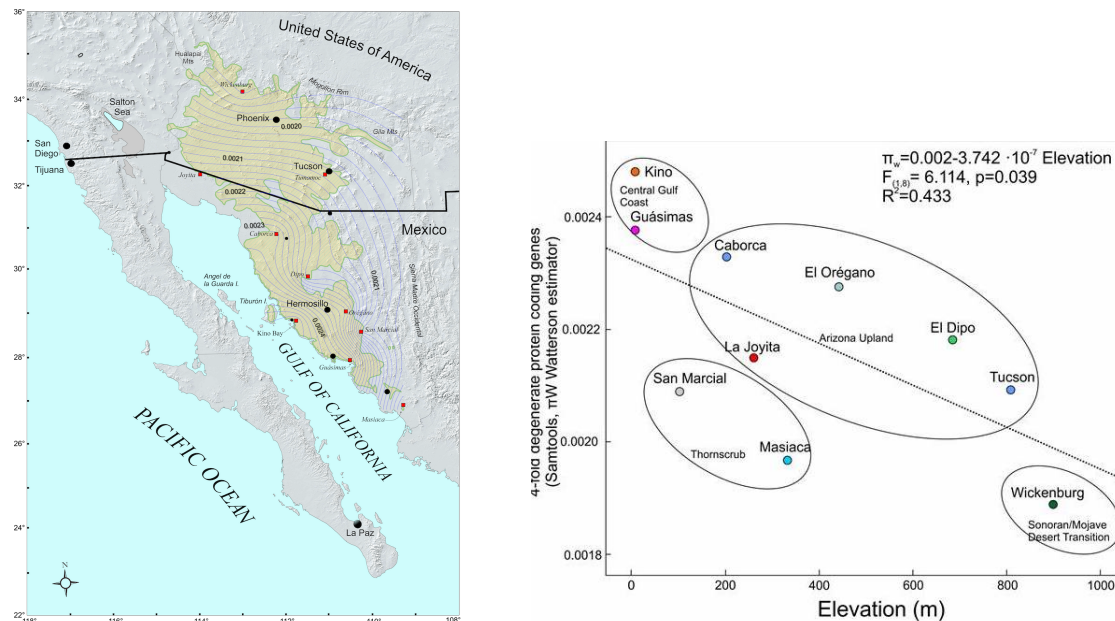381 lowest diversities, are populations near the southern and northern distribution margins

Fig. 3. a) Contour plot of estimated nucleotide diversity ($\pi_W$) in 4-fold degenerate sites in 10 populations of saguaro. Range of saguaro shown in yellow. b) Relationship between elevation and estimated nucleotide diversity ($\pi_W$).

respectively. The spatial distribution is similar for both estimators of nucleotide diversity. Using the Watterson estimator for the 4-fold degenerate sites in protein coding, the isoline of 0.0021 delimits well what we consider saguaro populations in good health and number. Although the interpolated surface is only derived from the genomic data, there is a significant negative linear trend with elevation (Fig. 3), and with the vegetation units described by Shreve (1951), displaying the highest genomic diversity in the Central Gulf Coast, intermediate values in the Arizona Upland and the thornscrub, and the lowest diversity in the transition to the Mojave Desert.

## *Phylogenetic data sets*

The three complete phylogenetic data sets ranged in size from 47511 and 106629 sites for 4-fold and 0-fold degenerate sites to nearly 3 million in the Large100k data set (Table 2). The fraction of genotypes called as heterozygous ranged from 17-20%, whereas the fraction of sites with missing data was 8-9%. Missing data was distributed unevenly

<sup>395</sup> among samples, with the two San Marcial samples having 40-54% missing genotypes and

<sup>396</sup> the Masiaca_M85 individual having 20%.

<sup>397</sup>      Analyses of pairwise genotype compatibility in the three complete data sets

<sup>398</sup> revealed a strong signal of local statistical dependence along the genome coordinate axis,

<sup>399</sup> presumably due to linkage (Fig. 4). The 4-fold data set lost dependence most quickly,

<sup>400</sup> within a $k_{\text{thin}}$ of 10 sites; the 0-fold data was about 20 sites; and the Large100k data about

<sup>401</sup> 100 sites. Note that these distances are in units of sites in the data matrices: two

<sup>402</sup> neighboring sites in the phylogenetic data set might be separated by a long coordinate

<sup>403</sup> distance along the scaffold because of intervening nonvariable sites.
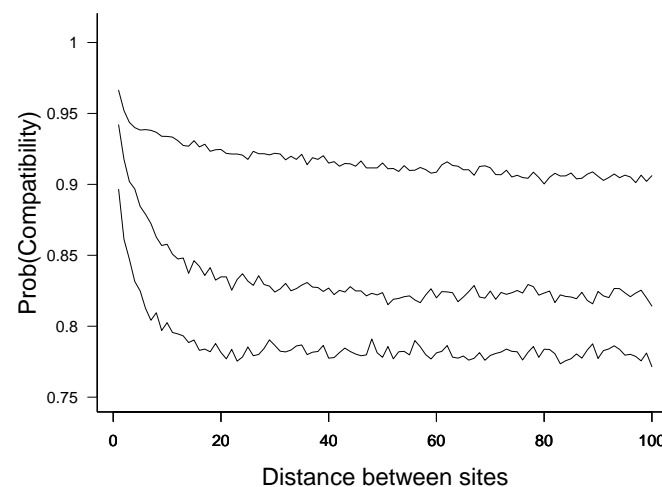


Fig. 4. Decay of pairwise genotype compatibility versus distance between sites in data sets (from top to bottom: Large100k; 0-fold; 4-fold).

<sup>404</sup>      Thinning by PLINK's algorithm for identifying sites in approximate linkage

<sup>405</sup> equilibrium produced three data sets having 27457, 13977, and 88480 sites for 0-fold,

<sup>406</sup> 4-fold, and Large100k data sets, respectively (Table 2). These all correspond to average

<sup>407</sup> thin distances that are smaller than found by our pairwise compatibility method and thus

<sup>408</sup> produced data sets with more sites each.

### Intraspecific Genotype Phylogenies

409

410        Trees inferred using PP and compatibility thinning in the three data sets were

411  highly congruent (Fig. 5). A northern clade of four populations, Caborca, La Joyita,

412  Tucson and Wickenburg, was consistently strongly supported, though La Joyita and

413  Caborca within this clade showed slightly lowered support for the monophyly of individual

414  populations. The remaining populations were all highly supported as monophyletic except

415  at the root. Rerooting of the dollop solution with our rerooting script often improved the

416  optimality score, but the root was consistently still placed within the San Marcial

417  population. The Kino Bay population, which has the highest nucleotide diversity and lies

418  at sea level on the Sea of Cortez to the west of San Marcial was nested further within the

419  genotype phylogeny. Trees inferred from the Large100k data differed slightly with respect

420  to Kino Bay, having it as sib group to the rest of the tree, except for San Marcial.
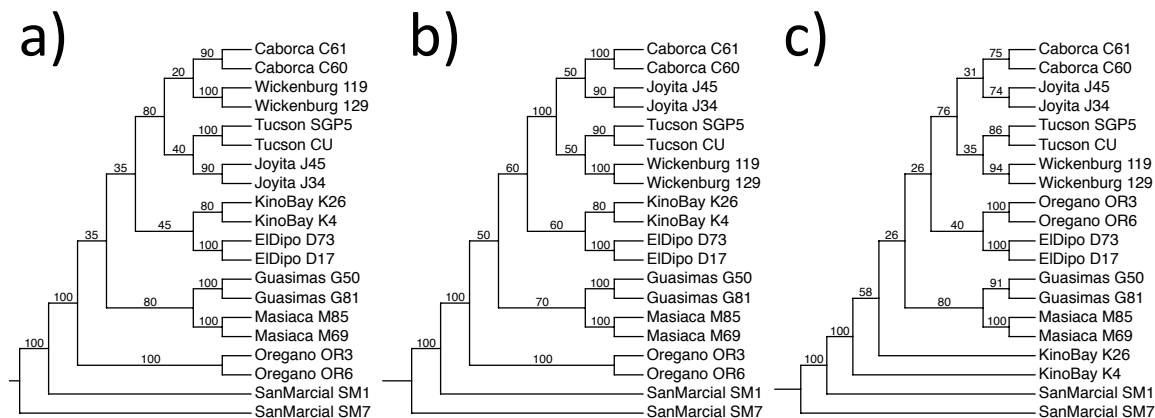


Fig. 5. Genotype trees reconstructed with PP in Phylip. These are majority rule trees of sets of $k_{\text{thin}}$ compatibility-thinned data sets for a) biallelic 0-fold sites ($k_{\text{thin}} = 20$); b) biallelic 4-fold sites ($k_{\text{thin}} = 10$); c) Large100k sites ($k_{\text{thin}} = 100$); Proportions on edges reflect the number of trees containing indicated bipartitions. Trees are rooted by the PP method.

421        Figure 6 assesses the autocorrelation between trees inferred by PP from the

422  compatibility-thinned data sets for the Large100k data. The Robinson-Foulds distance

423  between trees estimated at different lag distances appears to show no relationship to

424  relative genome coordinate distance, implying that there is no autocorrelation among the

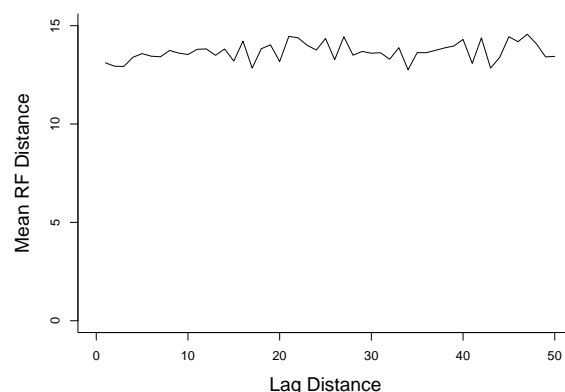425   PP trees inferred from the thinned data sets.



Fig. 6. Mean RF distances between trees estimated from $k_{\text{thin}} = 100$ thinned data sets subsampled from Large100k data set, $\{\hat{\Psi}_1, \hat{\Psi}_2, ..., \hat{\Psi}_{k_{\text{thin}}}\}$, at different lag distances. Lag distance, $\lambda$, refers to a pair of trees, $\{\hat{\Psi}_i, \hat{\Psi}_{i+\lambda}\}$.

426   Trees inferred with PP using the data sets thinned via PLINK (Fig. 7) all agreed

427   with the 0-fold and 4-fold trees from PP and compatibility thinning, and also were rooted

428   in San Marcial.



Fig. 7. Genotype trees reconstructed with PP in Phylip. These are bootstrap majority rule trees from the three data sets thinned to approximate linkage equilibrium in PLINK for a) biallelic 0-fold sites; b) biallelic 4-fold sites; c) Large100k sites. Trees are rooted by the PP method.

429   Downsampling the three data sets to assess the impact of missing sites resulted in

430   smaller data sets (Table 2), but with much more even distributions of missing data (Table

431   5). Because of their smaller size, the 0-fold and 4-fold downsamples were thinned to fewer

432   data sets: 6 and 3 data sets respectively, which maintains approximately the same thinned

₄₃₃ data set size as in the non-downsampled data. The Large100k data was thinned into 100

₄₃₄ data sets (a $k_{thin}$ value that should produce sites thinned well beyond the $k_{thin}$ distance in

₄₃₅ the original Large100k data set). The 0-fold and 4-fold majority rule trees (Fig. 8) were

₄₃₆ identical to those in Fig. 5. However, the Large100k analysis departed from this otherwise

₄₃₇ consistent picture, albeit with only moderate bootstrap values. It was rooted in the

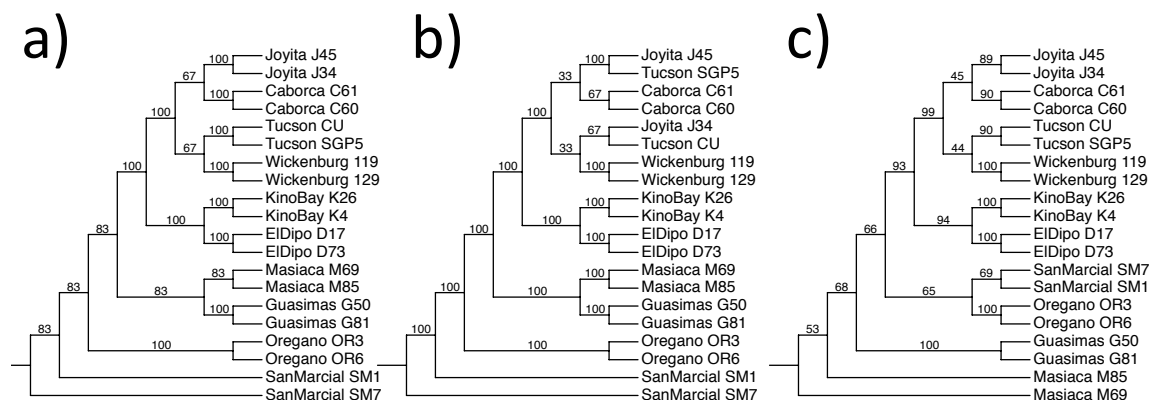₄₃₈ Masiaca population at the southern limit of the saguaro range.



Fig. 8. PP genotype trees based on downsampled data sets (see Table 5). Trees are majority rule trees of sets of $k_{thin}$ thinned data sets for a) biallelic 0-fold sites ($k_{thin} = 6$); b) biallelic 4-fold sites ($k_{thin} = 3$); c) Large100k sites ($k_{thin} = 100$); Proportions on edges reflect the number of trees containing indicated bipartitions. Trees are rooted by the PP method.

₄₃₉ Comparable trees constructed by SVDquartets are shown in Figure 9 for the two

₄₄₀ gene data sets. The unrooted topology is highly congruent between 0-fold and 4-fold data

₄₄₁ sets; all pairs of individuals within populations are supported as clades; and there remains

₄₄₂ some variability in the relationship of the northernmost populations. One minor difference

₄₄₃ with the PP trees is that SVDquartets moves the San Marcial population by one nearest

₄₄₄ neighbor interchange so it is slightly further from Orégano on the tree..

₄₄₅                                    DISCUSSION

₄₄₆ Inger (1967) published a phylogeny of frogs based on musculo-skeletal characters

₄₄₇ using code written by Felsenstein, implementing polymorphism parsimony for the first

₄₄₈ time (Felsenstein 1979). Felsenstein (1979) and Farris (1978) later formalized the problem
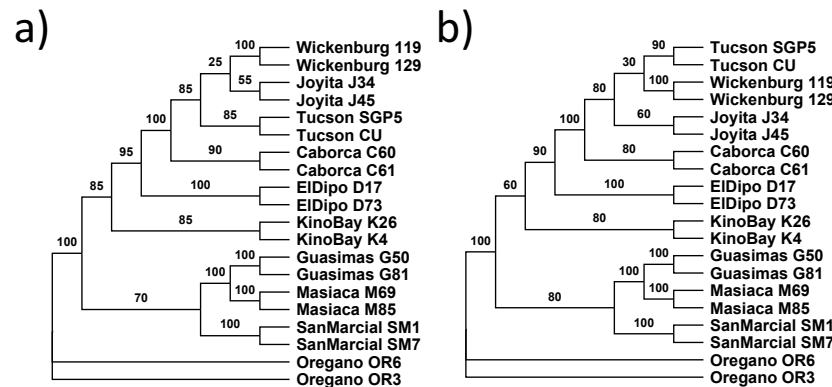
Fig. 9. Genotype phylogenies based on SVDquartets thinned by compatibility thinning as in PP analyses. Phylogenies are based on a) 0-fold, and b) 4-fold data sets. Proportions on edges reflect the number of trees among the thinned runs containing indicated bipartitions. Trees are unrooted.

⁴⁴⁹ statements and algorithms for assigning a PP score to a tree, both citing for inspiration

⁴⁵⁰ work on *Drosophila* inversion polymorphisms. The algorithm presented by Felsenstein is

⁴⁵¹ similar in spirit and runtime to other parsimony scoring methods and shares with them the

⁴⁵² potential to scale to large data sets, at least with heuristic tree search algorithms (though,

⁴⁵³ like MP and ML, PP is NP-complete: Day and Sankoff 1987). Despite the prevalence of

⁴⁵⁴ polymorphic character data in the real world, and the obvious connection to diploid

⁴⁵⁵ genotype data, PP has been used sparingly: e.g., for plant morphological data (Baum

⁴⁵⁶ 1983), and bird retrotransposons (Suh et al. 2015).

⁴⁵⁷ In this paper we showed that polymorphism parsimony can be placed in the

⁴⁵⁸ framework of minimizing deep coalescences in the genotypes' underlying allele tree. This is

⁴⁵⁹ intuitive on the one hand—the language of minimizing the extent of polymorphism on a

⁴⁶⁰ tree is similar to minimizing deep coalescences—but also somewhat surprising. One

⁴⁶¹ method takes only discrete character data as input, while the other seems much more

⁴⁶² directly related to a reconciliation between an allele tree and genotype tree. The criterion

⁴⁶³ of minimizing the number of polymorphic edges in the genotype tree is equivalent to

⁴⁶⁴ minimizing deep coalescences in the "best" allele tree that can be constructed from the

⁴⁶⁵ genotype data. An equivalent interpretation of these optimality criteria is that they

⁴⁶⁶ minimize the number of "fixation/loss" events (Farris 1978), that is, losses of one or the

467    other allele from a polymorphic ancestor.

468        Though intuitive, this equivalence is not automatic. It holds if the allele tree evolves

469    according to an infinite sites model. In this setting the equivalence of the two methods

470    places PP from the late 1970s into a rich framework of gene tree reconciliation, which dates

471    to about the same time (Goodman et al. 1979) but has led to a very large computational,

472    algorithmic and empirical literature (Page 1994; Maddison 1997; Nakhleh 2013).

473        One particularly attractive feature of polymorphism parsimony is that it provides

474    an estimate of the root, because its underlying asymmetrical character evolution model

475    penalizes rootings with deeper polymorphisms. For example, the tree,

476    $(0/0, (0/0, (0/1, 0/1)))$ has a PP score of two, but the rerooted tree, $(0/1, (0/1, (0/0, 0/0)))$

477    has a score of three. In the saguaro data this provided some evidence to dispute the

478    hypothesis that the population with highest nucleotide diversity in saguaro is also at the

479    root of the tree. PP shares this feature with some long standing methods, such as midpoint

480    rooting, UPGMA and likelihood inference using a molecular clock (Huelsenbeck et al.

481    2002; Felsenstein 2004), and some more recent methods of species tree inference, such as

482    BPP (Rannala and Yang 2017).

483        *Assumptions.—*  The PP method applied to genome-wide genotype data is likely

484    to be most appropriate when individual sampling is sparse relative to population

485    differentiation and isolation. For example, one individual sampled per population within a

486    species in which population differentiation is significant may balance the requirement of

487    limited gene flow, so that an individual is a stand-in for a population, with limited

488    sequence divergence, so that the infinite sites assumption is met. As sampling within

489    panmictic populations increases, we expect more unresolved nodes within populations on

490    the genotype tree, and lower confidence levels.

491        In the saguaro data we sampled two individuals per species, which allowed tests of

492    monophyly of the 10 populations. Had their been little phylogenetic structure between

493    populations, these support values should have been low, but most populations were

supported at 90% levels or above for a variety of data sets and treatments. The main exception was the rooting of the trees within the San Marcial population, which causes this population to be paraphyletic in several analyses.

The algorithm for minimizing deep coalescence forces homozygous alleles in an individual to coalesce at the leaves of the genotype tree. This is what a simple minimization criterion based on the DC score should do, and it is engineered by the reconcilation mapping function (Yu et al. 2011, Thm 6), but it is unlikely to be accurate for individuals that are closely related within a single population. Under the coalescent model every labeled history on the $2n$ alleles is equally probable (Xu and Yang 2016), and the chance of alleles of a genotype character pairing up as siblings within each individual on the allele tree is small. The parsimony rationale is to assume shallow coalescence unless there is evidence of deep coalescence based on the genotype data, but a prior based on the Kingman coalescent and random mating would disagree. This naturally places a premium on answering the basic question of how much the sampled individuals depart from random mating.

In general, the validity of an infinite sites model for our data is difficult to test, because homoplasy in the data can be caused by either multiple hits or discordant gene trees (hemiplasy: Avise and Robinson 2008) or both. The fraction of tri- or quadri-allelic variants was less than 0.5% in our three data sets, which implies that the rate of multiple hits is low.

*Thinning.*— Thinning SNP data to a subset of sites that are in linkage equilibrium and phylogenetically independent is recommended for a number of methods of SNP-based phylogenetic inference, including PP, SVDquartets and variants (Chifman and Kubatko 2014; Vachaspati and Warnow 2018), and SNAPP (Bryant et al. 2012). We used a method implemented in PLINK based on the correlation between genotype vectors at pairs of sites, and a new method based on pairwise phylogenetic compatibility, to thin data. Even though we used a high stringency in the PLINK analysis, it tended to retain

521    more sites than the phylogenetic assay. The resulting genotype trees on these larger data

522    sets were highly congruent with our results from compatibility-thinned data, and both sets

523    of results agreed about the root of the saguaro genotype phylogeny.

524         These results raise several questions for future work. First, as SNP data sets for

525    phylogenetic inference grow more dense, the size of thinned data sets will only increase to

526    a certain point, and the question of how to combine these data sets will have to be

527    addressed. Here we showed that there was little correlation between the trees constructed

528    from the different thinned data sets, and we proposed a simple consensus tree to serve as a

529    summary statistic, but it remains somewhat surprising that these trees are uncorrelated

530    when there is strong evidence that the underlying site data are correlated locally in the

531    genome. Second, once data sets are thinned enough, there remains the question of how

532    many thinned data sets to construct: fewer with more sites, or more with fewer sites.

533    Finally, there remains an open problem of how to identify true coalescent independent sites

534    rigorously in the face of recombination, sequencing error, and homoplasy.

535         *Haplotype based inference.—*   An alternative to thinning is to infer haplotype

536    blocks from all the SNPs (reviewed in Browning and Browning 2011; Gusfield 2014). In a

537    local block of nonrecombining sites, all SNPs should evolve on the same haplotype tree,

538    which motivates the Perfect Phylogeny Haplotype Problem (Gusfield 2002). For short

539    tracts and under the infinite sites assumption, all SNPs should be compatible with one

540    haplotype tree, or nearly so. In practice, sequencing error and homoplasy can limit the

541    applicability of this approach. More generally, haplotype trees that are not perfect

542    phylogenies can be sought (Sridhar et al. 2007), but the search space gets quite large. Not

543    only is there the usual search space across trees, but there is the additional exponential

544    growth in the number of alternative haplotypes for a given set of heterozygous SNPs. An

545    additional empirical problem evident in our saguaro data is that relatively low genetic

546    variation can mean that there is a long coordinate distance between SNPs, which leads to a

547    tradeoff between combining enough SNPs to build a reliable haplotype tree before reaching

548  recombination distances that allow independent coalescent genealogies (Springer and

549  Gatesy 2016). A general problem is to find a set of intervals that either partition or cover

550  the SNP sites within which there is a perfect haplotype tree, but between which the trees

551  may differ (Gramm et al. 2009; Wang 2013). Given a collection of such haplotype trees

552  across genomes for multiple individuals, we would still have the problem of how to integrate

553  them into a genotype tree. The classical MDC problem, using individuals as leaves of the

554  containing tree and a set of independent haplotype trees, would be one potential solution.

555      *Relationship to other methods.—*   Because the inner loop of tree search heuristics

556  for PP involves a linear time PP scoring algorithm (Felsenstein 1979), runtimes for our data

557  were quite managable even for the largest data sets examined here, some of which (after

558  thinning in various ways) had nearly 100,000 sites. No runs took more than about one hour

559  of CPU time on an HPC linux node, and the bootstrap and thinning replicates could be

560  trivially distributed on the cluster to make the analyses described here quite tractable.

561      Other computationally fast methods for inferring a population tree that scale to

562  data sets of the size considered here include TreeMix (Pickrell and Pritchard 2012), which

563  constructs a covariance matrix based on gene frequencies from SNPs and infers a tree on

564  which migration events can then be overlain; and SVDquartets (Chifman and Kubatko

565  2014; Vachaspati and Warnow 2018), which uses a non-parametric statistic based on the

566  multi-species coalescent for quartets and then combines the quartets to build a full tree.

567  The former assumes a simple allele frequency diffusion model and the latter a standard

568  DNA substitution model of the kind used in molecular phylogenetics, together with the

569  MSC. Both methods infer an unrooted tree. TreeMix can summarize genotype data as

570  allele frequencies within populations, but SVDquartets needs sequences to be entered for

571  each allele and then populations to be circumscribed. RevPoMo (Schrempf et al. 2016,

572  2019) is analogous to an allele frequency approach but models transitions between allele

573  count states instead. It returns an unrooted tree though a predecessor (De Maio et al.

574  2015) returned a rooted tree. Its authors "discourage ... using revPoMo on sequence data

575  where no population data is available yet" (Schrempf et al. 2016, p. 369) because of the

576  need for demographic parameters in the model.

577      Perhaps the most computationally ambitious current methods, BPP (Rannala and

578  Yang 2017), IMA-3 (Hey et al. 2018) and *Beast2 (Ogilvie et al. 2017), use a model-based

579  Bayesian approach to integrate over population genetics, demography and phylogenetics.

580  However, these appear to come at the cost of perhaps two or more orders of magnitude

581  slower running times (Hey et al. 2018; Ogilvie et al. 2017). These methods can return a

582  rooted tree by assuming a clock or the infinite sites assumption itself. Polymorphism

583  parsimony provides a rapid alternative for estimating a reasonable genotype or population

584  tree first, which can then be used in a two-step procedure to make more nuanced inferences

585  about demography, admixture and migration.

586      *Saguaro nucleotide diversity and phylogeography.*— Estimates of nucleotide

587  diversity for saguaro as a whole were near 0.0025, which is low among plants even

588  compared to some other long lived trees such as spruce, which has about 2.5 times the

589  diversity of saguaro (Chen et al. 2019). The efficacy of selection (Chen et al. 2017), defined

590  as the ratio of 0-fold to 4-fold degenerate site diversities, was nearly 0.49, which is also

591  unusually high among plants (c.f., 0.44 in one spruce species, Chen et al. 2019). This

592  indicates much higher relative genetic diversity for potentially fitness-related variants than

593  seen in other plants. These results may stem from low effective population sizes and/or

594  historical population crashes during glacial/interglacial periods. Irrespective of cause, the

595  0-fold phylogenetic data set we analyzed here had twice as many sites as the 4-fold

596  degenerate data set and evidently as much or more information about genotype phylogeny.

597      Genetic diversity was highest at sea level at Kino Bay, Sonora, Mexico, a bit south

598  of the approximate center of the species' range. Diversity dropped off quickly to the south

599  and east and more slowly to the north. The range of saguaro has been heavily influenced

600  by the ebb and flow of North American glaciation. Evidence from packrat midden remains

601  indicates, for example, that saguaros were absent in the US until recent reinvasions about

602  10,000 years ago (McAuliffe and Van Devender 1998). Given that the species split from its

603  closest relatives several million years ago (Copetti et al. 2017), it seems likely that glacial

604  refugia must have existed somewhere in the south of its current range or perhaps further

605  south.

606       Our PP analyses places the root of the saguaro genotype tree to the east or possibly

607  south of Kino Bay in either the San Marcial or Masiaca populations, which are at the

608  margins of the current range of saguaro. The latter two populations were favored

609  differentially by analyses that were thinned in different ways, suggesting that such thinning

610  choices can be significant. None of the analyses directly supported a rooting within the

611  Kino Bay population. San Marcial and Masiaca are in thornscrub vegetation and both

612  have substantially lower genetic diversity than Kino Bay.

613       Few simple genetic patterns are evident in the geographic distribution of saguaro

614  and other columnar cacti (Bustamante et al. 2016). This is particularly true for columnar

615  cacti in their northwestern range of distribution where frost limits these sensitive plants. In

616  the case of the saguaro, its range is constrained by elevation (towards the Sierra Madre

617  Occidental to the east and the Mogollon Rim to the north), westward by the Gulf of

618  California, and most likely by the very low precipitation of the Gran Desierto where

619  Sonora, California and Arizona meet (Albuquerque et al. 2018).

620       The pattern of higher diversity in the center of saguaro's range around Kino Bay

621  does generally support the so-called "center-periphery hypothesis" (Pironon et al. 2016)

622  but likely this is modulated by dynamic changes in the very recent past (Lázaro-Nogal

623  et al. 2017) and details of the demography and evolution of these populations. The

624  position of the southern populations of San Marcial and Masiaca as outgroups to the

625  remaining populations hints at possible crashes during the interglacial aggravated by

626  ongoing global climate change. These factors likely led to differentiation and genomic

627  diversity reduction through small effective population size and restricted local gene flow.

628       Albuquerque et al. (2018) found that the saguaro distribution is contracting. They

629   estimated a mean loss of almost 7% by 2050 under different climate change scenarios. That

630   contraction is occurring on the western edge of the range from western Arizona to the

631   southernmost extent in Mexico. Our rooted trees largely imply a basal grade of small,

632   isolated southern and southeastern populations giving rise to several distinct large,

633   thriving populations to the north and the northwest, which reinforces the idea of a refuge

634   near the southern boundaries of the actual distribution range at the lowest elevations. The

635   actual distribution of extant nucleotide diversity suggests a Pleistocene refuge in the

636   lowlands of the Gulf of California coast in southern Sonora, and the small range of

637   diversity estimates may indicate a rapid expansion northwards during the interglacial. If

638   we exclude the outlier populations of San Marcial and Masiaca and concentrate on more

639   central, numerically abundant, and/or the expanding, continuous populations in the

640   northern range, there is a hint of a latitudinal decrease in nucleotide diversity. Further

641   analyses on correlations between SNP variation and climatic, geographic and ecological

642   variables are necessary to understand these issues more fully.

## Acknowledgements

## References

647   Albuquerque, F., B. Benito, M. Á. M. Rodriguez, and C. Gray. 2018. Potential changes in

648       the distribution of *Carnegiea gigantea* under future scenarios. PeerJ 6:e5623.

649   Andrews, S. 2018. FastQC. Dowloaded from

650       https://www.bioinformatics.babraham.ac.uk/projects/fastqc/.

651   Avise, J. C. and T. J. Robinson. 2008. Hemiplasy: A new term in the lexicon of

652       phylogenetics. Systematic Biology 57:503–507.

653  Bansal, M. S., J. G. Burleigh, and O. Eulenstein. 2010. Efficient genome-scale phylogenetic

654      analysis under the duplication-loss and deep coalescence cost models. BMC

655      Bioinformatics 11:S42.

656  Baum, B. R. 1983. A phylogenetic analysis of the tribe Triticeae (Poaceae) based on

657      morphological characters of the genera. Canadian Journal of Botany-Revue Canadienne

658      De Botanique 61:518–535.

659  Bayzid, M. S. and T. Warnow. 2012. Estimating optimal species trees from incomplete

660      gene trees under deep coalescence. Journal of Computational Biology 19:591–605.

661  Bayzid, M. S. and T. Warnow. 2018. Gene tree parsimony for incomplete gene trees:

662      addressing true biological loss. Algorithms for Molecular Biology 13:1–12.

663  Bennett, M. and I. Leitch. 2012. Plant dna c-values database (release 6.0, dec. 2012).

664  Bolger, A. M., M. Lohse, and B. Usadel. 2014. Trimmomatic: a flexible trimmer for

665      Illumina sequence data. Bioinformatics 30:2114–2120.

666  Bravo, G. A., A. Antonelli, C. D. Bacon, K. Bartoszek, M. P. K. Blom, S. Huynh,

667      G. Jones, L. L. Knowles, S. Lamichhaney, T. Marcussen, H. Morlon, L. K. Nakhleh,

668      B. Oxelman, B. Pfeil, A. Schliep, N. Wahlberg, F. P. Werneck, J. Wiedenhoeft,

669      S. Willows-Munro, and S. V. Edwards. 2019. Embracing heterogeneity: coalescing the

670      Tree of Life and the future of phylogenomics. PeerJ 7.

671  Browning, S. R. and B. L. Browning. 2011. Haplotype phasing: existing methods and new

672      developments. Nature Reviews Genetics 12:703–714.

673  Bryant, D., R. Bouckaert, J. Felsenstein, N. A. Rosenberg, and A. RoyChoudhury. 2012.

674      Inferring species trees directly from biallelic genetic markers: Bypassing gene trees in a

675      full coalescent analysis. Molecular Biology and Evolution 29:1917–1932.

676  Buckler, E. S. and T. P. Holtsford. 1996. *Zea* systematics: Ribosomal ITS evidence.

677      Molecular Biology and Evolution 13:612–622.

678  Bustamante, E., A. Búrquez, E. Scheinvar, and L. E. Eguiarte. 2016. Population genetic

679     structure of a widespread bat-pollinated columnar cactus. PLOS One 11:e0152329.

680  Chang, C. C., C. C. Chow, L. Tellier, S. Vattikuti, S. M. Purcell, and J. J. Lee. 2015.

681     Second-generation PLINK: rising to the challenge of larger and richer datasets.

682     GigaScience 4:7.

683  Chen, J., S. Glemin, and M. Lascoux. 2017. Genetic diversity and the efficacy of purifying

684     selection across plant and animal species. Molecular Biology and Evolution

685     34:1417–1428.

686  Chen, J., L. L. Li, P. Milesi, G. Jansson, M. Berlin, B. Karlsson, J. Aleksic, G. G.

687     Vendramin, and M. Lascoux. 2019. Genomic data provide new insights on the

688     demographic history and the extent of recent material transfers in Norway spruce.

689     Evolutionary Applications 12:1539–1551.

690  Chifman, J. and L. Kubatko. 2014. Quartet inference from SNP data under the coalescent

691     model. Bioinformatics 30:3317–3324.

692  Copetti, D., A. Búrquez, E. Bustamante, J. L. M. Charboneau, K. L. Childs, L. E.

693     Eguiarte, S. Lee, T. L. Liu, M. M. McMahon, N. K. Whiteman, R. A. Wing, M. F.

694     Wojciechowski, and M. J. Sanderson. 2017. Extensive gene tree discordance and

695     hemiplasy shaped the genomes of North American columnar cacti. Proc. Natl. Acad.

696     Sci., USA 114:12003–12008.

697  Day, W. H. E. and D. Sankoff. 1987. Computational-complexity of inferring phylogenies

698     from chromosome inversion data. Journal of Theoretical Biology 124:213–218.

699  De Maio, N., D. Schrempf, and C. Kosiol. 2015. PoMo: An allele frequency-based approach

700     for species tree estimation. Systematic Biology 64:1018–1031.

701  Degnan, J. H. and N. A. Rosenberg. 2009. Gene tree discordance, phylogenetic inference

702     and the multispecies coalescent. Trends in Ecology and Evolution 24:332–340.

703  Degnan, J. H. and L. A. Salter. 2005. Gene tree distributions under the coalescent process.

704    Evolution 59:24–37.

705  Durvasula, A., A. Fulgione, R. M. Gutaker, S. I. Alacakaptan, P. J. Flood, C. Neto,

706    T. Tsuchimatsu, H. A. Burbano, F. X. Pico, C. Alonso-Blanco, and A. M. Hancock.

707    2017. African genomes illuminate the early history and transition to selfing in

708    *Arabidopsis thaliana.* Proceedings of the National Academy of Sciences of the United

709    States of America 114:5213–5218.

710  Farris, J. S. 1978. Inferring phylogenetic trees from chromosome inversion data. Systematic

711    Zoology 27:275–284.

712  Felsenstein, J. 1979. Alternative methods of phylogenetic inference and their

713    interrelationship. Systematic Zoology 28:49–62.

714  Felsenstein, J. 2004. Inferring Phylogenies. Sinauer Press, Sunderland, MA.

715  Felsenstein, J. 2005. PHYLIP (Phylogeny Inference Package) version 3.6.

716  Ferretti, L., E. Raineri, and S. Ramos-Onsins. 2012. Neutrality tests for sequences with

717    missing data. Genetics 191:1397–1401.

718  Goodman, M., J. Czelusniak, G. W. Moore, A. E. Romeroherrera, and G. Matsuda. 1979.

719    Fitting the gene lineage into its species lineage, a parsimony strategy illustrated by

720    cladograms constructed from globin sequences. Systematic Zoology 28:132–163.

721  Gramm, J., T. Hartman, T. Nierhoff, R. Sharan, and T. Tantau. 2009. On the complexity

722    of SNP block partitioning under the perfect phylogeny model. Discrete Mathematics

723    309:5610–5617.

724  Gusfield, D. 2002. Haplotyping as perfect phylogeny: conceptual framework and efficient

725    solutions. Pages 166–175 *in* RECOMB '02: Proceedings of the Sixth Annual

726    International Conference on Computational biology.

727 Gusfield, D. 2014. ReCombinatorics: The Algorithmics of Ancestral Recombination Graphs

728    and Explicit Phylogenetic Networks. MIT Press, Cambridge, MA.

729 Hein, J., M. H. Schierup, and C. Wiuf. 2005. Gene Genealogies, Variation and Evolution:

730    A Primer in Coalescent Theory. Oxford University Press, USA.

731 Hey, J., Y. J. Chung, A. Sethuraman, J. Lachance, S. Tishkoff, V. C. Sousa, and Y. Wang.

732    2018. Phylogeny estimation by integration over isolation with migration models.

733    Molecular Biology and Evolution 35:2805–2818.

734 Hudson, R. R. and N. L. Kaplan. 1985. Statistical properties of the number of

735    recombination events in the history of a sample of DNA-sequences. Genetics

736    111:147–164.

737 Huelsenbeck, J. P., J. P. Bollback, and A. M. Levine. 2002. Inferring the root of a

738    phylogenetic tree. Systematic Biology 51:32–43.

739 Inger, R. F. 1967. Development of a phylogeny of frogs. Evolution 21:369–384.

740 Junier, T. and E. M. Zdobnov. 2010. The Newick utilities: high-throughput phylogenetic

741    tree processing in the Unix shell. Bioinformatics 26:1669–1670.

742 Kimura, M. 1969. Number of heterozygous nucleotide sites maintained in a finite

743    population due to steady flux of mutations. Genetics 61:893–903.

744 Kingman, J. F. C. 1982. On the genealogy of large populations. Journal of Applied

745    Probability 19:27– 43.

746 Korneliussen, T. S., A. Albrechtsen, and R. Nielsen. 2014. ANGSD: Analysis of next

747    generation sequencing data. BMC Bioinformatics 15:356.

748 Langmead, B. and S. L. Salzberg. 2012. Fast gapped-read alignment with Bowtie 2. Nature

749    Methods 9:357–359.

750   Lázaro-Nogal, A., S. Matesanz, A. García-Fernández, A. Traveset, and F. Valladares. 2017.
751        Population size, center-periphery, and seed dispersers' effects on the genetic diversity
752        and population structure of the mediterranean relict shrub *Cneorum tricoccon*. Ecol
753        Evol. 7:7231–7242.

754   Lee, T. H., H. Guo, X. Y. Wang, C. Kim, and A. H. Paterson. 2014. SNPhylo: a pipeline
755        to construct a phylogenetic tree from huge SNP data. BMC Genomics 15:162.

756   Li, H., B. Handsaker, A. Wysoker, T. Fennell, J. Ruan, N. Homer, G. Marth, G. Abecasis,
757        R. Durbin, and P. Genome Project Data. 2009. The sequence alignment/map format
758        and SAMtools. Bioinformatics 25:2078–2079.

759   Liu, L., C. Anderson, D. Pearl, and S. Edwards. 2019. Modern phylogenomics: Building
760        phylogenetic trees using the multispecies coalescent model. Methods in Molecular
761        Biology 1910:211–239.

762   Ma, B., M. Li, and L. Zhang. 2001. From gene trees to species trees. SIAM J. Comput.
763        30:729–752.

764   Maddison, W. P. 1997. Gene trees in species trees. Systematic Biology 46:523–536.

765   Maddison, W. P. and L. L. Knowles. 2006. Inferring phylogeny despite incomplete lineage
766        sorting. Systematic Biology 55:21–30.

767   Maddison, W. P. and D. R. Maddison. 2000. MacClade 4: Analysis of phylogeny and
768        character evolution. Sinauer, Sunderland, MA.

769   Nakhleh, L. 2013. Computational approaches to species phylogeny inference and gene tree
770        reconciliation. Trends in Ecology and Evolution 28:719–728.

771   Ogilvie, H. A., R. R. Bouckaert, and A. J. Drummond. 2017. StarBEAST2 brings faster
772        species tree inference and accurate estimates of substitution rates. Molecular Biology
773        and Evolution 34:2101–2114.

Page, R. D. M. 1994. Maps between trees and cladistic analysis of historical associations among genes, organisms and areas. Systematic Biology 43:58–77.

Page, R. D. M. and M. A. Charleston. 1997. From gene to organismal phylogeny: Reconciled trees and the gene tree/species tree problem. Molecular Phylogenetics and Evolution 7:231–240.

Pickrell, J. K. and J. K. Pritchard. 2012. Inference of population splits and mixtures from genome-wide allele frequency data. PLOS Genetics 8:e1002967.

Pironon, S., G. Papuga, J. Villellas, A. L. Angert, M. B. García, and J. D. Thompson. 2016. Geographic variation in genetic and demographic performance: new insights from an old biogeographical paradigm. Biological Reviews (Cambridge) 92:1877–1909.

Potts, A. J., T. A. Hedderson, and G. W. Grimm. 2014. Constructing phylogenies in the presence of intra-individual site polymorphisms (2ISPs) with a focus on the nuclear ribosomal cistron. Systematic Biology 63:1–16.

Rannala, B. and Z. Yang. 2017. Efficient Bayesian species tree inference under the multispecies coalescent. Syst. Biol. 66:823–842.

Rheindt, F. E., M. K. Fujita, P. R. Wilton, and S. V. Edwards. 2014. Introgression and phenotypic assimilation in *Zimmerius* flycatchers (Tyrannidae): Population genetic and phylogenetic inferences from genome-wide SNPs. Systematic Biology 63:134–152.

Schmidt-Lebuhn, A. N., N. C. Aitken, and A. Chuah. 2017. Species trees from consensus single nucleotide polymorphism (SNP) data: Testing phylogenetic approaches with simulated and empirical data. Molecular Phylogenetics and Evolution 116:192–201.

Schrempf, D., B. Q. Minh, N. De Maio, A. von Haeseler, and C. Kosiol. 2016. Reversible polymorphism-aware phylogenetic models and their application to tree inference. Journal of Theoretical Biology 407:362–370.

798 Schrempf, D., B. Q. Minh, A. von Haeseler, and C. Kosiol. 2019. Polymorphism-aware

799 species trees with advanced mutation models, bootstrap, and rate heterogeneity.

800 Molecular Biology and Evolution 36:1294–1301.

801 Shreve, F. 1951. Vegetation and Flora of the Sonoran Desert vol. 1. Carnegie Institution,

802 Washington, DC.

803 Springer, M. S. and J. Gatesy. 2016. The gene tree delusion. Molecular Phylogenetics and

804 Evolution 94:1–33.

805 Sridhar, S., F. Lam, G. E. Blelloch, R. Ravi, and R. Schwartz. 2007. Direct maximum

806 parsimony phylogeny reconstruction from genotype data. BMC Bioinformatics 8:472.

807 Subramanian, S., U. Ramasamy, and D. Chen. 2019. VCF2PopTree: a client-side software

808 to construct population phylogeny from genome-wide SNPs. PeerJ 7:e8213.

809 Suh, A., L. Smeds, and H. Ellegren. 2015. The dynamics of incomplete lineage sorting

810 across the ancient adaptive radiation of neoavian birds. PLOS Biology 13:e1002224.

811 Swofford, D. L. 2002. PAUP*. Phylogenetic Analysis Using Parsimony (*and Other

812 Methods). 4.0 ed. Sinauer, Sunderland, MA.

813 Swofford, D. L. and S. H. Berlocher. 1987. Inferring evolutionary trees from gene-frequency

814 data under the principle of maximum parsimony. Systematic Zoology 36:293–325.

815 Than, C. and L. Nakhleh. 2009. Species tree inference by minimizing deep coalescences.

816 PLOS Computational Biology 5:e1000501.

817 Than, C. and L. Nakhleh. 2010. Inference of parsimonious species trees from multi-locus

818 data by minimizing deep coalescences book section 5, Pages 79–98. Wiley-Blackwell.

819 Tian, Y. and L. Kubatko. 2017. Rooting phylogenetic trees under the coalescent model

820 using site pattern probabilities. BMC Evolutionary Biology 17:263.

821  Vachaspati, P. and T. Warnow. 2018. SVDquest: Improving svdquartets species tree

822      estimation using exact optimization within a constrained search space. Molecular

823      Phylogenetics and Evolution 124:122–136.

824  Wang, J. R. 2013. Analysis and Visualization of Local Phylogenetic Structure within

825      Species. Thesis.

826  Wang, W. S., R. Mauleon, Z. Q. Hu, D. Chebotarov, S. S. Tai, Z. C. Wu, M. Li, T. Q.

827      Zheng, R. R. Fuentes, F. Zhang, L. Mansueto, D. Copetti, M. Sanciangco, K. C. Palis,

828      J. L. Xu, C. Sun, B. Y. Fu, H. L. Zhang, Y. M. Gao, X. Q. Zhao, F. Shen, X. Cui,

829      H. Yu, Z. C. Li, M. L. Chen, J. Detras, Y. L. Zhou, X. Y. Zhang, Y. Zhao, D. Kudrna,

830      C. C. Wang, R. Li, B. Jia, J. Y. Lu, X. C. He, Z. T. Dong, J. B. Xu, Y. H. Li, M. Wang,

831      J. X. Shi, J. Li, D. B. Zhang, S. Lee, W. S. Hu, A. Poliakov, I. Dubchak, V. J. Ulat,

832      F. N. Borja, J. R. Mendoza, J. Ali, J. Li, Q. Gao, Y. C. Niu, Z. Yue, M. E. B. Naredo,

833      J. Talag, X. Q. Wang, J. J. Li, X. D. Fang, Y. Yin, J. C. Glaszmann, J. W. Zhang, J. Y.

834      Li, R. S. Hamilton, R. A. Wing, J. Ruan, G. Y. Zhang, C. C. Wei, N. Alexandrov, K. L.

835      McNally, Z. K. Li, and H. Leung. 2018. Genomic variation in 3,010 diverse accessions of

836      Asian cultivated rice. Nature 557:43–49.

837  Xu, B. and Z. H. Yang. 2016. Challenges in species tree estimation under the multispecies

838      coalescent model. Genetics 204:1353–1368.

839  Xu, D., Y. Jaber, P. Pavlidis, and O. Gokcumen. 2017. VCFtoTree: a user-friendly tool to

840      construct locus-specific alignments and phylogenies from thousands of anthropologically

841      relevant genome sequences. BMC Bioinformatics 18:426.

842  Yu, Y., T. Warnow, and L. Nakhleh. 2011. Algorithms for MDC-based multi-locus

843      phylogeny inference: Beyond rooted binary gene trees on single alleles. Journal of

844      Computational Biology 18:1543–1559.

845  Zhang, L. X. 2011. From gene trees to species trees II: Species tree inference by minimizing

846  deep coalescence events. IEEE-ACM Transactions on Computational Biology and

847  Bioinformatics 8:1685–1691.

848  Zhao, Y. P., G. Y. Fan, P. P. Yin, S. Sun, N. Li, X. N. Hong, G. Hu, H. Zhang, F. M.

849  Zhang, J. D. Han, Y. J. Hao, Q. W. Xu, X. W. Yang, W. J. Xia, W. B. Chen, H. Y. Lin,

850  R. Zhang, J. Chen, X. M. Zheng, S. M. Y. Lee, J. Lee, K. Uehara, J. A. Wang, H. M.

851  Yang, C. X. Fu, X. Liu, X. Xu, and S. Ge. 2019. Resequencing 545 *Ginkgo* genomes

852  across the world reveals the evolutionary history of the living fossil. Nature

853  Communications 10:4201.

854  Zhu, J. F. and L. Nakhleh. 2018. Inference of species phylogenies from bi-allelic markers

855  using pseudo-likelihood. Bioinformatics 34:376–385.

856  ## Appendix

857  We assume each tree, $T$, is rooted, with edge set, $E(T)$, node set, $V(T)$, and leaf

858  set, $L(T) \subset V(T)$. Define the internal node set $\dot{V}(T) \subset V(T)$ to be all nodes with (in or

859  out)degree two or more. Each leaf $x \in L(T)$ has indegree one and is labeled from the set $\chi$.

860  Let $\mathcal{T}_\chi$ be the set of all such rooted trees. An internal node $v$ with outdegree one is called

861  *unary* and two is *binary*. If all $v \in \dot{V}(T)$ are binary, $T$ is *binary*.

862  For $T \in \mathcal{T}_\chi$, the most ancestral node having outdegree two or more, if it exists, is

863  called the *crown root*, $\rho_T$. A node with indegree zero and outdegree one, if it exists, is

864  called the *stem root*, $\rho'_T$. If $T$ has a crown root with indegree zero it is called a *crown tree*.

865  If $T$ has a stem root that is the parent of a crown root node, then $T$ is a *stem tree*, with a

866  stem edge, $e_s = (\rho'_T, \rho_T)$ (see Fig. 10).

867  If node $v'$ is an ancestor of $v$, we write $v' > v$. Define $d_T(v', v)$ as the number of

868  edges on the path between $v'$ and $v$, where $d_T(v, v) = 0$.

869  The set of leaves descended from $v$ is $C_T(v)$. The node of the most recent common

870  ancestor of a set of leaves, $A$, is $\text{MRCA}_T(A)$. For any binary $v \in \dot{V}(T)$, let $l(v)$, $r(v)$, and

871    $a(v)$ be the left child, right child and parent node of $v$, respectively, if these exist.

872         Define $T\|_U$ to be the minimal subtree of $T$ containing a set of nodes $U \subseteq V(T)$,

873    and $T|_U$ is obtained from $T\|_U$ by suppressing any unary nodes on $T\|_U$ (Ma et al. 2001)
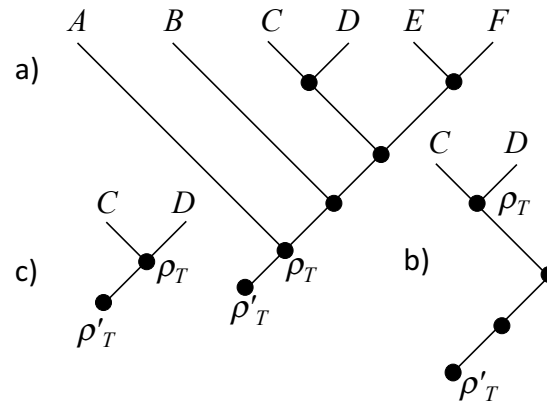
874    (See Fig. 10).



Fig. 10. Illustration of rooted stem tree, $T$, and two subtree operations. Let $U = \{C, D, \rho'_T\}$. a) Original rooted stem tree, (note notation for stem and crown root nodes); b) subtree $T\|_U$ (note retention of unary nodes), which is neither a stem tree or crown tree by our definitions; c) subtree $T|_U$ (note suppression of unary nodes), which is a stem tree.

875                    *Solution to PP Score Problem (Felsenstein 1979)*

876         For reference, we describe Felsenstein's (1979) two-pass algorithm to reconstruct

877    ancestral states under polymorphism parsimony, modified slightly to allow missing data.

878    Let the genotype tree, $\Psi$ be a rooted binary stem tree, and denote the genotype states of a

879    biallelic site, in which alleles are labeled 0 or 1, with an ordered pair of boolean variables

880    to indicate their presence or absence $(y^0, y^1), y \in \{+, -\}$. Thus, genotypes 0/0, 0/1, 1/1,

881    and ./. correspond to $(+, -), (++), (-+), (-, -)$, respectively.

882         More generally, let $y_v = (y_v^0, y_v^1)$ be a pair of boolean variables that indicates the

883    presence or absence of the respective alleles among any descendant leaves of $v$. If $v$ is a leaf

884    node, this is exactly the genotype of $v$. The algorithm first does a downpass over $v \in \dot{V}$,

885    assigning $y_v$ as follows:

$$y_v = (y^0_{l(v)} \vee y^0_{r(v)}, y^1_{l(v)} \vee y^1_{r(v)}), \tag{6}$$

where present/absent is treated as true/false for the logical "or" operator, $\vee$.

Next an uppass traversal over $v \in \dot{V}$ computes the actual optimal ancestral genotype states:

$$Y_v = (y^0_v, \mathrm{median}(Y^1_{a(v)}, y^1_{l(v)}, y^1_{r(v)})), \tag{7}$$

where "median$(b_1, b_2, b_3)$", is the most frequent state in three boolean variables. At $v = \rho_\Psi$, the state $Y^1_{a(v)}$ is initialized to '$-$' because $\rho'_\Psi$ has genotype of $(+, -)$, by assumption. The asymmetry in the two components of $Y_v$ arises because, by assumption, the 0 allele is always present in the ancestor of the root, but the 1 allele is not. Once the final internal states have been inferred, then the PP score, $c_{PP}(\Psi, G)$, is the number of nodes for which $Y_v = Y_{a(v)} = (+, +)$.

## *Solution to DC-G Score Problem*

*Preliminaries.—* In addition to the rooted binary stem genotype tree, $\Psi$, also now assume an underlying diploid rooted binary stem allele tree, $t$, for which there are at most two alleles for each leaf of $\Psi$. For any set of leaves of $t$, $W \subseteq L(t)$, let $\alpha(W)$ be the corresponding set of leaves of $\Psi$.

Define the *MRCA-mapping*, $\mathcal{M}$, from node $v$ in $t$ to a node in $\Psi$ as $\mathcal{M}_\Psi(v) = \mathrm{MRCA}_\Psi(\alpha(C_t(v)))$.

Define a *coalescent history* as a mapping, $h$, of all nodes in $t$ to nodes in $\Psi$. Informally, this describes how the allele tree is imbedded in the genotype tree. For edge, $e = (u, v) \in E(t)$, let $D_\Psi(e, h(t)) = d_\Psi(h(u), h(v))$ be the number of edges in $\Psi$ in the coalescent history's path from $h(u)$ to $h(v)$ (Fig. 11). Now, the *optimal coalescent history*,
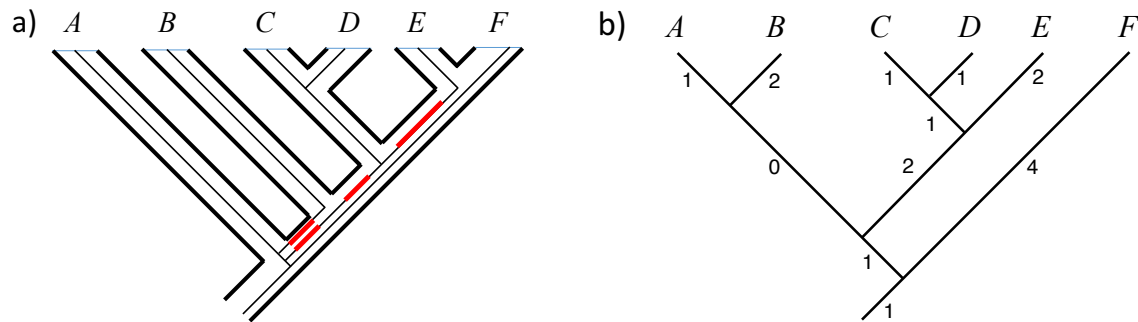
Fig. 11. a) Genotype tree, $\Psi$, with imbedded allele tree, $t$; extra lineages shown in red; b) Allele tree, $t$. Numbers next to edges are $D_\Psi(e, h(t))$, which count the number of edges in $\Psi$ in which this edge of $t$ is imbedded. Here $|\Psi|\|_{L(t)} = 11$ from (a), $\sum_e D_\Psi(e, h(t)) = 15$ from (b), and DC score is therefore 4.

$h^*(t)$, minimizes

$$\Lambda_\Psi(t) = \sum_{e \in E(t)} D_\Psi(e, h(t)) \tag{8}$$

over all $h(t)$. Although the number of distinct histories, $h(t)$, can be quite large (Degnan and Salter 2005), the optimal history for a crown tree is just given by $h^*(v) = \mathcal{M}_\Psi(v)$ for all $v \in \dot{V}(t)$ (Yu et al. 2011; Bayzid and Warnow 2012). For a stem allele tree, however, we need the following.

**Lemma 1** (Optimal coalescent history for stem allele tree). *Assume $g$ is a stem allele tree with crown tree, $g_c$, and stem edge $e_s = (\rho_g', \rho_g)$, and let $u$ in $\Psi$ be the node to which $\rho_g$ maps via $\mathcal{M}$. Further, constrain $\rho_g'$ to be present at some node, $u' \geqslant u$ of $\Psi$ so that there is a path on $\Psi$ between $u'$ and $u$ having $d(u', u) \geqslant 0$ edges. The optimal coalescent history of $g$ has two parts: that for $g_c$, which is determined by $\mathcal{M}$, as with any crown allele tree, and the history below this node down to $u'$, comprising a single edge of $g$. Thus,*

$\Lambda_\Psi(g) = \Lambda_\Psi(g_c) + d(u', u)$.

*Proof.* The optimal coalescent history of $g_c$ is given by $\mathcal{M}$ (Yu et al. 2011), which maps $\rho_g$ to $u$ on $\Psi$. With $e_s$ present, $\rho_g'$ is forced to be at $u' \geqslant u$, corresponding to $d(u', u)$ edges of $\Psi$ between $u'$ and $u$. By the definition of $\mathcal{M}$, the only way to modify the coalescent history of $g$ would be to move $\rho_g$ closer to the root of $\Psi$, decreasing $d(u', u)$. However, because $\rho_g$ has outdegree two, this would pull two allele tree lineages deeper in the genotype tree, and

<sub>923</sub> increase $\Lambda_\Psi(g_c)$ by *two* for each edge lost in $d(u', u)$. Thus, the optimal coalescent history

<sub>924</sub> of $g$ is given by the history for $g_c$ plus the path corresponding to $e_s$. $\qquad\square$

<sub>925</sub> There is some ambiguity in the literature (Than and Nakhleh 2010; Yu et al. 2011;

<sub>926</sub> Bayzid and Warnow 2012, 2018) (and in software implementations: Bayzid and Warnow

<sub>927</sub> 2012) about counting imbedded allele tree lineages and computing the DC score when the

<sub>928</sub> allele tree is missing from some leaves of the genotype tree: i.e., when $\alpha(L(t)) \subset L(\Psi)$. In

<sub>929</sub> that case, $t$ could be considered to be imbedded in either $\Psi|_{L(t)}$ or in $\Psi||_{L(t)}$ (Zhang 2011;

<sub>930</sub> Bayzid and Warnow 2012), and since the number of edges differs, so too could the DC

<sub>931</sub> score. We adopt the framework established in Bayzid and Warnow (2012) and Bayzid and

<sub>932</sub> Warnow (2018) which uses $\Psi||_{L(t)}$. This counts edges of $\Psi$ joined at unary nodes of $\Psi||_{L(t)}$,

<sub>933</sub> whereas $\Psi|_{L(t)}$ collapses these into a single edge.

<sub>934</sub> *DC score for stem tree with missing data.—* Now we are in a position to define the

<sub>935</sub> deep coalescence score in more general terms. The DC score, $c_{\mathrm{DC}}(\Psi, t)$, is the number of

<sub>936</sub> "extra" imbedded allele tree lineages above what is minimally necessary given the

<sub>937</sub> genotype data. Suppose leaf set $A \subset L(\Psi)$ has genotype data present, and suppose the

<sub>938</sub> stem root of allele tree $t$ is forced to be present at some node $u$ on $\Psi$, then let

<sub>939</sub> $U = A \cup \{u\}$. The minimum number of edges of $\Psi$ that must have at least one imbedded

<sub>940</sub> allele tree edge is $|E(\Psi||_U)|$, and therefore

$$c_{\mathrm{DC}}(\Psi, t) = \Lambda_\Psi(t) - |E(\Psi||_U)|. \tag{9}$$

<sub>941</sub> The last term depends only on the genotype tree and the data, not the unknown allele

<sub>942</sub> tree, $t$, so minimizing $\Lambda_\Psi(t)$ is equivalent to minimizing the number of deep coalescences.

<sub>943</sub> *Optimal allele trees.—* Now, we need to *find* an allele tree, $t$ that minimizes $\Lambda_\Psi(t)$

<sub>944</sub> for given $\Psi$ over all $t$ consistent with $G$ and the infinite alleles assumption. If there were no

<sub>945</sub> constraints on $t$, the lemma below could be used to find this tree easily. However, $t$ cannot

<sub>946</sub> be just any tree; it is constrained by the assumptions of the model so that: (i) it must

include two disjoint subtrees, $t^0$, and $t^1$, the leaves of which have these respective allele

states; (ii) the stem root node of $t^0$ must be present at the stem root of $\Psi$, because the 0

allele is ancestral; and (iii) the stem edge of $t^1$ (which is where the single $0 \rightarrow 1$ mutation

occurs on $t$), attaches to some edge of $t^0$. Denote the set of all allele trees satisfying

conditions (i)-(iii) as $\mathcal{T}_G^{01}$.

Because $\Lambda_\Psi(t)$ is additive over edges, its value for a tree comprising two subtrees

meeting at a single internal node is

$$\Lambda_\Psi(t) = \Lambda_\Psi(t^0) + \Lambda_\Psi(t^1) \tag{10}$$

The following lemma shows how to find an optimal subtree for each of these terms

separately. Because each of these terms has a lower bound determined only by the

genotype data and $\Psi$, these optimal subtrees will minimize the sum, and therefore the

assembled tree will also be the optimal.

**Lemma 2** (Optimal allele tree). *Given a set of allele tree leaf nodes, $W$, with $\alpha(W)$ the*

*corresponding leaves of $\Psi$, and $u = \mathcal{M}(\alpha(W))$, and given a node, $u' > u$ of $\Psi$, at which the*

*stem root of some allele tree, $g$, will be constrained to be present, then let $U = \alpha(W) \cup u'$.*

*For all $g \in \mathcal{T}_W$, $\Lambda_\Psi(g) \geqslant |E(\Psi||_U)|$. Moreover, there exists an optimal allele stem tree,*

*$g_{min}$, that achieves this lower bound, having the following properties: (i) for any leaf, $x$, of*

*$\Psi$, having two leaves in $W$, the two leaves form a cherry (a pair of siblings) in $g_{min}$, with a*

*parent we call $\xi(x)$; and (ii) the "backbone" subtree of $g_{min}$, obtained by replacing any such*

*cherries with their corresponding nodes, $\xi(x)$, has a topology given by $\Psi|_U$.*

*Proof.* Suppose $g_{min}$ is given as above. Each node $\xi(x)$ has two child edges, $e = (u, v)$ in

which both $u$ and $v$ map to node $x$ on $\Psi$, because they correspond to the same genotype

leaf, and thus $D_\Psi(e, h^*(g_{min})) = 0$ for both. Therefore, these cherries do not contribute to

$\Lambda_\Psi(g)$.

Because the remaining "backbone" subtree, $g'_{min}$, is the same as $\Psi|_U$ (i.e., cluster by

cluster), each edge of $g'_{min}$, $e_i = (u, v)$, corresponds to a unique path, $p_i$ in $\Psi||_U$ from $\mathcal{M}(u)$

972 to $\mathcal{M}(v)$. These paths together must partition $\Psi||_U$ edgewise and therefore

973 $\Lambda_\Psi(g_{min}) = |E(\Psi||_U)|$. Alternatively, note that for every node of $\Psi||_U$ with outdegree two,

974 $D_\Psi(e_i, h^*(g_{min})) = 1$ (the path has one edge). For any node with outdegree one, the path

975 has more than one edge, but in either case, these will be counted correctly in the value

976 $|E(\Psi||_U)|$.

977 $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad\square$

978 *Construction of $t_{min}$.—* Now we describe the construction of $t_{min}$, the optimal

979 $t \in \mathcal{T}_G^{01}$. Let $W^0$ and $W^1$ be the subsets of leaf alleles having the 0 and 1 states

980 respectively. The 0 allele must be present at the stem root of $\Psi$, by assumption, because it

981 is ancestral according to the model. In Lemma 2 we therefore have $U = \alpha(W^0) \cup \{\rho'_\Psi\}$, and

982 $t_{min}^0$ has the topology of $\Psi|_U$ (see example in Fig. 10).

983 The solution for $t_{min}^1$ is complicated slightly by the final disposition of the stem

984 edge of $t_{min}^1$, in particular, where it attaches to $t_{min}^0$. By this we mean also, because $t_{min}^0$ is

985 imbedded at its optimal coalescent history, to which node, $u''$, in $\Psi$ does it attach? The

986 node $\rho_{t_{min}^1}$ is placed at $u' = \mathcal{M}_\Psi(\rho_{t_{min}^1})$, and according to Lemma 1, constraining $t_{min}^1$ to be

987 present at $u''$ adds $d(u', u'')$ to the score of imbedded allele edges. Thus we should choose

988 the attachment point to minimize $d(u', u'')$. If $t_{min}^0$ and $u'$ overlap, then we can choose

989 $u' = u''$. Then the stem root of $t_{min}^1$ can attach to $t_{min}^0$ "just below" $u'$, effectively given the

990 stem edge a length of zero in keeping with $d(u', u'') = 0$. If they do not overlap, we choose

991 $u''$ to be as close as possible to $u'$ in the path from $u'$ to the root, as long as $u'' > u'$.

992 This construction for $u'' = u'$ accords well with the reconcilation framework. For

993 example, the MRCA mapping of both the stem and crown root of $t_{min}^1$ is $u'$, which means

994 we infer a duplication at the stem root node, where it joins $t_{min}^0$ (Goodman et al. 1979).

995 However, the "boundary" case of $u'' > u'$ is a bit more interesting, because it exposes some

996 of the assumptions of the polymorphism parsimony model. This case arises only when

997 there are no heterozygotes in the data and either all genotypes are trivially 0/0, or some

genotypes are 1/1 and those genotypes form a clade on $\Psi$ (if there is more than one 1/1).

No internal node is reconstructed as 0/1 under these circumstances, but the PP model

assumes that a 0/1 ancestor served as an intermediate. Moreover it also assumes a

duplication/coalescent event was necessary to evolve the 0/1 genotype. Essentially, the

model postulates a hidden polymorphic state in the tree despite none having been

observed in the data.

We can engineer this hidden state within a reconciliation framework by inserting a

node, $x$ and pendant edge, $e_x$, on the stem edge of $t_{min}^1$. This subdivides the original stem

edge into a parent edge, $e_a$ and child edge, $e_c$ that subtends the crown root node of $t_{min}^1$

(Fig. 12). Node $x$ represents the "duplication" event, and the origin of the new 1 allele

state occurs along $e_c$. The edge, $e_x$ terminates immediately and represents the "loss" of the

0-allele. Because of this, $e_x$ adds zero to $\Lambda_\Psi$. Both allele trees are present

simultaneously—briefly—indicating ephemeral polymorphism.



Fig. 12. Illustration of a deeper interpretation for constructing $t_{min}$, needed in the boundary case in which there is no overlap between allele subtrees, $t_{min}^0, t_{min}^1$. To maintain the assumptions of the PP model, we could insert a subtree (red dotted lines) inside the genotype tree edge in place of the stem edge of $t_{min}^1$. In this subtree there is a duplication, $x$, a mutation to 1 in one child edge, $e_c$, and an immediate loss of the other edge, $e_b$. This generates an ephemeral ancestral heterozygote (dotted circle).

In both cases, the construction of the two subtrees leads to an optimal allele tree

$_{1012}$ with minimum DC score, but what is this score? Any edge of $\Psi$ can have 0, 1, or 2

$_{1013}$ imbedded allele tree edges, but only edges with 2 contribute to the DC score. These are

$_{1014}$ edges in which both $t_{min}^0$ and $t_{min}^1$ are present, which is the subtree of $\Psi$ where the two

$_{1015}$ subtrees overlap, $\psi$. Thus, the DC-G score is:

$$c_{\text{DC-G}}(\Psi, G) = |E(\psi)|, \tag{11}$$

$_{1016}$ which makes clear that $c_{\text{DC-G}}(\Psi, G)$ is nonzero only if the two subtrees overlap.

$_{1017}$ <div align="center">*Proof of Theorem 1*</div>

$_{1018}$     We will show that the ancestral states reconstructed by PP, $Y_v = (Y_v^0, Y_v^1)$, are

$_{1019}$ equivalent to the presence or absence of alleles in an imbedded allele tree that matches $t_{min}$

$_{1020}$ for the DC-G computation. Thus the two methods produce the same score for a given

$_{1021}$ genotype tree and character.

$_{1022}$     The proof has two parts: (i) the PP downpass effectively traces out a tree for each

$_{1023}$ allele that is the minimal subtree of $\Psi$ for the leaves having that allele present and

$_{1024}$ extending down to the crown root (for the 0 allele this will be the final tree, $t_{min}^0$); (ii)

$_{1025}$ however, for the 1 allele, $t_{min}^1$ must be truncated toward the root so that it does not extend

$_{1026}$ below the last possible point of origin of the 1 allele. This is taken care of by the uppass

$_{1027}$ phase, which only affects the 1 allele.

$_{1028}$     To prove (i), note that for allele $i \in \{0, 1\}$ at node $v$, the downpass value, $y_v^i$ is set to

$_{1029}$ $+$ if any descendant nodes of $v$ are $+$. The downpass thus traces out the minimal subtree

$_{1030}$ of $+$ states that extends from the crown root of $\Psi$ to each leaf of $\Psi$ having the allele. For

$_{1031}$ the 0 allele this is the final set of states because the 0 allele is assumed to be present at the

$_{1032}$ stem root of $\Psi$, so $Y_v^0 = y_v^0$, and the tree implied by PP is exactly $t_{min}^0$ described above.

$_{1033}$     To prove (ii) is true for the 1 allele, consider the instances when the uppass might

$_{1034}$ *change* a downpass value $y_v^1$. From Eq. 7, the uppass value, $Y_v^1$, depends on states at its

$_{1035}$ three adjacent nodes:

$$Y_v^1 = \text{median}(Y_{a(v)}^1, y_{l(v)}^1, y_{r(v)}^1), \tag{12}$$

1036    There are 8 possible combinations of $+$ and $-$ arguments in this function, and

1037    enumeration of these shows that only two of them lead to cases in which $Y_v^1$ is different

1038    from $y_v^1$ (partly because of how $y_v^1$ is computed from its two child states): either

1039    $Y_{a(v)}^1 = -, y_{l(v)}^1 = +, y_{r(v)}^1 = -$ or $Y_{a(v)}^1 = -, y_{l(v)}^1 = -, y_{r(v)}^1 = +$. In both cases the change is

1040    from $y_v^1 = +$ to $Y_v^1 = -$. The uppass is initialized by setting the stem root state to $-$ (see

1041    Eq. 7), so that $Y_{a(v)}^1 = -$ when $v$ is the crown root. Thus, the uppass traversal *might* start

1042    out having to invoke one of these changes of state iff only one of its two children has the $+$

1043    state. If so, this will continue until the uppass reaches the most recent node that has both

1044    children with a $+$ state; that is, a node that *must* have the 1 allele (rather than merely

1045    *may* have it).

1046    This node of $\Psi$ is then the proper location of the crown root node of $t_{min}^1$, and the

1047    uppass is therefore equivalent to construction of the $t_{min}^1$ crown allele tree. Its stem edge

1048    attaches to $t_{min}^0$ as described earlier. Since the DC-G score depends only on the two

1049    subtrees and their intersection (Eq. 11), this means that the PP score and DC-G score are

1050    the same (see Fig. 13).

1051    □

## *Pairwise Genotype Compatibility*

1052

1053    Under the infinite sites model in the absence of recombination there should be

1054    blocks of neighboring loci that are compatible with each other (Hudson and Kaplan 1985).

1055    Recombination and sequencing error limits the lengths of such blocks. To estimate the

1056    distance between effectively unlinked, statistically independent sites, we assayed

1057    compatibility among genotypes of sites. For example, suppose two sites' genotypes for an

1058    individual are 0/1 and 1/1, then the possible haplotype assignments ($=$ "gametes") are 01

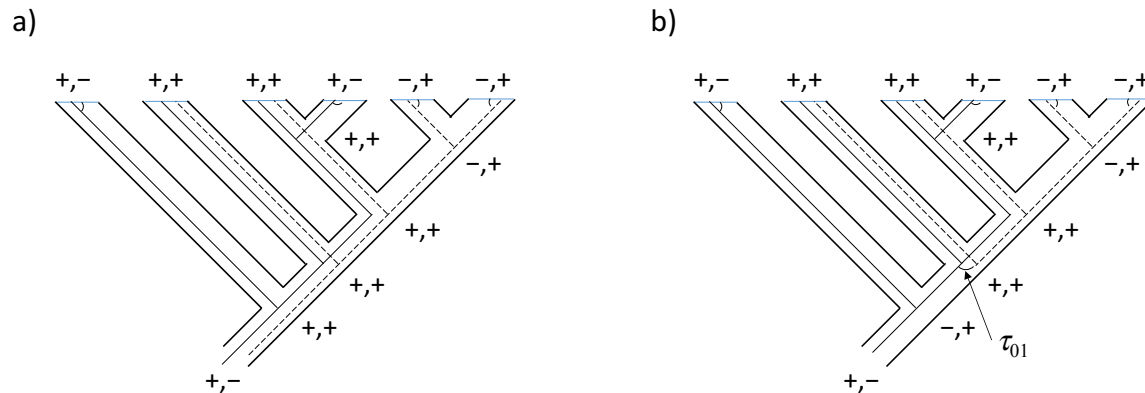1059    and 11, and each must be present. Thus, two of the four "gametes" (01 and 11) allowed by

Fig. 13. Example of equivalence between PP and DC-G ancestral state reconstructions. Solid line is allele tree $t_{min}^0$; dashed line is $t_{min}^1$. The stem edge of $t_{min}^1$ joins the two (here labeled $\tau_{01}$). a) After the downpass phase: initial allele subtrees for DC-G and the $y_v$ states for each node are shown. b) After the uppass phase: final allele tree for DC-G and final PP ancestral states, $Y_v$, are shown.

the four gamete test (Hudson and Kaplan 1985) must already be present in this one

individual's pair of genotypes. If the other two gamete types are also present in possible

haplotype assignments for any other individuals in the data, then this pair of sites is

genotype incompatible. The diagram in Fig. 14 is useful for enumerating the haplotype

(=gamete) types that are implied by a pair of genotypes. By checking off observed gamete

types among such genotype pairs for all individuals, one can determine if the two sites are

genotype-compatible (see Wang 2013, for further details).

This uses the unrooted definition of character compatibility so is an unrooted

version of genotype compatibility.

Fig. 14. Pairwise genotype compatibility. Genotype pairs at two sites are shown within boxes. Large pairs of numbers are allelic haplotypes implied by the genotype pairs. In some cases (dashed connectors), a genotype pair implies both of two haplotype pairs. If individuals in the data set collectively have haplotypes (="gametes") of all four types, the pair of sites is "genotype-incompatible". The 0/1-0/1 genotype pair is ambiguous and can go in either direction indicated, which means if three gamete types are already observed among the individuals, a 0/1-0/1 genotype pair will *not* cause genotype incompatibility.

Table 1. *Sequencing read statistics for samples*

| Sample | Population | Reads ($10^6$) | Bp ($10^9$) | Coverage (x) | Mean Length (bp) | Mean Quality Score |
|--------|-----------|------|------|------|------|------|
| Sample_C60 | Caborca | 87.18 | 11.64 | 8.3 | 133 | 33 |
| Sample_C61 | Caborca | 92.95 | 12.28 | 8.8 | 132 | 33 |
| Sample_D17 | El Dipo | 65.93 | 6.26 | 4.5 | 94 | 37 |
| Sample_D73 | El Dipo | 74.01 | 7.05 | 5.0 | 95 | 37 |
| Sample_G50 | Guasimas | 83.68 | 11.19 | 8.0 | 133 | 33 |
| Sample_G81 | Guasimas | 81.44 | 10.90 | 7.8 | 133 | 33 |
| Sample_K26 | Kino Bay | 87.85 | 11.78 | 8.4 | 134 | 33 |
| Sample_K4 | Kino Bay | 56.04 | 7.44 | 5.3 | 132 | 33 |
| Sample_J34 | La Joyita | 59.52 | 5.67 | 4.0 | 95 | 37 |
| Sample_J45 | La Joyita | 65.96 | 6.29 | 4.5 | 95 | 37 |
| Sample_M69 | Masiaca | 71.93 | 6.84 | 4.9 | 95 | 37 |
| Sample_M85 | Masiaca | 43.39 | 4.12 | 2.9 | 94 | 37 |
| SGP-OR3 | Oregano | 89.82 | 11.69 | 8.3 | 130 | 34 |
| SGP-OR6 | Oregano | 76.94 | 9.99 | 7.1 | 129 | 34 |
| SGP-SM1 | San Marcial | 28.90 | 3.78 | 2.7 | 130 | 34 |
| SGP-SM7 | San Marcial | 22.98 | 3.03 | 2.2 | 131 | 34 |
| SGP5_L004 | Tucson | 256.20 | 24.49 | 17.5 | 95 | 38 |
| Sample_CU | Tucson | 76.08 | 7.24 | 5.2 | 95 | 37 |
| Sample_119 | Wickenburg | 110.12 | 10.48 | 7.5 | 95 | 37 |
| Sample_129 | Wickenburg | 73.01 | 6.96 | 5.0 | 95 | 37 |

Table 2. *Variants, genotype statistics and phylogenetic data sets*

| Data set | 0-Fold Sites | 4-Fold Sites | Large 100k Scaffolds |
|----------|------|------|------|
| Total effective sites | 19563072 | 4592714 | 267653586 |
| Number of variants | 106998 | 47741 | 2871509 |
| Biallelic sites in complete matrix | 106629 | 47511 | 2858685 |
| Tri- or quad-allelic sites | 369 | 230 | 12824 |
| Fraction of tri- or quad-allelic sites | 0.0035 | 0.0048 | 0.0045 |
| Fraction of those sites homozygous | 0.72 | 0.72 | 0.74 |
| Fraction of those sites heterozygous | 0.19 | 0.19 | 0.17 |
| Fraction of those sites missing | 0.08 | 0.09 | 0.09 |
| Fraction of those sites with no missing data | 0.18 | 0.17 | 0.14 |
| Sites in each compatibility-thinned matrix | 5332 | 4752 | 28586 |
| Sites in PLINK-thinned matrix | 27457 | 13977 | 88480 |
| Sites in downsampled matrix | 31994 | 14367 | 755194 |
| Sites in each thinned downsampled matrix | 5333 | 4789 | 7552 |

Table 3. *Nucleotide diversity estimates[a]*

| Population | N | Large100k Scaffolds | | | Protein coding genes | | | | | | | | |
| | | | | | 0-fold degenerate sites | | | 4-fold degenerate sites | | | Efficacy[b] | |
| | | Variants | $\pi_T$ | $\pi_W$ | Variants | $\pi_T$ | $\pi_W$ | Variants | $\pi_T$ | $\pi_W$ | $\pi_T^0/\pi_T^4$ | $\pi_W^0/\pi_W^4$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Pooled | 20 | 2871509 | 0.00241 | 0.00258 | 106998 | 0.00124 | 0.00131 | 47741 | 0.00253 | 0.00250 | 0.490 | 0.524 |
| Caborca | 2 | 1077484 | 0.00224 | 0.00223 | 40336 | 0.00117 | 0.00115 | 19117 | 0.00237 | 0.00233 | 0.494 | 0.494 |
| El Dipo | 2 | 863443 | 0.00195 | 0.00194 | 34688 | 0.00108 | 0.00106 | 16645 | 0.00221 | 0.00218 | 0.489 | 0.486 |
| Guasimas | 2 | 1067561 | 0.00224 | 0.00223 | 40603 | 0.00118 | 0.00116 | 19418 | 0.00242 | 0.00238 | 0.488 | 0.487 |
| Kino Bay | 2 | 1087860 | 0.00239 | 0.00239 | 39744 | 0.00122 | 0.00121 | 18990 | 0.00251 | 0.00248 | 0.486 | 0.488 |
| La Joyita | 2 | 822958 | 0.00196 | 0.00194 | 33323 | 0.00106 | 0.00105 | 15958 | 0.00218 | 0.00215 | 0.486 | 0.488 |
| Masiaca | 2 | 625567 | 0.00172 | 0.00169 | 26248 | 0.00097 | 0.00095 | 12736 | 0.00202 | 0.00197 | 0.480 | 0.482 |
| El Oregano | 2 | 1026983 | 0.00222 | 0.00217 | 37817 | 0.00114 | 0.00111 | 18075 | 0.00233 | 0.00227 | 0.489 | 0.489 |
| San Marcial | 2 | 239978 | 0.00176 | 0.00173 | 11651 | 0.00106 | 0.00103 | 5457 | 0.00214 | 0.00208 | 0.495 | 0.495 |
| Tucson | 2 | 872562 | 0.00182 | 0.00182 | 36312 | 0.00105 | 0.00104 | 17162 | 0.00212 | 0.00209 | 0.495 | 0.498 |
| Wickenburg | 2 | 845475 | 0.00180 | 0.00179 | 32406 | 0.00096 | 0.00093 | 15389 | 0.00193 | 0.00189 | 0.497 | 0.492 |

[a] $\pi_T$ and $\pi_W$ are Tajima and Watterson estimators

[b] Superscripts denote degeneracy class

Table 4. *Nucleotide diversity estimates (Angsd)*

| Population | N | Large100k $\pi_T$ | Scaffolds $\pi_W$ |
|---|---|---|---|
| 10 pops(x2) | 20 | 0.00277 | 0.00325 |
| Caborca | 2 | 0.00235 | 0.00232 |
| El Dipo | 2 | 0.00203 | 0.00199 |
| Guasimas | 2 | 0.00243 | 0.00240 |
| Kino Bay | 2 | 0.00255 | 0.00253 |
| La Joyita | 2 | 0.00198 | 0.00194 |
| Masiaca | 2 | 0.00181 | 0.00174 |
| El Oregano | 2 | 0.00250 | 0.00245 |
| San Marcial | 2 | 0.00180 | 0.00173 |
| Tucson | 2 | 0.00184 | 0.00183 |
| Wickenburg | 2 | 0.00185 | 0.00181 |

Table 5. *Fractional missing data in downsampled datasets*

| Sample | 0-Fold Sites | 4-Fold Sites | Large100k |
|---|---|---|---|
| Caborca_C60 | 0.0092 | 0.0094 | 0.0054 |
| Caborca_C61 | 0.0160 | 0.0144 | 0.0078 |
| ElDipo_D17 | 0.0553 | 0.0575 | 0.0544 |
| ElDipo_D73 | 0.0298 | 0.0289 | 0.0395 |
| Guasimas_G50 | 0.0148 | 0.0123 | 0.0094 |
| Guasimas_G81 | 0.0161 | 0.0144 | 0.0109 |
| Joyita_J34 | 0.0728 | 0.0770 | 0.0850 |
| Joyita_J45 | 0.0366 | 0.0380 | 0.0536 |
| KinoBay_K26 | 0.0144 | 0.0107 | 0.0081 |
| KinoBay_K4 | 0.0640 | 0.0611 | 0.0500 |
| Masiaca_M69 | 0.0400 | 0.0397 | 0.0506 |
| Masiaca_M85 | 0.0519 | 0.0717 | 0.0588 |
| Oregano_OR3 | 0.0155 | 0.0174 | 0.0115 |
| Oregano_OR6 | 0.0262 | 0.0264 | 0.0196 |
| SanMarcial_SM1 | 0.0297 | 0.0301 | 0.0594 |
| SanMarcial_SM7 | 0.0416 | 0.0409 | 0.0542 |
| Tucson_CU | 0.0284 | 0.0298 | 0.0314 |
| Tucson_SGP5 | 0.0051 | 0.0036 | 0.0004 |
| Wickenburg_119 | 0.0090 | 0.0061 | 0.0074 |
| Wickenburg_129 | 0.0308 | 0.0322 | 0.0337 |