

# BONSAI: AN EFFICIENT METHOD FOR INFERRING LARGE HUMAN PEDIGREES FROM GENOTYPE DATA

Ethan M. Jewett\*, Kimberly F. McManus, William A. Freyman,  
the 23andMe Research Team, and Adam Auton.

23andMe, Inc. Sunnyvale, CA.

**Address correspondence to:** [ejewett@23andme.com](mailto:ejewett@23andme.com)

## 1. ABSTRACT

Pedigree inference from genotype data is a challenging problem, particularly when pedigrees are sparsely sampled and individuals may be distantly related to their closest genotyped relatives. We present a new method that infers small pedigrees of close relatives and then assembles them into larger pedigrees. To assemble large pedigrees, we introduce several new formulas and tools including a new likelihood for the degree separating two small pedigrees, a method for detecting individuals who share background identity-by-descent (IBD) that does not reflect recent common ancestry, and a method for identifying the ancestral branches through which distant relatives are connected. Our method also takes several new approaches that help to improve the accuracy and efficiency of pedigree inference. In particular, we incorporate age information directly into the likelihood rather than using ages only for consistency checks and we employ a heuristic branch-and-bound-like approach to more efficiently explore the space of possible pedigrees. Together, these approaches make it possible to construct large pedigrees that are challenging or intractable for current inference methods. The new method, Bonsai, is available at <https://github.com/23andMe/bonsaitree>.

## 2. INTRODUCTION

The ability to infer complex multi-generational pedigrees from genotype data has many applications ranging from genealogical research to the study of diseases. As human genotyping datasets continue to grow in size, there is increasing interest in computational methods that can reconstruct large pedigrees efficiently and accurately.

Although the problem of pedigree inference has been studied extensively, the majority of pedigree inference methods are designed for non-human species. A major challenge for pedigree reconstruction in non-human populations is that pairwise relationships can be difficult to infer with high accuracy, even when the degree of a relationship is small, because high quality genotype data may be unavailable. As a result, methods typically require that all or most individuals in a pedigree are sampled so that pedigrees can be assembled by connecting strings of parent-child, full-sibling, or half-sibling pairs (Almudevar 2003, Almudevar and Anderson 2012, Cowell 2009; 2013, Cussens et al. 2013, Wang 2004, Jones and Wang 2017, Kirkpatrick et al. 2011, Riester et al. 2009, Sheehan

## 2 BONSAI: AN EFFICIENT METHOD FOR INFERRING LARGE HUMAN PEDIGREES FROM GENOTYPE DATA

et al. 2014). Although it is possible to connect slightly more distant relationships (Huisman 2017, Anderson and Ng 2016), the majority of existing pedigree inference algorithms can be characterized as methods for either jointly or independently inferring pairwise parent-child pairs and full or half sibling sets, which are then consistent with a pedigree structure when assembled together.

In contrast to non-human pedigrees, genotype data for human populations is comparatively abundant and close relationships, such as parent-offspring or sibling pairs, can be inferred with a high degree of accuracy. The major challenge of pedigree inference in human populations is the fact that pedigrees are often sparsely sampled, with few genotyped sibling and parent pairs and few genotyped individuals beyond the most recent two or three generations. In human datasets, including direct-to-consumer genetic databases, genotyped individuals may have only a few genotyped relatives within a radius extending to first or second cousins and it is common for an individual's closest relative to be more distant than a second cousin. As a result, it is difficult to construct solid frameworks of close relatives and their genotyped ancestors into which other genotyped individuals can be placed.

There are currently two state-of-the-art methods for inferring complex human pedigrees from genotype data, both of which are maximum likelihood approaches that attempt to find a pedigree that maximizes the sum of log likelihoods of pairwise relationships, given observed patterns of identity-by-descent (IBD) sharing. The two methods differ primarily in the approaches they take to find the maximum likelihood pedigree.

The first and older method, PRIMUS (Staples et al. 2014), explores the space of possible pedigrees by starting with a seed individual and then iteratively adding individuals to the pedigree. Each time an individual is added, the method considers all possible positions that are consistent with the estimated pairwise relationships and the highest likelihood configuration is selected. When adding an individual to the pedigree, each pedigree at the previous step serves as a seed pedigree onto which the individual can be added in multiple ways. By constructing a large set of pedigrees in this way, the algorithm efficiently explores the space of pedigrees that are compatible with the estimated pairwise relationships.

In contrast to PRIMUS, the more recent CLAPPER method (Ko and Nielsen 2017) begins by connecting all individuals together into an initial guess of a pedigree. Then, at each subsequent step, the CLAPPER algorithm rearranges the relationships in the pedigree. This update step is done using a Markov chain Monte Carlo (MCMC) approach in which there are many different possible moves that can be made, such as adding or subtracting a degree of relatedness between two individuals, swapping the labels of two nodes, or pruning off an individual and their descendants and attaching them somewhere else.

The CLAPPER method is typically more accurate than PRIMUS (Ko and Nielsen 2017), whereas the PRIMUS approach can be faster than the MCMC approach used by CLAPPER. However, neither approach was designed to infer the large and sparse pedigrees that are common in direct-to-consumer genetic datasets where the degree of relationship separating a pair of genotyped individuals may be large, verging on degrees where individuals frequently share no detectable IBD. For such pedigrees, searching a broad pedigree space using the approach of PRIMUS or CLAPPER is computationally infeasible. Instead, it is useful to develop an inference approach that dramatically narrows the space

BONSAI: AN EFFICIENT METHOD FOR INFERRING LARGE HUMAN PEDIGREES FROM GENOTYPE DATA 3

73 of possible pedigrees, while being careful not to exclude the portion of the space containing the true  
74 pedigree.

75 Here, we introduce a new method, Bonsai, for inferring large and sparse pedigrees. To make  
76 inference efficient and accurate, we first infer small pedigrees of closely-related individuals using  
77 an approach that efficiently explores the space of possible pedigrees. This approach is similar to  
78 PRIMUS, but differs in key ways that make the search of the pedigree space both more efficient  
79 and more thorough. The small pedigrees are then assembled into larger pedigrees using several new  
80 techniques, including a generalized version of the DRUID method of Ramstetter et al. (2018), which  
81 allows our method to link distantly related individuals into large and sparsely sampled pedigrees.  
82 We refer to the first stage as “Small Bonsai” and to the second stage as “Big Bonsai” (Figure 1).  
83 We first describe the small and big Bonsai methods, then use both simulated and real data to  
84 investigate the performance of the methods and their components.

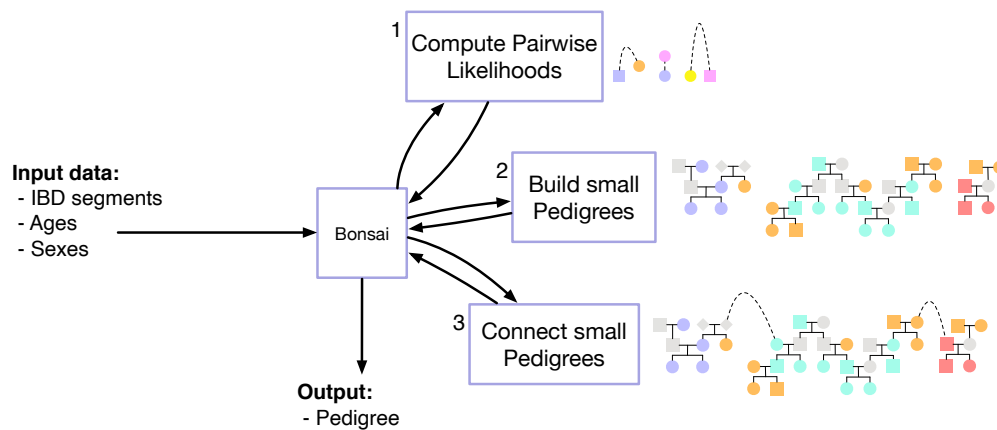


FIGURE 1. **Overview of the full Bonsai method.** Details of methods 1, 2, and 3 are presented in Algorithms 1, 2, and 4, respectively.

85

### 3. SUBJECTS AND METHODS

86 **3.1. Overview of the Bonsai method.** The Bonsai method is summarized in Figure 1. The  
87 input to the method consists of ages and sexes for a set of putatively-related individuals, along with  
88 IBD segments inferred between each pair of individuals. The method then proceeds through three  
89 stages in sequence.

90 First, the relationship between each pair of genotyped individuals is inferred using age and  
91 pairwise IBD data. The likelihoods of many other possible relationships are also computed and  
92 stored for each pair. Next, small pedigrees of closely-related individuals are inferred from these  
93 pairwise likelihoods. Finally, the inferred small pedigrees are assembled into large and sparse  
94 pedigrees.

95 Constructing small pedigrees and combining them together allows us to make use of information  
96 in small pedigree structures to improve the accuracy with which more distant relationships are

97 inferred. This approach allows us to more precisely infer the ancestral lineages through which small  
98 pedigrees are connected, the number of common ancestors shared by each pair of individuals, and  
99 segments of so-called background IBD that do not reflect recent ancestry. Each of these additional  
100 pieces of information makes it possible to proactively reduce the space of possible pedigrees that  
101 must be searched, making inference tractable for large and sparse pedigrees.

**3.2. Stage 1: inferring pairwise relationships.** The first stage of the Bonsai method is to infer the likelihoods of many possible relationships between each pair of putative relatives. To make the computation of the likelihood efficient without large sacrifices in accuracy, we use a composite likelihood that is the product of the likelihoods of different IBD summary statistics and the likelihoods of the pairwise age differences between the individuals. The genetic component  $\mathcal{L}_{\mathcal{R}}^g$  of the likelihood, computed from IBD, is multiplied by the age component  $\mathcal{L}_{\mathcal{R}}^a$  of the likelihood to obtain the final likelihood  $\mathcal{L}_{\mathcal{R}}$  of a given relationship type,  $\mathcal{R}$ :

$$\mathcal{L}_{\mathcal{R}} = \mathcal{L}_{\mathcal{R}}^g \mathcal{L}_{\mathcal{R}}^a. \quad (1)$$

102 The likelihood is composite, rather than exact, because we do not model the joint distribution of the  
103 IBD count and length summary statistics whose product is  $\mathcal{L}_{\mathcal{R}}^g$  and because there is an underlying  
104 joint distribution of IBD sharing and age difference that is not captured by the product of the two  
105 likelihoods  $\mathcal{L}_{\mathcal{R}}^g$  and  $\mathcal{L}_{\mathcal{R}}^a$ .

106 **3.2.1. Pairwise genetic likelihoods.** To compute the genetic component of the composite pairwise  
107 relationship likelihood, we consider regions of the genome shared identically by descent in a haploid  
108 fashion on just one chromatid in each individual, as well as regions shared IBD in a diploid fashion  
109 on both chromatids. We use the terms “IBD1 segment” and “IBD2 segment” to refer to regions of  
110 haploid and diploid IBD, respectively. The genetic component of the pairwise likelihood is computed  
111 using the total length of IBD1 segments, the total length of IBD2 segments, the total number of  
112 IBD1 segments, and the total number of IBD2 segments.

113 It is possible to compute the probability of an observed shared pattern of IBD analytically.  
114 However, in practice we find that error in IBD inference leads to differences between the empirical  
115 and analytical IBD distributions for each relationship type, especially for close relationships. Thus,  
116 we use likelihoods obtained as moment-fitted Poisson and Gaussian approximations of simulated  
117 distributions.

Let  $T_1^{i,j}$  and  $T_2^{i,j}$  be the total lengths of IBD 1 and 2, respectively for a pair of individuals,  $i$  and  $j$  and let  $C_1^{i,j}$  and  $C_2^{i,j}$  be the counts of the number of IBD 1 and 2 segments shared between two individuals. We follow the convention that uppercase variables  $T_1, T_2, C_1, C_2$ , etc. denote random variables and their lowercase counterparts,  $t_1, t_2, c_1, c_2$ , etc. denote their observed values. The genetic component of the composite likelihood for a given relationship type,  $\mathcal{R}$ , between a pair of individuals,  $i$  and  $j$ , is then computed as

$$\mathcal{L}_{\mathcal{R}}^g(i, j) \approx f_{\mathcal{R}}(t_1) f_{\mathcal{R}}(t_2) \mathbb{P}_{\mathcal{R}}(c_1) \mathbb{P}_{\mathcal{R}}(c_2), \quad (2)$$

118 where  $f_{\mathcal{R}}(t_1) \equiv f_{T_1^{i,j}}(t_1; \mathcal{R})$  is the probability density function of the sum of lengths of all IBD 1  
119 segments for a relationship of type  $\mathcal{R}$  and  $\mathbb{P}_{\mathcal{R}}(c_1) \equiv \mathbb{P}(C_1 = c_1; \mathcal{R})$  is the probability mass function

120 for the total number of segments of IBD1 for a relationship of type  $\mathcal{R}$ . The quantities  $f_{\mathcal{R}}(t_2)$  and  
121  $\mathbb{P}_{\mathcal{R}}(c_2)$  are defined analogously for segments of IBD 2.

122 In Equation (2), the quantities  $f_{\mathcal{R}}(t_1)$  and  $f_{\mathcal{R}}(t_2)$  are modeled as Gaussian distributions and the  
123 distributions  $\mathbb{P}_{\mathcal{R}}(c_1)$  and  $\mathbb{P}_{\mathcal{R}}(c_2)$  are Poisson with means given by the expected numbers of IBD1  
124 and IBD2 segments, respectively between two individuals of relationship type  $\mathcal{R}$ . The mean and  
125 variance of  $T_i^{\mathcal{R}}$ , and the mean of  $C_i^{\mathcal{R}}$  for a relationship of type  $\mathcal{R}$  were obtained empirically using  
126 simulations. Details of the simulations used to obtain these moments are provided in Section 3.6.4.

127 **3.2.2. Pairwise age likelihoods.** The pairwise age likelihood for a given relationship type,  $\mathcal{R}$ , was  
128 obtained by moment-fitting a Gaussian distribution to the differences between the ages of 23andMe  
129 customers who self-reported to be of relationship type  $\mathcal{R}$  (Figure 2). We required that the self-  
130 reported relationship between each pair of individuals could be verified through a string of inferred  
131 parent-child or full-sibling relationships. For example, a self-reported first-cousin relationship  
132 between individuals  $i$  and  $j$  was verified if  $i$  and  $j$  each had inferred parents in the 23andMe  
133 database, and if these parents in turn had the same pair of inferred parents, or were inferred to be  
134 full siblings.

For two customers,  $i$  and  $j$ , with ages  $a_i$  and  $a_j$ , the age component of the likelihood for relationship  
type  $\mathcal{R}$  was modeled as a Gaussian distribution with the empirical mean and variance:

$$\mathcal{L}_{\mathcal{R}}^a(i, j) = \frac{e^{-[(a_i - a_j) - \mu_a^{\mathcal{R}}]^2 / 2(\sigma_a^{\mathcal{R}})^2}}{\sigma_a^{\mathcal{R}} \sqrt{2\pi}}. \quad (3)$$

135 In Equation (3),  $\mu_a^{\mathcal{R}}$  and  $\sigma_a^{\mathcal{R}}$  are the moment-fitted mean and standard deviation of the empirical  
136 age difference for all pairs of customers who reported the pairwise relationship,  $\mathcal{R}$ . Note that the  
137 probability  $\mathcal{L}_{\mathcal{R}}^a(i, j)$  is not symmetrical in the ages,  $a_i$  and  $a_j$ . This is useful for determining the  
138 directionality of the relationship between two people, such as parent-child or nephew-aunt when age  
139 information is available.

**3.3. The likelihood of a pedigree.** The composite likelihood,  $\mathcal{L}_{\mathcal{P}}$ , of a pedigree  $\mathcal{P}$  is computed  
as the product of genetic and age likelihoods (Equation 1) for all pairs of individuals in the pedigree,

$$\mathcal{L}_{\mathcal{P}} = \prod_{i, j \in \mathcal{P}} \mathcal{L}_{\mathcal{R}}^g(i, j) \mathcal{L}_{\mathcal{R}}^a(i, j). \quad (4)$$

140 where  $\mathcal{R}$  is the relationship between  $i$  and  $j$  implied by the pedigree structure. This likelihood is  
141 efficiently computed as each new individual is added to the pedigree. By doing so, we can inductively  
142 extend the existing relationships of the parents and/or children of the newly-added person to obtain  
143 the relationships of the new person to all existing individuals in the pedigree. We then add the log  
144 likelihoods of each of these new pairwise relationships to the log likelihood of the pedigree without  
145 the new individual.

146 **3.4. The “Small” Bonsai method.** To construct a pedigree from pairwise likelihoods, the Small  
147 Bonsai method begins by placing a focal individual by itself in the pedigree (Figure 3). This  
148 focal individual is typically the person with the closest average degree of relationship to all other  
149 individuals in the putatively-related set, but any individual can be chosen. At each subsequent  
150 step of the Small Bonsai algorithm, the next individual to be placed is chosen to be the unplaced

6 BONSAI: AN EFFICIENT METHOD FOR INFERRING LARGE HUMAN PEDIGREES FROM GENOTYPE DATA

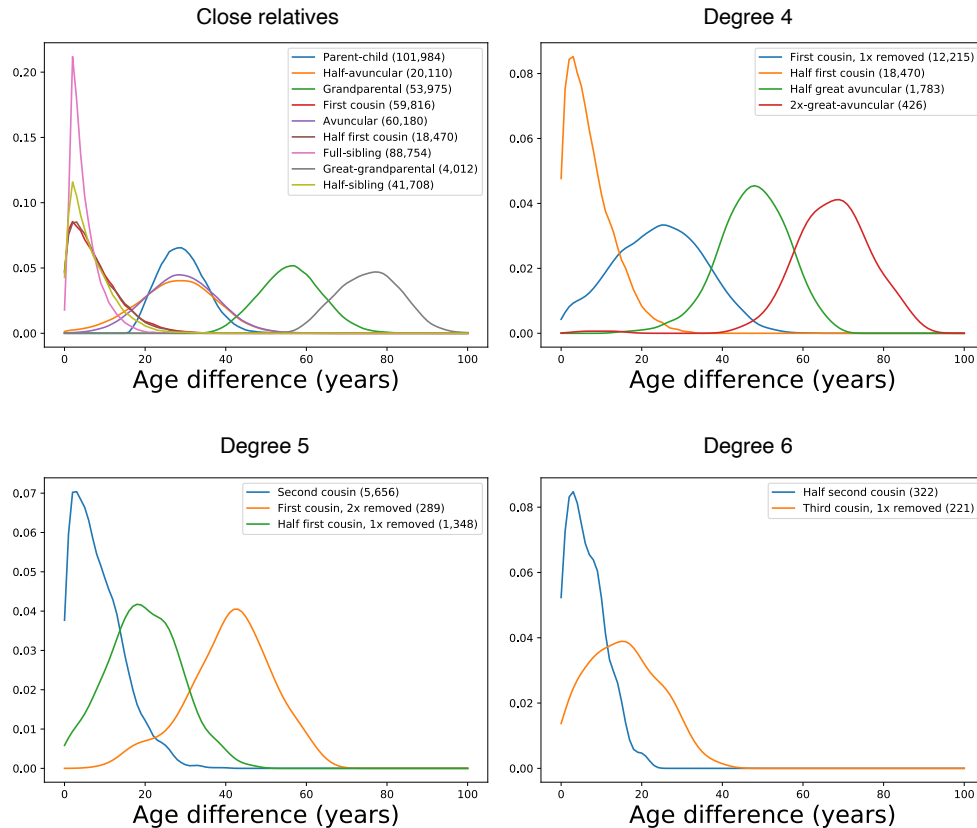


FIGURE 2. **Empirical age difference distributions.** Kernel densities for the absolute difference in age between a pair of relatives of a given type. The number of pairs of each type used in the analysis is given in parentheses. Different panels show relationships of different degrees.

151 individual with the closest inferred degree of relationship with one of the individuals already placed  
 152 in the pedigree, where ties are broken by the total amount of IBD shared. Because each pair of  
 153 individuals has many possible relationships, we determine the order in which individuals are added  
 154 using the most likely pairwise relationship for each pair.

155 The next individual to be placed is considered in all ways that are consistent with the most  
 156 likely inferred pairwise relationships to individuals already placed. In particular, for a user-specified  
 157 parameter  $r$ , we consider the top  $r$  most likely pairwise relationships between the new individual  
 158 and their closest relative in the set of placed individuals and we place the individual in all ways  
 159 that are compatible with each of these  $r$  most-likely relationships. When adding an individual to  
 160 the pedigree, we must not only add them in all possible ways to a particular pedigree, we also add  
 161 them in all  $r$  ways to all high-likelihood pedigrees that were formed at the previous step.

162 When two or more pedigrees formed by adding an individual would be topologically identical, we  
 163 only construct one of the pedigrees. For example, in the second row of Figure 3, because the sexes  
 164 of the parents are unknown and there are no placed relatives except the focal individual that can  
 165 be used for triangulation, adding an avuncular relative through the right parent is topologically

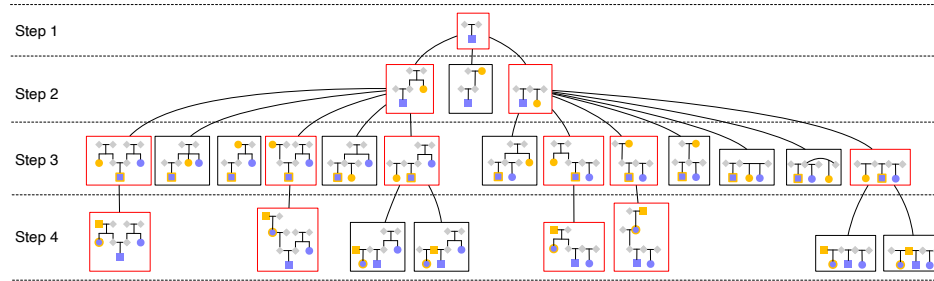


FIGURE 3. **The Small Bonsai method.** An example of the sequence of steps for building a small pedigree is shown. The sequence proceeds from top to bottom in the figure. The  $i$ th row of pedigrees in rectangles represents the  $i$ th step of the Small Bonsai algorithm in which the  $i$ th individual is added to a pedigree. Red boxes indicate pedigrees that are retained and carried forward to the next step. Black boxes indicate pedigrees with low likelihoods that are discarded.

166 identical to adding them through the left parent. Therefore, we only build one of these pedigrees  
167 (the one on the far left of the second row).

168 To avoid a rapid expansion in the number of pedigrees at each step, we employ a heuristic  
169 branch-and-bound-like procedure in which we discard each pedigree at the end of each step that is  
170 very unlikely, compared with the most likely pedigree. In particular, we discard all pedigrees whose  
171 likelihoods are less than a fraction  $f_\ell$  of the likelihood of the most likely pedigree, where a pedigree's  
172 likelihood is the product over the likelihoods of the pairwise relationships implied by the pedigree  
173 (Section 3.3). In practice, when individuals are closely related, there are only a few pedigrees that  
174 have high likelihoods and the rest can be discarded. As a result, the likelihood threshold has a low  
175 impact on accuracy while serving to dramatically speed up pedigree building.

176 This heuristic branch-and-branch-and-bound-like procedure is repeated until no un-  
177 placed individual has a pairwise point-estimated degree that is within a user-specified degree  $d$   
178 of any placed individual. At this point, the Small Bonsai algorithm is terminated. If unplaced  
179 individuals remain, a new focal individual is chosen from among the unplaced individuals and the  
180 Small Bonsai algorithm is applied again. The Small Bonsai algorithm is applied repeatedly, choosing  
181 a new focal individual each time, until all individuals have been placed into some pedigree.

182 Figure 3 shows an example sequence for constructing a pedigree using the Small Bonsai method.  
183 In the first row of the figure, a focal individual (shaded yellow square) is placed into a pedigree  
184 on their own. Grey diamonds indicate their parents, whose sexes are unspecified. In the second  
185 row, the unplaced individual with the closest degree of relationship to the placed individual, is  
186 placed into the pedigree (yellow circle). The new individual is placed in all ways that are consistent  
187 with the top  $r$  most likely relationships inferred in the pairwise relationship inference step (Section  
188 3.2). Here, we have chosen  $r = 3$ . These  $r = 3$  most likely relationships happen to be “avuncular,”  
189 “grandparental,” and “half-sibling” in the example shown. This is the “branch” step of the heuristic  
190 branch-and-bound-like procedure.

191 Before placing the next individual, we evaluate the likelihood of each pedigree, computed as the  
192 product of pairwise likelihoods of the relationships induced by the pedigree. We retain only those

193 pedigrees whose likelihoods are at least a fraction  $f_\ell$  of the likelihood of the most likely pedigree.  
194 This is the “bound” step of the heuristic branch-and-bound-like procedure.

195 In the third row of the diagram, the unplaced individual (yellow circle) with the closest degree  
196 of relationship to a placed individual is added to all pedigrees that were carried forward from the  
197 previous step. The new individual is added to each pedigree in all ways that are consistent with  
198 the top  $r$  most likely relationships to their closest placed relative (purple square with a yellow  
199 boundary). Again, these relationships happen to be “avuncular,” “grandparent,” and “half-sibling”  
200 in the example. We then perform the “bound” step, retaining only those pedigrees whose likelihoods  
201 are at least a fraction  $f_\ell$  of the likelihood of the most likely pedigree.

202 In the fourth row, we show one final iteration of the procedure. Again, the unplaced individual  
203 (yellow circle) is added in all ways that are consistent with the top  $r$  most likely pairwise point  
204 estimated relationships with their closest relative (purple circle with yellow a boundary). In this  
205 case the most likely point-estimated relationship happens to be “parental.” Because parent-child  
206 relationships are inferred with near certainty, we have only placed the individual as a parent in the  
207 diagram, omitting the next 2 most likely relationships which will be considerably less likely.

### 208 3.5. The “Big” Bonsai method.

209 3.5.1. *Overview of the Big Bonsai method.* When building a pedigree containing distantly-related  
210 individuals, the Small Bonsai method is first applied repeatedly to build sets of small non-overlapping  
211 pedigrees. The union of individuals in these small pedigrees is equal to the set of individuals in the  
212 full pedigree. The Big Bonsai method is then applied to combine the small pedigrees together, one  
213 pair at a time, with the two pedigrees sharing the most total IBD combined at each step.

214 The Big Bonsai method relies on several new methods we introduce that are useful for different  
215 aspects of combining pedigrees together. The first method is a generalized version of the DRUID  
216 estimator (Ramstetter et al. 2018) for inferring the degree of relatedness separating the common  
217 ancestors of two small pedigrees. The DRUID estimator was derived for specific pedigree structures,  
218 such as a set of siblings and their avuncular relatives connected to another such pedigree through  
219 the common grandparental ancestors of the two pedigrees. Here, we generalize the DRUID estimator  
220 to any pair of outbred pedigrees and, in Appendix 6.3, we further generalize the DRUID estimator  
221 to the case in which two pedigrees are connected through two individuals who are not the common  
222 ancestors of their respective pedigrees.

223 The second tool we introduce is an approximation of the likelihood of the degree separating  
224 two pedigrees, as a function of the total IBD shared between the two pedigrees. This likelihood,  
225 which was inspired by the DRUID estimator, makes it possible to evaluate the relative likelihoods  
226 of different degrees separating two pedigrees in addition to obtaining a point estimate of the degree.

227 The third tool we introduce is a new test for detecting segments of background IBD. Background  
228 IBD segments are regions of the genome that are shared identically-by-state (IBS) and which did  
229 not arise by transmission from a single shared common ancestor. Instead, these segments arose  
230 because of demographic or evolutionary processes, such as a population bottleneck. They are long  
231 regions of IBS with hidden recombination events and they can provide misleading information about



TABLE 1. **Variable definitions.**

Variable	Definition
$\mathcal{R}$	A specific relationship type (e.g., parent-child).
$\mathcal{L}_{\mathcal{R}}$	Likelihood of relationship $\mathcal{R}$ .
$\mathcal{L}_{\mathcal{R}}^g$	Genetic component of the likelihood of relationship $\mathcal{R}$ .
$\mathcal{L}_{\mathcal{R}}^a$	Age component of the likelihood of relationship $\mathcal{R}$ .
$C_i$	Count of segments of IBD of type $i$ , where $i \in \{1, 2\}$ .
$T_i$	Total length of IBD of type $i$ , where $i \in \{1, 2\}$ .
$\mu_a^{\mathcal{R}}$	Mean age difference for two individuals with relationship $\mathcal{R}$ .
$\sigma_a^{\mathcal{R}}$	Standard deviation of the age difference for two individuals with relationship $\mathcal{R}$ .
$a_i$	Age of individual $i$ .
$\mathcal{P}$	A pedigree.
$\mathcal{N}$ (or $\mathcal{N}_A, \mathcal{N}_{\mathcal{P}}, \mathcal{N}_S$ )	A set of individuals (corresponding to ancestor $A$ , pedigree $\mathcal{P}$ , or pedigree set $S$ ).
$A$ (or $A_{\mathcal{N}}, A_i$ )	A specific common ancestor (of $\mathcal{N}, \mathcal{N}_i$ ).
$\mathcal{A}$ (or $\mathcal{A}_{\mathcal{N}}, \mathcal{A}_i$ )	Set of ancestors (of $\mathcal{N}, \mathcal{N}_i$ ).
$\Lambda$ (or $\Lambda_i$ )	Induced subtree relating a set of nodes $\mathcal{N}$ or $\mathcal{N}_i$ .
$d_{i,j}$	True genetic degree separating individuals $i$ and $j$ .
$d_L(i, j)$	Maximum likelihood estimate of the degree between individuals $i$ and $j$ .
$d_D(i, j)$	Generalized DRUID estimate of the degree between individuals $i$ and $j$ .
$G$	Set of common ancestors connecting two individuals or pedigrees.
$O_i$	The event that IBD is observed in individual $i$ .
$p_{i,0}$	The probability that an allele is not observed at a specified locus in individual $i$ .
$p_{i,1}$	The probability that an allele is observed at a specified locus in individual $i$ .
$T_{i,j}$	The total length of IBD observed between sets $\mathcal{N}_i$ and $\mathcal{N}_j$ .
$L_{i,j}$	Length of a single merged segment of IBD observed between sets $\mathcal{N}_i$ and $\mathcal{N}_j$ .
$\mathcal{I}$	The event that an ancestrally transmitted allele is shared IBD between sets $\mathcal{N}_i$ and $\mathcal{N}_j$ .
$L_{genome}$	Length of the genome in centimorgans.
$C_i$	The set of children of node $i$ .
$r$	Number of most likely relationships considered for each new individual in Small Bonsai.
$f_{\ell}$	Fraction of likelihood of most likely pedigree below which we discard a pedigree.
$(d_1, d_2, n)$	Degree tuple of the form (up, down, number common ancs) (Ko and Nielsen 2017).
<i>Degree</i>	Defined as $d_1 + d_2 - n + 1$ .

232 the degree of relationship between a pair of individuals. Background IBD segments can lead to  
 233 mis-inferred pedigrees, particularly when pedigrees are sparsely genotyped.

234 The fourth tool we introduce is a method for determining the correct ancestral lineages through  
 235 which two or more pedigrees are connected. This approach relies on detecting overlapping IBD  
 236 segments that are inconsistent with certain lineage combinations.

237 We also derive a recursive formula for computing the probability of an observed presence-absence  
 238 pattern of an ancestrally transmitted allele in their descendants. This formula is useful for developing  
 239 the generalized DRUID estimator and the likelihoods for degree estimation and background IBD  
 240 detection.

241 Together, the new tools we introduce can be used to identify the ancestors through which two  
 242 small pedigrees are connected, infer the degree separating the two ancestors, and identify and discard  
 243 individuals whose IBD sharing patterns appear to be background IBD. By using these inference

10BONSAI: AN EFFICIENT METHOD FOR INFERRING LARGE HUMAN PEDIGREES FROM GENOTYPE DATA

244 tools to identify highly-likely ways of connecting pedigrees, the space of possible pedigrees can be  
 245 considerably reduced. We now describe each of these approaches in detail.

246 3.5.2. *The probability of a presence-absence pattern of an ancestral allele.* Consider two pedigrees  
 247  $\mathcal{P}_1$  and  $\mathcal{P}_2$  of genotyped individuals,  $\mathcal{N}_1$  and  $\mathcal{N}_2$ , related through a common ancestor (or pair of  
 248 ancestors),  $G$  (Figure 4). Let  $A_1$  be the common ancestor of  $\mathcal{N}_1$  in  $\mathcal{P}_1$  and let  $A_2$  be the common  
 249 ancestor of  $\mathcal{N}_2$  in  $\mathcal{P}_2$ .

250 Consider an allele transmitted from one chromatid in  $G$  to its descendants. We begin by deriving  
 251 the probability of the observed pattern of presence and absence of the ancestral allele among  
 252 descendants of  $A_1$  and  $A_2$ . Let  $d_{A_1,G}$  and  $d_{A_2,G}$  be the degrees separating  $A_1$  and  $A_2$  from the set  
 253 of most recent common ancestors,  $G$ , of the pedigree.  $G$  corresponds to two individuals if  $A_1$  and  
 254  $A_2$  are descended from an ancestral couple and  $G$  corresponds to a single common ancestor if  $A_1$   
 255 and  $A_2$  are descended from a pair of half siblings. We do not consider cases of endogamy, where  $G$   
 256 corresponds to more than one ancestor other than a mate pair.

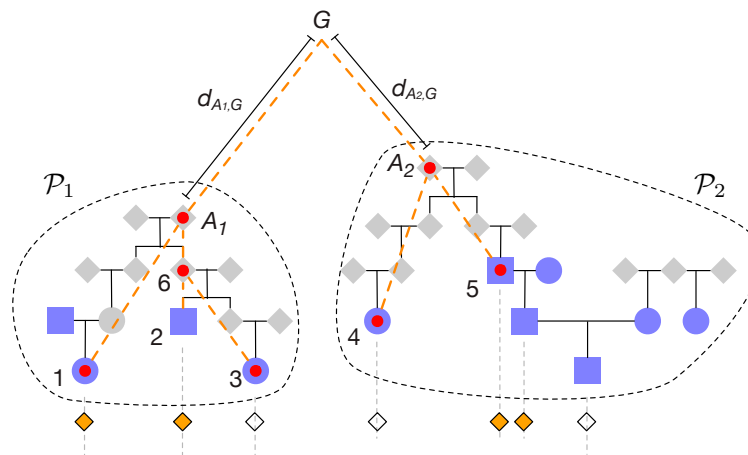


FIGURE 4. **Example of an observed pattern of presence and absence of an ancestral allele.** Genotyped individuals are shaded in purple. Filled and empty diamonds below indicate the presence or absence of the allele from  $G$ . Red dots on purple genotyped individuals indicate the set of genotyped individuals with no direct genotyped ancestors. Red dots on grey ungenotyped individuals indicate the most recent common ancestors transmitting the segments to the genotyped individuals. Dashed orange lines indicate the paths by which the allele is transmitted from common ancestor  $G$ . The number of meioses separating  $A_1$  and  $A_2$  from a common ancestor,  $G$ , are  $d_{A_1,G}$  and  $d_{A_2,G}$ .

257 Figure 4 shows a presence-absence pattern of an inherited allele among genotyped individuals in  
 258 the two small pedigrees  $\mathcal{P}_1$  and  $\mathcal{P}_2$ . The probability of the observed presence and absence pattern  
 259 can be computed recursively by conditioning on whether the allele was observed in the ancestor of  
 260 each individual. This approach is similar to Felsenstein's tree pruning algorithm (Felsenstein 1981).

261 Let  $O_i$  be a random variable describing the event that a copy of the allele is transmitted to  
 262 descendant  $i$  and is observed. We set  $O_i = 1$  if the allele is observed in individual  $i$  and  $O_i = 0$  if it

263 is not observed. The probabilities  $\mathbb{P}(O_i = 0)$  and  $\mathbb{P}(O_i = 1)$  can be computed by conditioning on  
264 whether the allele in  $G$  was observed in the node of the induced subtree immediately ancestral to  $i$ .

Defining

$$p_{i,0} \equiv \mathbb{P}(O_i = 0), p_{i,1} \equiv \mathbb{P}(O_i = 1), \quad (5)$$

we show in Appendix 6.1 that the probabilities can be computed using the recursion

$$\begin{aligned} p_{i,0} &= [1 - 2^{-d_{i,a(i)}}]p_{a(i),1} + p_{a(i),0}, \\ p_{i,1} &= 2^{-d_{i,a(i)}}p_{a(i),1}, \end{aligned} \quad (6)$$

265 with the base conditions  $p_{g,0} = 0$  and  $p_{g,1} = 1$  for each allelic copy,  $g$ , in  $G$ . The probability of an  
266 observed IBD sharing pattern  $\{O_1, \dots, O_k\}$  across  $k$  leaf nodes can be computed recursively using  
267 Equation (6).

268 3.5.3. *The generalized DRUID estimator.* Ramstetter et al. (2018) developed a method for inferring  
269 degrees of relatedness among distant relatives. The method addresses the problem that the amount  
270 of IBD shared between two individuals decreases exponentially with their degree of relatedness,  
271 resulting in very little information for inferring degrees between distant relatives. In fact, there can  
272 be a non-negligible probability that distant relatives will share no IBD segments at all, especially if  
273 information contained in short IBD segments is discarded to reduce the rate of false positive IBD  
274 segments.

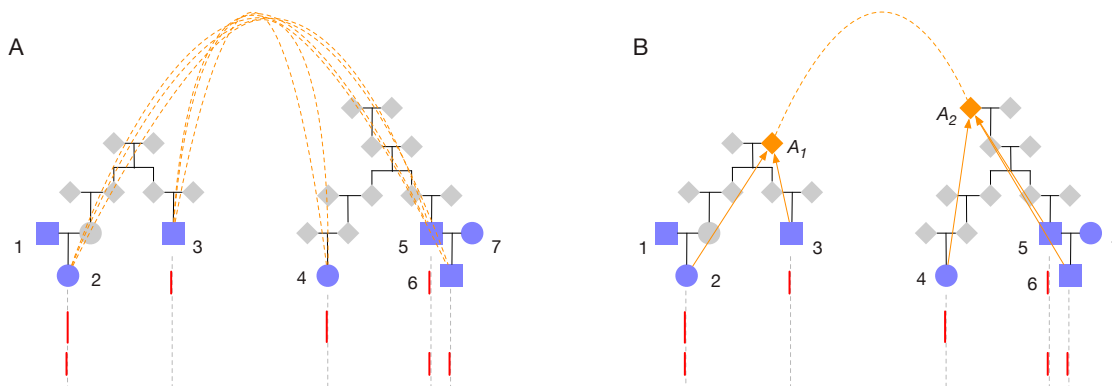
275 Because two genealogically-related individuals may share little or no IBD, it is helpful to leverage  
276 IBD segments shared among close relatives of the two individuals when inferring their degree of  
277 relatedness. Figure 5 illustrates the utility of considering IBD segments among groups of individuals  
278 rather than pairwise IBD when the degree of relatedness is not small. In particular, individuals 3  
279 and 4 in Figure 5 share no IBD segments. Thus, one cannot infer their degree of relatedness without  
280 additional information. However, if close relatives of 3 and 4 do share IBD with one another, and if  
281 pedigrees can be inferred relating these close relatives to 3 and 4, then we can use the IBD in these  
282 relatives to estimate the degree of relationship between 3 and 4.

283 Leveraging IBD shared by close relatives has the effect of increasing the amount of available data  
284 for inferring pairwise relationships. Ramstetter et al. (2018) demonstrated that an approach based  
285 on summarizing IBD from close relatives can greatly improve the accuracy of estimates of distant  
286 degrees of relatedness compared with the approach of computing a composite likelihood over all  
287 pairs of individuals (Staples et al. 2016). These two approaches are shown in Figure 5.

288 Let  $\mathcal{N}_1$  and  $\mathcal{N}_2$  be two sets of genotyped individuals; for example, sets  $\mathcal{N}_1 = \{2, 3\}$  and  $\mathcal{N}_2 =$   
289  $\{4, 5, 6\}$  in Figure 5B. Let  $A_1$  and  $A_2$  be the most recent common ancestors of  $\mathcal{N}_1$  and  $\mathcal{N}_2$ , respectively  
290 and let  $d(A_1, A_2)$  denote the degree between  $A_1$  and  $A_2$ . The DRUID estimator of  $d(A_1, A_2)$  derived  
291 by Ramstetter et al. is obtained by first merging all IBD segments observed between  $\mathcal{N}_1$  and  $\mathcal{N}_2$ .  
292 The total merged IBD is then converted into a point estimate of the amount of IBD shared between  
293 the common ancestor  $A_1$  and the common ancestor  $A_2$ .

294 The amount of IBD shared between  $A_1$  and  $A_2$  is estimated by considering the fraction  $\varphi_1$  of  
295 the genome of  $A_1$  that is passed on to its genotyped descendants in  $\mathcal{N}_1$  and the fraction  $\varphi_2$  of  
296 the genome of  $A_2$  that is passed on to its genotyped descendants in  $\mathcal{N}_2$ . If  $IBD(A_1, A_2)$  is the

12BONSAI: AN EFFICIENT METHOD FOR INFERRING LARGE HUMAN PEDIGREES FROM GENOTYPE DATA



**FIGURE 5. Leveraging IBD from close relatives to infer the degree of relatedness between individuals.** Each panel in the figure shows a comparison between two pedigrees. Purple-shading indicates individuals who have been genotyped. Red lines indicate IBD shared between a genotyped individual in the pedigree containing 1, 2, and 3 and a genotyped individual in the pedigree containing 4, 5, 6, and 7. Orange dashed lines in Panel (A) indicate pairwise degrees of relatedness among all cross-pedigree pairs. The orange dashed line in Panel (B) indicates the degree of relatedness between the common ancestor  $A_1$  of the IBD-carrying individuals in the left pedigree and the common ancestor  $A_2$  of the IBD-carrying individuals in the right pedigree. In Panel (A), pairwise IBD is summarized to infer the degree separating the two pedigrees. In Panel (B), IBD information in the descendants of  $A_1$  and  $A_2$  is summarized to infer the degree of relatedness among these two common ancestors. The approach in Panel (A) is taken by the PADRE method (Staples et al. 2016). The approach in Panel (B) is taken by the DRUID method (Ramstetter et al. 2018).

297 amount of IBD shared between  $A_1$  and  $A_2$ , then the expected amount shared between  $\mathcal{N}_1$  and  
 298  $\mathcal{N}_2$  is  $IBD(\mathcal{N}_1, \mathcal{N}_2) = \varphi_1 \varphi_2 IBD(A_1, A_2)$ . Solving for  $IBD(A_1, A_2)$  yields a point estimator of  
 299  $IBD(A_1, A_2)$  in terms of the observed quantity  $IBD(\mathcal{N}_1, \mathcal{N}_2)$ .

300 Ramstetter et al. (2018) derived formulas for  $\varphi_1$  and  $\varphi_2$  for specific pedigree configurations, such  
 301 as sets of siblings or siblings together with avuncular relatives. Here, we generalize the DRUID  
 302 estimator to arbitrary outbred pedigrees and further generalize the method to include the case  
 303 in which  $A_1$  is descended from a descendant  $A_2$  who is ancestral to only a subset of genotyped  
 304 individuals in  $\mathcal{N}_2$ .

The fraction  $\varphi_i$  of the genome of  $A_i$  that is passed on to some descendant in  $\mathcal{N}_i$  can be computed as

$$\varphi_i = 1 - \prod_{n \in \mathcal{N}_i} p_{n,0}, \quad (7)$$

where the quantities  $p_{n,0}$  are computed recursively using Equation (6). Thus, an estimate of the amount of IBD shared between  $A_1$  and  $A_2$  is

$$\widehat{IBD}(A_1, A_2) = \frac{IBD(\mathcal{N}_1, \mathcal{N}_2)}{\varphi_1 \varphi_2}. \quad (8)$$

Using the expression  $\hat{\phi} = \widehat{IBD}(A_1, A_2)/4L_{genome}$  for the kinship coefficient when all IBD is of type 1, we obtain the generalized DRUID estimator

$$d_D(A_1, A_2) = d : \frac{1}{2^{d+3/2}} \leq \frac{IBD(\mathcal{N}_1, \mathcal{N}_2)}{4\varphi_1\varphi_2L_{genome}} < \frac{1}{2^{d+1/2}}, \quad (9)$$

where the bounds come from Manichaikul et al. (2010) and are the ones used for the DRUID estimator presented in Ramstetter et al. (2018).

In Appendix 6.3, we demonstrate how the DRUID estimator can be further generalized to the case in which  $A_1$  is descended from one of the individuals in  $\mathcal{N}_2$ , or from an internal node of the induced subtree that is a descendant of  $A_2$ . Thus, we obtain a version of the DRUID estimator that can be applied to general outbred pedigrees.

3.5.4. *The likelihood of the degree of relatedness among groups of individuals.* Using the DRUID principle, we can develop a likelihood estimator of the pairwise degree of relatedness between the common ancestors  $A_1$  and  $A_2$ , given the observed total IBD  $T_{1,2}$  between the genotyped descendants of  $A_1$  and  $A_2$ .

Consider again the scenario depicted in Figure 4 in which two sets of genotyped individuals,  $\mathcal{N}_1$  and  $\mathcal{N}_2$ , are related through a common ancestor or pair of ancestors,  $G$ . The probability that a given allele from  $G$  is observed IBD between  $\mathcal{N}_1$  and  $\mathcal{N}_2$  can be obtained by conditioning on the events that it is observed in  $A_1$  and  $A_2$ . Let  $\mathcal{I}$  denote the event that the allele is observed IBD. Then

$$\begin{aligned} \mathbb{P}(\mathcal{I}) &= \varphi_1\mathbb{P}(O_{A_1} = 1)\varphi_2\mathbb{P}(O_{A_2} = 1) \\ &= \varphi_1\varphi_22^{-(d_{A_1,G}+d_{A_2,G})}, \end{aligned} \quad (10)$$

where  $\varphi_i$  is computed using Equation (7).

If  $A_1$  and  $A_2$  had exactly one common ancestor with one allele to transmit, then Equation (10) would be the fraction of the genome in which we expect to find some segment shared IBD between some member of  $\mathcal{N}_1$  and some member of  $\mathcal{N}_2$ . However, we must now account for the fact that each common ancestor of  $A_1$  and  $A_2$  in  $G$  carries two allelic copies and that there can be either one or two such common ancestors.

Let  $|G|$  denote the number of common ancestors of  $A_1$  and  $A_2$ , each of which carries two alleles at the locus of interest. Let  $\mathcal{I}^c$  denote the complement of event  $\mathcal{I}$ , i.e., the event that a specific allele from  $G$  is not observed IBD. Thus, we have

$$\mathbb{P}(\mathcal{I}^c) = 1 - \mathbb{P}(\mathcal{I}). \quad (11)$$

Then the probability that none of the  $2|G|$  alleles is observed IBD is  $\mathbb{P}(\mathcal{I}^c)^{2|G|}$ , and the probability that at least one of the alleles is observed is  $1 - \mathbb{P}(\mathcal{I}^c)^{2|G|}$ .

We can use the probability of observing an allele IBD to obtain an approximate likelihood of the total length  $T_{1,2}$  of IBD observed between descendants of  $A_1$  and  $A_2$ . The mean of this distribution is simply the expected length of the genome in a state of IBD between the two clades, which is

$$E[T_{1,2}] = (1 - \mathbb{P}(\mathcal{I}^c)^{2|G|})L_{genome}. \quad (12)$$

An approximation of the variance of  $T_{1,2}$  is derived in Section 6.2 and is given by

$$\text{Var}(T_{1,2}) \approx (1 - \mathbb{P}(\mathcal{I}^c)^{2|G|})L_{genome} \frac{E[L_{1,2}^2]}{E[L_{1,2}]}, \quad (13)$$

326 where  $L_{1,2}$  is the length of any given IBD segment between  $A_1$  and  $A_2$  formed by merging all IBD  
 327 segments between leaf nodes in  $A_1$  and  $A_2$  that overlap one another. The moments  $E[L_{1,2}^m]$  are  
 328 derived in Appendix 6.2 and can be computed using Equation (28) or (29).

If the segments,  $L_{1,2}$  were each exponentially distributed, then  $T_{1,2}$  would have a gamma distribution. In practice, a gamma distribution is an accurate approximation for the distribution of  $T_{1,2}$ , given that the length  $T_{1,2}$  is greater than zero. Thus, we can approximate the distribution of  $T_{1,2}$  by

$$T_{1,2}|T_{1,2} > 0 \sim \text{Gamma}(k_{1,2}, \theta_{1,2}),$$

where  $k_{1,2}$  and  $\theta_{1,2}$  are found by matching the mean and variance of the gamma distribution with  $E[T_{1,2}]$  and  $\text{Var}(T_{1,2})$ . Thus, we obtain

$$T_{1,2}|T_{1,2} > 0 \sim \text{Gamma}\left(\frac{E[L_{1,2}]^2}{\text{Var}(L_{1,2})}, \frac{\text{Var}(L_{1,2})}{E[L_{1,2}]}\right), \quad (14)$$

329 where  $E[L_{1,2}]$  and  $E[L_{1,2}^2]$  are given by Equation (29).

If every IBD segment has some length, we can assume that  $T_{1,2}$  is only identically zero when there are no IBD segments. The distribution of the number of segments can be modeled as a Poisson random variable with mean  $E[N_{1,2}]$  equal to the expected number  $N_{1,2}$  of merged segments shared between  $\mathcal{N}_1$  and  $\mathcal{N}_2$ . The probability that there are no segments is then  $e^{-E[N_{1,2}]}$ . Thus, we have the approximation

$$f_{T_{1,2}}(t_{1,2}) \approx \begin{cases} \frac{t_{1,2}^{k-1}}{\Gamma(k)\theta^k} e^{-t_{1,2}/\theta} (1 - e^{-E[N_{1,2}]}) & \text{if } t_{1,2} > 0 \\ e^{-E[N_{1,2}]} & \text{if } t_{1,2} = 0. \end{cases}, \quad (15)$$

330 where  $k = E[L_{1,2}]^2/\text{Var}(L_{1,2})$ ,  $\theta = \text{Var}(L_{1,2})/E[L_{1,2}]$  and  $E[N_{1,2}]$  is given in Equation (25). Figure  
 331 S2 shows analytical values computed using Equations (12) and (13) compared to empirical values  
 332 from simulations. Figure S3 shows the approximate analytical distribution computed using Equation  
 333 (15) compared to the empirical distribution computed from simulations. Although the gamma  
 334 distribution in Equation (15) provides a good fit to the empirical distribution, a Gaussian distribution  
 335 can be more robust in practice because the gamma approximation is slightly underdispersed compared  
 336 with the true distribution. In practice, we use the Gaussian distribution for inference.

A maximum likelihood estimator of the degree between  $A_1$  and  $A_2$  can be obtained by determining the degree  $d_L(A_1, A_2)$  between  $A_1$  and  $A_2$  for which value of the distribution in Equation (15) is maximized. This gives the maximum likelihood estimator

$$d_L(A_1, A_2) = \arg \max_d f_{T_{1,2}}(t_{1,2}; d). \quad (16)$$

337 3.5.5. *Likelihoods for identifying background IBD.* Individuals with no recent relationship can share  
 338 small segments of IBD by chance, especially in populations with recent or severe bottlenecks. This  
 339 kind of IBD is referred to as background IBD and it poses a considerable challenge to accurate  
 340 pedigree inference.

BONSAI: AN EFFICIENT METHOD FOR INFERRING LARGE HUMAN PEDIGREES FROM GENOTYPE DATA<sup>15</sup>

341 Previous methods have addressed background IBD by various approaches. For example, the  
342 authors of the ERS method (Huff et al. 2011) presented an approach for modeling the distribution  
343 of background IBD among unrelated individuals and then performing a likelihood ratio test to  
344 determine whether the IBD shared between a new pair of individuals was significantly different from  
345 background.

346 Power for detecting background IBD can be increased by comparing sets of individuals rather than  
347 pairs of individuals, leveraging the information inherent in previously-inferred pedigree structures.  
348 As we demonstrate, such an approach makes it possible to detect background IBD between sets of  
349 individuals without prior knowledge of the distribution of background IBD. This is useful because it  
350 can be challenging to know a priori the expected amount of background IBD between a given pair  
351 of individuals.

352 We take an approach to identifying background IBD in which we consider the information  
353 contained in IBD sharing patterns across multiple individuals to determine when IBD is background  
354 and when it is due to true recent ancestry. In particular, we consider the problem in which all of the  
355 IBD observed in an individual is either background IBD, or true IBD due to a recent relationship.

356 To illustrate the approach, consider the IBD sharing pattern shown in Figure 6. Individuals 3 and  
357 4 share relatively large amounts of IBD with 5 and 6, compared with the amount shared between  
358  $\{1, 2\}$  and  $\{5, 6\}$ . If 1 and 2 were much more distantly related to 5 and 6 than 3 and 4, we might  
359 not consider the amount of IBD they share with 5 and 6 to be unusually small. However, because 1,  
360 2, 3, and 4 have similar degrees of relatedness to 5 and 6, the amount of IBD shared by 1 and 2  
361 appears to be unusually low. If we can say that the amount of IBD shared below node 7 is smaller  
362 than expected by chance, then we can assume that the IBD observed in 1 and 2 is background IBD  
363 and remove these nodes from consideration when connecting the left and right pedigrees.

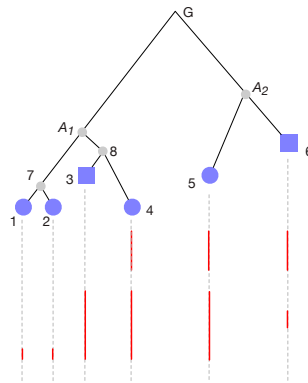


FIGURE 6. **Detecting background IBD.** Genotyped individuals are shaded in purple. Vertical red lines indicate IBD segments shared between the genotyped descendants of  $A_1$  and the genotyped descendants of  $A_2$ .

364 We test for background IBD in practice through a series of hypothesis tests. Given that IBD is  
365 observed between two sets of nodes,  $\mathcal{N}_1$  and  $\mathcal{N}_2$ , we take the putative common ancestors  $A_1$  and  
366  $A_2$  through which the IBD was inherited to be the most recent common ancestors of  $\mathcal{N}_1$  and  $\mathcal{N}_2$ ,  
367 respectively. We then consider each of the descendant nodes,  $c$ , immediately below  $A_1$  in turn (e.g.,

368 7 and 8 in Figure 6) and we ask whether the amount of observed IBD below the node is much  
369 lower or higher than expected by chance, given the degree between  $A_1$  and  $A_2$  inferred using all the  
370 descendant nodes below  $A_1$ , excluding  $c$ .

371 If we assume that some individuals in  $\mathcal{N}_1$  are related to some individuals in  $\mathcal{N}_2$ , then on average  
372 the observed IBD will represent true IBD, plus background IBD. The individuals sharing the greatest  
373 amount of IBD, relative to their genealogical positions, are likely to be the truly-related individuals.  
374 When testing for background IBD, we assume that the individuals sharing the greatest amount of  
375 IBD are truly related and we test for background IBD only in the individuals sharing less IBD.  
376 Thus, the child node  $c^*$  of  $A_1$  with the greatest IBD sharing with  $\mathcal{N}_2$  is exempt from our test.

377 We drop all nodes that reject the null hypothesis of this test and re-set the ancestral node to  
378 be the common ancestor of all remaining IBD-carrying nodes. For example, if we detected that  
379 the clade below node 7 in Figure 6 had much lower IBD than expected by chance, we would drop  
380 node 7 and its descendants from consideration and set the true common ancestor relating the two  
381 pedigrees to be node 8. We iteratively repeat this procedure until no nodes are dropped. We then  
382 repeat the procedure for the nodes immediately below  $A_2$ .

383 Let  $\mathcal{C}_n$  denote the set of children of node  $n$ . To test whether the IBD observed below a child  
384 node  $c \in \mathcal{C}_n$  is background IBD, we establish an approximation of the null hypothesis  $H_0$  that the  
385 observed IBD below node  $c$  is real and we ask whether this hypothesis is rejected in favor of the  
386 alternative hypothesis  $H_1$  that the IBD is background.

Under  $H_0$ , we assume that the degree  $d_{H_0}(A_1, A_2)$  between  $A_1$  and  $A_2$  is the maximum likelihood estimate  $d_{H_0}(A_1, A_2) = d_L(A_1, A_2 \setminus c)$ , or the generalized DRUID estimate  $d_{H_0}(A_1, A_2) = d_D(A_1, A_2 \setminus c)$  ignoring clade  $c$ . We then perform the following test

Reject  $H_0$  at level  $\alpha$  if:

$$\begin{aligned} & \mathbb{P}(T_{c,A_2} \leq t_{c,A_2}; d_{H_0}(A_1, A_2)) < \alpha/2, \\ \text{or} & \quad \mathbb{P}(T_{c,A_2} \geq t_{c,A_2}; d_{H_0}(A_1, A_2)) < \alpha/2. \end{aligned} \tag{17}$$

387 where  $T_{c,A_2}$  is the random variable describing the amount of IBD between descendants of  $c$  and  
388 descendants of  $A_2$  with observed value  $t_{c,A_2}$ . The distribution of  $T_{c,A_2}$  is given by Equation (15). It  
389 is reasonable to be conservative when dropping background IBD so that true relationships are called  
390 as background IBD only a small fraction of the time. Thus, in practice, we take  $\alpha$  to be small, such  
391 as  $\alpha = 10^{-3}$ .

392 3.5.6. *Determining the ancestral branches through which to connect pedigrees.* One difficulty in  
393 constructing large pedigrees is determining the ancestors through which two sets of genotyped  
394 individuals are related. A simple fundamental question is whether two lineages are both on the  
395 maternal side of an individual, both on the paternal side, or on opposite parental sides. Without  
396 genotyped parents, the side through which a lineage passes can be difficult to determine, although  
397 sex chromosomes and mitochondrial haplotypes can be used to resolve the parent of origin in some  
398 cases.

399 We consider the problem of inferring whether two distant sets of relatives are related through  
400 the same parent of a focal individual, or through different parents. The scenario we consider is



BONSAI: AN EFFICIENT METHOD FOR INFERRING LARGE HUMAN PEDIGREES FROM GENOTYPE DATA<sup>17</sup>

401 illustrated in Figure 7. The amount of IBD shared among the red and purple pedigrees in Figure 7  
402 is uninformative about whether they are related through the same parent. Even if the purple and  
403 red pedigrees in Figure 7 shared no IBD, they could still be related to individual 1 through the  
404 same parent by passing through different grandparents. However, if the red and purple pedigrees  
405 are related to the focal individual 1 through the same parent, the IBD segments the purple pedigree  
406 shares with individual 1 cannot spatially overlap with the segments the red pedigree shares with  
407 individual 1. This is because two overlapping segments would have undergone recombination in the  
408 parent (i.e., individual 10). The result will either be a spliced segment (Figure 7), or the replacement  
409 of one segment by the other with possible reduction in segment size.

410 In the Big Bonsai method, when there are multiple possible grandparents of a common ancestor  
411 through which we can connect a focal set of nodes  $\mathcal{N}$  in a focal pedigree  $\mathcal{P}$  to two distantly-related  
412 pedigrees  $\mathcal{P}_1$  and  $\mathcal{P}_2$ , we examine whether the IBD segments between  $\mathcal{P}_1$  and  $\mathcal{N}$  overlap with the  
413 IBD segments between  $\mathcal{P}_2$  and  $\mathcal{N}$ . The efficacy of checking segment overlaps is discussed in Section  
414 4.3 using simulated data.

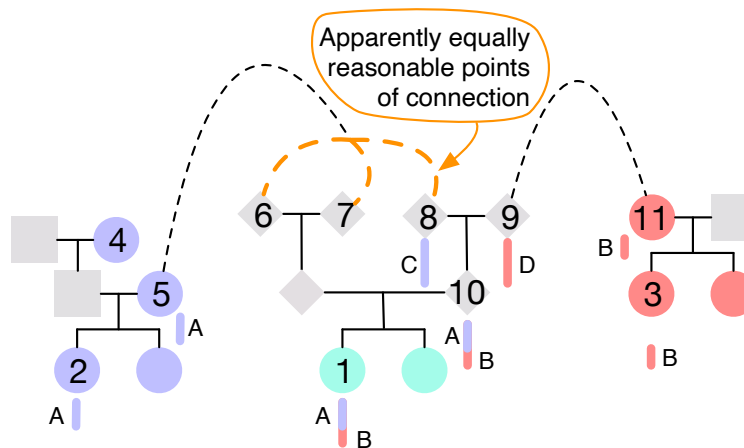


FIGURE 7. **Determining the parental side of distant relatives.** Individual 1 in the cyan pedigree shares segment *A* IBD with individuals 2 and 5 in the purple pedigree and they share segment *B* IBD with individuals 3 and 11 in the red pedigree. If the lineage connecting individual 1 to the purple pedigree passes through ancestor 8 and the lineage connecting individual 1 to the red pedigree passes through individual 9, then the ranges of segments *A* and *B* cannot overlap because individual 10 only transmits one recombined haplotype to individual 1. Observing abutting segments *A* and *B* is evidence that the cyan pedigree is connected to the purple and red pedigrees through the same parent. Observing spatially overlapping segments *A* and *B* is evidence that the purple and red pedigrees are connected through different parents of individual 1. In the absence of segment overlaps and splicing information, the orange dashed lines indicate equally reasonable ways to connect the purple and cyan pedigrees.

415 3.5.7. *Summary of the Big Bonsai algorithm.* We combine the tools in Sections 3.5.2 – 3.5.6 to obtain  
416 the Big Bonsai method presented in Algorithm 4. The input for the Big Bonsai method consists of  
417 small pedigrees inferred using the Small Bonsai method. It assembles these small pedigrees into

418 a large and sparsely-sampled pedigree by iteratively combining the two pedigrees that share the  
419 greatest total length of IBD until all pedigrees have been agglomerated into a single pedigree, or  
420 discarded because they cannot be combined in a reasonable way.

421 We assume that a pair of pedigrees,  $\mathcal{P}_1$  and  $\mathcal{P}_2$ , can only be combined in ways that connect  
422 individuals who share IBD. When combining two pedigrees, the Big Bonsai method identifies the  
423 sets,  $\mathcal{N}_1$  and  $\mathcal{N}_2$  of genotyped nodes in each pedigree that share at least one IBD segment with an  
424 individual in the other pedigree. If the set  $\mathcal{N}_i$  does not have at least one common ancestor, we find  
425 the set  $\tilde{\mathcal{A}}_i$  of most recent ancestral nodes whose descendants comprise  $\mathcal{N}_i$ . The pair of ancestors  
426  $A_1 \in \tilde{\mathcal{A}}_1$  and  $A_2 \in \tilde{\mathcal{A}}_2$  whose descendants share the greatest total length of IBD is then determined  
427 and we redefine  $\mathcal{N}_1$  and  $\mathcal{N}_2$  to be the genotyped descendants of  $A_1$  and  $A_2$ , respectively.

428 Our objective is to identify pairs of individuals through which  $\mathcal{P}_1$  and  $\mathcal{P}_2$  can be connected in  
429 such a way that all individuals in  $\mathcal{N}_1$  are related to all individuals in  $\mathcal{N}_2$ . This is accomplished  
430 if and only if the sets  $\mathcal{N}_1$  and  $\mathcal{N}_2$  share at least one common ancestor. Sets  $\mathcal{N}_1$  and  $\mathcal{N}_2$  will be  
431 connected through a common ancestor if their respective common ancestors,  $A_1$  and  $A_2$ , share a  
432 common ancestor or if  $A_1$  is descended from any individual in  $\mathcal{N}_2$  or from any ancestor on the  
433 induced subtree  $\Lambda_2$  of pedigree  $\mathcal{P}_2$  relating  $\mathcal{N}_2$  to one another. Similarly, sets  $\mathcal{N}_1$  and  $\mathcal{N}_2$  will have a  
434 common ancestor if  $A_2$  is descended from any individual in  $\mathcal{N}_1$  or from any ancestor on the induced  
435 tree  $\Lambda_1$  of pedigree  $\mathcal{P}_1$  relating  $\mathcal{N}_1$ .

436 We present a generalized DRUID estimator in Appendix 6.3 for connecting pedigrees through  
437 individuals  $A$  who are not common ancestors of  $\mathcal{N}_1$  or  $\mathcal{N}_2$ . However, connecting pedigrees  $\mathcal{P}_1$  and  
438  $\mathcal{P}_2$  through all possible pairs can be computationally inefficient. Instead, we can accept a certain  
439 loss in accuracy and allow pedigrees to be connected only through common ancestors. We find  
440 that this approach works well in practice, generating pedigrees that are nearly as accurate as those  
441 constructed by connecting  $\mathcal{P}_1$  and  $\mathcal{P}_2$  in all possible ways.

442 Let  $A_1$  be a most recent common ancestor of  $\mathcal{N}_1$  and let  $A_2$  be a most recent common ancestor  
443 of  $\mathcal{N}_2$ . For each pair of possible ancestors  $(A_1, A_2)$ , we compute the generalized DRUID estimate  
444  $d_D(A_1, A_2)$  of the degree using Equation (9). We then perform the test for background IBD described  
445 in Section 3.5.5, which potentially results in a new pair of common ancestors  $A'_1$  and  $A'_2$  whose  
446 descendants do not share detectable background IBD. If the pair  $(A'_1, A'_2)$  differs from the original  
447 pair  $(A_1, A_2)$ , we replace  $A_1$  and  $A_2$  with  $A'_1$  and  $A'_2$  and recompute the generalized DRUID estimate  
448  $d_D(A_1, A_2)$ . At the end of these steps, we have a set of possible ancestral pairs through which  $\mathcal{P}_1$   
449 and  $\mathcal{P}_2$  can be connected, along with point estimates,  $d_D(A_1, A_2)$ , of the total degree separating  
450 each pair.

451 It remains to evaluate the likelihood of each pair and degree. Following the notation of Ko and  
452 Nielsen (2017), denote the relationship between a pair of individuals  $A_1$  and  $A_2$  with common  
453 ancestor (or ancestral pair)  $G$  by  $(d_1, d_2, n)$ , where  $d_1$  is the number of meiotic events separating  
454  $A_1$  from  $G$ ,  $d_2$  is the number of meiotic events separating  $A_2$  from  $G$ , and  $n = |G|$  is the number  
455 of common ancestors. For a given estimate  $d_D(A_1, A_2)$  of the degree between  $A_1$  and  $A_2$  and a  
456 number of common ancestors  $n$ , we consider all relationship types  $(d_1, d_2, n)$  corresponding to degree  
457  $d_D(A_1, A_2)$ ; in other words, we consider all relationship types such that  $d_1 + d_2 = d_D(A_1, A_2) + n - 1$ .

458 For a given pair of ancestors  $A_1$  and  $A_2$ , and for each relationship  $(d_1, d_2, n)$ , we connect  $A_1$  to  $A_2$   
459 through all such relationships and we evaluate the composite likelihoods of the resulting pedigrees  
460 computed using Equation (4). All pedigrees whose likelihoods are at least a fraction  $f_\ell$  of that of  
461 the most likely pedigree are stored and the rest are discarded. We also apply the test in Section  
462 3.5.6 for incompatible ancestral lineages to each retained pedigree and we retain only those pairs  
463 that pass the test.

464 Here, we have considered the procedure for combining two pedigrees  $\mathcal{P}_1$  and  $\mathcal{P}_2$ . However, the  
465 output of the Small Bonsai method is a set of high-likelihood pedigrees  $S$  and the input to the  
466 Big Bonsai method is a list  $\vec{S} = [S_1, \dots, S_K]$  of such sets. Let  $\mathcal{N}_S$  denote the genotyped node set  
467 corresponding to the pedigree set  $S$ ; in other words,  $\mathcal{N}_S$  is the genotyped node set of every pedigree  
468  $\mathcal{P} \in S$ . If  $\mathcal{N}$  is the set of genotyped nodes in the full pedigree, then  $\bigcup_{i=1}^K \mathcal{N}_{S_i} = \mathcal{N}$ .

469 At each step of the Big Bonsai method, we compare each pair of genotyped sets  $\mathcal{N}_{S_i}$  and  $\mathcal{N}_{S_j}$   
470 ( $1 \leq i, j \leq K$ ) to determine the pair with the greatest shared total amount of IBD. Here, the total  
471 amount of IBD is the total length of IBD obtained by merging the segments shared between all  
472 pairs of individuals  $(u, v) \in \mathcal{N}_{S_i} \otimes \mathcal{N}_{S_j}$ . We then identify the subsets  $\mathcal{N}_i \subseteq \mathcal{N}_{S_i}$  and  $\mathcal{N}_j \subseteq \mathcal{N}_{S_j}$  that  
473 share IBD and we combine each pair of pedigrees  $(\mathcal{P}_i, \mathcal{P}_j) \in S_i \otimes S_j$  through all pairs of possible  
474 most recent common ancestors of  $\mathcal{N}_i$  and  $\mathcal{N}_j$ . The full algorithm is presented in Algorithm 4.

475 It is possible to mis-infer relationships early in the process of pedigree building that lead to  
476 conflicts several steps later in the process. The downstream effects of a misplaced individual can  
477 be difficult to predict and prevent without a bird's-eye view of the pedigree, but misplaced pairs  
478 of relatives can often be detected after the pedigree is built. In practice, we include a final step  
479 in the pedigree building process to detect internal inconsistencies by comparing the final pairwise  
480 relationships implied by the pedigree structure to the initial pairwise likelihood predictions. When  
481 the inferred relationships have low pairwise likelihoods, we rebuild the pedigree, iteratively expanding  
482 the number of pedigrees that are retained at each step to increase the chances that the correct  
483 pedigree is explored. We also correct pairwise point estimates that are likely to be incorrect when  
484 viewed in the context of a fully-built pedigree before attempting to re-infer the pedigree.

485 Putting together the point estimator, the Small Bonsai method, and the Big Bonsai method, we  
486 obtain the full Bonsai method shown in Figure 1. Outlines of the three primary stages of Bonsai  
487 are shown in Algorithms 1, 2, and 4. The Bonsai method performs these stages in series.

488 **3.6. Subjects and simulations.** Our empirical analyses are based on simulated data, as well as a  
489 dataset comprised of the pedigrees of 23andMe research participants. All simulations and analyses  
490 that used real genotype data were performed using individuals consented for research at 23andMe.

491 **3.6.1. Overview of simulations.** Simulations were carried out using two different general approaches.  
492 In one approach, no genotype or customer data were used and IBD segments were known with  
493 certainty, their positions and lengths being recorded during the simulation process. In the second  
494 simulation approach, the full-genome genotypes of research-consented 23andMe customers were used  
495 for the pedigree founders and genotypes were simulated for individuals in all subsequent generations  
496 through cross-over events. IBD segments were then inferred between each pair of individuals using

---

**Algorithm 1 Pairwise likelihoods and point estimates.** Compute the likelihood of many different relationships between a pair of individuals,  $i$  and  $j$ , and obtain a point estimate of the relationship between  $i$  and  $j$ .

**Input:** Ages  $a_i, a_j$ , pairwise total lengths  $t_1^{i,j}$  and  $t_2^{i,j}$  if IBD1 and IBD2, and pairwise counts  $c_1^{i,j}$  and  $c_2^{i,j}$  of IBD1 and IBD2 segments.

---

```
 $\vec{\mathcal{R}} \leftarrow$  List of relationships at which to evaluate the likelihood
 $L =$  Initialize dictionary mapping relationships to likelihoods
 $\mathcal{L}_{\max} = -\infty$ 
for  $\mathcal{R} \in \vec{\mathcal{R}}$  do
  Compute  $\mathcal{L}_{\mathcal{R}}^g$  using Equation (2)
  Compute  $\mathcal{L}_{\mathcal{R}}^a$  using Equation (3)
  Compute  $\mathcal{L}_{\mathcal{R}} = \mathcal{L}_{\mathcal{R}}^a \mathcal{L}_{\mathcal{R}}^g$  as in Equation (1)
   $L[\mathcal{R}] = \mathcal{L}_{\mathcal{R}}$ 
  if  $\mathcal{L}_{\mathcal{R}} > \mathcal{L}_{\max}$  then
     $\mathcal{L}_{\max} = \mathcal{L}_{\mathcal{R}}$ 
     $\hat{\mathcal{R}} = \mathcal{R}$  # Estimated pairwise relationship
  end if
end for
return  $\hat{\mathcal{R}}, L$ 
```

---

497 an in-house method for inferring IBD on unphased data (Henn et al. 2012), which is similar to that  
498 of Seidman et al. (2020)

499 In all simulations, the number of cross-over events in each meiosis was drawn such that the  
500 expected number of events was one per 100 cM and the locations of cross-overs were sampled  
501 uniformly along chromosomes.

502 3.6.2. *Validated real pedigrees.* To evaluate Bonsai on true pedigrees, we constructed 204 validated  
503 pedigrees for individuals in the 23andMe database. By considering pedigrees in which all individuals  
504 were genotyped, we were able to construct each pedigree with a high degree of certainty by connecting  
505 parent-child pairs inferred using Algorithm 1. To ensure that the true pedigree was known with  
506 certainty, we considered quartets of genotyped customers with at least two full-sibling children  
507 and two parents. We identified pedigrees in which each individual was connected to every other  
508 individual through a chain composed of these building blocks. We further restricted our attention  
509 to pedigrees that spanned at least three generations with at least one pair of first cousins.

510 Pedigrees identified in this way allowed us to know the true pedigree structure because parent-  
511 offspring and full-sibling pairs can be inferred with nearly perfect accuracy and the quartet structure  
512 allows us to further confirm each inferred relationship using the other pairs in the quartet. In  
513 particular, we required that each sibling pair had inferred child-parent relationships with the same  
514 two parents using Algorithm 1. We also required the self-reported ages of both parents to be at  
515 least 17 years older than the self-reported ages of the children.

516 3.6.3. *Self-reported pedigrees.* The Family Tree feature provided by 23andMe allows users to edit  
517 and validate relationships in their pedigrees. We considered a set of such pedigrees where users had  
518 either verified or changed relationships, indicating that they knew the correct relationships for at  
519 least a subset of individuals in the pedigree. We considered only individuals in these pedigrees who

---

**Algorithm 2 Small Bonsai algorithm.** Infer a Small Bonsai pedigree.

---

**Input:**

- `like_dict[id1][id2][L]` # Dictionary mapping pairs of IDs to likelihood dictionaries  $L$  produced by Algorithm 1
- `max_deg` # Max degree between any placed pair.
- `max_append_types` # Max number of ways to add a new person to a pedigree.
- `fl` # Fraction of least to most likely pedigree likelihoods.
- `focal_id` # ID of focal individual.

```

 $\mathcal{P}$  = Initialize pedigree with focal_id
 $U$  = Initialize set of (“unplaced”) individuals not in  $\mathcal{P}$ 
 $S = \{\mathcal{P}\}$  # Set of pedigrees built so far
while  $|U| > 0$  and  $\min_{u \in U, p \in \mathcal{P}} \{\hat{d}_{u,p}\} \leq \text{max\_deg}$  do
     $(u, p) = \arg \min_{u \in U, p \in \mathcal{P}} \{\hat{d}_{u,p}\}$  # Unplaced individual closest to any placed individual
     $S'$  = Initialize empty set of pedigrees built on this step
    for  $\mathcal{P} \in S$  do
         $\vec{\mathcal{R}} = \text{reverseargsort}(\text{like\_dict}[p][u])$  # Sort by highest to lowest likelihood relationship
        for  $\mathcal{R}$  in  $\vec{\mathcal{R}}[0 : \text{max\_append\_types}]$  do
            for  $\rho \equiv \mathcal{R}$  do # Relationships  $\rho$  consistent with  $\mathcal{R}$ 
                 $\mathcal{P}' = \mathcal{P}$  with  $u$  placed in relationship  $\rho$ , relative to  $p$ 
                 $S' = S' \cup \{\mathcal{P}'\}$ 
            end for
        end for
    end for
     $\mathcal{P}^* = \arg \max_{\mathcal{P} \in S'} (\mathcal{L}(\mathcal{P}))$  # Most likely pedigree
     $S = \{\mathcal{P} \in S' : \mathcal{L}(\mathcal{P}) \geq f_l \mathcal{L}(\mathcal{P}^*)\}$ 
end while
return  $S$ 

```

---

**Algorithm 3 Detect background IBD.** Detect whether the IBD observed in one of the clades directly descended from  $A_1$  in pedigree  $\mathcal{P}_1$  carries background IBD relative to the descendants of  $A_2$  in pedigree  $\mathcal{P}_2$ .

---

```

function DROPBACKGROUND( $A_1, A_2$ )
     $d_D(A_1, A_2) = \text{Generalized DRUID estimated degree}$ 
     $\Lambda_1(A_1) = \text{Induced subtree below } A_1 \text{ connecting } \mathcal{N}_{A_1}$ 
     $\mathcal{N} = \mathcal{N}_{A_1}$ 
     $c^* = \arg \max_{c \in \text{Children}(A_1)} \{T_{c,A_2}\}$  # Clade sharing most IBD with  $A_2$ .
    for  $c \in \text{Children}(A_1) \setminus c^*$  do
        if  $\text{Reject } H_0(T_{c,A_2})$  then
             $\mathcal{N} = \mathcal{N} \setminus \mathcal{N}_c$ 
        end if
    end for
    if  $A_1 = A_{\mathcal{N}}$  then
        Return  $A_1$ 
    else
        Return DropBackground( $A_{\mathcal{N}}, A_2$ )
    end if
end function

```

---

---

**Algorithm 4 Big Bonsai algorithm.** Combine Small Bonsai pedigrees into a Big Bonsai pedigree.

---

```

# Infer small pedigrees
U = Initialize set of unplaced individuals
S→ = Initialize empty list of sets of pedigrees
while |U| > 0 do
  focal_id = arg minu ∈ U (  $\frac{1}{|U \setminus u|} \sum_{u' \in U \setminus u} \hat{d}_{u,u'}$  ) # ID with closest mean degree to all other IDs
  S = SmallBonsai(focal_id) # Infer a set S of likely pedigrees for focal_id
  S→.append(S)
end while

# Combine small pedigrees
while length(S→) > 0 do
  S1, S2 = arg maxS1, S2 ∈ S→ (TNS1, NS2) # Pedigree sets with greatest shared total IBD
  N1~ = Subset of N1 related to N2
  N2~ = Subset of N2 related to N1
  S = Initialize empty set of pedigrees
  for P1 ∈ S1, P2 ∈ S2 do
    AN1~ = CommonAncestors(N1~) # Set of common ancestors of N1~
    AN2~ = CommonAncestors(N2~) # Set of common ancestors of N2~
    for A1 ∈ AN1~, A2 ∈ AN2~ do
      dD(A1, A2) = Infer generalized DRUID estimate between A1 and A2
      A1' = DropBackground(A1, A2) # Algorithm 3
      A2' = DropBackground(A2, A1) # Algorithm 3
      if OverlapIBD(A1', A2') then
        Continue # Ignore pairs (A1', A2') that fail the overlap conflict test in Section 3.5.6
      end if
      dD(A1', A2') = Infer generalized DRUID estimate between A1' and A2'
      for |G| = 1, 2 do # Number of common ancestors G of A1' and A2'
        d = dD(A1', A2') + |G| - 1
        for d1 = 0, ..., d do
          d2 = d - d1
          ρ = (d1, d2, |G|) # relationship
          P1,2 = New pedigree from connecting A1' to A2' through relationship ρ
          S.add(P1,2)
        end for
      end for
    end for
  end for
  P* = arg maxP ∈ S (L(P))
  S = {P ∈ S : L(P) ≥ fℓL(P*)}
end for
S1 = S
S→.delete(S2)
end while
return S→

```

---

520 were consented for research and re-built the pedigree using only the subset of research-consented  
521 individuals. The inferred relationships in the pedigree could then be compared with the user-verified  
522 relationships.

523 *3.6.4. Simulations for fitting empirical pairwise genetic likelihood distributions.* The distribution of  
524 the total length of IBD1 and IBD2, the distribution of lengths of IBD1 and IBD2 segments, and the  
525 distribution of the total counts of IBD1 and IBD2 segments for a specified relationship type  $\mathcal{R}$  were  
526 obtained by simulating full genomes for 100 pairs of individuals of the relationship type. For each  
527 simulation replicate, a pedigree was specified containing the relationship of interest and cross-over  
528 events were simulated within the pedigree.

529 Over the 100 replicates, we computed the mean  $\mu_Q$  and standard deviation  $\sigma_Q$  of the quantities  
530  $Q = T_1, T_2, C_1$ , and  $C_2$  where  $T_1$  is the total genome-wide length of IBD1,  $T_2$  is the total genome-  
531 wide length of IBD2,  $C_1$  is the total genome-wide count of IBD1 segments, and  $C_2$  is the total  
532 genome-wide count of IBD2 segments.

533 *3.6.5. Large simulated pedigrees.* The 204 validated customer pedigrees described in Section 3.6.2  
534 are small enough that the Small Bonsai method is capable of building them without resorting  
535 to the Big Bonsai method. To evaluate the Big Bonsai method, we required considerably larger  
536 pedigrees whose structures were known with certainty. Although many pedigrees for 23andMe  
537 research-consented customers are large, the relationships within them are typically not known with  
538 certainty. Therefore, we simulated large pedigrees to evaluate the Big Bonsai method.

539 Exact IBD was simulated for pedigrees with a depth of five generations by choosing a focal  
540 individual and building the “cone” of ancestors comprised of two parents, four grandparents, eight  
541 great-grandparents, and sixteen great-great-grandparents. For each individual in the ancestral cone,  
542 a second spouse was added with probability 0.2. Then, two children were created for every pair  
543 of spouses in the pedigree. Two children were repeatedly sampled for every spouse pair with no  
544 children until the generation with the focal individual was reached. An example of a pedigree  
545 generated by this approach is shown in Supplemental Figure S1.

546 *3.6.6. Simulated pedigrees for testing degree inference.* The approach for simulating pedigrees for  
547 degree inference was similar to that in Section 3.6.5; however, the pedigree structure was different.  
548 For these pedigrees, we were interested in inferring the degree between a pair of common ancestors  
549  $A_1$  and  $A_2$ , given IBD observed between their descendants  $\mathcal{N}_1$  and  $\mathcal{N}_2$ .

550 For this analysis, we created two identical small pedigrees  $\mathcal{P}_1$  and  $\mathcal{P}_2$ . Each small pedigree had  
551 the same structure comprised of the common ancestor  $A_1$  or  $A_2$ , their spouse, their two children,  
552 and four grandchildren, where the grandchildren were comprised of two children for each child of  
553  $A_1$  or  $A_2$ . The ancestors  $A_1$  and  $A_2$  were then connected by degree  $d(A_1, A_2)$  through a pair of  
554 common ancestors, where the degree  $d$  varied from 1 to 10.

555

## 4. RESULTS

556 We considered both simulated and real data to investigate the performance of the small and big  
557 Bonsai methods and their components.

24BONSAI: AN EFFICIENT METHOD FOR INFERRING LARGE HUMAN PEDIGREES FROM GENOTYPE DATA

558 **4.1. Degree estimation.** To evaluate the accuracy of degree inference using the likelihood estimator  
 559 (Equation 16) and the generalized DRUID estimator (Equation 9), we applied these estimators to  
 560 infer the degree between common ancestors  $A_1$  and  $A_2$  of two small pedigrees  $\mathcal{P}_1$  and  $\mathcal{P}_2$  (Section  
 561 3.6.6). Figure 8 shows the accuracy of the likelihood estimator  $d_L$  and the generalized DRUID  
 562 estimator  $d_D$  for inferring the degree  $d$ , conditional on the event that any IBD at all was observed  
 563 between the leaf nodes in  $\mathcal{P}_1$  and  $\mathcal{P}_2$ . From Figure 8 it can be seen that both the maximum likelihood  
 564 estimator  $d_L$  and the generalized DRUID estimator  $d_D$  have similar accuracies for inferring the  
 565 degree  $d$ . Moreover, the DRUID estimate is nearly identical to the maximum likelihood estimate,  
 566 which is important in practice because it implies that connecting two pedigrees through the degree  
 567 inferred by DRUID results in a pedigree that is approximately the maximum likelihood pedigree.  
 568 This result can dramatically speed up pedigree inference and, in practice, we use the generalized  
 569 DRUID estimator for inferring the degree of separation between two small pedigrees.

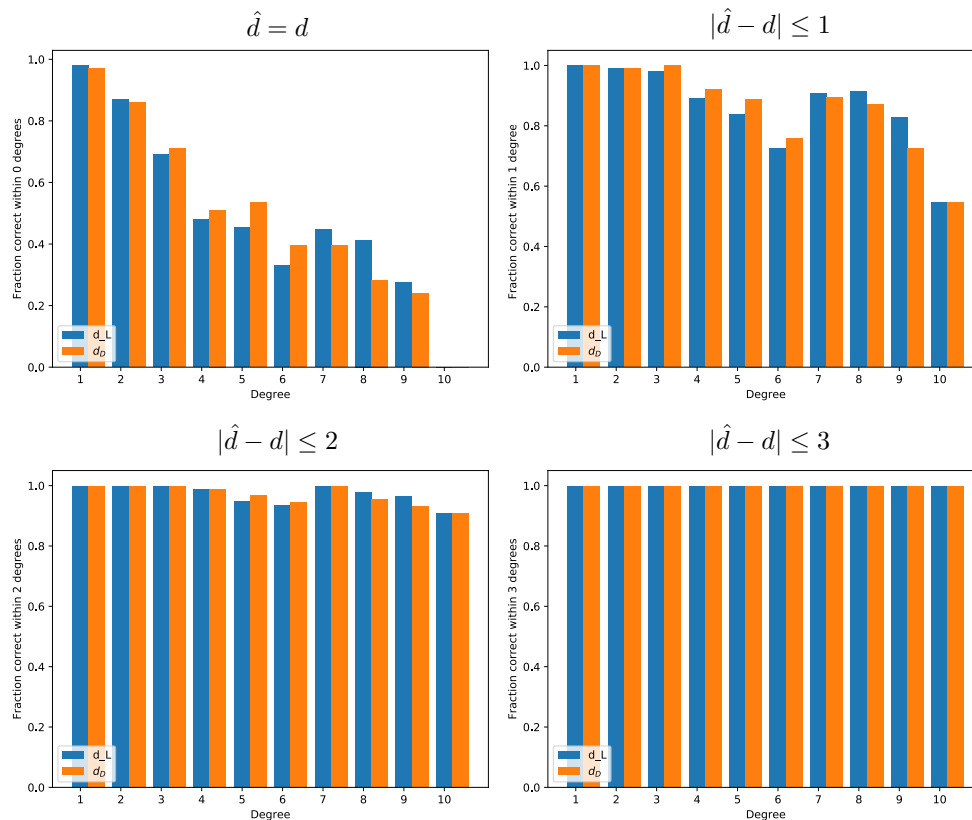


FIGURE 8. The accuracy of the likelihood method (Equation 16) and the generalized DRUID method (Equation 9) for inferring the degree between a pair of common ancestors. The accuracy of the estimate is shown for four different tolerances: exactly equal to the true degree, within one degree of the true degree, within two degrees of the true degree, and within three degrees of the true degree.

570 **4.2. Background IBD detection.** To evaluate the efficacy of the test in Equation (17) for  
 571 detecting background IBD, we simulated pedigrees comprised of three small pedigrees,  $\mathcal{P}_1$ ,  $\mathcal{P}_2$ , and



572  $\mathcal{P}_3$ , connected together (Figure 9). In one set of simulations, pedigree  $\mathcal{P}_1$  was related only to  $\mathcal{P}_2$   
573 and not  $\mathcal{P}_3$  (Figure 9, Scenario 1). This allowed us to simulate background IBD among all pairs  
574 of individuals and then attempt to detect it. In another set of simulations, pedigree  $\mathcal{P}_1$  was truly  
575 related to all other individuals (Figure 9, Scenario 2). This second set of simulations allowed us to  
576 evaluate the rate at which background IBD was detected even when there was true IBD between  $\mathcal{P}_1$   
577 and  $\mathcal{P}_3$  as well as background IBD. Note that in all simulations, all pairs shared a nonzero expected  
578 amount of background IBD so that even truly related individuals carried additional background  
579 IBD.

580 For all pedigrees, we simulated background IBD between each pair of individuals by sampling  
581 the number of background IBD segments from a Poisson distribution with mean  $\lambda_{bgd} = 0.05, 0.5, 1,$   
582 or 5. We then sampled the length  $L$  of each observed background segment from a thresholded  
583 exponential distribution with mean 7 cM and minimum length of 5 cM. A minimum of 5 cM was  
584 chosen because, in practice, small segments can be difficult to infer and it is a common practice to  
585 employ a minimum cutoff on the length of IBD segments to reduce the rate of false positives (Huff  
586 et al. 2011).

587 The rates of background IBD we tested corresponded to values spanning the empirically observed  
588 range of background segment counts in broad human populations in the 23andMe database. When  
589 considering only segments greater than 5 cM in length, the average number of background IBD  
590 segments between a pair of individuals is between 0.01 and 0.02 for most human populations.  
591 However, for populations with historical bottlenecks, the expected background IBD count can be  
592 closer to  $\lambda_{bgd} = 5$ .

593 Figure 9A shows the fraction of times the null hypothesis  $H_0$  in Equation (17) was rejected when  
594 individuals shared an average of 0.05 background IBD segments. Blue bars correspond to simulation  
595 replicates in which all IBD shared between  $\mathcal{P}_1$  and  $\mathcal{P}_3$  was background (Scenario 1) and orange  
596 bars correspond to simulation replicates in which background IBD between  $\mathcal{P}_1$  and  $\mathcal{P}_3$  existed in  
597 addition to true IBD.

598 Each bar in Figure 9A was calculated using 50 simulated pedigrees. From Figures 9A and 9B, it  
599 can be seen that for a level of background IBD consistent with the majority of human populations,  
600 the test correctly identified background IBD a large fraction of the time. Moreover, the test typically  
601 did not detect background IBD when there was also true IBD in addition to background IBD.

602 For high levels of background IBD consistent with populations with severe bottlenecks, it was  
603 much more difficult to detect background IBD (Figures 9C and 9D). This was especially true when  
604 the internal branch length  $b$  was long and background IBD dominated true IBD.

605 This is problematic because the detection of background IBD is particularly important in these  
606 populations. However, the amount of background IBD is controllable to some degree by establishing  
607 a threshold for the minimum length of an IBD segment to be included in an analysis. The higher  
608 the threshold, the fewer the number of false positive segments and the longer the length of the  
609 background IBD segments that are not filtered by the threshold. By using this threshold when  
610 detecting background IBD, it is possible to increase the power for detecting background IBD. This  
611 can be seen in Figures 9E-G, which show the same simulations shown in 9B-D, but discarding all  
612 segments shorter than 10 cM.

26BONSAI: AN EFFICIENT METHOD FOR INFERRING LARGE HUMAN PEDIGREES FROM GENOTYPE DATA

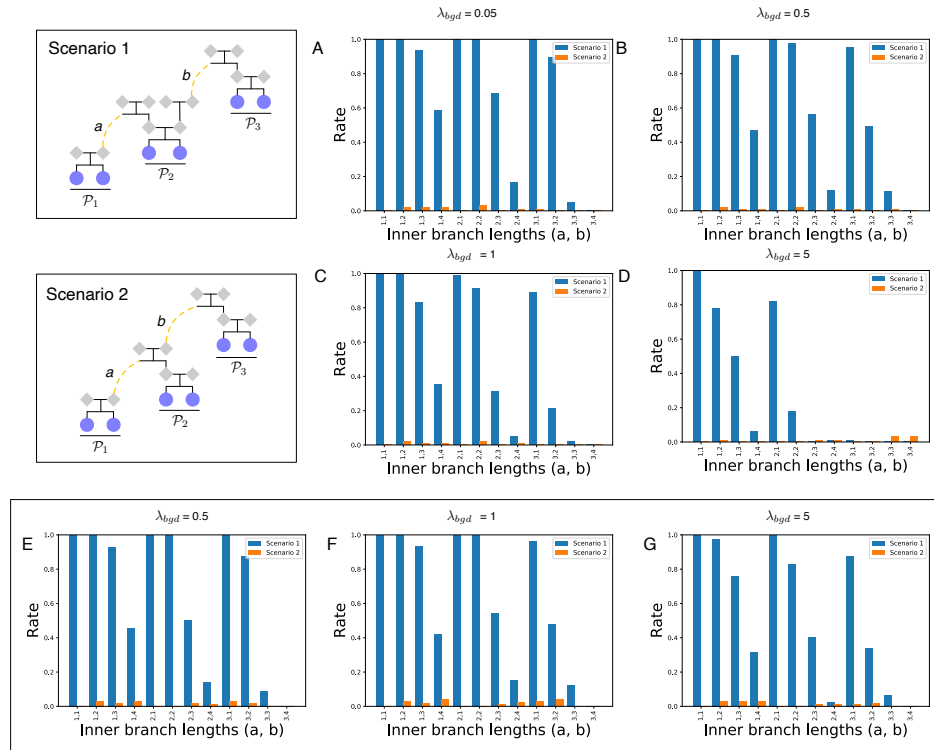


FIGURE 9. **Evaluating the use of the test in Equation (17) for detecting background IBD.** Scenario 1 shows a pedigree structure in which any IBD observed between  $\mathcal{P}_1$  and  $\mathcal{P}_3$  is background IBD. Scenario 2 shows a pedigree structure in which IBD shared between  $\mathcal{P}_1$  and  $\mathcal{P}_3$  comprises both true and background IBD. Branch lengths  $a$  and  $b$  were variable. (A)-(D) Rates for detecting background IBD under scenarios 1 and 2. (E)-(G) Same as (B)-(D), but using only segments at least 10 cM in length. Plots are shown for  $\alpha = 10^{-3}$ .

613 **4.3. Segment overlap detection.** We evaluated the degree to which overlapping IBD segments  
 614 can be informative about the ancestors through which two pedigrees are connected using the  
 615 large simulated pedigrees described in Section 3.6.5. For each pedigree, we considered the four  
 616 grandparents of a focal individual and the leaves descended from all lineages extending up from each  
 617 of the four grandparents. In the example large pedigree shown in Figure S1, the focal individual is  
 618 one of the yellow leaf nodes and the clades corresponding to the four leaf sets are colored in green,  
 619 cyan, red, and magenta.

620 For a pair of leaf sets related to the focal individual through an ancestral couple, we expect to see  
 621 no overlap in the IBD segments shared with the focal individual. For a pair of leaf sets related to the  
 622 focal individual through two grandparents who are not a couple, we expect to observe overlapping  
 623 segments occasionally.

624 Figure 3 shows the rate at which segments from one leaf set overlapped segments from another  
 625 leaf set by more than a fraction  $f$  of the total IBD observed between the two leaf sets, combined,  
 626 for  $f \in \{0.01, 0.05, 0.1, 0.15, 0.2\}$ .

BONSAI: AN EFFICIENT METHOD FOR INFERRING LARGE HUMAN PEDIGREES FROM GENOTYPE DATA<sup>27</sup>

627 Let  $i$  denote the focal individual. For leaf sets  $\mathcal{N}_1$  and  $\mathcal{N}_2$  with total amounts of IBD to the  
 628 focal individual denoted by  $T_{i,\mathcal{N}_1}$  and  $T_{i,\mathcal{N}_2}$ , let  $T_{i,\mathcal{N}_1 \cup \mathcal{N}_2}$  denote the total length of merged segments  
 629 between focal individual  $i$  and either set. We recorded an overlap in segments if the following  
 630 relationship was satisfied:  $T_{i,\mathcal{N}_1} + T_{i,\mathcal{N}_2} - T_{i,\mathcal{N}_1 \cup \mathcal{N}_2} > fT_{i,\mathcal{N}_1 \cup \mathcal{N}_2}$ .

631 Figure 10 indicates that even with few sampled leaves from each leaf set, it is possible to  
 632 detect overlapping IBD segments a large fraction of the time when the leaves are related through  
 633 grandparents who are not a couple. Each bar in Figure 3 was computed using 100 pedigrees, each  
 634 with four pairs of leaf sets related to individual 1 through a pair of grandparents who were not a  
 635 couple. Only IBD segments greater than 5 cM in length were considered.

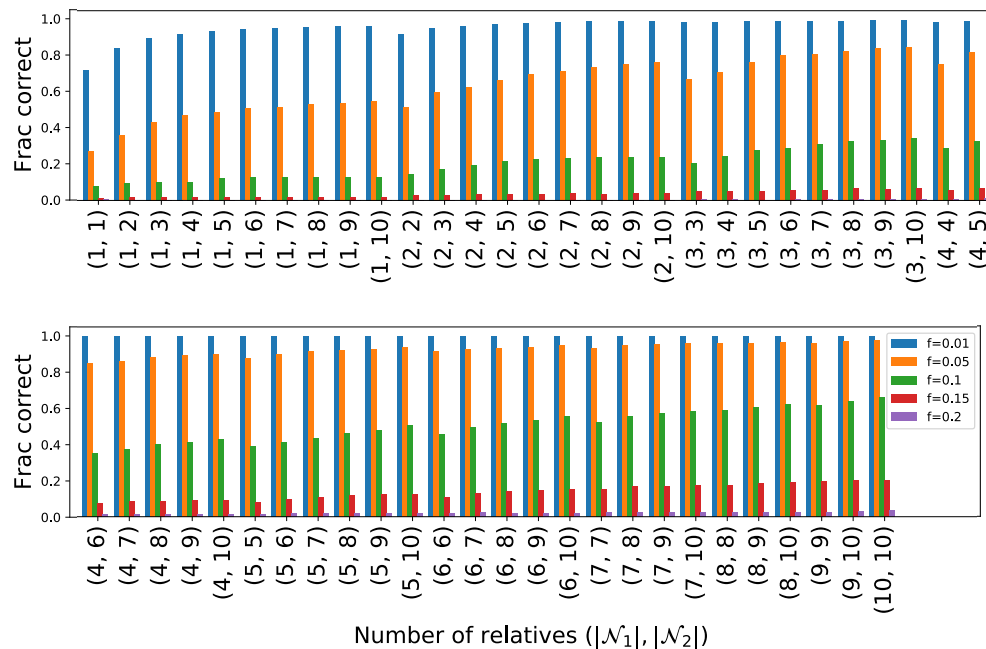


FIGURE 10. **The probability of observing an IBD segment overlap.** The plot shows the probability of observing an overlap of at least fraction  $f$  ( $f = 0.01, 0.05, 0.1, 0.2$ ) among segments shared IBD between the focal individual and sets of leaves related to the focal individual through ancestors who are not a couple. IBD segments were simulated for large pedigrees like that shown in Figure S1. IBD was computed between the focal individual and the leaf nodes of each of the four clades related to the focal individual through each of the four grandparents (colored green, cyan, red, and magenta in Figure S1). An observed IBD segment overlap was evidence that the lineages were related to the focal individual through a pair of ancestors who were not a couple.

636 **4.4. Timing and accuracy of Small Bonsai, compared with PRIMUS.** To evaluate the  
 637 accuracy and runtime of Bonsai in comparison with PRIMUS, we applied PRIMUS and Bonsai to a  
 638 set of 204 pedigrees comprised of research-consented 23andMe customers (Section 3.6.2) in which  
 639 all individuals were genotyped and for which the true pedigree was known with a high degree of  
 640 certainty.

641 Pedigrees in which all individuals have been genotyped are simple to infer by connecting together  
642 first degree relatives. The difficulty is in constructing pedigrees in which only a small fraction of  
643 individuals have been genotyped. Therefore, to evaluate the accuracy of Bonsai, we subsampled the  
644 validated pedigrees and performed inference using the subset of individuals, ignoring the remaining  
645 individuals. The resulting pedigree could then be compared to the subgraph of the true pedigree  
646 corresponding to the subsampled individuals to determine the accuracy of the inference.

647 We subsampled each pedigree to 10, 20, 30, 40, 50, 60, 70, 80, 90, or 100% of its members with a  
648 minimum of at least two individuals sampled per pedigree. Figure 11 shows the degree to which  
649 each method recovered each relationship type. From Figure 11, it can be seen that the Bonsai  
650 algorithm achieved improved accuracy for inferring relationships, and that when at least 50% of  
651 individuals were sampled, Bonsai inferred the correct pedigree with near perfect accuracy.

652 We also compared the runtime of the Bonsai method to the runtime of PRIMUS for the same set  
653 of pedigrees described in Section 3.6.2. Figure 12 shows the runtime for Small Bonsai compared to  
654 the runtime for PRIMUS for different percentages of sampled lineages from each of the pedigrees.  
655 The runtimes for Bonsai and PRIMUS were similar when few or many individuals were sampled,  
656 although Bonsai was slightly faster. In these regimes, pedigrees were fast to construct because  
657 individuals could be connected through close relatives, which were inferred with high confidence.  
658 However, when the number of individuals sampled was moderate and there were several possible  
659 pedigree configurations with high likelihoods, the Bonsai method was significantly faster than  
660 PRIMUS.

661 **4.5. Timing and accuracy of the Big Bonsai method.** Reconstruction of large pedigrees using  
662 the Small Bonsai method can be computationally challenging due to a quickly-expanding state  
663 space of possible pedigrees. Figure 13 shows timing and accuracy results for reconstructing large  
664 five-generation pedigrees simulated using the approach described in Section 3.6.5. For these analyses,  
665 we were interested in the ability of Big Bonsai to infer pedigrees that were realistic representations  
666 of direct-to-consumer genetic data. In realistic pedigrees, individuals beyond the most recent two  
667 generations may not have been sampled. When sampling individuals for the pedigree, we sampled  
668 individuals only in the most recent two generations who were not founders. This provided a pool of  
669 approximately 130 individuals who could be sampled out of a total of at least 261 individuals in  
670 each pedigree, including pedigree founders.

671 To evaluate the ability of the Big Bonsai method to reconstruct pedigrees with sparsely sampled  
672 individuals, we further subsampled 10, 20, 30, 40, 50, 60, 70, 80, 90, or 100% of the approximately  
673 130 non-founder individuals in the most recent two generations. Sampling 10% of these individuals  
674 corresponds to sampling approximately 5% of all individuals in the full pedigree and sampling 100%  
675 of these individuals corresponds to sampling approximately 50% of all individuals in the pedigree  
676 overall. Our sampling scheme presents a further challenge to pedigree reconstruction because the  
677 samples did not contain ancestral individuals who could provide additional information about the  
678 degrees of distant relationships.

679 From Figure 13A, it can be seen that the runtime is on the order of several seconds per pedigree,  
680 even though pedigree sizes were large. Bonsai built pedigrees with over one hundred sampled  
681 individuals in tens of seconds.

BONSAI: AN EFFICIENT METHOD FOR INFERRING LARGE HUMAN PEDIGREES FROM GENOTYPE DATA<sup>29</sup>

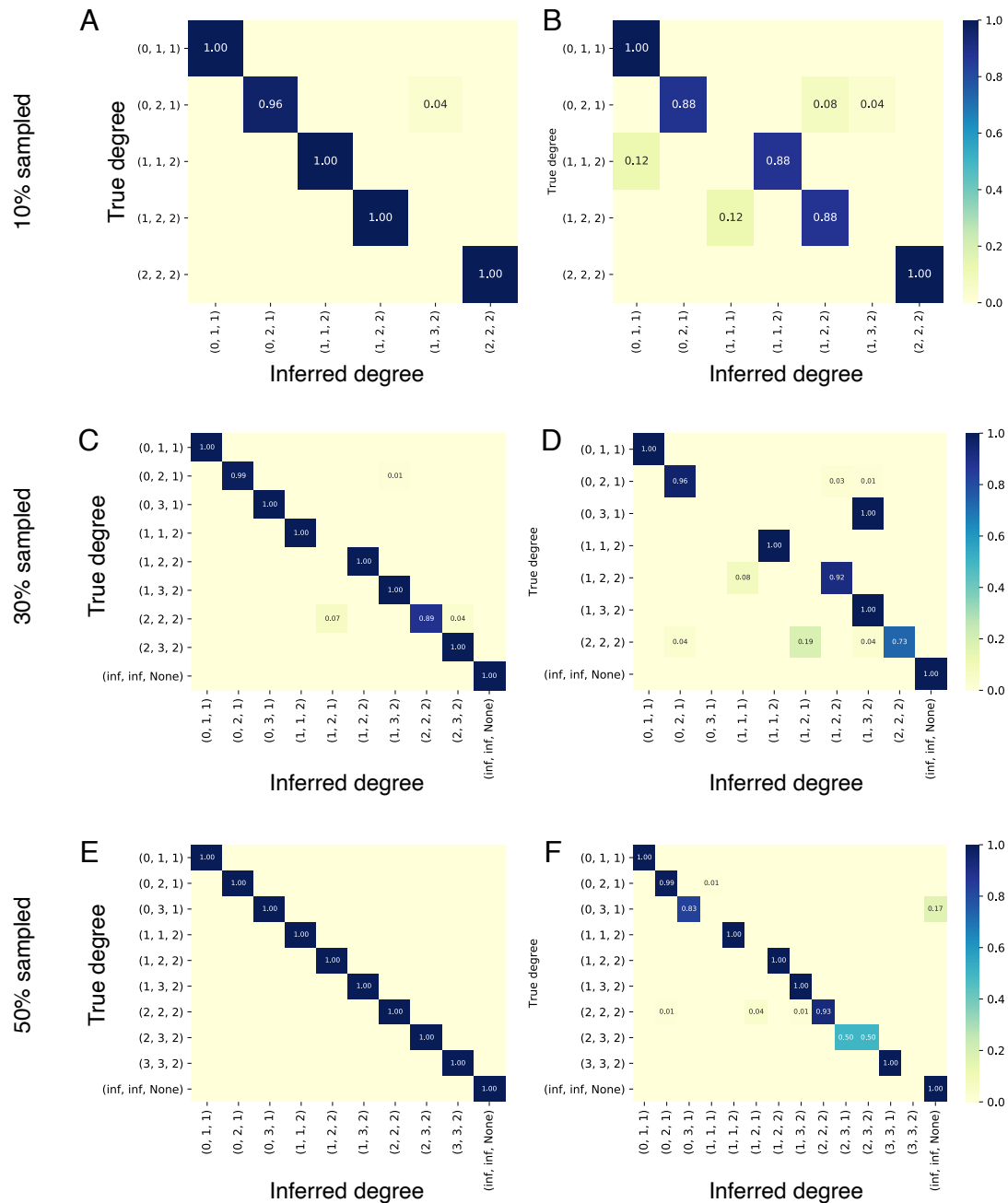


FIGURE 11. **Comparison of Bonsai with PRIMUS.** Panel A shows the accuracy of Bonsai for inferring different relationships when 10% of individuals were sampled from each pedigree. The relationship type between a pair of individuals  $i$  and  $j$  is indicated as a tuple of the form  $(d_{i,G}, d_{j,G}, |G|)$  following the notation of Ko and Nielsen (2017). Panel B shows the accuracy for PRIMUS applied to the same individuals with the same pairwise likelihoods as Panel A. Panels C and D compare Bonsai with PRIMUS when 30% of individuals were sampled and panels E and F compare Bonsai with PRIMUS when 50% of individuals were sampled. Accuracy of Bonsai was perfect for 60-100% of lineages, although PRIMUS continued to mis-infer some relationships.

### 30BONSAI: AN EFFICIENT METHOD FOR INFERRING LARGE HUMAN PEDIGREES FROM GENOTYPE DATA

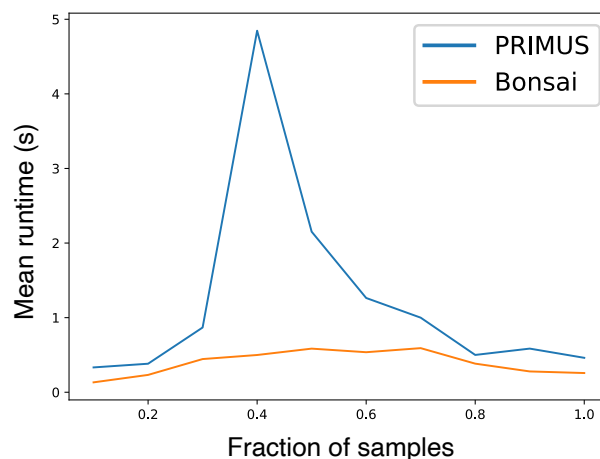


FIGURE 12. **Comparison of runtime between Bonsai and PRIMUS.** Runtime was evaluated using 204 pedigrees of 23andMe research participants that were known with a high degree of certainty. Because every individual in each pedigree was genotyped, we sampled 10, 20, 30, 40, 50, 60, 70, 80, 90, or 100% of their members uniformly at random without replacement. The subsampled individuals were then used to reconstruct the pedigrees using PRIMUS and Bonsai using the same IBD, age, and sex data. The x-axis in the figure is the fraction of sampled individuals in each of the 204 pedigrees that were used for pedigree inference. The y-axis is the average time in seconds required to reconstruct a pedigree.

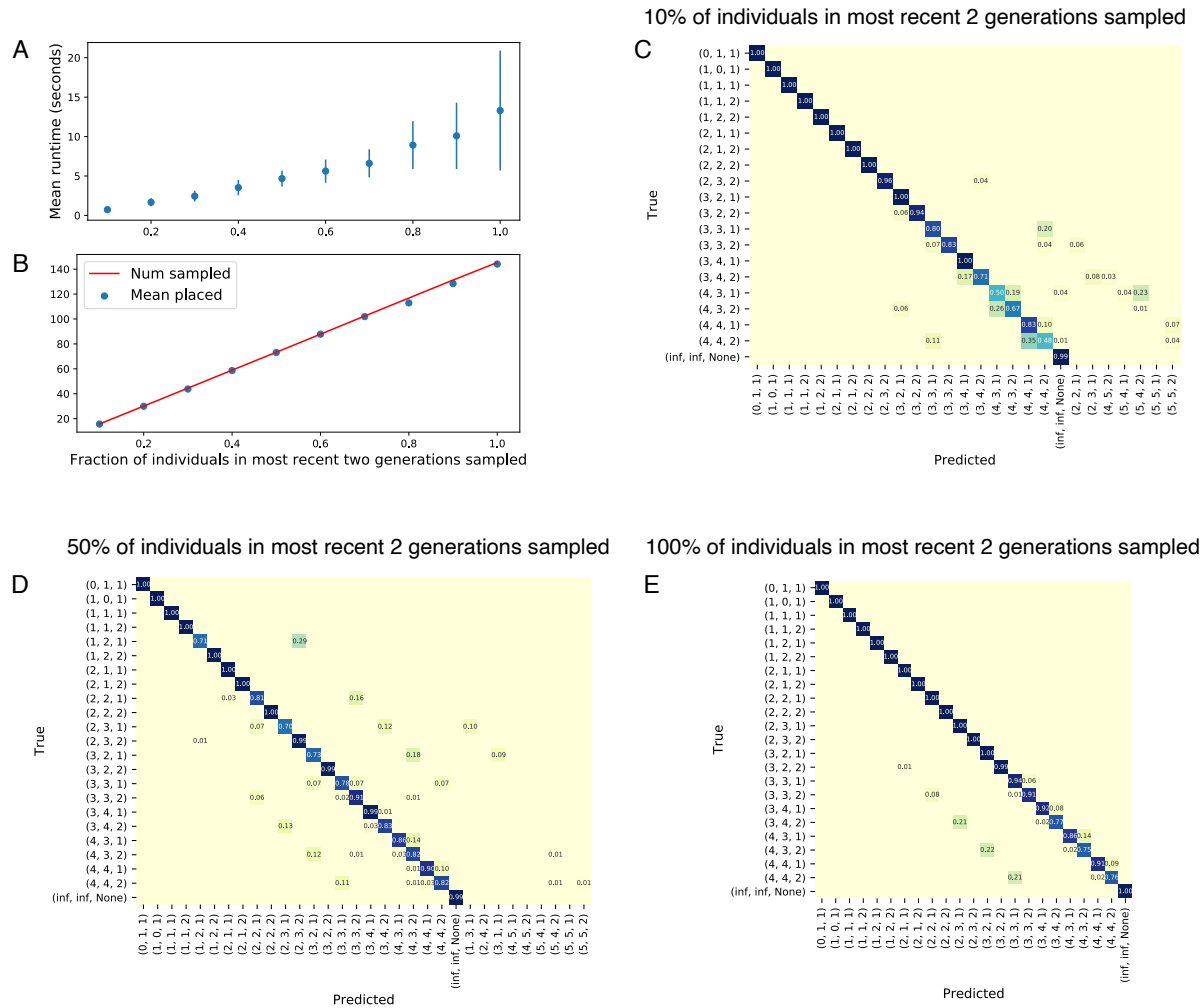
682 The Big Bonsai method is designed to drop small pedigrees from consideration, rather than  
683 combining them with the other pedigrees when an inconsistency is detected. This can occur, for  
684 example, if the small pedigree is inferred with a very unlikely relationship despite re-running with  
685 parameter values that search a broader pedigree space and attempting to correct relationships that  
686 are judged to be inaccurate. Figure 13B indicates that the fraction of times individuals or small  
687 pedigrees were dropped was small, as the number of placed individuals was typically very close to  
688 the number of sampled individuals.

689 Figures 13C-D show the accuracy for inferring large pedigrees when different fractions of indi-  
690 viduals were sampled. Close relationships were typically reconstructed accurately, whereas distant  
691 relationships were more challenging, yet still generally accurate especially when the fraction of  
692 sampled individuals was high.

693 Note that, because the ages of individuals in the pedigree conformed to average age differences  
694 between generations, it was sometimes possible to distinguish distant half relationships from distant  
695 full relationships. For example, a pair of individuals of the same age related by four degrees of  
696 separation is more likely to be a pair of half first cousins, rather than a full first cousin once removed.  
697 Half relationships are likely to be more challenging to infer in practice, given that age differences  
698 may differ from expectation.

699 **4.6. Reconstruction of self-reported pedigrees using Big Bonsai.** We also compared rela-  
700 tionships inferred by Bonsai with self-reported relationships using 265 pedigrees for which the

BONSAI: AN EFFICIENT METHOD FOR INFERRING LARGE HUMAN PEDIGREES FROM GENOTYPE DATA<sup>31</sup>



**FIGURE 13. Timing and accuracy of the Big Bonsai method.** Large pedigrees were simulated with a depth of five generations and two offspring per pair as described in Section 3.6.5. To capture the sparsity of pedigrees observed in direct-to-consumer customer data, we sampled only a fraction of individuals in each pedigree and used these as the genotyped individuals to infer the pedigree. Individuals were only sampled from the most recent two generations because ancestors are often ungenotyped in human data and because inference is more challenging when genotyped ancestors are unavailable. (A) Runtime for Big Bonsai as a function of the fraction of sampled individuals in the most recent two generations. (B) The number of sampled individuals in each pedigree and the mean number placed, averaged across 100 replicates. (C)-(E) The fraction of pairs with a given relationship type that were inferred to have each other relationship type. Tuples  $(d_i, G, d_j, G, |G|)$  indicate a specific relationship type between individuals  $i$  and  $j$  using the notation of Ko and Nielsen (2017): (up, down, number of common ancestors). The tuple (inf,inf,None) indicates unrelated individuals.

701 relationships between two or more individuals had been self-reported by the focal individual for  
702 whom the pedigree was built (Section 3.6.3).

703 Figure 14 shows the correspondence of each inferred relationship type with the self-reported  
704 relationship type. The plots show the fraction of times the self-reported and inferred relationships  
705 agreed exactly in that their relationship tuples (up, down, number of ancestors) were the same. The  
706 plots also show the fraction of times the relationships agreed in degree, the fraction of times the  
707 relationships agreed within one degree, and the fraction of times relationships agreed within two  
708 degrees.

709 The inferred and self-reported relationships typically agreed for close relationships up to first  
710 cousins. However, the inferred relationship often differed from the self-reported relationship for  
711 distant relationship types, and occasionally for relatives as close as siblings or parents. For parent-  
712 child and full sibling pairs, it is possible to check whether the self-reported relationship is correct  
713 because the IBD sharing patterns for these relationships are very distinct from other relationship  
714 types. It is of interest to note that in all but one case in which the inferred and self-reported  
715 relationships differed for a parent-child or full sibling pair, the self-reported relationship was, in fact,  
716 incorrect due to impossible levels of shared IBD. In these cases, it was frequently the case that a  
717 self-reported parent-child pair shared no IBD, or that a self-reported full sibling pair shared no IBD2  
718 and instead had an IBD sharing pattern that was more consistent with a half sibling or a cousin. In  
719 only one case was the self-reported relationship type consistent with the IBD sharing pattern, and  
720 in this case one individual had a self-reported age much greater than 100 years, leading to a strong  
721 contribution from the age component of the likelihood and an incorrectly inferred relationship type.

722 For distant relationships, we observed greater disparities between the self-reported and inferred  
723 values. However, the inferred degree was often within one or two degrees of the self-reported  
724 relationship, even for relationships as distant as seventh degree or higher in some cases. Moreover,  
725 relationships for which the self-reported and inferred degrees differed by more than two degrees  
726 typically had few self-reported pairs (Figure 14). This relatively high accuracy for distant relationship  
727 degree is consistent with our analysis of the accuracy of the generalized DRUID estimator.

728

## 5. DISCUSSION

729 We have presented a method for inferring large pedigrees quickly and accurately, even when the  
730 fraction of genotyped individuals in a pedigree is low and the distance between an individual and  
731 their closest relative can be moderate or large. Our method has three component algorithms that are  
732 applied in sequence: 1) a method to infer the likelihoods of pairwise relationships between each pair  
733 of individuals using both age and IBD data, 2) a method for inferring pedigrees of small-to-moderate  
734 size, and 3) a novel method for combining small pedigrees together into large and sparsely-sampled  
735 pedigrees.

736 Our Small Bonsai algorithm efficiently explores the space of possible pedigrees using a constructive  
737 approach. This approach is similar to that of PRIMUS (Staples et al. 2014), but it employs several  
738 new features that make it more efficient and more accurate than PRIMUS, including incorporating  
739 ages directly into the likelihoods, expanding the set of pedigrees that are explored, and introducing  
740 a branch-and-bound-like method for exploring the space of pedigrees more efficiently.



BONSAI: AN EFFICIENT METHOD FOR INFERRING LARGE HUMAN PEDIGREES FROM GENOTYPE DATA33

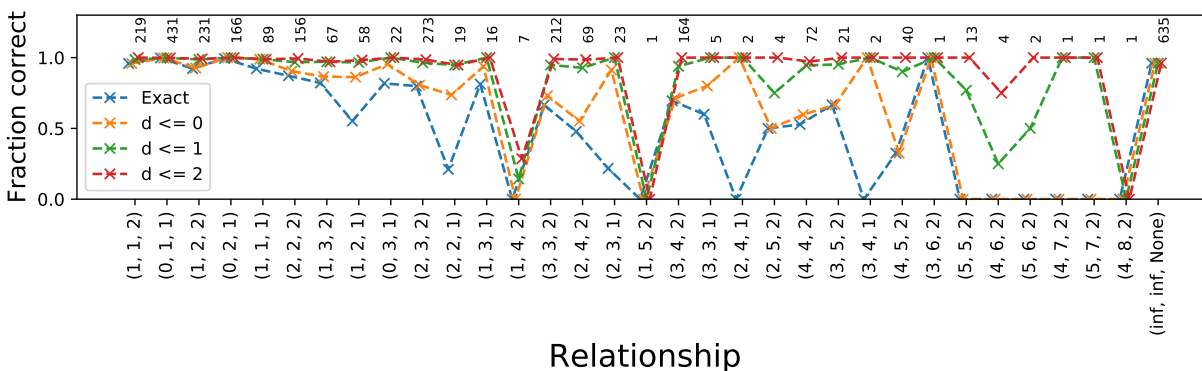


FIGURE 14. **Comparison with self-reported pedigrees.** Comparison of predicted relationships with self-reported relationships. Blue markers show the fraction of relationship pairs for which the inferred and self-reported relationships agreed exactly. The orange, green, and red markers show the fraction of pairs for which the degrees of the inferred and self-reported relationships differed by at most 0, 1, or 2 degrees, respectively. The number of pairs for each relationship is shown above the curves. Dashed lines are included to improve visibility.

741 Although the new methodological approaches implemented in the Small Bonsai method provide a  
 742 pedigree inference algorithm with improved accuracy and performance, the primary novelty of the  
 743 Bonsai method is in the Big Bonsai algorithm, which combines small pedigrees together into large  
 744 and sparsely-sampled pedigrees. This algorithm makes it possible to construct pedigrees that are  
 745 much bigger than the maximum size that can be constructed by current approaches.

746 The construction of large and sparse pedigrees requires a fundamentally different approach from  
 747 combining individuals one at a time as is done in PRIMUS, or searching a broad pedigree space by  
 748 rearranging pedigrees as is done in CLAPPER. Because the space of possible pedigrees is large, it  
 749 is useful to proactively and dramatically narrow the set of possible pedigrees to include only the  
 750 pedigrees with the highest likelihoods.

751 Combining small pedigrees together into large and sparse pedigrees makes it possible to leverage  
 752 information in the previously-inferred small pedigrees to identify the most likely ways in which the  
 753 small pedigrees can be connected together. Leveraging information across small pedigrees allows  
 754 us to more accurately infer the degree of relatedness between two small pedigrees and to identify  
 755 background IBD.

756 We have introduced three new tools for combining pedigrees together. First, we have generalized  
 757 the DRUID method of Ramstetter et al. (2018) to apply to general outbred pedigrees, rather  
 758 than specific pedigree structures. We have also extended the method to allow pedigrees to be  
 759 connected through pairs of individuals who are not common ancestors. We have also shown that the  
 760 generalized DRUID estimate is nearly identical to the maximum likelihood estimate. Thus, rather  
 761 than exploring multiple ways of connecting two pedigrees and selecting the most likely pedigree, we  
 762 can simply connect the two pedigrees through the DRUID point estimate and achieve nearly the  
 763 same result, greatly speeding up the inference process.

764 The second tool we have introduced is an approximate likelihood for the degree separating the  
765 common ancestors of two pedigrees as a function of the total length of IBD shared by the pedigrees.  
766 This likelihood is used as the foundation for our method for testing whether the IBD shared between  
767 two sets of individuals is the result of a true relationship, or whether the IBD is background IBD.  
768 Our approach obviates the need to infer the population or family-level distribution of IBD, which  
769 is useful because the expected amount of background IBD between a pair of individuals can be  
770 challenging to know in advance. By testing IBD between groups of individuals rather than pairs, we  
771 also reduce problems with multiple testing.

772 Although we intentionally did not incorporate the population or family-level distribution of  
773 background IBD into our background IBD detection method, one can imagine a method that  
774 combines our approach with such a distribution to improve the power for detecting background  
775 IBD when the population or family-level distribution of background IBD is known.

776 Finally, we have introduced a method for determining when the connection of pedigrees through  
777 certain ancestral branches is inconsistent with patterns of IBD overlap. This method makes it  
778 possible to assign two pedigrees to the correct parental sides of a focal individual in a focal pedigree.  
779 Using only information contained in pairwise IBD sharing, these inconsistent pedigrees would not  
780 be detected; pedigrees formed by connecting two pedigrees through incompatible grandparental  
781 lineages would appear to have the same likelihood as the true pedigree. Our approach achieves high  
782 accuracy even when few relatives on each parental side have been sampled.

783 Compared to previous methods for inferring complex human pedigrees, the Bonsai method yields  
784 improvements in both accuracy and computational efficiency and makes it possible to build pedigrees  
785 that are considerably larger than those that were possible before. The speed of pedigree building  
786 depends on the complexity of the pedigree, the proportion of individuals who are genotyped, and the  
787 distribution of these individuals throughout the generations of the pedigree. As a result, it can be  
788 difficult to characterize the runtime of Bonsai relative to other methods. However, in a comparison  
789 of runtime on 204 real-world pedigrees, Bonsai was always faster than the current fastest method  
790 PRIMUS. For large complicated pedigrees, Bonsai built pedigrees in a matter of seconds that took  
791 hours or which did not complete when built with PRIMUS or the Small Bonsai method alone.

792 Although we have presented an approach based on IBD segment overlaps for partitioning sets  
793 of distant relatives into their respective parental sides, relative to a focal individual or clade, it is  
794 likely that additional resolution could be gained by using IBD detected on sex chromosomes. At  
795 present, the Bonsai method uses only autosomal IBD to avoid considering the sexes of ancestral  
796 individuals along the paths connecting each pair of individuals when computing the likelihood of  
797 their relationship. Increased accuracy can also be obtained by using SNP-level information in our  
798 test of IBD overlap, such as opposite homozygotes, instead of IBD segments, as overlaps often occur  
799 between segments that are too short to be identified by existing IBD methods.

800 There is also potential to improve close relationship estimates by using phasing information.  
801 Williams et al. (2020) have demonstrated that half-sibling, avuncular, and grandparental relation-  
802 ships, which have been difficult to differentiate historically due to the fact that the total amount of  
803 IBD is the same for each of these relationship types, can be differentiated by making use of long-range  
804 phasing information. Phased IBD estimates, obtained from programs such as the PhasedIBD method

805 of Freyman et al. (2020), could provide a considerable boost in accuracy for close relationships.  
806 Improved close relationships would lead to improved distant relationships due to the fact that the  
807 small pedigree structures being connected would be more accurate. The PhasedIBD method of  
808 Freyman et al. could also improve distant relationship estimates through more accurate inference of  
809 short IBD segments.

810 Although our method for detecting background IBD is able to distinguish background IBD from  
811 true IBD when the level of background IBD is low, the approach struggles when there is a significant  
812 quantity of background IBD. In such cases, other approaches for accounting for background IBD  
813 when inferring relationships can be used. One approach is to detect the amount of “self” IBD shared  
814 between homologous chromosomes in each individual in a pedigree. Assuming that all individuals  
815 in the pedigree come from the same population, the amount of self IBD provides an expected level  
816 of background IBD sharing between two haplotypes that can then be subtracted from each pairwise  
817 relationship of which the individual is a member. We find that this approach improves pedigree  
818 inference accuracy in practice.

819 The approach of using self IBD to adjust pairwise IBD estimates can also be an effective approach  
820 when inferring pedigrees with recent consanguinity. The current Bonsai method assumes that  
821 pedigrees are graphs without cycles. However, it is possible to include cycles when adding new  
822 individuals to the pedigree if individuals have substantial self IBD and their relationships with  
823 others are indicative of recent consanguinity. This approach can be used together with distributions  
824 that are specifically trained on relationships with consanguinity.

825 Approaches for inferring pedigrees in the context of background IBD and consanguinity are  
826 important for improving pedigree inference in all human populations. Although the theoretically  
827 maximal accuracy with which a pedigree can be inferred differs across human populations due to  
828 differences in demographic histories, it is likely that improvements in accuracy can be attained for  
829 all populations through improved methodology, such as the improvement of pairwise relationship  
830 inference by methods such as deep-learning trained in specific populations, the inclusion of additional  
831 consanguineous relationship types, and the inclusion of additional genetic information from sex  
832 chromosomes and mitochondrial DNA. By nature, pedigree inference is a complicated problem  
833 requiring methods that can handle a wide variety of pedigree structures and input data. However,  
834 our results show that the inference of large and sparse human pedigrees is tractable, and that  
835 accuracy will continue to increase as pedigrees become increasingly densely sampled.

836

## 6. APPENDIX

837 **6.1. The probability of a pattern of IBD.** Consider the induced subtree in a pedigree relating  
838 a set of genotyped individuals. This tree is shown with dashed red lines in Figure 4 with nodes of  
839 the tree indicated with red dots. Let  $a(i)$  denote the direct ancestral node of node  $i$  in this tree.  
840 For example, in the tree in Figure 4, we have  $a(1) = A_1$ ,  $a(6) = A_1$ ,  $a(2) = 6$ ,  $a(3) = 6$ ,  $a(4) = A_2$ ,  
841  $a(5) = A_2$ ,  $a(A_1) = G$  and  $a(A_2) = G$ .

Assuming that all IBD segments are observed, we have

$$\mathbb{P}(O_i = 1) = \mathbb{P}(O_i = 1 | O_{a(i)} = 1) \mathbb{P}(O_{a(i)} = 1) + \mathbb{P}(O_i = 1 | O_{a(i)} = 0) \mathbb{P}(O_{a(i)} = 0)$$

$$\begin{aligned}
 &= \mathbb{P}(O_i = 1 | O_{a(i)} = 1) \mathbb{P}(O_{a(i)} = 1) \\
 &= 2^{-d_{i,a(i)}} \mathbb{P}(O_{a(i)} = 1),
 \end{aligned} \tag{18}$$

where  $d_{i,a(i)}$  is the number of meioses separating individual  $i$  from their ancestor  $a(i)$ . Similarly, we have

$$\begin{aligned}
 \mathbb{P}(O_i = 0) &= \mathbb{P}(O_i = 0 | O_{a(i)} = 1) \mathbb{P}(O_{a(i)} = 1) + \mathbb{P}(O_i = 0 | O_{a(i)} = 0) \mathbb{P}(O_{a(i)} = 0) \\
 &= [1 - 2^{-d_{i,a(i)}}] \mathbb{P}(O_{a(i)} = 1) + \mathbb{P}(O_{a(i)} = 0).
 \end{aligned} \tag{19}$$

842 In the final lines of Equations (18) and (19), we have used the fact that the probability that an  
 843 allelic copy is transmitted in one meiosis is  $1/2$ .

Equations (18) and (19) establish a recursion for computing the probability of an observed presence and absence pattern from a given ancestral allelic copy at a single base of the genome. Defining

$$p_{i,0} \equiv \mathbb{P}(O_i = 0), p_{i,1} \equiv \mathbb{P}(O_i = 1),$$

we can express the recursion compactly as

$$\begin{aligned}
 p_{i,0} &= [1 - 2^{-d_{i,a(i)}}] p_{a(i),1} + p_{a(i),0}, \\
 p_{i,1} &= 2^{-d_{i,a(i)}} p_{a(i),1},
 \end{aligned}$$

844 with the base conditions  $p_{g,0} = 0$  and  $p_{g,1} = 1$  for each chromatid,  $g$ , in  $G$ . The probability of an  
 845 observed IBD sharing pattern  $\{O_1, \dots, O_k\}$  across  $k$  leaf nodes can be computed recursively using  
 846 Equation (6).

847 **6.2. Approximating the variance of  $T_{1,2}$ .** Here, we derive an approximation of the variance of  
 848 the total length,  $T_{1,2}$ , of IBD shared across the genotyped descendants of two acenstral individuals,  $A_1$   
 849 and  $A_2$ . When we encounter a patch of IBD at a locus, the length of the patch can be approximated  
 850 as the maximum length of  $|\mathcal{N}_1| \times |\mathcal{N}_2|$  different IBD segments, where  $\mathcal{N}_i$  is the set of genotyped  
 851 nodes below ancestor  $A_i$  at locus  $m$  in which the IBD segment is observed. This approximation  
 852 comes from conceptualizing IBD sharing among the  $|\mathcal{N}_1|$  IBD segment carrying descendants of  $A_1$   
 853 and the  $|\mathcal{N}_2|$  IBD segment carrying descendants of  $A_2$  as  $|\mathcal{N}_1| \times |\mathcal{N}_2|$  independent segments with  
 854 a single point at which all segments overlap. The length of the merged segment to one side of  
 855 this focal point then has a distribution given by the maximum of  $|\mathcal{N}_1| \times |\mathcal{N}_2|$  exponential random  
 856 variables whose means depend on the degree of separation between the corresponding pairs of leaf  
 857 individuals. To simplify matters, we assume that the length of the full merged overlapping segment  
 858 (not just to the left or right) is exponentially distributed.

859 This approximation is an oversimplification of the IBD sharing pattern because the segments are  
 860 not truly independent and need not overlap at a single point. Moreover, under this approximation,  
 861 the length of the merged segment would be the maximum over sums of identically distributed  
 862 random variables, representing the sum of the length of a segment to the right of the center point and  
 863 the length of the segment to the left. However, we are not overly concerned with these drawbacks of  
 864 the conceptualization because our main goal is to obtain an accurate, yet simple approximation  
 865 of the variance of the distribution. We also assume that no member of  $\mathcal{N}_i$  is the direct ancestor

866 of another member of the set, which holds in practice if we drop all individuals from  $\mathcal{N}_i$  who are  
867 descended from others.

868 The length,  $\ell_{i,j}$ , of an IBD segment between leaf nodes  $i$  and  $j$  is can be modeled as an exponentially  
869 distributed random variable with mean length  $\mu_{ij} = L_{genome}/d_{i,j}R$ , where  $d_{i,j}$  is the degree of  
870 relationship between them and  $R$  is the expected number of recombination events, genome wide, in  
871 one meiosis (Huff et al. 2011). When the length of the genome is expressed in centimorgans (cM), the  
872 expected number of recombination events in the genome is  $L_{genome}/100$ . Thus, the expected length  
873 in cM of an IBD segment between individuals  $i$  and  $j$  separated by  $d_{i,j}$  meioses is  $\mu_{ij} = 100/d_{i,j}$ .

874 Let  $L_{1,2}$  denote a random variable describing the length of the segment formed by merging all  
875 segments at a given locus  $m$  between descendants of  $A_1$  and  $A_2$ . If the lengths of all segments at  
876 this locus were independent, their merged length in our conceptualization would have a distribution  
877 given by the maximum over independent exponentially distributed random variables with means  
878  $\{\mu_{i,j}\}_{i \in \mathcal{N}_1, j \in \mathcal{N}_2}$ .

If the leaf nodes with observed IBD at the marker are  $\mathcal{N}_1$  and  $\mathcal{N}_2$ , then we have  $L_{1,2} = \max(\{\ell_{i,j}\}_{i \in \mathcal{N}_1, j \in \mathcal{N}_2})$ . Under this condition, the cumulative density function (CDF)  $F_L(\ell; \mathcal{N}_1, \mathcal{N}_2)$  of  $L$  is

$$\begin{aligned}
 F_{L_{1,2}}(\ell; \mathcal{N}_1, \mathcal{N}_2) &= \mathbb{P}(L_{1,2} < \ell; \mathcal{N}_1, \mathcal{N}_2) \\
 &= \mathbb{P}(\ell_{i,j} < \ell, \text{ for } i \in \mathcal{N}_1, j \in \mathcal{N}_2) \\
 &= \prod_{i \in \mathcal{N}_1} \prod_{j \in \mathcal{N}_2} \mathbb{P}(\ell_{i,j} < \ell) \\
 &= \prod_{i \in \mathcal{N}_1} \prod_{j \in \mathcal{N}_2} (1 - e^{-\lambda_{i,j}\ell}) \\
 &= 1 - \sum_{i \in \mathcal{N}_1, j \in \mathcal{N}_2} e^{-\lambda_{i,j}\ell} + \sum_{i,u \in \mathcal{N}_1, j,v \in \mathcal{N}_2} e^{-(\lambda_{i,j} + \lambda_{u,v})\ell} (1 - \delta_{(i,j),(u,v)}) \\
 &\quad - \sum_{i,u,w \in \mathcal{N}_1, j,v,z \in \mathcal{N}_2} e^{-(\lambda_{i,j} + \lambda_{u,v} + \lambda_{z,w})\ell} (1 - \delta_{(i,j),(u,v)})(1 - \delta_{(i,j),(z,w)})(1 - \delta_{(u,v),(z,w)}) + \dots,
 \end{aligned} \tag{20}$$

879 where  $\lambda_{i,j} = 1/\mu_{i,j} = d_{i,j}/100$  and  $\delta_{(a,b),(c,d)}$  is the Kronecker delta between tuples  $(a, b)$  and  $(c, d)$ ,  
880 which is equal to one when  $(a, b) = (c, d)$  and zero, otherwise.

The sets  $\mathcal{N}_1$  and  $\mathcal{N}_2$  are, themselves, random variables. Summing over all sets  $\mathcal{N}_1$  and  $\mathcal{N}_2$ , we have

$$F_{L_{1,2}}(\ell) = \sum_{\mathcal{N}_1, \mathcal{N}_2} F_{L_{1,2}}(\ell; \mathcal{N}_1, \mathcal{N}_2) \mathbb{P}(\mathcal{N}_1) \mathbb{P}(\mathcal{N}_2), \tag{22}$$

881 where the probabilities  $\mathbb{P}(\mathcal{N}_1)$  and  $\mathbb{P}(\mathcal{N}_2)$  are probabilities of observing IBD in the sets of leaf nodes  
882 below  $A_1$  and  $A_2$  and can be approximated using the recursion in Equation (6).

883 Over the length of the genome, the number  $N_{1,2}$  of IBD segments between the descendants of  
884  $A_1$  and  $A_2$  is approximately Poisson distributed with mean  $(1 - \mathbb{P}(\mathcal{I}^c)^{2|G|})L_{genome}/E[L_{1,2}]$ . This  
885 rate comes from the fact that the average total amount of the genome in a patch of IBD is

886  $(1 - \mathbb{P}(\mathcal{I}^c)^{2|G|})L_{genome}$  while the average length of any given segment is  $E[L_{1,2}]$ . When the lengths  
 887 of IBD are short and far apart, which they are when the degree between  $A_1$  and  $A_2$  is large, this is  
 888 a reasonable approximation. This is precisely the regime in which the distribution in Equation (15)  
 889 is most useful.

The total length  $T_{1,2}$  of merged IBD among the descendants of  $A_1$  and  $A_2$  is

$$T_{1,2} = \sum_{n=1}^{N_{1,2}} L_{1,2}. \quad (23)$$

We can derive the variance of  $T_{1,2}$  using the law of total variance as

$$\begin{aligned} \text{Var}(T_{1,2}) &= \mathbb{E}[\text{Var}(T_{1,2}|N_{1,2})] + \text{Var}(\mathbb{E}[T_{1,2}|N_{1,2}]) \\ &= \mathbb{E}[N_{1,2}\text{Var}(L_{1,2})] + \text{Var}(N_{1,2}\mathbb{E}[L_{1,2}]) \\ &= \mathbb{E}[N_{1,2}]\text{Var}(L_{1,2}) + \text{Var}(N_{1,2})\mathbb{E}[L_{1,2}]^2. \end{aligned} \quad (24)$$

Note that because  $N_{1,2} \sim \text{Poisson}((1 - \mathbb{P}(\mathcal{I}^c)^{2|G|})L_{genome}/\mathbb{E}[L_{1,2}])$ , we have

$$\mathbb{E}[N_{1,2}] = \text{Var}(N_{1,2}) = (1 - \mathbb{P}(\mathcal{I}^c)^{2|G|})L_{genome}/\mathbb{E}[L_{1,2}]. \quad (25)$$

So Equation (24) simplifies to

$$\begin{aligned} \text{Var}(T_{1,2}) &= \frac{(1 - \mathbb{P}(\mathcal{I}^c)^{2|G|})L_{genome}}{\mathbb{E}[L_{1,2}]} [\text{Var}(L_{1,2}) + \mathbb{E}[L_{1,2}]^2] \\ &= (1 - \mathbb{P}(\mathcal{I}^c)^{2|G|})L_{genome} \frac{\mathbb{E}[L_{1,2}^2]}{\mathbb{E}[L_{1,2}]}, \end{aligned} \quad (26)$$

890 where we have used the fact that  $\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2$ .

It remains to find  $\mathbb{E}[L_{1,2}]$  and  $\mathbb{E}[L_{1,2}^2]$ . Using the CDF of  $L_{1,2}$  in Equation (22) and the fact that  $\mathbb{E}[X^m] = m! \int_{\mathbb{R}} x^{m-1} [1 - F_X(x)] dx$ , we have

$$\begin{aligned} \mathbb{E}_{\mathcal{N}_1, \mathcal{N}_2}[L_{1,2}^m] &= m! \int_{\ell=0}^{\infty} x^{m-1} [1 - F_{L_{1,2}}(\ell; \mathcal{N}_1, \mathcal{N}_2)] d\ell \\ &= \sum_{i \in \mathcal{N}_1, j \in \mathcal{N}_2} \int_{\ell=0}^{\infty} m! \ell^{m-1} e^{-\lambda_{i,j} \ell} d\ell - \sum_{i, u \in \mathcal{N}_1, j, v \in \mathcal{N}_2} \int_{\ell=0}^{\infty} m! \ell^{m-1} e^{-(\lambda_{i,j} + \lambda_{u,v}) \ell} d\ell \\ &\quad + \sum_{i, u, w \in \mathcal{N}_1, j, v, z \in \mathcal{N}_2} \int_{\ell=0}^{\infty} m! \ell^{m-1} e^{-(\lambda_{i,j} + \lambda_{u,v} + \lambda_{z,w}) \ell} d\ell + \dots \\ &= \sum_{i \in \mathcal{N}_1, j \in \mathcal{N}_2} \frac{m!}{\lambda_{i,j}^m} - \sum_{i, u \in \mathcal{N}_1, j, v \in \mathcal{N}_2} \frac{m!}{(\lambda_{i,j} + \lambda_{u,v})^m} \\ &\quad + \sum_{i, u, w \in \mathcal{N}_1, j, v, z \in \mathcal{N}_2} \frac{m!}{(\lambda_{i,j} + \lambda_{u,v} + \lambda_{z,w})^m} + \dots \end{aligned} \quad (27)$$

891 where the integrals in Equation (27) can be evaluated by noting that they are essentially expressions  
 892 for the moments of exponential random variables with parameters  $\lambda_i$ ,  $(\lambda_i + \lambda_j)$ ,  $(\lambda_i + \lambda_j + \lambda_k)$ , etc.

Thus, we can use Equation (27) to compute

$$\mathbb{E}[L_{1,2}^m] = \sum_{\mathcal{N}_1, \mathcal{N}_2} \mathbb{E}_{\mathcal{N}_1, \mathcal{N}_2}[L_{1,2}^m] \mathbb{P}(\mathcal{N}_1, \mathcal{N}_2), \quad (28)$$

893 where  $\mathbb{P}(\mathcal{N}_1, \mathcal{N}_2)$  is the probability of observing IBD segments at the leaves  $\mathcal{N}_1$  and  $\mathcal{N}_2$ , and is  
894 approximated using the recursion in Equation (6). We then plug Equation (28) in to obtain the  
895 variance of  $T_{1,2}$  in Equation (26).

In practice, it is too computationally demanding to compute the sums in Equation (28) because the terms  $\mathbb{E}_{\mathcal{N}_1, \mathcal{N}_2}[L_{1,2}]$ ,  $\mathbb{E}_{\mathcal{N}_1, \mathcal{N}_2}[L_{1,2}^2]$ , and  $\mathbb{P}(\mathcal{N}_1, \mathcal{N}_2)$  are not fast to compute in large quantities. However, the probabilities  $\mathbb{P}(\mathcal{N}_1, \mathcal{N}_2)$  can be computed quickly enough, allowing us to find the most likely sets of leaf nodes,  $\hat{\mathcal{N}}_1$  and  $\hat{\mathcal{N}}_2$ , with observed IBD. Thus, in practice we use an approximation in which we assume that the most likely IBD pattern has been observed and we compute

$$\mathbb{E}[L_{1,2}^m] \approx \mathbb{E}_{\hat{\mathcal{N}}_1, \hat{\mathcal{N}}_2}[L_{1,2}^m]. \quad (29)$$

896 The assumption used in this approximation is that most patterns of observed IBD at the leaves are  
897 unlikely compared with the most likely patterns and that most likely patterns of IBD will yield  
898 similar moments  $\mathbb{E}[L_{1,2}^m]$ .

899 **6.3. Re-rooting the DRUID estimator.** In some scenarios, the common ancestors,  $A_1$  and  $A_2$ ,  
900 of sets of individuals  $\mathcal{N}_1$  and  $\mathcal{N}_2$ , may not be related through a common ancestor or ancestral pair of  
901 both  $A_1$  and  $A_2$ . In particular  $A_2$  can be the direct descendant of  $A_1$ , or vice versa. This scenario,  
902 along with the scenario treated in Section 3.5.3 in which  $\mathcal{N}_1$  and  $\mathcal{N}_2$  are connected through their  
903 common ancestors, covers all possible ways in which  $\mathcal{N}_1$  and  $\mathcal{N}_2$  can be connected such that they  
904 are mutually related, i.e., so that they share a common ancestor.

905 We now describe an approach for computing the generalized DRUID estimate when  $A_2$  is  
906 descended from an individual  $A$  who is the common ancestor of only a subset of  $\mathcal{N}_1$ . We consider  $A$   
907 to be any node ancestral to some node in  $\mathcal{N}_1$ , including any member of  $\mathcal{N}_1$  itself.

908 Let  $\Lambda_1(A_1)$  denote the induced subtree in pedigree  $\mathcal{P}_1$  that relates  $A_1$  and their descendants  
909  $\mathcal{N}_1$ . To obtain the generalized DRUID estimate when  $A_2$  is descended from  $A$ , we re-root the tree  
910  $\Lambda_1(A_1)$  at  $A$  to obtain a re-rooted tree  $\tilde{\Lambda}_1(A)$ . We then compute the generalized DRUID estimate  
911 from Section 3.5.3 using the re-rooted tree  $\tilde{\Lambda}_1(A)$ . The estimate between  $A$  and  $A_2$  obtained using  
912 Equation (9) applied to  $\tilde{\Lambda}_1(A)$  and  $\Lambda_2(A_2)$  is then the number of meioses separating  $A$  and  $A_2$ .

913 The one complication is that  $A_2$  can be descended from both  $A$  and a spouse  $A'$ , who is also an  
914 ancestor of one or more of  $A$ 's genotyped descendants  $\mathcal{N}_A$ . In this case,  $A_2$  is more closely related  
915 to  $\mathcal{N}_A$  than to  $\mathcal{N}_1 \setminus \mathcal{N}_A$  by one degree. We solve this problem by representing the clades of shared  
916 descendants *twice* on the re-rooted tree, obtaining a multi-labeled tree (Figure 15). In contrast, if  
917  $A_2$  is descended from  $A$  and a spouse  $A''$  who is not ancestral to any genotyped descendants, we do  
918 not duplicate the descendants of  $A$  on the tree.

## 919 7. ACKNOWLEDGEMENTS

920 We would like to thank the employees and research participants of 23andMe who made this  
921 research possible. We would also like to thank Amy Williams and Ying Qiao for helpful discussions.  
922 Members of the 23andMe Research Team are Michelle Agee, Stella Aslibekyan, Elizabeth Babalola,  
923 Robert K. Bell, Jessica Bielenberg, Katarzyna Bryc, Emily Bullis, Briana Cameron, Daniella Coker,  
924 Gabriel Cuellar Partida, Devika Dhamija, Sayantan Das, Sarah L. Elson, Teresa Filshtein, Kipper  
925 Fletez-Brant, Pierre Fontanillas, Pooja M. Gandhi, Karl Heilbron, Barry Hicks, David A. Hinds,

SAIBONSAI: AN EFFICIENT METHOD FOR INFERRING LARGE HUMAN PEDIGREES FROM GENOTYPE DATA

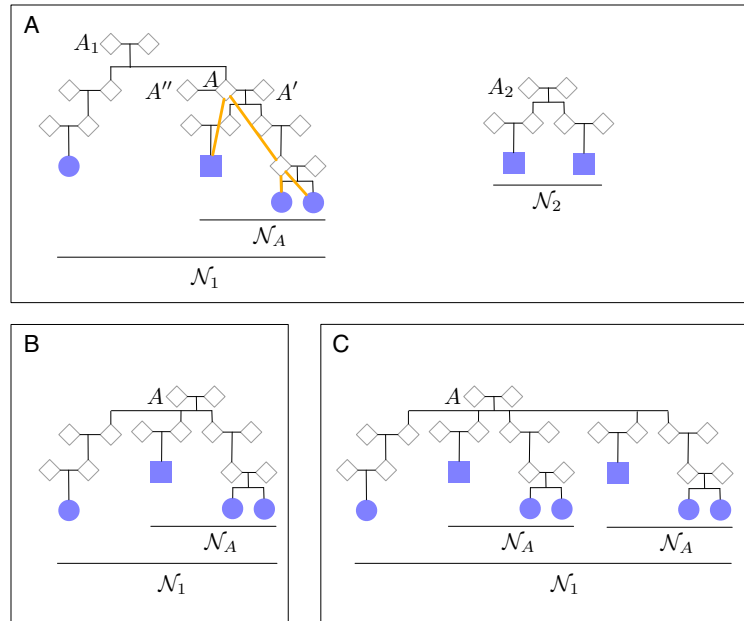


FIGURE 15. **Re-rooting the DRUID estimator.** (A) The pedigrees relating  $\mathcal{N}_1$  and  $\mathcal{N}_2$ , respectively. The induced subtree relating the descendants of internal node  $A$  is shown in orange. (B) The re-rooted tree when  $A_2$  is descended from  $A$  and  $A''$ . (C) The re-rooted tree when  $A_2$  is descended from  $A$  and  $A'$ .

926 Karen E. Huber, Yunxuan Jiang, Aaron Kleinman, Katelyn Kukar, Keng-Han Lin, Maya Lowe,  
 927 Marie K. Luff, Jennifer C. McCreight, Matthew H. McIntyre, Steven J. Micheletti, Meghan E.  
 928 Moreno, Joanna L. Mountain, Sahar V. Mozaffari, Priyanka Nandakumar, Elizabeth S. Noblin, Jared  
 929 O'Connell, Aaron A. Petrakovitz, G. David Poznik, Anjali J. Shastri, Janie F. Shelton, Jingchunzi  
 930 Shi, Suyash Shringarpure, Chao Tian, Vinh Tran, Joyce Y. Tung, Xin Wang, Wei Wang, Catherine  
 931 H. Weldon, and Peter Wilton.

932 8. DECLARATIONS OF INTERESTS

933 The authors are employees of 23andMe, Inc., and hold stock or stock options in 23andMe.

934 9. WEB RESOURCES

935 The Bonsai code is available at <https://github.com/23andMe/bonsaitree>.

936 10. SUPPLEMENTAL FIGURES





**FIGURE S1. Example of a large pedigree for testing the Big Bonsai method.** 100 such pedigrees were simulated using the approach described in Section 3.6.5. The focal individual is labeled “1.”

S430NSAI: AN EFFICIENT METHOD FOR INFERRING LARGE HUMAN PEDIGREES FROM GENOTYPE DATA

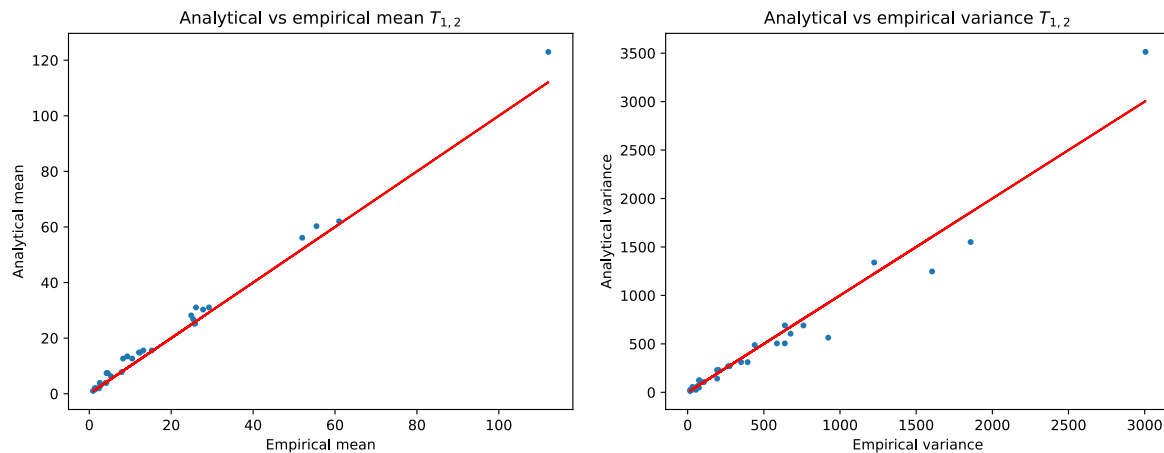


FIGURE S2. **Analytical versus empirical mean and variance of  $T_{1,2}$ .** Analytical means were computed using Equation (12) and analytical variances were computed using Equation (13). Empirical means and variances were computed using simulated pedigrees comprised of two small pedigrees  $\mathcal{P}_1$  and  $\mathcal{P}_2$  connected through either one or two common ancestors. Pedigree  $\mathcal{P}_1$  had a randomly generated structure simulated by starting with the pair of root individuals and their two children. At each subsequent generation, each leaf node had a probability  $1/2$  of having one offspring and  $1/2$  of having two offspring. Pedigree  $\mathcal{P}_1$  was extended down from the root nodes until the total number of leaves was  $|\mathcal{N}_1|$ . Pedigree  $\mathcal{P}_2$  was simulated in the same way, but independently of  $\mathcal{P}_1$ . One common ancestor  $A_1$  of  $\mathcal{P}_1$  was then connected to one common ancestor  $A_2$  of  $\mathcal{P}_2$  through either one or two common ancestors,  $G$ . The degrees  $d_{A_1,G}$  and  $d_{A_2,G}$  were set to either 3 or 5 and the number of common ancestors  $|G|$  was either 1 or 2. Thus, the genealogical degrees between  $A_1$  and  $A_2$  were in the set  $\{5, 6, 7, 8, 9, 10\}$ . The number of leaves in each pedigree was either  $\mathcal{N}_i = 2$  or  $\mathcal{N}_i = 4$ . We ran 10 simulation replicates for each configuration of  $d_{A_1,G}$ ,  $d_{A_2,G}$ ,  $|\mathcal{N}_1|$ ,  $|\mathcal{N}_2|$ , and  $|G|$ .

937

REFERENCES

- 938 A. Almudevar. A simulated annealing algorithm for maximum likelihood pedigree reconstruction.  
939 *Theor. Popul. Biol.*, 63:63–75, 2003.
- 940 A. Almudevar and E.C. Anderson. A new version of prt software for sibling groups reconstruction  
941 with comments regarding several issues in the sibling reconstruction problem. *Mol. Ecol. Resour.*,  
942 12:164–178, 2012.
- 943 E.C. Anderson and T.C. Ng. Bayesian pedigree inference with small numbers of single nucleotide  
944 polymorphisms via a factor-graph representation. *Theor. Popul. Biol.*, 107:39–51, 2016.
- 945 R.G. Cowell. Efficient maximum likelihood pedigree reconstruction. *Theor. Popul. Biol.*, 76:285–291,  
946 2009.
- 947 R.G. Cowell. A simple greedy algorithm for reconstructing pedigrees. *Theor. Popul. Biol.*, 83:55–63,  
948 2013.
- 949 J. Cussens, M. Bartlett, E.M. Jones, and N.A. Sheehan. Maximum likelihood pedigree reconstruction  
950 using integer linear programming. *Genet. Epid.*, 37:69–83, 2013.

BONSAI: AN EFFICIENT METHOD FOR INFERRING LARGE HUMAN PEDIGREES FROM GENOTYPE DATA 343

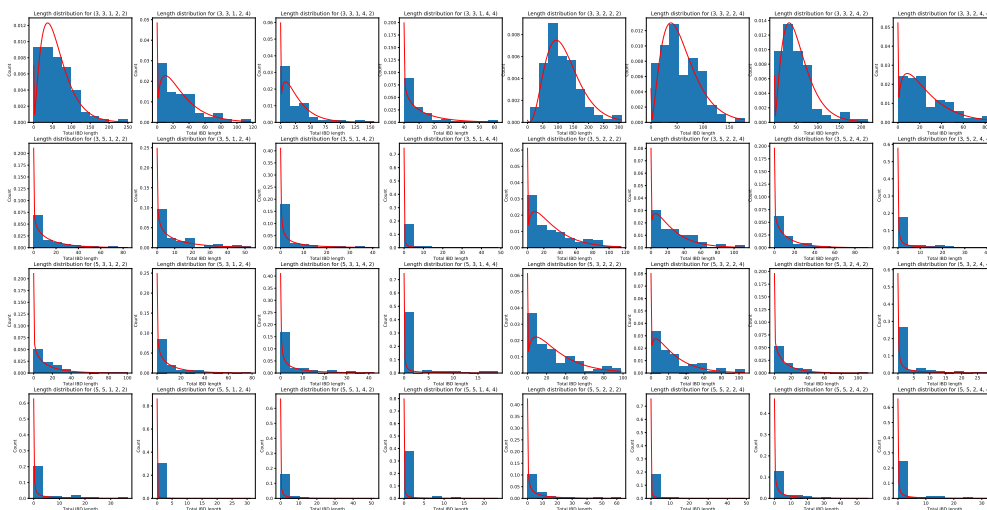


FIGURE S3. **Analytical versus empirical distributions of  $T_{1,2}$ .** Analytical distributions were computed using Equation (15). Empirical distributions were simulated in the manner described in Figure S2.

- 951 J. Felsenstein. Evolutionary trees from dna sequences: A maximum likelihood approach. *J. Mol.*  
 952 *Evol.*, 17:368–376, 1981.
- 953 W.A. Freyman, K.F. McManus, Shringarpure S.S., E.M. Jewett, K. Bryc, 23andMe Research Team,  
 954 and A. Auton. Fast and robust identity-by-descent inference with the templated positional  
 955 burrows-wheeler transform. *Mol Biol Evol.*, 2020. doi: 10.1093/molbev/msaa328.
- 956 B.M. Henn, L. Hon, J.M. Macpherson, N. Eriksson, S. Saxonov, I. Pe’er, and J.L. Mountain. Cryptic  
 957 distant relatives are common in both isolated and cosmopolitan genetic samples. *PLOS One.*, 7:  
 958 e34267, 2012.
- 959 C.D. Huff, D.J. Witherspoon, T.S. Simonson, J. Xing, W.S. Watkins, Y. Zhang, T.M. Tuohy,  
 960 D.W. Neklason, R.W. Burt, S.L. Guthery, S.R. Woodward, and L.B. Jorde. Maximum-likelihood  
 961 estimation of recent shared ancestry (ersa). *Genome Research*, 21:768–774, 2011.
- 962 J. Huisman. Pedigree reconstruction from SNP data: parentage assignment, sibship clustering and  
 963 beyond. *Mol. Ecol.*, 17:1009–1024, 2017.
- 964 O.R. Jones and J. Wang. COLONY: a program for parentage and sibship inference from multilocus  
 965 genotype data. *Mol. Ecol. Resour.*, 10:551–555, 2017.
- 966 B. Kirkpatrick, S.C. Li, R.M. Karp, and E. Halperin. Pedigree reconstruction using identity by  
 967 descent. *J. Comp. Biol.*, 18:1481–1493, 2011.
- 968 A. Ko and R. Nielsen. Composite likelihood method for inferring local pedigrees. *PLOS Genet.*, 13:  
 969 e1006963, 2017.
- 970 J.C. Manichaikul, A. and Mychaleckyj, S.S. Rich, K. Daly, M. Sale, and W.M. Chen. Robust  
 971 relationship inference in genome-wide association studies. *Bioinformatics*, 26:2867–2873, 2010.

SAHONSAI: AN EFFICIENT METHOD FOR INFERRING LARGE HUMAN PEDIGREES FROM GENOTYPE DATA

- 972 M.D. Ramstetter, S.A. Shenoy, T.D. Dyer, D.M. Lehman, J.E. Curran, R. Duggirala, J. Blangero,  
973 J.G. Mezey, and A.L. Williams. Inferring identical-by-descent sharing of sample ancestors promotes  
974 high-resolution relative detection. *Am. J. Hum. Genet.*, 103:30–44, 2018.
- 975 M. Riester, P.F. Stadler, and K. Klemm. FRANz: reconstruction of wild multi-generation pedigrees.  
976 *Bioinformatics*, 25:2134–2139, 2009.
- 977 D.N. Seidman, S.A. Shenoy, M. Kim, R. Babu, I.G. Woods, T.D. Dyer, D.M. Lehman, J.E. Curran,  
978 R. Duggirala, J. Blangero, and A.L. Williams. Rapid, phase-free detection of long identity-by-  
979 descent segments enables effective relationship classification. *Am. J. Hum. Genet.*, 106:453–466,  
980 2020.
- 981 N.A. Sheehan, M. Bartlett, and J. Cussens. Improved maximum likelihood reconstruction of complex  
982 multi-generational pedigrees. *Theor. Popul. Biol.*, 97:11–19, 2014.
- 983 J. Staples, D. Qiao, M.H. Cho, E.K. Silverman, University of Washington Center for Mendelian Ge-  
984 nomics, D.A. Nickerson, and J.E. Below. PRIMUS: Rapid reconstruction of pedigrees from  
985 genome-wide estimates of identity by descent. *Am. J. Hum. Genet.*, 95:553–564, 2014.
- 986 J. Staples, D.J. Witherspoon, L.B. Jorde, D.A. Nickerson, University of Washington Center for  
987 Mendelian Genomics, J.E. Below, and C.D. Huff. PADRE: Pedigree-aware distant-relationship  
988 estimation. *Am. J. Hum. Genet.*, 0:<https://doi.org/10.1101/2020.02.25.965376>, 2016.
- 989 J. Wang. Sibship reconstruction from genetic data with typing errors. *Genetics*, 166:1963–1979,  
990 2004.
- 991 C.M. Williams, B. Scelza, C.R. Gignoux, and B.M. Henn. A rapid, accurate approach to inferring  
992 pedigrees in endogamous populations. *bioRxiv*, 99:154–162, 2020.