

1 **SNAPPY: Single Nucleotide Assignment of Phylogenetic Parameters on the Y chromosome**

2 **Alissa L. Severson^{1,*}, Jonathan A. Shortt^{2,*}, Fernando L. Mendez¹, Genevieve L. Wojcik¹,**

3 **Carlos D. Bustamante¹, Christopher R. Gignoux^{2,3}**

4 ¹Department of Genetics, Stanford University School of Medicine, Stanford, CA 94305, USA;

5 ²Colorado Center for Personalized Medicine, University of Colorado Anschutz Medical Campus,

6 Aurora, CO 80045, USA

7 *Contributed equally

8 ³To whom correspondence should be addressed

9 **Abstract**

10 *Summary*

11 The assignment of Y chromosome data to related clusters, or haplogroups, is a common

12 application in human population genetics. To enable this at scale, we developed SNAPPY.

13 SNAPPY is a software program used to assign Y-chromosome phylogeny-informed haplotypes

14 using dense genotype data. The program efficiently tests all haplotypes in a provided Y-

15 chromosome database to find the haplogroup that is best supported by the input genotypes.

16 Importantly, the method considers both the amount of support for the specific haplogroup, as

17 well as its ancestral haplogroups via parsimony. This accounts for the underlying genealogy the

18 haplotypes represent, strengthening the accuracy of the assignments. SNAPPY is fast, scalable,

19 and uses standard file formats, making it easy to integrate into analytical pipelines.

20 *Availability and Implementation*

21 The program is implemented in python. The program, a user manual, haplotype databases, and
22 test datasets are available for download at github.com/chrisgene/snappy.

23 *Contact*

24 Jonathan.shortt@ucdenver.edu, Chris.gignoux@ucdenver.edu

25 **Introduction**

26 Analyses of Y-chromosome haplogroups have yielded important insights into human
27 migration history (Bergstrom, et al., 2016; Chiaroni, et al., 2009), cultural customs (Seielstad, et
28 al., 1998), and ancestral population sizes (Karmin, et al., 2015). These insights are possible
29 because of the two unique characteristics of the Y-chromosome: it is passed directly from father
30 to son, and much of the Y-chromosome does not recombine with any other chromosome; this
31 leaves a long tract of unbroken DNA that serves as faithful record of evolution of the
32 chromosome throughout human history. Constant accumulation of variation within the
33 chromosome over human history created new combinations of variants (haplotypes), which trace
34 patrilineal descent.

35 Despite its use in ancestry, efficient software tools that are able to rapidly assign Y-
36 chromosome haplotypes from dense genotype data into clusters of related haplotypes
37 (haplogroups) at population scales are limited in scope. Here we describe SNAPPY (Single
38 Nucleotide Assignment of Phylogenetic Parameters on the Y-chromosome), a tool to assign Y-
39 chromosome haplogroups using dense genotyping data. Haplogroup assignment is based on the
40 well-established polymorphisms along the Y-chromosome with haplogroup information

41 maintained by the International Society of Genetic Genealogists (ISOGG, isogg.org) and others
42 [e.g., (Karafet, et al., 2015; Poznik, et al., 2016)], using only phylogenetically-informative alleles
43 to determine which haplogroup has the highest support, thus avoiding complications related to
44 the reversion of alleles to ancestral states. Importantly, the method is able to identify haplogroups
45 from both leaf and interior nodes of the Y-chromosome tree. Here we briefly outline the
46 implementation of the program, and present Y-chromosome haplotype assignments from The
47 1000 Genomes Project (Genomes Project, et al., 2015) data to validate the algorithm.

48 **Implementation**

49 *Dependencies*

50 SNAPPY is implemented in python requiring only numpy in addition to the standard python
51 library. In addition, SNAPPY makes use of plink (Chang, et al., 2015) for initial file format
52 conversions from either array-based formats (e.g., .ped/.bed) or from vcf.

53 *Algorithm Description*

54 SNAPPY leverages the Y-chromosome phylogeny and a database of haplotype-informative
55 SNPs derived from ISOGG and other sources to store nested haplotypes in memory-efficient
56 dictionaries. Genotypes from individual samples are stored as a list of dictionaries keyed by
57 chromosome position. Currently the tree is optimized for sites on the commonly-used Illumina
58 Multi-Ethnic Global Array (pagestudy.org/mega), however tree files can easily be generated
59 from relevant sources. To compute haplogroup scores for each of the 565 possible Y-
60 chromosome haplogroups present in our Y-chromosome haplogroup reference library, SNAPPY
61 looks up a sample's genotypes—stored in a dictionary—and counts the proportion of matching
62 alleles that are at haplogroup-informative positions.

63 SNAPPY determines a score for each haplotype using the number of haplogroup-informative
64 derived alleles adjusted by the sample's number of non-missing informative positions. Note that
65 a particular haplogroup's score uses alleles from both its own informative positions as well as its
66 ancestral nodes on the tree. This ensures that haplogroup assignments take into account the full
67 phylogenetic structure of the Y-chromosome, and enables the creation of easily traversable data
68 structures that eliminate redundant storage of informative positions due to the highly
69 parsimonious nature of variants on the Y (Poznik, et al., 2016). Haplogroup scores for every
70 haplogroup are stored in a two-dimensional numpy array to allow for efficient storage and quick
71 processing.

72 Finally, haplogroup assignments are made by evaluating the scores of nodes of each
73 branch, starting at the most distal node of the branch that has both a score above a user-defined
74 threshold and no descendant haplotypes with a score that exceeds the user-defined threshold, and
75 then working towards the root of the tree. This ensures that all nodes are considered potential
76 haplogroups for the sample, even if they are not terminal nodes on the tree. SNAPPY makes its
77 haplogroup assignment based on the highest scoring node or the deepest node with a score higher
78 than a user-defined threshold. This algorithm ensures that the deepest haplogroup with sufficient
79 support is assigned. In addition to reporting the single haplogroup with the most support for each
80 sample, SNAPPY also reports each haplogroup that has a score greater than a user-defined
81 minimum score so that haplogroup assignments can be adjusted or investigated where necessary.

82 **Validation and Testing**

83 *Data Sources*

84 We downloaded a list of Y-chromosome variants found on the Multi-Ethnic Genotyping Array
85 (MEGA) (Bien, et al., 2016) (a list of variants found on the MEGA can be found at
86 <https://pagestudy.org/index.php/multi-ethnic-genotyping-array>). Because the MEGA Y content
87 was designed to be ancestry informative, the SNPs included are well-distributed throughout the
88 Y-chromosome tree, which allows for accurate and precise assignment of haplogroups. MEGA
89 variant positions were converted from GRCh38.p7 to GRCh37 using NCBI's genome remapper
90 tool (<https://www.ncbi.nlm.nih.gov/genome/tools/remap>) in preparation for extracting MEGA
91 positions from Y-chromosome Phase 3 of The 1000 Genomes Project (Genomes Project, et al.,
92 2015) data.

93 For reference, we downloaded 1,233 Y-chromosome genotypes from Phase 3 of The
94 1000 Genomes Project from NCBI (Genomes Project, et al., 2015) ([ftp://ftp-](ftp://ftp-trace.ncbi.nih.gov/1000genomes/ftp)
95 [trace.ncbi.nih.gov/1000genomes/ftp](ftp://ftp-trace.ncbi.nih.gov/1000genomes/ftp)). The vcf was filtered to contain only the variants that are
96 present on MEGA using plink (Chang, et al., 2015). This yielded a final total of 2,366 variants.

97 Similarly, we provide Y-chromosome haplogroup tree information and informative SNPs,
98 available in the SNAPPY distribution found at github.com/chrisgene/snappy. This information is
99 largely consistent with phylogenetic information present on ISOGG.

100 *Results*

101 To assess the speed and accuracy of SNAPPY, we tested 1,233 males from Phase 3 of The 1000
102 Genomes Project (Genomes Project, et al., 2015). The program is computationally efficient: for a
103 set of 1,233 samples genotyped at 2,366 SNPs, it runs in an average of 4.96 seconds on a single

104 2.3 GHz core of an Intel Xeon processor, using 278 Mb of RAM. SNAPPY scales linearly
105 indicating good performance for larger data sets.

106 Y-chromosome haplotypes for 1000 Genomes males have been previously well-
107 characterized (Poznik, et al., 2016), and we used this characterization to assess the accuracy of
108 SNAPPY on genotype data compared to full sequences. We found that most individuals had top
109 haplogroup scores >95%, (Supplemental Figure 1), correctly predicting over >99% of major
110 haplogroup assignments for all individuals, with minor differences in fine-grained haplogroup
111 designations, even given topological differences between our tree and prior examples
112 (Supplemental Table 1). The three individuals with discordant major haplogroup assignments
113 were assigned by SNAPPY to the haplogroup P1 rather than the closely-related reference
114 assignment of Q1a. We note that the P1 and Q haplogroups had the same score (0.958) in these
115 individuals, but the P1 haplogroup was chosen because it resides deeper in the tree than Q. These
116 inconsistencies are expected to resolve with increased genotype density or tree topologies with
117 improved resolution and can be easily adjudicated on manual inspection. A fourth individual was
118 correctly assigned to the A0 haplogroup consistent with no derived haplogroup with support
119 above our recommended 60% minimum match. Some differences between the sets are to be
120 expected because the number of variants between the two datasets differ substantially (>60K for
121 the reference set vs. 2,399 for our tests), and different Y-chromosome phylogenies (*de novo*
122 reconstruction vs. ISOGG download) may result in subtle yet important differences between
123 branch points. Nevertheless, the ability to recover major haplogroups from individuals
124 representing nearly all major branches of the Y haplogroup tree indicates that SNAPPY, and
125 informative array-genotyped sites, are robust and unbiased in ancestry assignment.

126 **Conclusions**

127 The ability to determine Y-chromosome haplogroups is of broad interest in many lines of
128 research from population genetics to anthropology to medicine. Here we have introduced and
129 demonstrated SNAPPY, a method to rapidly and accurately assign Y-chromosome haplogroups
130 to population-scale data sets using dense genotyping. SNAPPY is user-friendly, requiring input
131 in common plink format and just a single command to run. Additionally, SNAPPY is flexible,
132 with several user-defined parameters, the ability to use user-curated y-haplogroup trees, and to
133 manually interpret haplogroup assignments.

134 **Funding Information**

135 This work was supported in part by National Institutes of Health grant number T32HG00044 to
136 CRG. The content is solely the responsibility of the authors and does not necessarily represent
137 the official views of the National Institutes of Health.

138 **Acknowledgments**

139 We gratefully acknowledge the community in ISOGG for maintaining and providing Y-
140 chromosome haplogroup trees, David Poznik for assistance with the original phylogeny, and
141 Peter Underhill for general discussions.

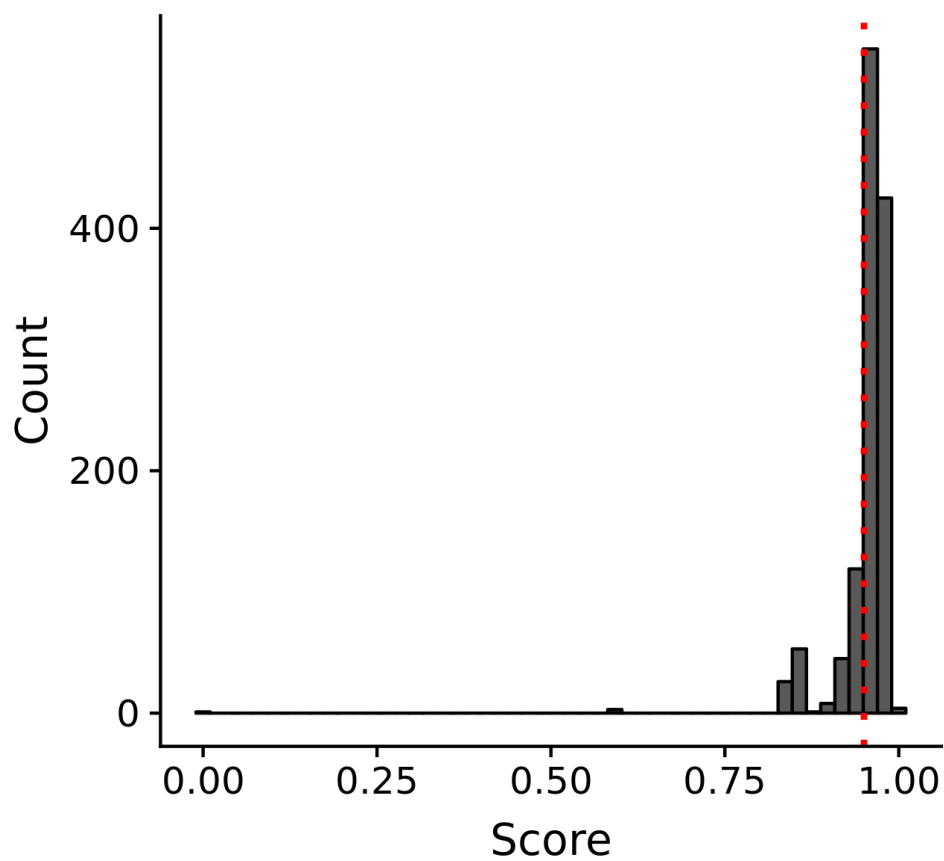
142 **References**

143 Bergstrom, A., *et al.* Deep Roots for Aboriginal Australian Y Chromosomes. *Curr Biol*
144 2016;26(6):809-813.
145 Bien, S.A., *et al.* Strategies for Enriching Variant Coverage in Candidate Disease Loci on a
146 Multiethnic Genotyping Array. *PLoS One* 2016;11(12):e0167758.
147 Chang, C.C., *et al.* Second-generation PLINK: rising to the challenge of larger and richer
148 datasets. *Gigascience* 2015;4:7.

- 149 Chiaroni, J., Underhill, P.A. and Cavalli-Sforza, L.L. Y chromosome diversity, human
150 expansion, drift, and cultural evolution. *Proc Natl Acad Sci U S A* 2009;106(48):20174-20179.
151 Genomes Project, C., *et al.* A global reference for human genetic variation. *Nature*
152 2015;526(7571):68-74.
153 Karafet, T.M., *et al.* Improved phylogenetic resolution and rapid diversification of Y-
154 chromosome haplogroup K-M526 in Southeast Asia. *Eur J Hum Genet* 2015;23(3):369-373.
155 Karmin, M., *et al.* A recent bottleneck of Y chromosome diversity coincides with a global
156 change in culture. *Genome Res* 2015;25(4):459-466.
157 Poznik, G.D., *et al.* Punctuated bursts in human male demography inferred from 1,244
158 worldwide Y-chromosome sequences. *Nat Genet* 2016;48(6):593-599.
159 Seielstad, M.T., Minch, E. and Cavalli-Sforza, L.L. Genetic evidence for a higher female
160 migration rate in humans. *Nat Genet* 1998;20(3):278-280.

161

162



Supplemental Figure 1 Distribution of SNAPPY-assigned haplogroup scores.

Histogram of haplogroup scores assigned by SNAPPY. A dashed red line indicates 95%.

