

# yHaplo™ | Software Manual

G. David Poznik  
23andMe

Last updated: April 2, 2019

## Table of Contents

<b>1</b>	<b>Introduction.....</b>	<b>2</b>
<b>2</b>	<b>Downloading and running yHaplo .....</b>	<b>2</b>
<b>3</b>	<b>Input files .....</b>	<b>4</b>
3.1	Primary tree structure .....	4
3.2	Phylogenetically informative SNPs.....	4
3.3	Quality control for phylogenetically informative SNPs.....	4
3.4	Representative Markers .....	4
3.5	Sample genotype data.....	5
<b>4</b>	<b>Genotype data.....</b>	<b>6</b>
4.1	Sample-major formats .....	6
4.1.1	Genos .....	6
4.1.2	Ped/Map .....	6
4.1.3	23andMe internal storage.....	6
4.2	Variant-major formats .....	7
4.2.1	VCF .....	7
4.2.2	23andMe text download.....	7
<b>5</b>	<b>Haplogroup output files.....</b>	<b>8</b>
5.1	Haplogroup calls.....	8
5.2	Counts of ancestral- and derived-allele genotypes.....	8
5.3	Haplogroup paths .....	9
5.4	Derived-allele genotypes observed .....	10
5.5	Derived-allele SNP details .....	10
5.6	Ancestral-allele genotypes observed .....	11
5.7	Ancestral-allele SNP details .....	12
<b>6</b>	<b>Additional output files .....</b>	<b>13</b>
6.1	Phylogenetically informative SNPs.....	13
6.1.1	Post-QC SNPs.....	13
6.1.2	Unique post-QC SNPs .....	13
6.1.3	Markers dropped from consideration.....	13
6.2	Trees and traversals .....	13
6.2.1	Trees.....	13
6.2.2	Traversals .....	14
<b>7</b>	<b>References.....</b>	<b>14</b>

## 1 Introduction

yHaplo identifies the Y-chromosome haplogroup of each male in a sample of one to millions. It does not rely on any particular genotyping modality or platform, and it is robust to missing data, genotype errors, mutation recurrence, and other complications. Although full sequences yield the most granular haplogroup classifications, genotyping arrays can yield reliable calls, provided a reasonable number of phylogenetically informative variants has been assayed.

Briefly, haplogroup calling involves two steps. The program first builds a representation of the Y-chromosome phylogeny by reading its primary structure from (Newick-formatted) text and then importing phylogenetically informative SNPs from the ISOGG database<sup>1</sup>, storing each SNP within a specific node and growing the tree as necessary. It then traverses the tree for each individual, identifying for each the path of derived alleles leading to a haplogroup designation.

To learn more about the algorithm, please see our bioRxiv pre-print<sup>2</sup>:

Poznik GD. 2016. Identifying Y-chromosome haplogroups in arbitrarily large samples of sequenced or genotyped men. bioRxiv doi: 10.1101/088716

<http://biorxiv.org/content/early/2016/11/19/088716>

## 2 Downloading and running yHaplo

yHaplo is available at the 23andMe code repository for non-commercial use pursuant to the terms of the non-exclusive license agreement included with the software distribution:

<https://github.com/23andMe/yhaplo>

To clone the repository, issue the following command from the directory in which yHaplo is to be stored (henceforth `pathTo/`):

```
git clone git@github.com:23andMe/yhaplo.git
```

Alternatively, users can download yHaplo from the repository web site linked above. To do so, click the “Clone or download” link, and select “Download ZIP.” Unzip the file, and move/rename the directory to “`pathTo/yhaplo`.”

Add the yHaplo directory to your path by, for example, adding the following line your `.bashrc` or `.bash_profile` file:

```
PATH=pathTo/yhaplo:$PATH
```

To test run the software, issue the following command from any directory, again substituting the full path for “`pathTo`”:

```
callHaplogroups.py -i pathTo/yhaplo/data/1000Y.subset.genos.txt
```

The program will first read and process a standard set of input files (**section 3**). Next, it will detect from the genotype file's name extension (`.genos.txt`) that the data to be analyzed are in a simple sample-major text format (**section 4.1.1**). It will read these data and identify the Y-chromosome haplogroup for each individual in the sample.

The software can operate on other common genotype formats (**section 0**) and produce several helpful auxiliary output files (**section 5**). To see the full list of command-line options, some of which are described in this manual, issue the following command:

```
callHaplogroups.py -h
```

The most useful auxiliary output files are elicited with the following options:

```
callHaplogroups.py -i projectName.genos.txt -c -hp -ds -dsd -as -asd
```

To test on a larger dataset, you could download the Y-chromosome genotypes for all 1,244 males from phase 3 of the 1000 Genomes Project<sup>3</sup>:

```
ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/supporting/chrY/  
ALL.chrY_10Mbp_mask.glia_freebayes_maxLikGT_siteQC.20130502.60555_biallelic_snps.vcf.gz
```

For clarity, rename this file to `1000Y.all.vcf.gz`, then run yHaplo as above:

```
callHaplogroups.py -i 1000Y.all.vcf.gz
```

## 3 Input files

### 3.1 Primary tree structure

This Newick-formatted text file codifies the primary structure of the Y-chromosome tree.

```
input/y.tree.primary.DATE.nwk
```

### 3.2 Phylogenetically informative SNPs

This table includes details for ~20,000 phylogenetically informative SNPs. The contents were curated by the International Society of Genetic Genealogy (ISOGG) and scraped from their website<sup>1</sup> on the date specified in the filename.

```
input/isogg.DATE.txt
```

### 3.3 Quality control for phylogenetically informative SNPs

Several files are necessary to correct a number of glitches in the file described above. We encourage users to augment these files in light of any further inconsistencies they may find.

Two files correct physical coordinates and ancestral- and derived-allele designations, respectively. These files have the same format as the post-QC SNPs file described in **section 5.1.1**:

```
input/isogg.correct.coordinate.txt  
input/isogg.correct.polarize.txt
```

Three files, also formatted as in **section 5.1.1**, indicate SNPs to drop in light of inconsistencies observed in test data:

```
input/isogg.omit.bad.txt  
input/isogg.omit.bad.23andMe.txt  
input/isogg.omit.branch.conflict.txt
```

This file lists physical coordinates of multiallelic sites to be excluded:

```
input/isogg.multiallelic.txt
```

This file lists SNPs with known, shared, recurrences not reflected in the ISOGG table. Though the software does not use this file, we provide it for completeness.

```
input/isogg.split.txt
```

### 3.4 Representative Markers

These two files list the names of markers deemed representative of corresponding haplogroups:

```
input/representative.SNPs.isogg.YEArTree.txt  
input/representative.SNPs.additional.txt
```

The first was generated by parsing the ISOGG tree of the given year, and the second is a manually edited supplementary list. yHaplo ignores the first column of each file.

### 3.5 Sample genotype data

This file bears the Y-chromosome genotypes of 35 globally diverse males at ~60,000 SNPs:

```
data/1000Y.subset.genos.txt
```

These data are derived from full sequences generated by the 1000 Genomes Project<sup>3</sup>, and the file format is described in **section 0**.

## 4 Genotype data

yHaplo can read genotype data in both variant-major and sample-major orientations. Since haplogroups are properties of individuals, it is more memory efficient to use a sample-major format, but the program will run on variant-major data as long as the full data matrix can fit in memory. This should be the case for sample sizes on the order of ten thousand or fewer individuals, but to call haplogroups in samples on the order of hundreds of thousands or millions, it is probably best to transpose or split the data before running.

### 4.1 Sample-major formats

#### 4.1.1 Genos

To run on sample-major text data, use: `-i projectName.genos.txt`

The expected format is:

Row 0	Physical coordinates (GRCh37)
Column 0	Sample IDs. The first element is ignored.
Cell ( $i, j$ )	Genotype for the individual indicated by the $i^{\text{th}}$ row of column 0 at the position indicated by the $j^{\text{th}}$ , column of row 0. Genotypes are encoded as single characters from the set { A, C, G, T, . }, with "." representing an unobserved value.

#### 4.1.2 Ped/Map

In this common sample-major format, the `.ped` file has six identifier columns and  $2m$  allele columns, one pair for each of  $m$  markers. The `.map` file has  $m$  rows, one for each marker, ordered as in the `.ped` file. For data in plink<sup>4</sup> format (`.bed`, `.bim`, `.fam`), one can export to ped/map format with the following command:

```
plink --bfile projectName --recode --out projectName
```

To then convert the ped/map data to genos format, use the utility script included with yHaplo:

```
convert2genos.py projectName.ped
```

#### 4.1.3 23andMe internal storage

yHaplo can load data directly from 23andMe's internal data storage. The corresponding options will only work within the 23andMe research environment.

To run on all male customers, use: `-a`

To run on a subset of research IDs, use: `-i projectName.resid.txt`

## 4.2 Variant-major formats

### 4.2.1 VCF

To run on VCF<sup>5</sup>-formatted data, use: `-i projectName.vcf`  
Or, if the VCF is gzipped: `-i projectName.vcf.gz`

**Note:** When calling genotypes from sequencing data, please be sure to **emit all confident sites**, not just those at which alternative alleles were observed. Ref/alt status is unimportant to yHaplo, but ancestral/derived status is, and the reference sequence contains many derived alleles. yHaplo will not be happy if you discard these valuable data by generating “variants-only” VCFs. To limit compute time and file size, you could safely restrict to positions in:

```
output/isogg.snps.unique.DATE.txt,
```

as these are the only SNPs yHaplo considers. To generate this file, just run:

```
callHaplogroups.py
```

with no arguments.

### 4.2.2 23andMe text download

The utility script referred to in **section 4.1.2** can also convert single-sample text files downloaded from 23andMe. Just rename the file to `sampleID.23andMe.txt` and run:

```
convert2genos.py sampleID.23andMe.txt
```

This file can have header lines whose first characters are “#” or “rsid,” and its data columns are assumed to be:

1. SNP identifier. The script ignores this column.
2. Chromosome. The script retains only those rows with values of “24” or “Y.”
3. Physical coordinate. Reference human assembly build 37 (GRCh37) is assumed.
4. Genotype. The script drops rows with values not in { A, C, G, T, D, I }.

A fifth column bearing a “second allele” is allowed. The script ignores any row for which the value of the fifth column differs from that of the fourth.

## 5 Haplogroup output files

yHaplo generates a log file and a haplogroup calls file. Users also have the option to generate any of several helpful auxiliary output files. The `-o` option specifies the output directory, which is `output/` by default.

### 5.1 Haplogroup calls

**File name:** `output/haplogroups.projectName.txt`

#### Columns

1. ID
2. Haplogroup short form, with the name of a SNP observed in the derived state
3. Haplogroup short form, with the name of a representative SNP
4. Haplogroup long form, using Y-Chromosome Consortium nomenclature<sup>6</sup>
5. Prior haplogroup call (if imported)
6. Discordance indicator (if prior call imported): “.” = concordant, “\*” = discordant

#### Example line

```
HG01791 R-M167 R-M167 R1b1a2a1a2a1b1a1 R1b1a2a1a2a1b1a1 .
```

#### Important note on the distinction between columns 2 and 3

The SNP names of columns 2 and 3 may differ. yHaplo picks for each branch of the tree a single SNP to represent the equivalence class of SNPs associated with the branch. Column 3 uses this representative SNP, but the individual may not have been genotyped for it. Furthermore, it is actually possible that the individual has been genotyped for the SNP, and an ancestral allele was observed. This can happen when the ISOGG representation of the tree is incomplete, and the set of SNPs associated with a branch is not a true equivalence class. For example, if there are 10 SNPs associated with a given branch, and we observe 5 ancestral alleles and 5 derived alleles, the true haplogroup is one that shares some portion of the given branch before diverging from it, and the representative SNP may have arisen after this divergence. In contrast, the SNP cited in column 2 will necessarily have been observed in the derived state for this individual, although the SNP itself may be less commonly known.

### 5.2 Counts of ancestral- and derived-allele genotypes

This optional output file records for each individual counts of ancestral- and derived-allele genotypes encountered along each branch visited in the search path. The file contains one block for each individual, with the last line indicating the final haplogroup call.

**Command-line option:** `-c, --ancDerCounts`

**File name:** `output/counts.ancDer.projectName.txt`

#### Columns



### *Body of each block*

1. ID
2. Branch label (haplogroup)
3. Number of ancestral-allele genotypes observed
4. Number of derived-allele genotypes observed

### *Last line of each block*

1. ID
2. YCC haplogroup
3. Prior haplogroup call (if imported)
4. Pipe separator (|)
5. Haplogroup short form, with the name of a SNP observed in the derived state
6. Haplogroup short form, with the name of a representative SNP

### **Example block**

```
NA18960  A0          30    0
NA18960  A1          0    12
NA18960  A1a        15    0
NA18960  A1b         0    30
NA18960  BT         0   333
NA18960  B           1    0
NA18960  CT         0   191
NA18960  DE         0    26
NA18960  E          95    0
NA18960  D           0    36
NA18960  D1b         0     5
NA18960  D1b2       29    0
NA18960  D1b1a       1     0
NA18960  D1b1b       1     0
NA18960  D1b1d       0     1
NA18960  D1b1a2      2     0
NA18960  D1b1d1     29    0
NA18960  D1b1d D1b1d | D-CTS6609 D-CTS6609
```

## **5.3 Haplogroup paths**

This optional output file provides a compact summary of the path through the tree leading to the each individual's haplogroup call. It lists each branch on which derived-allele genotypes were observed and counts thereof. This information is a subset of that described in **section 5.2**.

**Command-line option:** `-hp, --haplogroupPaths`

**File name:** `output/paths.projectName.txt`

### **Columns**

1. ID
2. YCC haplogroup

3. Haplogroup short form, with the name of a representative SNP
4. Pipe separator (|)
- 5+. Branch labels (haplogroups) and numbers of derived-allele genotypes, colon-separated

### Example line

```
NA18960 D1b1d D-CTS6609 | A1:12 A1b:30 BT:333 CT:191 DE:26 D:36 D1b:5 D1b1d:1
```

## 5.4 Derived-allele genotypes observed

This optional output file includes for each individual a list of SNPs observed in the derived state on the path leading to the individual's lineage. This information is an expanded version of that described in **section 5.3**.

**Command-line option:** `-ds, --derSNPs`

**File name:** `output/derived.snps.projectName.txt`

### Columns

1. ID
2. YCC haplogroup
3. Haplogroup short form, with the name of a representative SNP
4. Pipe separator (|)
- 5+. SNPs, each with haplogroup and common name, separated by a colon

### Example line (with many SNPs omitted for clarity)

```
NA18960 D1b1d D-CTS6609 | A1:P305 ... A1:L1112 A1b:P108 ... A1b:Z17900 BT:M42 ...
    BT:Z17390 CT:M168 ... CT:Y1528 DE:M145 ... DE:PF1833 D:M174 ... D:IMS-JST021355
    D1b:M55 D1b:M64.1 D1b:M179 D1b:M359.1 D1b:P37.1 D1b1d:CTS6609
```

## 5.5 Derived-allele SNP details

This optional output file includes detailed information about each SNP whose derived allele was observed on the path through the tree. Each block is essentially the transpose one row of the file described in **section 5.4**, with additional information for each SNP.

**Command-line option:** `-dsd, --derSNPsDetail`

**File name:** `output/derived.snps.detail.projectName.txt`

### Columns

*First line of each block*

As in the haplogroup calls file (**section 5.1**).

*Body of each block*

1. ID
2. SNP name
3. Branch label (haplogroup)
4. Physical coordinate (GRCh37)
5. Mutation type: ancestral allele and derived allele, separated by a two-character arrow (->)

**Example block** (with many SNPs omitted for clarity)

```

NA18960 D-CTS6609 D-CTS6609 D1b1d D1b1d .
NA18960 P305 A1 2710154 A->G
...
NA18960 L1112 A1 8466995 G->A
NA18960 P108 A1b 15426248 C->T
...
NA18960 Z17900 A1b 23572106 G->A
NA18960 M42 BT 21866840 A->T
...
NA18960 Z17390 BT 28746408 G->C
NA18960 M168 CT 14813991 C->T
...
NA18960 Y1528 CT 15932327 A->G
NA18960 M145 DE 21717208 C->T
...
NA18960 PF1833 DE 22014732 C->G
NA18960 M174 D 14954280 T->C
...
NA18960 IMS-JST021355 D 2828425 A->G
NA18960 M55 D1b 21872738 T->C
NA18960 M64.1 D1b 21903383 A->G
NA18960 M179 D1b 14838700 C->T
NA18960 M359.1 D1b 14491671 T->C
NA18960 P37.1 D1b 14491684 T->C
NA18960 CTS6609 D1b1d 17008697 A->T

```

## 5.6 Ancestral-allele genotypes observed

This optional output file includes for each individual a list of SNPs encountered in the ancestral state during the search. This information is an expanded version of that described in **section 5.5**, and it complements the file described in **section 5.4**.

**Command-line option:** `-as, --ancSNPs`

**File name:** `output/ancestral.snps.projectName.txt`

### Columns

As in the derived-allele genotypes file (**section 5.4**).

**Example line** (with many SNPs omitted for clarity)

```

NA18960 D1b1d D-CTS6609 | A0:V148 ... A0:L1055 A1a:M31 ... A1a:V215 B:M181 E:M40
... E:CTS3337 D1b2:CTS131 ... D1b2:Z17171 D1b1a:M125 D1b1b:M151 D1b1a2:CTS107
D1b1a2:IMS-JST022457 D1b1d1:CTS504 ... D1b1d1:Z14880

```

## 5.7 Ancestral-allele SNP details

This optional output file includes detailed information about each SNP whose ancestral allele was encountered in the search. Each block is essentially the transpose one row of the file described in **section 5.6**, with additional information for each SNP. The file complements the one described in **section 5.5**.

**Command-line option:** `-asd, --ancSNPsDetail`

**File name:** `output/ancestral.snps.detail.projectName.txt`

### Columns

As in the derived-allele SNP details file (**section 5.5**).

### Example block (with many SNPs omitted for clarity)

```
NA18960 D-CTS6609 D-CTS6609 D1b1d D1b1d .
NA18960 V148 A0 6788191 G->A
...
NA18960 L1055 A0 21622096 G->A
NA18960 M31 A1a 21739754 G->C
...
NA18960 V215 A1a 2868499 T->G
NA18960 M181 B 14851554 T->C
NA18960 M40 E 2663943 C->T
...
NA18960 CTS3337 E 14845090 C->T
NA18960 CTS131 D1b2 2750950 G->A
...
NA18960 Z17171 D1b2 9813700 A->G
NA18960 M125 D1b1a 21930287 T->C
NA18960 M151 D1b1b 21892634 G->A
NA18960 CTS107 D1b1a2 2731887 C->T
NA18960 IMS-JST022457 D1b1a2 24464597 G->C
NA18960 CTS504 D1b1d1 6804892 T->C
...
NA18960 Z14880 D1b1d1 22589060 G->T
```

## 6 Additional output files

### 6.1 Phylogenetically informative SNPs

Upon reading and cleaning the raw set of phylogenetically informative SNPs (**section 3**), the program emits three files.

#### 6.1.1 Post-QC SNPs

This file lists all SNPs considered for haplogroup classification. An individual SNP can have multiple entries—one for each common name associated with it.

**File name:** `output/isogg.snps.cleaned.DATE.txt`

#### Columns

1. SNP name
2. Branch label (haplogroup)
3. Physical coordinate (GRCh37)
4. Mutation type: ancestral allele and derived allele, separated by a two-character arrow (->)

#### 6.1.2 Unique post-QC SNPs

This file is similar to the post-QC SNPs file (**section 6.1.1**), but with one entry per SNP.

**File name:** `output/isogg.snps.unique.DATE.txt`

**Format:** As in **section 6.1.1**, with a fifth column:  
a comma-separated list of all common names for the SNP

#### 6.1.3 Markers dropped from consideration

This file lists all markers dropped from consideration due to inconsistencies observed in test data.

**File name:** `output/isogg.snps.dropped.DATE.txt`

**Format:** As in **section 6.1.1**.

## 6.2 Trees and traversals

### 6.2.1 Trees

Upon building the tree, the program writes four Newick-formatted files:

```
output/y.tree.ycc.DATE.nwk
output/y.tree.ycc.hgSNP.DATE.nwk
output/y.tree.aligned.ycc.DATE.nwk
output/y.tree.aligned.hgSNP.DATE.nwk
```

Each details the topology of the tree, with the latter two, labeled “aligned,” including artificial branch lengths to align the leaf nodes when plotting. Branches are labeled with YCC-style haplogroup names in the files labeled “ycc” and with representative-SNP-style haplogroup names in the files labeled “hgSNP.” To visualize the final tree in plain text, users can plot it with the provided utility script, which wraps a Biopython<sup>7</sup> function:

```
plotTree.py -n output/y.tree.aligned.hgSNP.DATE.nwk
```

Alternatively, they can plot the tree graphically by importing to third-party tree plotting software.

### 6.2.2 Traversals

For an even simpler text-based visualization of the tree, use one of the traversal options:

```
-b, --breadthFirst    write bread-first traversal
-d, --depthFirst      write depth-first (pre-order) traversal
```

These yield:

```
y.tree.bf.traversal.DATE.txt
y.tree.df.traversal.DATE.txt
```

## 7 References

1. ISOGG. International Society of Genetic Genealogy. (2016). at <<http://www.isogg.org/>>
2. Poznik, G. D. Identifying Y-chromosome haplogroups in arbitrarily large samples of sequenced or genotyped men. *bioRxiv* (2016). doi:10.1101/088716
3. Poznik, G. D. *et al.* Punctuated bursts in human male demography inferred from 1,244 worldwide Y-chromosome sequences. *Nat. Genet.* **48**, 593–9 (2016).
4. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* **81**, 559–75 (2007).
5. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158 (2011).
6. The Y Chromosome Consortium. A nomenclature system for the tree of human Y-Chromosomal binary haplogroups. *Genome Res.* **12**, 339–348 (2002).
7. Cock, P. J. A. *et al.* Biopython: freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics* **25**, 1422–1423 (2009).