

**A peer-reviewed version of this preprint was published in PeerJ on 25 June 2018.**

[View the peer-reviewed version](https://peerj.com/articles/5149) (peerj.com/articles/5149), which is the preferred citable publication unless you specifically need to cite this preprint.

Smieszek S, Mitchell SL, Farber-Eger EH, Veatch OJ, Wheeler NR, Goodloe RJ, Wells QS, Murdock DG, Crawford DC. 2018. Hi-MC: a novel method for high-throughput mitochondrial haplogroup classification. PeerJ 6:e5149 <https://doi.org/10.7717/peerj.5149>

# Hi-MC: A novel method for high-throughput mitochondrial haplogroup classification

Sabrina L Mitchell <sup>1</sup>, Eric H Farber-Eger <sup>2</sup>, Olivia J Veatch <sup>3</sup>, Robert J Goodloe <sup>1</sup>, Quinn S Wells <sup>4,5</sup>, Deborah G Murdock <sup>6</sup>, Dana C Crawford <sup>Corresp. 7, 8</sup>

<sup>1</sup> Center for Human Genetics Research, Vanderbilt University Medical Center, Nashville, TN, United States

<sup>2</sup> Vanderbilt Institute for Clinical and Translational Research, Vanderbilt University Medical Center, Nashville, TN, United States

<sup>3</sup> Department of Neurology, Vanderbilt University Medical Center, Nashville, TN, United States

<sup>4</sup> Department of Medicine, Vanderbilt University Medical Center, Nashville, TN, United States

<sup>5</sup> Department of Pharmacology, Vanderbilt University, Nashville, TN, United States

<sup>6</sup> Center for Mitochondrial and Epigenomic Medicine, Children's Hospital of Philadelphia, Philadelphia, PA, United States

<sup>7</sup> Population and Quantitative Health Sciences, Case Western Reserve University, Cleveland, OH, United States

<sup>8</sup> Institute for Computational Biology, Case Western Reserve University, Cleveland, OH, United States

Corresponding Author: Dana C Crawford

Email address: [dana.crawford@case.edu](mailto:dana.crawford@case.edu)

Effective approaches for assessing mitochondrial DNA (mtDNA) variation are important to multiple scientific disciplines. Mitochondrial haplogroups characterize branch points in the phylogeny of mtDNA. Several tools exist for mitochondrial haplogroup classification. However, most require full or partial mtDNA sequence which is often cost prohibitive for studies with large sample sizes. The purpose of this study was to develop Hi-MC, a high-throughput method for mitochondrial haplogroup classification that is cost effective and applicable to large sample sizes making mitochondrial analysis more accessible in genetic studies. Using rigorous selection criteria, we defined and validated a custom panel of mtDNA single nucleotide polymorphisms (SNPs) that allows for accurate classification of European, African, and Native American mitochondrial haplogroups at broad resolution with minimal genotyping and cost. We demonstrate that Hi-MC performs well in samples of European, African, and Native American ancestries, and that Hi-MC performs comparably to a commonly used classifier. Implementation as a software package in R enables users to download and run the program locally, grants greater flexibility in the number of samples that can be run, and allows for easy expansion in future revisions. The source code is freely available at <https://github.com/vserch/himc>.

1 **Hi-MC: A novel method for high-throughput mitochondrial haplogroup classification**

2 **Running title: Hi-MC for mitochondrial haplogroup classification**

3 Sabrina L. Mitchell<sup>1</sup>, Eric H. Farber-Eger<sup>2</sup>, Olivia J. Veatch<sup>3</sup>, Robert J. Goodloe<sup>1</sup>, Quinn S.

4 Wells<sup>4,5</sup>, Deborah G. Murdock<sup>6</sup>, Dana C. Crawford<sup>7</sup>

5 <sup>1</sup>Center for Human Genetics Research, Vanderbilt University Medical Center, Nashville, TN

6 37232, USA

7 <sup>2</sup>Vanderbilt Institute for Clinical and Translational Research, Vanderbilt University, Nashville, TN

8 37232, USA

9 <sup>3</sup>Department of Neurology, Vanderbilt University Medical Center, Nashville, TN 37232, USA

10 <sup>4</sup>Department of Medicine, Vanderbilt University Medical Center, Nashville, TN 37232, USA

11 <sup>5</sup>Department of Pharmacology, Vanderbilt University Medical Center, Nashville, TN 37232, USA

12 <sup>6</sup>Center for Mitochondrial and Epigenomic Medicine, Children's Hospital of Philadelphia,

13 Philadelphia, PA 19104, USA

14 <sup>7</sup>Department of Population and Quantitative Health Sciences, Institute for Computational

15 Biology, Case Western Reserve University, Cleveland OH 44106, USA

16 **Corresponding Author**

17 Dana C. Crawford, PhD

18 2103 Cornell Road, Wolstein Research Building, Suite 2-527

19 Case Western Reserve University

20 Cleveland, OH 44106

21 Telephone: (216) 368-5546

22 Email: [dana.crawford@case.edu](mailto:dana.crawford@case.edu)

23 **Key words:** mitochondrial haplogroups, classification, genetic variation, R

24 **Abstract**

25 Effective approaches for assessing mitochondrial DNA (mtDNA) variation are important to  
26 multiple scientific disciplines. Mitochondrial haplogroups characterize branch points in the  
27 phylogeny of mtDNA. Several tools exist for mitochondrial haplogroup classification. However,  
28 most require full or partial mtDNA sequence which is often cost prohibitive for studies with large  
29 sample sizes. The purpose of this study was to develop Hi-MC, a high-throughput method for  
30 mitochondrial haplogroup classification that is cost effective and applicable to large sample sizes  
31 making mitochondrial analysis more accessible in genetic studies. Using rigorous selection  
32 criteria, we defined and validated a custom panel of mtDNA single nucleotide polymorphisms  
33 (SNPs) that allows for accurate classification of European, African, and Native American  
34 mitochondrial haplogroups at broad resolution with minimal genotyping and cost. We  
35 demonstrate that Hi-MC performs well in samples of European, African, and Native American  
36 ancestries, and that Hi-MC performs comparably to a commonly used classifier. Implementation  
37 as a software package in R enables users to download and run the program locally, grants greater  
38 flexibility in the number of samples that can be run, and allows for easy expansion in future  
39 revisions. The source code is freely available at <https://github.com/vserch/himc>.

## 40 Introduction

41 Human mitochondrial DNA (mtDNA) consists of a double-stranded, circular chromosome that  
42 spans 16,529 base pairs and encodes 22 transfer RNAs, 2 ribosomal RNAs, and 13 proteins that  
43 are part of the oxidative phosphorylation enzyme complexes. Compared with nuclear DNA,  
44 unique characteristics of mtDNA include uniparental (i.e. matrilineal) inheritance, lack of  
45 recombination, high copy number, and a high mutation rate. These characteristics make mtDNA a  
46 powerful tool for investigations in multiple disciplines including population and medical  
47 genetics, molecular anthropology, and forensics<sup>1</sup>. Strong evidence exists supporting the  
48 involvement of mtDNA variation in human disease phenotypes, underscoring the importance of  
49 integrating the mitochondrial genome in genetic association studies. Evidence includes the  
50 association of mtDNA single nucleotide polymorphisms (SNPs) and mitochondrial haplogroups  
51 with a number of phenotypes encompassing cancer, neurologic, ocular, cardiovascular, and  
52 metabolic traits<sup>2-7</sup>.

53 Mitochondrial haplogroups are collections of similar combinations of mtDNA SNPs  
54 inherited from a common ancestor. These haplogroups are formed via the sequential  
55 accumulation of mutations through the maternal lineage. As a result of population migration,  
56 distinct mitochondrial haplogroups are associated with different continental ancestries including  
57 African, European, Native American, Asian, and Oceanic<sup>4,8,9</sup>, allowing for accurate classification  
58 of maternal genetic ancestry in large datasets using a small subset of mitochondrial markers.

59 Currently, several methods are available for mitochondrial haplogroup classification  
60 including Haplogrep, HaploFind, MitoTool, HmtDB, MToolBox, and Phy-mer<sup>10-17</sup>. While these  
61 methods are powerful tools for mtDNA sequence analysis, including classification of  
62 mitochondrial haplogroups, most require full or partial mtDNA sequence, and some are limited in  
63 the number of samples that can be processed at once. To address limitations of existing methods  
64 we developed a high-throughput method for automated mitochondrial haplogroup classification

65 that can accommodate large sample sizes with SNP data recorded in the widely used pedigree  
66 (PED/MAP) file format.

67 Using a custom panel of mitochondrial SNPs we constructed a reduced mitochondrial  
68 phylogenetic tree, and developed an algorithm (Hi-MC) for broad classification of European,  
69 African, and Native American mitochondrial haplogroups. After employing Hi-MC, we  
70 determined mitochondrial haplogroup classifications of samples from the International HapMap  
71 Project<sup>18-20</sup>. To evaluate the performance of the algorithm we compared Hi-MC mitochondrial  
72 haplogroup classifications with those previously reported by HapMap and with classifications  
73 generated via Haplogrep, the most widely used web-based application for mitochondrial  
74 haplogroup classification. As expected, given the mitochondrial SNPs included in the custom  
75 panel, Hi-MC performs well on samples of European, African, and Native American ancestry, but  
76 does not perform as well resolving mitochondrial haplogroup in samples of Asian ancestry.  
77 Although Hi-MC does not yet resolve mitochondrial haplogroups for all populations, it provides  
78 a user-friendly method for high-throughput classification and is provided in an R software  
79 package that can be easily expanded in future revisions to capture additional mitochondrial  
80 haplogroups.

## 81 **Materials and methods**

### 82 *Algorithm*

83 The algorithm input is a list of mitochondrial SNP genotypes for each individual DNA sample,  
84 and the output is haplogroup classification. The Cambridge reference sequence (rCRS) is used to  
85 specify SNP positions. PhyloTree, a comprehensive phylogenetic tree of human mtDNA variation  
86 displaying relationships between mitochondrial haplogroups<sup>21</sup>, was used as a reference to create a  
87 reduced tree of 46 common haplogroups as presented in Mitchell et al<sup>22</sup>. This reduced  
88 classification tree was converted into a node-based tree structure. Each haplogroup node has a list  
89 of associated SNPs, a parent node, and zero or more child nodes. The SNPs associated with a

90 node define which SNP genotypes a subject must possess to belong to the corresponding  
91 haplogroup. Classification into a haplogroup also requires a subject to recursively meet the  
92 definition for the parent haplogroup. Haplogroups that require the reversion to the ancestral  
93 genotype (e.g.10398A to 10398G) are accommodated by adding a second hierarchy of required  
94 SNP genotypes.

95         The algorithm determines the appropriate haplogroup in a two-step process (Figure 1). In  
96 the first step, the algorithm passes mitochondrial SNP genotype data for each subject into the root  
97 node of the tree. The algorithm checks the list of SNP genotypes against those required by the  
98 root node. If the array meets the criteria for the parent node, this haplogroup is added to an  
99 accumulator. The algorithm then passes the list of SNP genotypes to each of the child nodes  
100 connected to that parent node until the tree is exhausted. Next, the algorithm ranks the list of  
101 haplogroups in the accumulator according to their distance from the root node. Any haplogroup  
102 with a path length less than that of the haplogroup with the longest path length is dropped. The  
103 remaining haplogroups, along with their path from the root node to the end node, are returned as  
104 a result.

### 105 ***Implementation***

106 The algorithm is implemented as a package in R<sup>23</sup> [<https://github.com/vserch/himc>]. Data input is  
107 standard PED/MAP formatted files. The output is an R dataframe object that includes subject IDs  
108 with a corresponding haplogroup classification and the path through the tree from root node to  
109 final classification. The output can easily be exported directly to a CSV file or text file. For  
110 further details on use of the Hi-MC package in R visit [www.icompbio.net](http://www.icompbio.net).

### 111 ***Mitochondrial SNP Selection***

112 The SNPs were selected for broad classification of European, African, and Native American  
113 mitochondrial haplogroup lineages as previously described<sup>22</sup>. Briefly, SNPs were chosen using  
114 Phylotree<sup>21</sup> and an extensive literature search for prior studies related to mitochondrial

115 haplogroup classification<sup>24-27</sup>. Preference was given to those SNPs that appear only once in  
116 Phylotree since such SNPs are specific to a single haplogroup. Sixty-three SNPs were selected,  
117 the majority of which are located in the coding region of the mitochondrial genome. Three  
118 Sequenom genotyping assay pools including all of these SNPs were designed using the  
119 MassARRAY software<sup>22</sup>. As described in Mitchell et al<sup>22</sup>, the custom SNP panel was genotyped  
120 in the National Health and Nutrition Examination Surveys (NHANES) accessed by the  
121 Epidemiologic Architecture for Genes Linked to Environment (EAGLE)<sup>28</sup>, a study site of the  
122 Population Architecture using Genomics and Epidemiology (PAGE) I study<sup>29</sup>. The Vanderbilt  
123 University Institutional Review Board determined that EAGLE was “non-human” subjects  
124 research.

### 125 *Application of Hi-MC*

126 To evaluate the performance of Hi-MC for mitochondrial haplogroup classification we genotyped  
127 the custom SNP panel in, and applied the algorithm to, HapMap Phase I and Phase III samples.  
128 We selected HapMap samples for the present study as HapMap samples were the preferred  
129 reference samples for individual study sites including this study as part of the larger PAGE I  
130 study<sup>29</sup>. The populations from HapMap Phase I included: individuals of Northern and Western  
131 European ancestry from the Centre d'Etude du Polymorphisme Humain samples collected in  
132 Utah, USA (CEU, n=90), Yoruba from Ibadan, Nigeria (YRI, n=90), Japanese in Tokyo, Japan  
133 (JPT, n=45), and Han Chinese in Beijing, China (CHB, n=45). The HapMap Phase III samples  
134 used in this study included only those of Mexican ancestry from Los Angeles, California (MXL,  
135 n=90). The International HapMap Consortium reported mitochondrial haplogroup classifications  
136 for the CEU, YRI, CHB, and JPT Phase I HapMap samples<sup>20</sup>; however, mitochondrial haplogroup  
137 classifications for the Phase III MXL samples have not been previously reported.

138 We genotyped the custom SNP panel in the CEU, YRI, and CHB/JPT Phase I HapMap  
139 samples and in the MXL samples from Phase III. Briefly, aliquots of DNA from HapMap CEU,



140 YRI, CHB/JPT, and MXL samples were obtained from the Coriell repository. SNPs were  
141 genotyped via the Agena Biosciences (formerly Sequenom) iPLEX® Gold MassArray platform.  
142 Multiplex primer extension was performed, and extension products were analyzed by MALDI-  
143 TOF mass spectrometry<sup>30</sup>.

144 SNP genotyping efficiency was set to greater than or equal to 0.90. The hypervariable  
145 region SNP mt16189 did not meet this threshold and was dropped from the analysis. Additionally,  
146 SNP mt9540 was excluded from the analysis due to poor genotyping efficiency. We determined  
147 that the primers for SNP mt9540 lacked specificity, consistent with the amplification of nuclear  
148 insertions of mitochondrial origin (NumtS) common in the human genome<sup>31</sup>. Therefore, SNP  
149 mt9540 is not included in the algorithm for classification. The final list of custom panel SNPs  
150 used to classify mitochondrial haplogroups is given in Supplementary Table 1.

151 Using genotype data from the custom SNP panel we employed Hi-MC and Haplogrep to  
152 determine mitochondrial haplogroup classifications in the HapMap samples. Although there are  
153 several tools available from which to compare Hi-MC, we selected Haplogrep for comparison  
154 given it is the most widely used tool to date with >180 citations in the peer-reviewed literature.  
155 We then compared the Hi-MC mitochondrial haplogroup classifications to the HapMap-reported  
156 classifications for Phase I samples<sup>20</sup>. We also compared Hi-MC haplogroup classifications to  
157 Haplogrep-based haplogroup classifications for both Phase I and Phase III HapMap samples. We  
158 calculated percent concordance for each comparison. Classifications were considered concordant  
159 if they were in the same haplogroup, even if one classification method resulted in finer resolution.  
160 For example, if one method classified a sample as A2 and another method classified the same  
161 sample as A2x, such classifications were considered concordant. Differences in the resolution of  
162 haplogroup classifications were not unexpected given differences in underlying methodology and  
163 the number of SNPs used for classification. The HapMap classifications were generated using  
164 more mitochondrial SNP genotypes compared to the reduced number of SNPs necessary to use

165 Hi-MC. HapMap Phase I sample data includes genotypes for 214 mitochondrial SNPs, 49 of  
166 which overlap with the custom SNP panel genotyped in this study (Supplementary Table 2).  
167 Additionally, Hi-MC uses a reduced tree for classification while Haplogrep employs all of  
168 Phylotree which can result in finer sub-haplogroup resolution.

169 To resolve discordant classifications, possibly due to missing key SNP genotypes, we  
170 used the publicly available Phase I HapMap mitochondrial SNP genotype data to determine the  
171 mitochondrial haplogroup classification via Haplogrep. If the classification returned from  
172 Haplogrep was concordant with the HapMap-reported classification, then we considered the  
173 discordance resolved, as it was likely due to missing SNP genotypes necessary for accurate  
174 haplogroup classification by Hi-MC.

## 175 **Results**

### 176 *CEU and YRI populations*

177 Overall, concordance between Hi-MC and both HapMap and Haplogrep was high for the CEU  
178 and YRI populations. Among the CEU samples mitochondrial haplogroup classifications were  
179 100% concordant between Hi-MC and HapMap, as well as between Hi-MC and Haplogrep  
180 (Table 1). In the YRI samples, concordance between Hi-MC and HapMap was 96.3% (Table 1).  
181 Among the YRI samples, three classifications were discordant between Hi-MC and HapMap, one  
182 classification was discordant between Hi-MC and Haplogrep, and four classifications were  
183 discordant between Haplogrep and HapMap. The three samples that were discordant between Hi-  
184 MC and HapMap were also discordant between Haplogrep and HapMap.

185 Among the eleven YRI samples that were either discordant or unclassified seven were  
186 resolved. These samples were missing many SNP genotypes and/or crucial haplogroup-defining  
187 SNPs in our genotype data which likely accounts for the discordance. The four YRI samples for  
188 which discordance could not be resolved (Y024-NA18861, Y024-NA18663, Y043-NA19137,  
189 and Y043-NA19139) were classified as ‘L1’ by HapMap, but were classified as ‘L0a’ by

190 Haplogrep using HapMap-generated genotype data. The ‘L0’ classification is consistent with the  
191 classification obtained via Hi-MC and Haplogrep when using genotypes from our custom SNP  
192 panel. In the HapMap genotype data, all of these samples have eight of the ten SNP genotypes  
193 that define haplogroup ‘L0’, suggesting that ‘L0’ is the correct classification.

#### 194 ***CHB/JPT populations***

195 Compared with the CEU and YRI populations, we observed less concordance among the  
196 CHB/JPT samples. Between Hi-MC and HapMap-reported classifications, 37 (41.6%) were  
197 concordant at the haplogroup level and 31 (34.8%) were considered concordant at the macro-  
198 haplogroup level. Concordance at the macro-haplogroup level is defined as appropriate macro-  
199 haplogroup classification in the absence of sub-haplogroup defining SNP genotype data. For  
200 example, consider that haplogroup E is a sub-haplogroup of the macro-haplogroup M. Genotypes  
201 for SNPs that define haplogroup E were not included in the custom SNP panel; therefore,  
202 individuals classified as haplogroup E by HapMap, but classified as haplogroup M by Hi-MC  
203 were considered concordant at the macro-haplogroup level. There were 21 (23.6%) discordant  
204 classifications among the CHB/JPT samples. These results were not unexpected given that the  
205 SNPs included on the custom panel do not capture all Asian-specific haplogroup lineages. Among  
206 the 21 CHB/JPT samples that were discordant, two samples were resolved at the haplogroup level  
207 and five samples were resolved at the macro-haplogroup level. The remaining samples with  
208 discordant classifications could not be resolved.

#### 209 ***Determination of mitochondrial haplogroups in HapMap Phase III samples of Mexican*** 210 ***ancestry***

211 The mitochondrial haplogroups for the samples of Mexican ancestry from HapMap Phase III  
212 have not been previously reported. Samples in this data set include 30 trios of Mexican ancestry  
213 from Los Angeles, CA. We applied Hi-MC to determine mitochondrial haplogroups in these  
214 samples and characterized the distribution of mitochondrial haplogroups among the MXL. Due to

215 matrilineal inheritance of mtDNA, offspring have the same mitochondrial haplogroup as their  
216 mother; therefore, offspring were excluded when calculating the frequency distribution of  
217 mitochondrial haplogroups. One additional sample was excluded from frequency calculations due  
218 to poor genotyping efficiency. Overall in the MXL samples, 84.8% of mitochondrial haplogroups  
219 identified were of Native American ancestry and 15.3% were of European ancestry (Table 2). The  
220 distribution of haplogroups in the HapMap MXL samples is similar to the distribution of  
221 haplogroups observed in Mexican Americans ascertained for the National Health and Nutrition  
222 Examination Surveys (NHANES)<sup>22</sup>.

223       To further evaluate the performance of Hi-MC, we compared the Hi-MC mitochondrial  
224 haplogroup classifications of MXL samples to Haplogrep-based classifications. Percent  
225 concordance between Hi-MC and Haplogrep for classification of the MXL samples was 98.9%.  
226 There was one sample out of 89 with a discordant mitochondrial haplogroup classification. This  
227 sample was missing the haplogroup H-defining SNP genotype therefore Hi-MC was unable to  
228 classify the sample beyond haplogroup 'HV.' Haplogrep classified this sample as H1c1b. For this  
229 individual the classifications differ between Hi-MC and Haplogrep due to differences in  
230 methodology.

## 231 **Discussion**

232 Using a custom panel of mitochondrial SNPs that we previously applied to participants in the  
233 NHANES data sets<sup>22</sup>, we developed Hi-MC, a method for high-throughput classification of  
234 European, African, and Native American mitochondrial haplogroup lineages. We evaluated the  
235 performance of Hi-MC, and with genotype data from the custom SNP panel, demonstrate that Hi-  
236 MC performs comparably to the widely-used tool Haplogrep. While Haplogrep is an excellent  
237 tool for mitochondrial haplogroup classification that accepts either sequence or SNP genotype  
238 data, it was developed primarily for sequence level data. The ability to alternatively genotype a

239 relatively small number of SNPs ( $n=63$ ) allows for rapid haplogroup classification in a large  
240 number of genetic samples.

241 Mitochondrial SNPs captured by standard genotyping arrays vary widely, and often the  
242 SNPs on these arrays are not informative for haplogroup determination. Hi-MC uses a defined  
243 panel of mitochondrial SNPs for classification of mitochondrial haplogroups. This defined panel  
244 of SNPs eliminates the need for investigators to spend time identifying appropriate SNPs for  
245 mitochondrial haplogroup classification. Additionally, the relatively small number of SNPs in the  
246 custom panel makes Hi-MC particularly useful for large data sets where full mitochondrial  
247 genome sequencing is not practical. As examples, approaches like Hi-MC promise to be of use to  
248 large biobank and cohort efforts such as Million Veteran Program<sup>32</sup> and the UK Biobank<sup>33</sup>, both  
249 of which continue to rely on cost-effective array-based assays rather than cost-prohibitive  
250 sequencing to generate genome-wide and mitochondrial data on hundreds of thousands to a  
251 million participants.

252 Hi-MC employs the commonly used PED/MAP file format as the input. There are a  
253 number of software programs that make use of the PED/MAP format, including PLINK<sup>34</sup> which  
254 is widely used for analyzing genotypic data. Thus, in contrast to Haplogrep, many Hi-MC users  
255 will not have to reformat data prior to use. Additionally, Hi-MC is an R-based software package  
256 that can be downloaded and run locally allowing for memory limits that are dependent on the  
257 machine where R is being run, thus granting greater flexibility in the number of samples that can  
258 be processed at once. Once samples have been classified using Hi-MC, figures or tables  
259 displaying haplogroup frequencies can be easily generated via other R packages such as  
260 `ggplot2`<sup>35</sup>.

261 We determined that Hi-MC performs well with samples of European, African, and Native  
262 American descent. However, because many Asian-specific haplogroups are not captured by the  
263 custom SNP panel it does not perform as well on samples of Asian maternal lineage. While

264 progress has been made in characterizing the phylogeny of Asian mtDNA<sup>36,37</sup>, in general, the  
265 Asian branches of the mitochondrial phylogenetic tree are not as well-defined as other parts of  
266 the tree. Thus, compared to other ancestries, classifying Asian lineage haplogroups continues to  
267 be more challenging. As more mtDNA sequences are obtained from individuals of Asian descent  
268 the phylogeny of mitochondrial genetic variation will be better understood. Future versions of Hi-  
269 MC will be updated to incorporate additional knowledge regarding subjects of Asian descent.

270 We applied Hi-MC to the HapMap Phase III MXL samples as the mitochondrial  
271 haplogroups for these participants have not been previously reported. The haplogroup distribution  
272 observed in the HapMap Phase III MXL samples is somewhat similar to the recently reported  
273 Haplogrep2-generated distribution for the MXL samples sequenced as part of the 1000 Genomes  
274 Project<sup>38</sup>. In this newer reference dataset, the most common reported haplogroup is A (25%)  
275 followed by B (15%) and C (9%)<sup>38</sup> compared with a higher A (A2) frequency in the present study  
276 (39%; Table 2). Overall, the distribution of Native American and European haplogroups in the  
277 MXL samples from HapMap Phase III is similar to the distribution observed in the NHANES  
278 Mexican American samples<sup>22</sup>. No African lineage mitochondrial haplogroups were identified  
279 among the HapMap MXL samples. This differs from the NHANES Mexican Americans in which  
280 4.4% had mitochondrial haplogroups of African ancestry<sup>22</sup>. The lack of African haplogroups in  
281 the HapMap MXL samples is likely due to the small sample size and the regional ascertainment  
282 of these samples. While the NHANES samples were collected from across the United States, the  
283 HapMap Phase III MXL samples were ascertained solely from Los Angeles, CA, therefore are  
284 likely to be more homogeneous.

285 While there are several benefits to Hi-MC, there are some limitations. Currently, Hi-MC  
286 employs a reduced mitochondrial phylogenetic tree for classification. As a result, it is currently  
287 limited to classification of the major haplogroups of European, African, and Native American  
288 lineages, and requires that SNPs from the described custom panel be genotyped. While this panel

289 was customized for populations expected for the PAGE I study, it is notable that several SNPs in  
290 this panel (MT1736, MT2092, MT3552, MT4883, MT10400, MT11177, MT11251, MT11719,  
291 MT12007, MT12308, MT12705, MT13368, MT14766) overlap with previously published  
292 panels<sup>1,39</sup>, suggesting the potential for both greater resolution and generalizability in future  
293 extensions of Hi-MC. Additionally, because the method relies on a limited number of SNPs, it is  
294 not very robust to missing genotype data and it has the ability to classify mitochondrial  
295 haplogroups at a broad level, but currently cannot capture sub-haplogroups at finer resolution. As  
296 such, in instances where sequence level data is available another method for mitochondrial  
297 haplogroup classification, such as Haplogrep, would be more appropriate.

298         Despite these limitations, Hi-MC offers several advantages including a defined panel of  
299 mitochondrial SNPs that is used in conjunction with the software for mitochondrial haplogroup  
300 classification. Hi-MC utilizes PED/MAP files for a user-friendly input file format, saving time  
301 and reducing opportunities for errors to be incorporated into the data. Also, Hi-MC is  
302 implemented in the commonly used statistical software environment R allowing for classification  
303 of relatively large sample sizes, as well as the ability to easily utilize other available R packages  
304 for visualization of results.

### 305 ***Conclusions***

306 We have developed a custom SNP panel and algorithm for mitochondrial haplogroup  
307 classification. The algorithm, Hi-MC is implemented in R and makes use of PED/MAP file  
308 format for data input. We evaluated the performance of Hi-MC and demonstrate that  
309 classifications are comparable to the widely-used tool Haplogrep. Hi-MC offers an algorithm that  
310 leverages a validated mtDNA SNP panel for mitochondrial haplogroup classification and is  
311 particularly valuable for studies in which sequencing is not feasible.

### 312 **Conflict of Interest**

313 The authors declare no conflict of interest.

314 **Acknowledgements**

315 Special thanks to Paxton Baker, MS, Melissa Allen, Ping Mayo, MS, and Nathalie Schnetz-  
316 Boutaud, PhD for their work in genotyping these samples. This work was supported in part by  
317 National Institutes of Health grant [U01 HG004798] and associated American Recovery and  
318 Reinvestment Act (ARRA) supplements.

319 **Figure 1: Hi-MC algorithm structure**

320 Input for the algorithm is a list of sample IDs and corresponding SNP genotype data in pedigree  
321 (PED/MAP) format. These genotypes are recursively analyzed through a node-based tree  
322 structure. Each successive genotype classification is passed on to the Accumulator. They are then  
323 ranked according to specificity [longer path through the tree -> more SNPs checked -> more  
324 specific], with the most specific haplogroup as the final output. MRCA = most recent common  
325 ancestor





326 **References**

- 327 1. Chaitanya L, van Oven M, Weiler N, Harteveld J, Wirken L, Sijen T, et al. Developmental  
328 validation of mitochondrial DNA genotyping assays for adept matrilineal inference of  
329 biogeographic ancestry at a continental level. *Forensic Science International: Genetics*.  
330 2014;11(Supplement C):39-51.
- 331 2. van der Walt JM, Dementieva YA, Martin ER, Scott WK, Nicodemus KK, Kroner CC, et  
332 al. Analysis of European mitochondrial haplogroups with Alzheimer disease risk. *Neuroscience*  
333 *Letters*. 2004;365(1):28-32.
- 334 3. K S, Jalali S, Scaria V, Bhardwaj A. MitoLSDB: A Comprehensive Resource to Study  
335 Genotype to Phenotype Correlations in Human Mitochondrial DNA Variations. *PLOS ONE*.  
336 2013;8(4):e60066.
- 337 4. Wallace DC. Bioenergetics in human evolution and disease: implications for the origins of  
338 biological complexity and the missing genetic variation of common diseases. *Philosophical*  
339 *Transactions of the Royal Society B: Biological Sciences*. 2013;368(1622).
- 340 5. Mitchell S, Hall J, Goodloe R, Boston J, Farber-Eger E, Pendergrass S, et al. Investigating  
341 the relationship between mitochondrial genetic variation and cardiovascular-related traits to  
342 develop a framework for mitochondrial phenome-wide association studies. *BioData Mining*.  
343 2014;7(1):6.
- 344 6. Hudson G, Gomez-Duran A, Wilson IJ, Chinnery PF. Recent Mitochondrial DNA  
345 Mutations Increase the Risk of Developing Common Late-Onset Human Diseases. *PLOS*  
346 *Genetics*. 2014;10(5):e1004369.
- 347 7. Fetterman Jessica L, Zelickson Blake R, Johnson Larry W, Moellering Douglas R,  
348 Westbrook David G, Pompilius M, et al. Mitochondrial genetic background modulates  
349 bioenergetics and susceptibility to acute cardiac volume overload. *Biochemical Journal*.  
350 2013;455(2):157-67.

- 351 8. Maca-Meyer N, Gonzalez AM, Larruga JM, Flores C, Cabrera VM. Major genomic  
352 mitochondrial lineages delineate early human expansions. *BMC Genet.* 2001;2:13.
- 353 9. Forster P. Ice Ages and the mitochondrial DNA chronology of human dispersals: a review.  
354 *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences.*  
355 2004;359(1442):255-64.
- 356 10. Kloss-Brandstätter A, Pacher D, Schönherr S, Weissensteiner H, Binna R, Specht G, et al.  
357 HaploGrep: a fast and reliable algorithm for automatic classification of mitochondrial DNA  
358 haplogroups. *Human Mutation.* 2011;32(1):25-32.
- 359 11. Weissensteiner H, Pacher D, Kloss-Brandstätter A, Forer L, Specht G, Bandelt HJ, et al.  
360 HaploGrep 2: mitochondrial haplogroup classification in the era of high-throughput sequencing.  
361 *Nucleic Acids Res.* 2016;44(W1):W58-W63.
- 362 12. Fan L, Yao Y-G. MitoTool: A web server for the analysis and retrieval of human  
363 mitochondrial DNA sequence variations. *Mitochondrion.* 2011;11(2):351-6.
- 364 13. Rubino F, Piredda R, Calabrese FM, Simone D, Lang M, Calabrese C, et al. HmtDB, a  
365 genomic resource for mitochondrion-based human variability studies. *Nucleic Acids Res.*  
366 2012;40(Database issue):D1150-D9.
- 367 14. Fan L, Yao Y-G. An update to MitoTool: Using a new scoring system for faster mtDNA  
368 haplogroup determination. *Mitochondrion.* 2013;13(4):360-3.
- 369 15. Vianello D, Sevini F, Castellani G, Lomartire L, Capri M, Franceschi C. HAPLOFIND: A  
370 New Method for High-Throughput mtDNA Haplogroup Assignment. *Human Mutation.*  
371 2013;34(9):1189-94.
- 372 16. Navarro-Gomez D, Leipzig J, Shen L, Lott M, Stassen AP, Wallace DC, et al. Phy-Mer: a  
373 novel alignment-free and reference-independent mitochondrial haplogroup classifier.  
374 *Bioinformatics.* 2015;31(8):1310-2.

- 375 17. Calabrese C, Simone D, Diroma MA, Santorsola M, Gutta C, Gasparre G, et al.  
376 MToolBox: a highly automated pipeline for heteroplasmy annotation and prioritization analysis  
377 of human mitochondrial variants in high-throughput sequencing. *Bioinformatics*.  
378 2014;30(21):3115-7.
- 379 18. The International HapMap Project. *Nature*. 2003;426(6968):789-96.
- 380 19. Consortium IH. Integrating common and rare genetic variation in diverse human  
381 populations. *Nature*. 2010;467(7311):52-8.
- 382 20. Consortium TIH. A haplotype map of the human genome. *Nature*. 2005;437(7063):1299-  
383 320.
- 384 21. van Oven M, Kayser M. Updated comprehensive phylogenetic tree of global human  
385 mitochondrial DNA variation. *Human Mutation*. 2009;30(2):E386-E94.
- 386 22. Mitchell SL, Goodloe R, Brown-Gentry K, Pendergrass SA, Murdock DG, Crawford DC.  
387 Characterization of mitochondrial haplogroups in a large population-based sample from the  
388 United States. *Hum Genet*. 2014;133(7):861-8.
- 389 23. R: A language and environment for statistical computing. Vienna, Austria: R Foundation  
390 for Statistical Computing; 2013 2013.
- 391 24. Herrstadt C, Elson JL, Fahy E, Preston G, Turnbull DM, Anderson C, et al. Reduced-  
392 Median-Network Analysis of Complete Mitochondrial DNA Coding-Region Sequences for the  
393 Major African, Asian, and European Haplogroups. *The American Journal of Human Genetics*.  
394 2002;70(5):1152-71.
- 395 25. van der Walt JM, Nicodemus KK, Martin ER, Scott WK, Nance MA, Watts RL, et al.  
396 Mitochondrial Polymorphisms Significantly Reduce the Risk of Parkinson Disease. *The*  
397 *American Journal of Human Genetics*. 2003;72(4):804-11.

- 398 26. Poole Jason C, Procaccio V, Brandon Martin C, Merrick G, Wallace Douglas C. Multiplex  
399 analysis of mitochondrial DNA pathogenic and polymorphic sequence variants. *Biol Chem*.  
400 2010;391(10):1115-30.
- 401 27. Paneto GG, Köhnemann S, Martins JA, Cicarelli RMB, Pfeiffer H. A single multiplex  
402 PCR and SNaPshot minisequencing reaction of 42 SNPs to classify admixture populations into  
403 mitochondrial DNA haplogroups. *Mitochondrion*. 2011;11(2):296-302.
- 404 28. Crawford DC, Goodloe R, Farber-Eger E, Boston J, Pendergrass SA, Haines JL, et al.  
405 Leveraging epidemiologic and clinical collections for genomic studies of complex traits. *Human*  
406 *Heredity*. 2015;79(3-4):137-46.
- 407 29. Matisse TC, Ambite JL, Buyske S, Carlson CS, Cole SA, Crawford DC, et al. The Next  
408 PAGE in Understanding Complex Traits: Design for the Analysis of Population Architecture  
409 Using Genetics and Epidemiology (PAGE) Study. *American Journal of Epidemiology*.  
410 2011;174(7):849-59.
- 411 30. Tang K, Oeth P, Kammerer S, Denissenko MF, Ekblom J, Jurinke C, et al. Mining disease  
412 susceptibility genes through SNP analyses and expression profiling using MALDI-TOF mass  
413 spectrometry. *J Proteome Res*. 2004;3(2):218-27.
- 414 31. Hazkani-Covo E, Zeller RM, Martin W. Molecular Poltergeists: Mitochondrial DNA  
415 Copies (numts) in Sequenced Nuclear Genomes. *PLOS Genetics*. 2010;6(2):e1000834.
- 416 32. Gaziano JM, Concato J, Brophy M, Fiore L, Pyarajan S, Breeling J, et al. Million Veteran  
417 Program: A mega-biobank to study genetic influences on health and disease. *Journal of Clinical*  
418 *Epidemiology*. 2016.
- 419 33. Sudlow C, Gallacher J, Allen N, Beral V, Burton P, Danesh J, et al. UK Biobank: An Open  
420 Access Resource for Identifying the Causes of a Wide Range of Complex Diseases of Middle and  
421 Old Age. *PLoS Med*. 2015;12(3):e1001779.

- 422 34. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, et al. PLINK: a  
423 tool set for whole-genome association and population-based linkage analysis. *Am J Hum Genet.*  
424 2007;81(3):559-75.
- 425 35. Wickham H. ggplot2: Elegant graphics for data analysis. New York: Springer-Verlag;  
426 2009.
- 427 36. Kivisild T, Tolk HV, Parik J, Wang Y, Papiha SS, Bandelt HJ, et al. The emerging limbs  
428 and twigs of the East Asian mtDNA tree. *Mol Biol Evol.* 2002;19(10):1737-51.
- 429 37. Kong QP, Bandelt HJ, Sun C, Yao YG, Salas A, Achilli A, et al. Updating the East Asian  
430 mtDNA phylogeny: a prerequisite for the identification of pathogenic mutations. *Hum Mol Genet.*  
431 2006;15(13):2076-86.
- 432 38. Rishishwar L, Jordan IK. Implications of human evolution and admixture for  
433 mitochondrial replacement therapy. *BMC Genomics.* 2017;18(1):140.
- 434 39. van Oven M, Vermeulen M, Kayser M. Multiplex genotyping system for efficient  
435 inference of matrilineal genetic ancestry with continental resolution. *Investigative Genetics.*  
436 2011;2(1):6.

**Table 1** (on next page)

Table 1.

Percent concordance in CEU and YRI populations for pair-wise comparisons of mitochondrial haplogroup classifications

**Table 1:** Percent concordance in CEU and YRI populations for pair-wise comparisons of mitochondrial haplogroup classifications

	<b>CEU (n=86*)</b>	<b>YRI (n=82*)</b>
<b>Hi-MC vs HapMap</b>	100%	96.3%
<b>Hi-MC vs Haplogrep</b>	100%	98.8%
<b>Haplogrep vs HapMap</b>	100%	95.1%

\*Due to missing genotypes at key haplogroup-defining SNPs four CEU and eight YRI samples were excluded from the percent concordance calculations.



**Table 2** (on next page)

Table 2. Distribution of mitochondrial haplogroups in the HapMap Phase III samples of Mexican ancestry in Los Angeles, CA

**Table 2:** Distribution of mitochondrial haplogroups in the HapMap Phase III samples of Mexican ancestry in Los Angeles, CA

<b>Mitochondrial Haplogroup</b>	<b>Number (%)</b>
<b>Native American</b>	
<b>A2</b>	23 (39.0%)
<b>B2</b>	11 (18.6%)
<b>C</b>	9 (15.3%)
<b>D1</b>	7 (11.9%)
<b>European</b>	
<b>H</b>	3 (5.1%)
<b>H/V</b>	2 (3.4%)
<b>U</b>	2 (3.4%)
<b>V</b>	1 (1.7%)
<b>W</b>	1 (1.7%)

Given that the mitochondrial haplogroup of the offspring is the same as that of the mother, offspring were excluded when determining the frequency distribution of haplogroups. One sample was excluded from frequency calculations due to missing genotype data (n=59).

**Figure 1**(on next page)

## Hi-MC algorithm structure

Input for the algorithm is a list of sample IDs and corresponding SNP genotype data in pedigree (PED/MAP) format. These genotypes are recursively analyzed through a node-based tree structure. Each successive genotype classification is passed on to the Accumulator. They are then ranked according to specificity [longer path through the tree -> more SNPs checked -> more specific], with the most specific haplogroup as the final output. MRCA = most recent common ancestor

