# Sequencing of Full Mitochondrial Genomes for NIST Population Samples

Kevin M. Kiesler[1], K.S. Andreaggi[2,3], J.D. Ring[2,3], C.R. Taylor[2,3], C.K. Marshall[2,3], and P.M. Vallone[1]

[1]U.S. National Institute of Standards and Technology, 100 Bureau Drive, Gaithersburg, MD 20899-8314, USA
[2]SNA International, Alexandria, VA, USA
[3]Armed Forces Medical Examiner System's Armed Forces DNA Identification Lab (AFMES-AFDIL), Dover, DE, USA

Email: Kevin.Kiesler@nist.gov

National Institute of Standards and Technology
Technology Administration, U.S. Department of Commerce

Poster # P146

The U.S. National Institute of Standards and Technology (NIST) has undertaken a population sequencing study of 607 mitochondrial genomes (mtGenomes). The goal of the study is to increase the availability of high-quality whole mtGenome haplotypes for haplotype frequency estimations, allowing the increased discrimination power of the full mtGenome to be harnessed. Whole mtGenome sequencing is enabled by Next Generation Sequencing (NGS) workflows which are vastly simplified compared to Sanger-type sequencing, which is typically limited to the mtDNA Control Region, or portions thereof, due to the labor required to produce mtDNA sequence data. Samples used in this study were grouped into categories of self-reported race: African American, U.S. Caucasian, and U.S. Hispanic individuals, representing three major U.S. populations. Assessment of resolution for mtGenome relative to control region data are presented. Previously characterized markers for this sample set include Y-chromosome Short Tandem Repeats (Y-STRs) and Ancestry Informative Single Nucleotide Polymorphisms (AISNPs). Comparisons of ancestry inference from mtGenome, AISNP, Y-STR data, and self-reporting are presented.

## Methods – mtGenome Sequencing

- **Full mitochondrial genome (mtGenome) sequences obtained for 607 NIST population samples**
  - Self-reported: African American (n = 236), U.S. Caucasians (n = 247), and U.S. Hispanics (n = 124)
- **Long PCR amplification of full mitochondrial genome (mtGenome)**
  - Primers & PCR conditions from Fendt et al. [1]
- **Sequenced on next generation instrument**
  - KAPA Hyperplus library procedure from Ring et al. [2]
  - MiSeq FGx instrument using MiSeq Reagent Kit v3-600
- **mtGenome data curated using CLC Bio Genomics Workbench – AQME tool [3]**
  - Data review and curation performed by AFMES-AFDIL staff
  - Minimum heteroplasmy variant frequency 5 %

## Ancestry inference using mtGenome, autosomal SNPs, and Y STRs

- **Full mitochondrial haplogroups assigned by CLC Bio Genomics Workbench – AQME Tool using EMPOP database values**
  - Self-reported: African American (n = 236), U.S. Caucasians (n = 247), and U.S. Hispanics (n = 124)
- **Autosomal ancestry-informative SNPs from ForenSeq**
  - Libraries prepared with Verogen ForenSeq Signature DNA Prep Kit with DNA Primer Mix B (DPMB)
  - Sequenced on MiSeq FGx using MiSeq FGx Reagent Kit
  - Ancestry informative SNPs from Kidd et. al. [4]
  - Ancestry predictions from SNP data made using FROG-kb [5]
- **Y-STR profiles from ForenSeq and Promega PPY-23**
  - From NIST 1036 population sequencing [6] and ForenSeq DNA Signature Kit (unpublished)
  - Y Hgs were estimated using the NevGen Haplogroup Predictor (desktop version)
    - https://www.nevgen.org/
  - Y haplogroups categorized by continent according to Kivisild [8] and Lao [9]

Figure 1: Mitochondrial genome showing control region with amplification strategy.
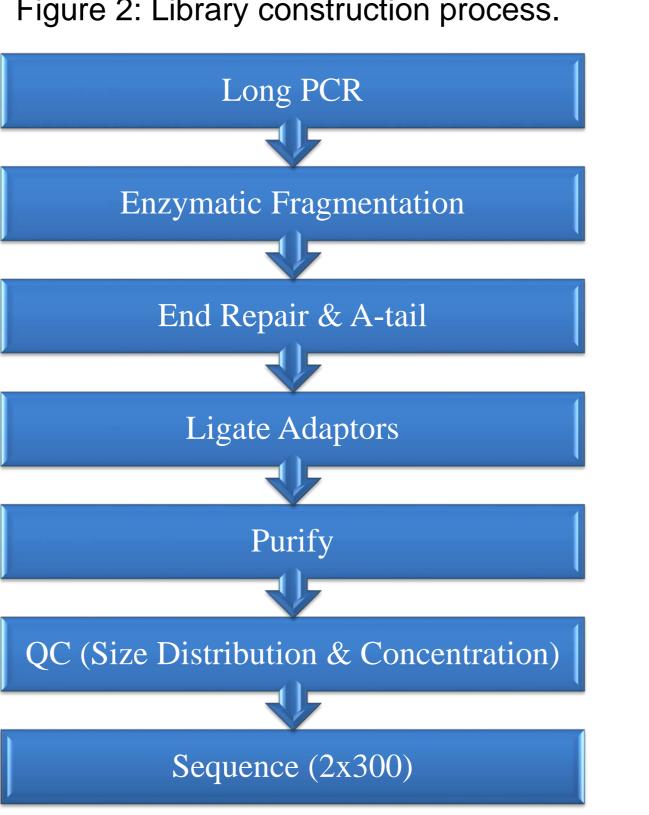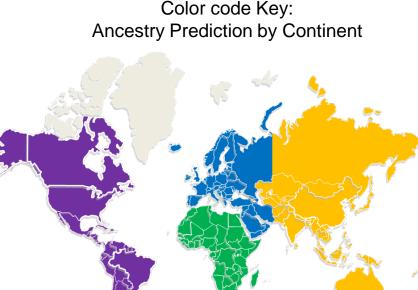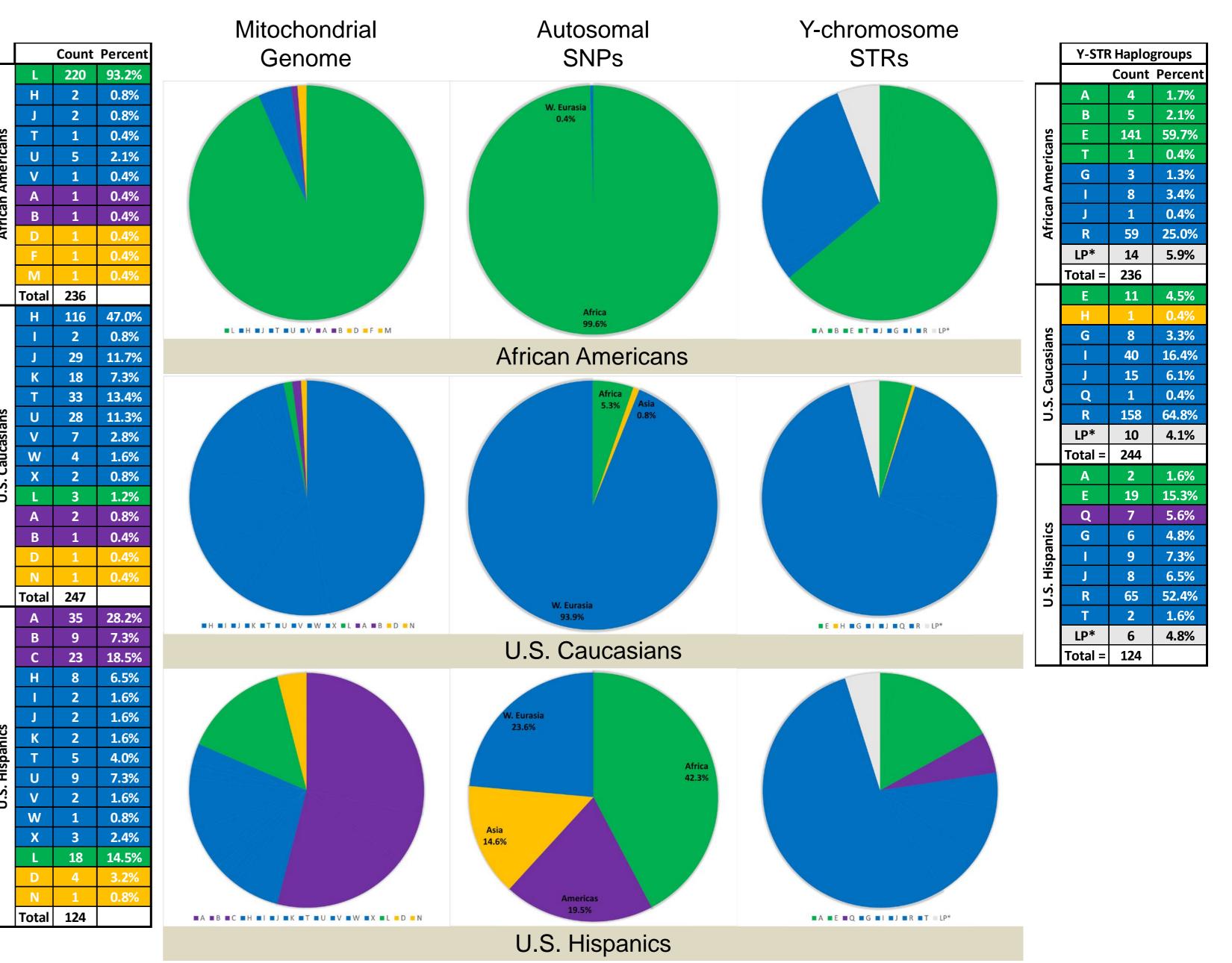
Figure 2: Library construction process.

Figure 3: Ancestry prediction for three marker types for the three populations studied. Predictions were grouped into broad continental origins: African (green), European (blue), Asian (Orange), and Americas (purple); for mitochondrial haplogroups according to Rishishwar [7] and for Y-STR haplogroups according to Kivisild [8] and Lao [9]. For SNP data the top metapopulation predicted with The Standalone FROG-kb Ancestry Inference Batch Likelihood Computation Tool found on GitHub: https://github.com/haseenaR/FROGAncestryCalc was used. Location of the metapopulation was converted into continental ancestry groups. Mitochondrial and Y-chromosome subclades were reduced to highest "single letter" level for simplicity of display in this figure.
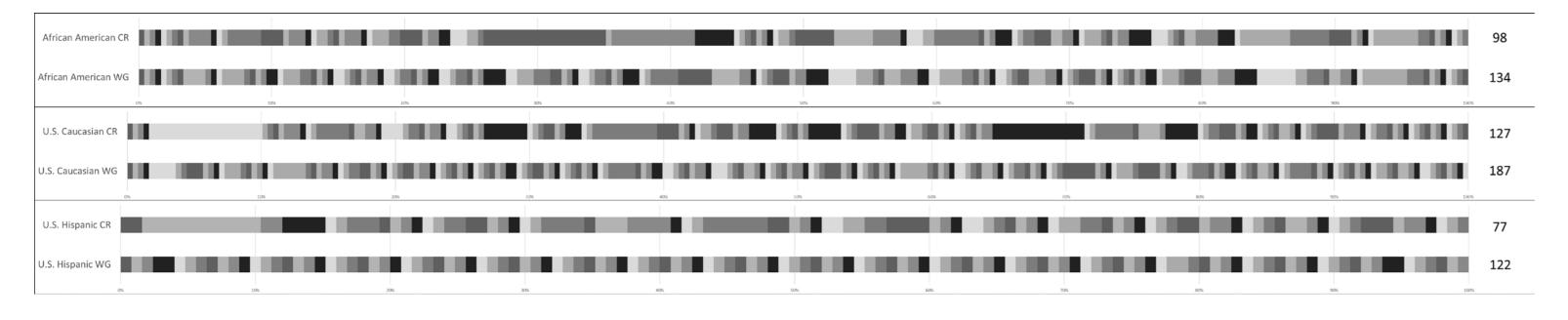
Table 1: Comparison of shared profiles from Control Region (nt 16,024 to 16,569 & 1 to 576) vs whole mtGenome.

| Number of Profiles Matching | All Samples Control Region | All Samples Whole Genome | African American Control Region | African American Whole Genome | U.S. Caucasian Control Region | U.S. Caucasian Whole Genome | U.S. Hispanic Control Region | U.S. Hispanic Whole Genome |
|---|---|---|---|---|---|---|---|---|
| 11 | 1 | 0 | | | 1 | 0 | | |
| 6 | 1 | 0 | | | | | 1 | 0 |
| 5 | 1 | 1 | | | 1 | 0 | 1 | 1 |
| 4 | 5 | 0 | 1 | 0 | 4 | 0 | | |
| 3 | 7 | 3 | 4 | 1 | 3 | 2 | | |
| 2 | 28 | 11 | 13 | 4 | 8 | 2 | 7 | 5 |
| Unique | 488 (88 %) | 571 (94 %) | 194 (82 %) | 225 (95 %) | 190 (77 %) | 237 (96 %) | 104 (84 %) | 109 (88 %) |
| Total # | 607 | 607 | 236 | 236 | 247 | 247 | 124 | 124 |

Figure 2: Comparison of haplogroups derived from Control Region vs whole mtGenome for three U.S. populations with number of haplogroups shown at end of each bar graph.

- African American CR — 98
- African American WG — 134
- U.S. Caucasian CR — 127
- U.S. Caucasian WG — 187
- U.S. Hispanic CR — 77
- U.S. Hispanic WG — 122

### Discussion - Whole mtGenome vs Control Region

- Indels at positions 309, 315, 573, and 16,193 were not considered differences in profiles.
- Point heteroplasmy was considered a difference between two profiles.
  - Heteroplasmy was observed 209 times in these 607 mtGenomes.
- Overall 83 additional profiles could be uniquely identified in this population set when using whole mtGenome (WG) analysis vs control region (CR), an increase of 17 % within the study.
- The Caucasian population benefitted the most from whole mtGenome analysis with an increase in unique profiles of 25 %, followed by African Americans with a 16 % increase, and Hispanics with a 5 % increase.
- What happened to the most common shared CR profiles when WG is used?
  - African American group of 4 → group of 3 + 1 unique profile
  - U.S. Caucasian group of 11 → 11 unique profiles
  - U.S. Hispanic group of 6 → group of 5 + 1 unique profile
- Whole mtGenome sequencing enables deeper subclade classification vs control region (see Figure 2).

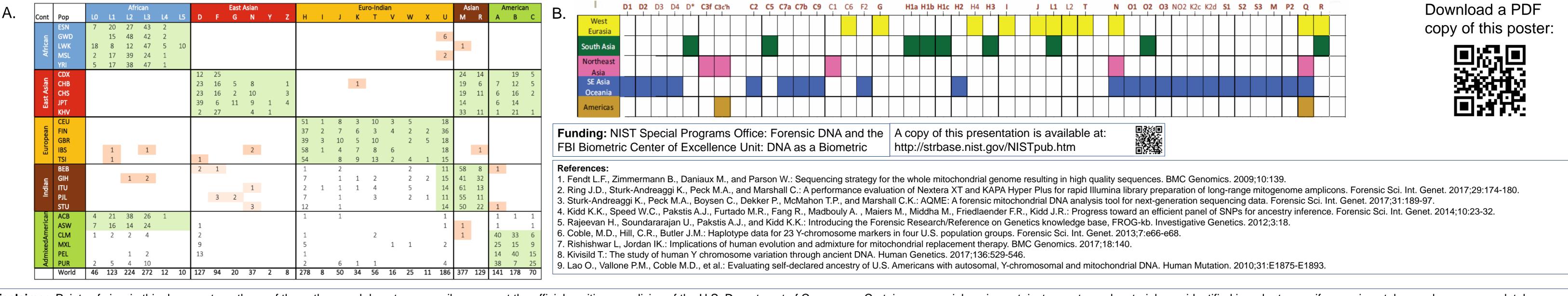### Discussion – Ancestry Inference relative to self-identified ancestry

- For the U.S. Caucasian samples all three marker systems gave similar predictions.
  - Categorized as "West Eurasian": mtGenome = 96.8 %, AISNP = 93.8 %, Y-STR = 91.0 %
- For the U.S. African American samples, mtDNA and AISNP inference performed similarly while Y-STR inference performed differently.
  - Categorized as "African": mtGenome = 93.3 %, AISNP = 99.6 %, Y-STR = 64.0 %
- For U.S. Hispanic samples continental categorizations indicate weaker correlation to self identified ancestry.
  - Categorized as "European" or "Americas": mtGenome = 81.5 %, AISNP = 43.1 %, Y-STR = 78.2 %

*LP = low probability Y-haplogroup prediction from STR data. Minimum probability threshold of 0.10 used for NevGen haplogroup predictor.
**See poster P145 for a larger population scale Y-STR analysis.

### Conclusions:

- As expected, greater discrimination is observed from Whole mtGenome sequence compared to the Control Region sequence.
  - Gains in unique profiles: African American = 16 %, U.S. Caucasian = 25 %, and U.S. Hispanics = 5 %
- Whole mtGenome sequence may be used to differentiate some individuals when control region mtDNA profiles do not.
- Ancestry prediction with mitochondrial DNA works similarly to Y-STR or nuclear DNA SNPs.
  - Differences in mtDNA and Y-chromosome ancestry proportions may be due to bias in admixture history.
  - Y STR haplogroup inferences from the NevGen Haplogroup Predictor are informative, but additional Y-STR or Y-SNP typing may be needed to support improved paternal lineage predictions.
  - Ancestry inference in groups such as U.S. Hispanics are more variable, due to complex admixture and natural history of such populations.
- There is a limitation in the Bayesian Y haplogroup predictor, NevGen, in that it always predicts a result, even when a subclade may not be present in the reference database or support for a prediction is weak. This could result in false positives.
- Ancestry is affected by human migration, both ancient and modern, with no clearly defined boundaries for haplogroups. Many haplogroups are represented on multiple continents at varying frequencies. This could lead to misclassification at the "continent" level.
- Future: finish the mitochondrial for the remaining NIST population samples (n ≈ 662) and submit to EMPOP to increase the number whole mitochondrial genomes in the database.

Figure 4: Haplogroups ascribed to continent for (A) mitochondrial DNA according to [7] and (B) Y-STR haplogroup predicted by NEVGEN Y-haplogroup prediction tool according to [8].

Download a PDF copy of this poster:

References:
1. Fendt L.F., Zimmermann B., Daniaux M., and Parson W.: Sequencing strategy for the whole mitochondrial genome resulting in high quality sequences. BMC Genomics. 2009;10:139.
2. Ring J.D., Sturk-Andreaggi K., Peck M.A., and Marshall C.: A performance evaluation of Nextera XT and KAPA Hyper Plus for rapid Illumina library preparation of long-range mitogenome amplicons. Forensic Sci. Int. Genet. 2017;31:174-180.
3. Sturk-Andreaggi K., Peck M.A., Boyser C., Dekker P., McMahon T.P., and Marshall C.K.: AQME: A forensic mitochondrial DNA analysis tool for next-generation sequencing data. Forensic Sci. Int. Genet. 2017;31:189-97.
4. Kidd K.K., Speed W.C., Pakstis A.J., Furtado M.R., Fang R., Madbouly A., Maiers M., Middha M., Friedlaender F.R., Kidd J.R.: Progress toward an efficient panel of SNPs for ancestry inference. Forensic Sci. Int. Genet. 2014;10:23-32.
5. Rajeevan H., Soundararajan U., Pakstis A.J., and Kidd K.K.: Introducing the Forensic Research/Reference on Genetics knowledge base, FROG-kb. Investigative Genetics. 2012;3:18.
6. Coble, M.D., Hill, C.R., Butler J.M.: Haplotype data for 23 Y-chromosome markers in four U.S. population groups. Forensic Sci. Int. Genet. 2013;7:e66-e68.
7. Rishishwar L, Jordan IK.: Implications of human evolution and admixture for mitochondrial replacement therapy. BMC Genomics. 2017;18:140.
8. Kivisild T.: The study of human Y chromosome variation through ancient DNA. Human Genetics. 2017;136:529-546.
9. Lao O., Vallone P.M., Coble M.D., et al.: Evaluating self-declared ancestry of U.S. Americans with autosomal, Y-chromosomal and mitochondrial DNA. Human Mutation. 2010;31:E1875-E1893.