



**HAL**  
open science

**Predicting haplogroups using a versatile machine learning program (PredYMaLe) on a new mutationally balanced 32 Y-STR multiplex (CombYplex): unlocking the full potential of the human STR mutation rate spectrum to estimate forensic parameters**

Caroline Bouakaze, Franklin Delehelle, Nancy Sáenz-Oyhéréguay, Andreia Moreira, Stéphanie Schiavinato, Myriam Croze, Solène Delon, Cesar Fortes-Lima, Morgane Gibert, Louis Bujan, et al.

► **To cite this version:**

Caroline Bouakaze, Franklin Delehelle, Nancy Sáenz-Oyhéréguay, Andreia Moreira, Stéphanie Schiavinato, et al.. Predicting haplogroups using a versatile machine learning program (PredYMaLe) on a new mutationally balanced 32 Y-STR multiplex (CombYplex): unlocking the full potential of the human STR mutation rate spectrum to estimate forensic parameters. *Forensic Science International: Genetics* , 2020, 48, pp.102342. 10.1016/j.fsigen.2020.102342 . hal-02906055

**HAL Id: hal-02906055**

**<https://hal.science/hal-02906055>**

Submitted on 6 Nov 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Journal Pre-proof

Predicting haplogroups using a versatile machine learning program (PredYMaLe) on a new mutationally balanced 32 Y-STR multiplex (CombYplex): unlocking the full potential of the human STR mutation rate spectrum to estimate forensic parameters

Caroline Bouakaze, Franklin Delehelle, Nancy Saenz-Oyh  r  guy, St  phanie Schiavinato, Andreia Moreira, Myriam Croze, Sol  ne Delon, Cesar Fortes-Lima, Morgane Gibert, Louis Bujan, Eric Huyghe, Gil Bellis, Rosario Calderon, Candela Lucia Hern  ndez, Efr  n Avenda  o-Tamayo, Gabriel Bedoya, Antonio Salas, St  phane Mazi  res, Jacques Charioni, Florence Migot-Nabias, Andres Ruiz-Linares, Jean-Michel Dugoujon, Catherine Th  ves, Catherine Mollereau-Manaute, Camille No  s, Nicolas Poulet, Turi King, Maria Eugenia D'Amato, Patricia Balaesque

PII: S1872-4973(20)30115-0

DOI: <https://doi.org/10.1016/j.fsigen.2020.102342>

Reference: FSIGEN 102342

To appear in: *Forensic Science International: Genetics*

Received Date: 19 December 2019

Revised Date: 10 June 2020

Accepted Date: 11 June 2020

Please cite this article as: Bouakaze C, Delehelle F, Saenz-Oyh  r  guy N, Schiavinato S, Moreira A, Croze M, Delon S, Fortes-Lima C, Gibert M, Bujan L, Huyghe E, Bellis G, Calderon R, Hern  ndez CL, Avenda  o-Tamayo E, Bedoya G, Salas A, Mazi  res S, Charioni J, Migot-Nabias F, Ruiz-Linares A, Dugoujon J-Michel, Th  ves C, Mollereau-Manaute C, No  s C, Poulet N, King T, D'Amato ME, Balaesque P, Predicting haplogroups using a versatile machine learning program (PredYMaLe) on a new mutationally balanced 32 Y-STR multiplex (CombYplex): unlocking the full potential of the human STR mutation rate spectrum to

estimate forensic parameters, *Forensic Science International: Genetics* (2020),  
doi: <https://doi.org/10.1016/j.fsigen.2020.102342>

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2020 Published by Elsevier.

**Predicting haplogroups using a versatile machine learning program (PredYMaLe) on a new mutationally balanced 32 Y-STR multiplex (CombYplex): unlocking the full potential of the human STR mutation rate spectrum to estimate forensic parameters.**

Caroline Bouakaze<sup>1,16</sup>, Franklin Delehelle<sup>1\*</sup>, Nancy Saenz-Oyhéréguy<sup>1\*</sup>, Stéphanie Schiavinato<sup>1</sup>, Andreia Moreira<sup>1</sup>, Myriam Croze<sup>1, 17</sup>, Solène Delon<sup>1</sup>, Cesar Fortes-Lima<sup>1,18</sup>, Morgane Gibert<sup>1</sup>, Louis Bujan<sup>2</sup>, Eric Huyghe<sup>2</sup>, Gil Bellis<sup>3</sup>, Rosario Calderon<sup>4</sup>, Candela Lucia Hernández<sup>4</sup>, Efrén Avendaño-Tamayo<sup>5</sup>, Gabriel Bedoya<sup>6</sup>, Antonio Salas<sup>7</sup>, Stéphane Mazières<sup>8</sup>, Jacques Charioni<sup>8,9</sup>, Florence Migot-Nabias<sup>10</sup>, Andres Ruiz-Linares<sup>8,11</sup>, Jean-Michel Dugoujon<sup>1</sup>, Catherine Thèves<sup>1</sup>, Catherine Mollereau-Manaute<sup>1</sup>, Camille Noûs<sup>12</sup>, Nicolas Poulet<sup>13</sup>, Turi King<sup>14</sup>, Maria Eugenia D'Amato<sup>15</sup> and Patricia Balaresque<sup>1</sup>

\* These authors contributed equally to the work.

<sup>1</sup>Laboratoire d'Anthropologie Moléculaire et Imagerie de Synthèse (AMIS), UMR5288 – CNRS & Université Toulouse III, 37 allées Jules Guesde, 31073 Toulouse Cedex 3, France

<sup>2</sup>Equipe d'accueil EA3694, Hôpital Paule de Viguié, 330 Avenue de Grande Bretagne, TSA 70034, 31059 Toulouse Cedex 9, France

<sup>3</sup>INED Institut National d'Etudes Démographiques, 133 Boulevard Davout, 75980 Paris cedex 20, France

<sup>4</sup>Department of Biodiversity, Ecology and Evolution, Faculty of Biology, Complutense University. 28040 Madrid, Spain

<sup>5</sup>Grupo de Ciencias Básicas Aplicadas del Tecnológico de Antioquia, Tecnológico de Antioquia, Institución Universitaria, Medellín 050034, Colombia

<sup>6</sup>GENMOL (Genética Molecular), Instituto de Biología, Universidad de Antioquia Medellín Colombia

<sup>7</sup>Unidade de Xenética, Instituto de Ciencias Forenses (INCIFOR), Facultade de Medicina, Universidade de Santiago de Compostela, and GenPoB Research Group, Instituto de Investigaciones

Sanitarias (IDIS), Hospital Clínico Universitario de Santiago (SERGAS), Galicia, Spain

<sup>8</sup>Aix Marseille Univ, CNRS, EFS, ADES, Marseille, France

<sup>9</sup>Etablissement Français du Sang PACA Corse, Marseille, France

10. Université de Paris, MERIT, IRD, F-75006, Paris, France
11. Ministry of Education Key Laboratory of Contemporary Anthropology, School of Life Sciences, Fudan University, Shanghai, China
12. Laboratoire Cogitamus, CNRS & Université Toulouse III, 31000 Toulouse, France
13. Pôle écohydraulique AFB-IMT, allée du Pr Camille Soula, 31400 Toulouse, France
14. Department of Genetics, University of Leicester, Leicester, United Kingdom
7. Forensic DNA Laboratory, Department of Biotechnology, Faculty of Natural Sciences, University of Western Cape, Cape Town, South Africa
8. *Present address:* Institut National de Police Scientifique, Laboratoire de Police Scientifique de Lyon, 31 Avenue Franklin Roosevelt, 69134 Ecully Cedex, France
9. *Present address:* Division of EcoScience, Ewha Womans University, Seoul
10. *Present address:* Sub-department of Human Evolution, Department of Organismal Biology, Evolutionary Biology Centre, Uppsala University, Norbyvagen 18C, SE-752 36 Uppsala, Sweden

\*Corresponding author:

**Patricia Balaesque**

CNRS & University of Toulouse III (UMR 5288)

Laboratoire d'Anthropobiologie Moléculaire et Imagerie de Synthèse 37, allées Jules Guesde, 31073 TOULOUSE France

**Phone number:** + 33 (0) 5 61 14 55 04

**E-mail address:** patricia.balaesque@univ-tlse3.fr

## HIGHLIGHTS

- 32 Y-STR well-balanced mutation rate (CombYplex) and machine-learning program (PredYMaLe)
- Y-STR-based haplogroup prediction
- Best predictions using SVM and Random Forest classifiers
- Assignment accuracy scores (or prediction scores) using SVM: 97%
- Heterogeneous haplogroup predictions among classes
- Potential effects of small sample sizes and gene conversion

## ABSTRACT

We developed a new mutationally well-balanced 32 Y-STR multiplex (**CombYplex**) together with a machine learning (ML) program **PredYMaLe** to assess the impact of STR mutability on haplogroup prediction, while respecting forensic community criteria (high DC/HD). We designed CombYplex around two sub-panels M1 and M2 characterized by average and high-mutation STR panels. Using these two sub-panels, we tested how our program PredYmaLe reacts to mutability when considering basal branches and, moving down, terminal branches. We tested first the discrimination capacity of CombYplex on 996 human samples using various forensic and statistical parameters and showed that its resolution is sufficient to separate haplogroup classes. In parallel, PredYMaLe was designed and used to test whether a ML approach can predict haplogroup classes from Y-STR profiles. Applied to our kit, SVM and Random Forest classifiers perform very well (average 97%), better than Neural Network (average 91%) and Bayesian methods (<90%). We observe

heterogeneity in haplogroup assignment accuracy among classes, with most haplogroups having high prediction scores (99-100%) and two (E1b1b and G) having lower scores (67%). The small sample sizes of these classes explain the high tendency to misclassify the Y-profiles of these haplogroups; results were measurably improved as soon as more training data were added. We provide evidence that our ML approach is a robust method to accurately predict haplogroups when it is combined with a sufficient number of markers, well-balanced mutation rate Y-STR panels, and large ML training sets. Further research on confounding factors (such as gene conversion) and ideal STR panels in regard to the branches analysed can be developed to help classifiers further optimize prediction scores.

**Keywords:** Y-STR, machine learning, assignment accuracy and haplogroup prediction (hg prediction), incremental mutation rates

**Running title:** Y-chromosome study: combined use of a 32 Y-STRs multiplex and machine learning methods for haplogroup prediction

## INTRODUCTION

The Y-chromosome has been extensively used to identify male individuals in forensic communities (Kayser, 2017) and to reconstruct the family and evolutionary history of paternal lineages in geneticists (Jobling and Tyler-Smith, 2003) and genealogists communities (Calafell and Larmuseau, 2017). Questions related to the latter research topic are diverse and to address them on the Y-chromosome which is characterized by a low genetic diversity in human species, it can be advantageous to capture not only long-term but also short-term genomic information. It would help to optimally study not only the biogeographic informativeness of Y-haplotypes (Pardo-Seco et al., 2019) but also Y-specific migration paths and social structure, surname diffusion, paternal history of royal family members, and paternal lineage diffusion (Gill et al., 1994; Austerlitz and Heyer, 1998; King et al., 2014, 2007; Bowden et al., 2008; Chaix et al., 2008; Heyer et al., 2009, 2015; King and Jobling, 2009a, 2009b; Verdu et al., 2010; Martinez-Cadenas et al., 2016; Calafell and Larmuseau, 2017). But whatever the objectives and the technics used, the key problem remains the same: finding a good equilibrium between the resolution needs (markers and mutation rates) and the costs involved. Retrieving long-term genomic information has classically been completed using Y-SNaPshot analyses (for a review on Y-SNP typing see Sobrino, Brión and Carracedo, 2005), and very recently by using massively parallel sequencing (Ralf et al., 2019). Retrieving short-term genomic information has mainly consisted in Y-STR profiling in accessing the maximum of STRs variants and

polymorphism either by (i) designing Y-STR multiplexes including highly mutable markers to better discriminate closely related individuals (Purps et al., 2014; Gopinath et al., 2016) or (ii) by sequencing and extracting length-based Y-STR polymorphism STR loci from Next Generation Sequencing technologies as implemented in STRait Razor (Warshauer et al., 2013) to get rid of the excess of variants. To access short and long-term information while diminishing costs, some studies have chosen to generate high resolution Y-STR data and to use previously developed tools to predict haplogroup classes (Young et al., 2011; Mirabal et al., 2010, Šehović et al., 2017, Jannuzzi et al., 2020). Among these methods, Neural Network-based models (Felix Immanuel, 2013; <http://www.y-str.org/>) and Bayesian-allele frequency approaches (Athey 2006) were the first to have been developed, although ML approaches have been also tested (Schlecht et al., 2008) (see Supplementary data 1 for a review). However, the large bias in haplogroup prediction error (Jannuzzi et al., 2020) has urged the development of ready-to-use predictive tools, while considering more carefully the impact of STR mutation rates. The human Y-STR mutation rate spectrum is wide with a 1000 to 10000-fold of magnitude. Although this represents a powerful source of variation for designing tools in forensic genetics (from molecular to computational-based types), it is currently poorly explored.

## SUPPLEMENTARY DATA 1

In this paper, we assessed whether a well-balanced STR multiplex, associated with machine learning (ML) approaches can efficiently predict haplogroups, while still providing the high Discrimination Capacity (DC) index required in forensic genetics. We designed a 32 Y-STR-typing kit "CombYplex" around two panels of STRs (M1 and M2) mutating at various rates (selected from  $3.85 \times 10^{-04}$  to  $1.45 \times 10^{-02}$  mutation/locus/generation) to test the impact STR mutability on Hg prediction. Then, we designed "PredYMaLe" (Predicting Y-lineages using ML models), a program that includes various ML approaches to predict Y-haplogroup classes from a set of Y-STR markers.

First, for the CombYplex design, we assembled and typed a panel of 996 male individuals from three continents (Africa, Europe, and South America) available in our collections; we tested the discrimination power of CombYplex by computing both classic forensic and statistics parameters, e.g. Haplotype Diversity (HD), Discrimination Capacity (DC), Principal Component Analysis (PCA) and Linear Discriminant Analysis (LDA). Second, we tested whether the ML approaches implemented in the PredYMaLe program could efficiently predict the haplogroup lineages. We used a sub-panel of 503 chromosomes on four panels of STRs (the full 32-STR CombYplex, the Y-filer, and the CombYplex\_M1 and M2 only) for which haplogroup data were available. We evaluated the impact of STR-assembly on assignment accuracy, by considering first seven main Hg classes (considered as basal Y-tree branches) and then 12 detailed Hg classes (including E-subdivided terminal-like branches) to test the impact of Hg subdivision. Although not all haplogroup lineages could be tested in this article, the wide range of coalescence ages associated with the Hgs tested here (from 5 KYA for R1b1a1a2a1a2a1b1a1-M167 to 45 KYA for Hgs I-170 or J-M304 (Kivisild et al. 2017) should give a good preview of the prediction scores for comparable clades existing in the Y-tree and of the associated divergence between the relative

haplotypes. Our results showed that: (i) the full and well-balanced STR profiles (CombYplex or Y-filer) give the best prediction scores using the SVM and Random Forest classifiers, whereas Neural Network or Bayesian approaches, the most currently used methods for Hg prediction, fall short; (ii) PredYMaLe and CombYplex can predict haplogroup classes with an average assignment accuracy of 97% using Support Vector Machines (SVM) and Random forest classifiers, but classifiers are sensitive to STR panel composition, STR number, and training dataset size. These results can be used in the future to design well-balanced STR panels with a high number of markers, featuring high discrimination capacity and accurate predictions of haplogroup lineages with appropriate ML methods.

## MATERIALS AND METHODS

### 1. Database of Y-STR characteristics

For 220 Y-STRs, we collected information on Y-STR molecular characteristics, mutation rates, and polymorphisms for humans. This database is available in Supplementary data 2.

## SUPPLEMENTARY DATA 2

### 2. Selecting Y-STRs and constructing multiplexes: CombYplex M1 and M2

We selected a set of 32 Y-STRs from our database to construct two complementary multiplexes: one with average-mutating markers (M1) and one with high-mutating markers (M2). These markers were chosen to be polymorphic and to have the simplest molecular structure as possible (see Table 1). M1 includes the following **18 Y-STRs**: DYS485, DYS588, DYS502, DYS461, DYS638, DYS643, DYS587, DYS575, DYS578, DYS632, DYS508, DYS640, DYS511, DYS577, DYS556, DYS517, DYS565 and DYS538. Their mutation rates range from  $3.85 \times 10^{-04}$  to  $3.21 \times 10^{-03}$  mutation/locus/generation. Their molecular structures, primers and conditions are detailed in Table 1 and an example of a M1\_CombYplex profile is proposed in Figure 1a.

## TABLE 1, FIGURE 1a

M2 includes a **sex-testing assay** (derived from Cadamuro et al., 2015) and the following **14 Y-STRs**: Y-GATA-A10, DYS570, DYS549, DYS460, DYS442, DYS510, DYS541, DYS576, DYS513, DYS458, DYS481, DYS612, DYS444, and DYS533. These markers were chosen to be highly polymorphic and to have the simplest molecular structure

as possible; however, when STR with pure molecular structures could not be selected, we compromised between a simple structure and high STR mutation rate (e.g. DYS612 and DYS533). Their mutation rates range from  $3.32 \times 10^{-03}$  to  $1.45 \times 10^{-02}$  mutation/locus/generation. Their molecular structures, primers and conditions are detailed in Table 1 and an example of a M2\_CombYplex profile is proposed in Figure 1b.

TABLE 1, FIGURE 1b

The multiplexes were designed using the shortest amplicons as possible, with a maximum size of 356 bp for DYS533 (M2). They were designed to be used independently or combined, according to the degree of resolution required. The cost of a full CombYplex reaction (32 Y-linked STRs + three sex-typing markers) is only 4.3 € (in France and based on public prices for all the reagents), and one of the assets of this multiplex in regard to its resolution. This tool was developed on an ABI Prism 3730 DNA Analyzer 48-capillary array system (Life Technologies), due to contextual and logistic reasons, but its design strategy can be transposed to Next Generation Sequencing.

### 3. Population samples

Samples, available from collaborations and internal collections, were obtained from healthy human volunteers with consent forms. They were extracted from various substrates including saliva and whole blood. A total of **996** samples were used in this study (plus one male control, one female control and 1 *AZFc* deleted Y-chromosome male to control for deletion) and genotyped with the CombYplex kit. This dataset includes **six native West African** populations: three populations from Benin: 59 Bariba (Parakou region), 47 Yoruba (Ketou region), and 68 Fon (Cotonou and Ouidah regions), two populations from Ivory Coast:

47 Ahizi (Nigui-Saff region) and 37 Yacouba (Danané region), and one population from Mali: 13 Bwa (Segou region), **three native South African** populations (97 Xhosa, 90 Zulu, and 33 Tswana), **three admixed African-descendant** populations (52 French Guyana and Suriname Noir Marron, 56 Ketou-Yoruba, 35 Brazil - Rio de Janeiro, 20 Colombia), **one native American** population (6 Palikour), and **11 European populations** (30 Spain Barcelona, 19 Spain Galicia, 24 Spain Granada, 25 Spain Huelva, 46 France Loire-Atlantique, 50 France Vendée, 21 France Sarthe, 30 France Maine and Loire, 81 France Ariège-Pyrénées, and 57 France Haute-Garonne).

#### 4. Analysis of grouped samples

DNA samples were grouped based on two criteria: geographic ("GEO" sample) and phylogenetic ("HAPLO" sample).

In the "**GEO sample**" the geographic location of individuals is based on two generations of residence. All the 996 male individuals are included in this sample, to evaluate forensic parameters and control the discrimination power of the sample.

The "**HAPLO sample**", a haplogroup-based sample, is a subset of the GEO sample, used to evaluate haplogroup predictions with ML methods. It includes 503 individuals for whom Y-SNP haplogroup and Y-filer profiles were also available. Since many studies have already tested the added value of PPY23 and Y-filer plus, we did not type these these additional products due to the costs involved. We used Y-filer, a mutationally relatively balanced Y-STR kit for which we already had data in our database. We removed DYS385a/b and analysed only 15-STRs from the Y-filer panel since we have found evidence of conversion and outlier alleles in previous work (Balaesque et al., 2014). Eight main Hgs were first considered to calculate forensics parameters (E1a, E1E1a, E1b1b, F, G, I, J, R1b1a1a2). However, haplogroup classes represented by a very low number of individuals were not included in the

subsequent ML analyses (7 individuals in Hg F-M213\*/F-M89\*, and 2 individuals in Hg E1b1b1b1a-M81 included in E1b1b for 12-classes analyses): 7-Main and 12 detailed classes were considered in ML-analyses. Hg G and E1b1b had the lowest sample sizes, with 9 and 12 individuals respectively; we kept these Hgs in the 7-main classes to test the potential impact of a low number of individuals. The results for these two Hgs will have to be considered carefully due to the effect of small training sets reported in the ML literature.

First, the HAPLO sample was used to test the efficiency of CombYplex using classic forensics parameters (Haplotype diversity, Gene Diversity, Discrimination Capacity and Match Probability) and to test whether CombYplex could discriminate haplogroup classes using discriminant analyses (PCA). Second, it was used to test whether haplogroups could be predicted from the full 32 Y-STR, from the M1 and M2 only (lower number of markers and contrasted mutation rate), or from the Y-filer Y-STR profiles using an ML program. The HAPLO sub-sample includes six European populations (n=201; 26 Spain Barcelona, 14 Spain Galicia, 19 Spain Granada, 22 Spain Huelva, 64 France Pyrenees, 56 France Haute-Garonne), five native African populations (n=191; 52 Benin Parakou Bariba, 60 Benin Cotonou Fon, 36 Ivory-Coast Ahizi, 30 Ivory-Coast Yacouba, 13 Mali Bwa), and five admixed African-descendant populations (n=111; 8 French Guyana Aluku, 50 Ketou-Yoruba, 27 Noir-Marron, 12 Brazil-Rio de Janeiro, and 14 Colombia).

## 5. DNA extraction

The DNA extraction method was chosen according to the sample substrate. DNA was extracted from: (i) **whole blood**, using the QiaAmp DNA Blood mini-kit (Qiagen), (ii) **serum**, using the i-genomic DNA Blood mini-kit (Euromedex), and (iii) **saliva**, using the OG-300 Oragene DNA Self-Collection Kit (DNA Genotek) following the respective manufacturer's

instructions. The quantity and quality of DNA extracted was estimated using a NanoDrop Spectrophotometer 2000C (LabTech).

## **6. PCR amplification conditions: CombYplex M1 and M2**

CombYplex M1 and M2 were amplified in a reaction volume of 12.5  $\mu$ l with 6.25  $\mu$ l of QIAGEN Multiplex PCR Plus Kit (Qiagen), 1.25  $\mu$ l Q-solution (Qiagen), 4  $\mu$ l of the CombYplex M1 or M2 primer mix (see Tables 1a and 1b for concentrations) and 5 ng of DNA template (limit of detection tested: 2-2.5 ng). Thermal cycling was conducted on a GeneAmp PCR System 2700 (Applied Biosystems) using the following conditions: 95°C for 5 min; 30 cycles: 95°C for 30 sec, 62°C for 90 sec, 72°C for 30 sec; 68°C for 30 min, 10°C hold. To ensure that the resultant PCR amplicons were A-tailed (thereby avoiding the split peak phenomenon when visualized), a 2  $\mu$ l reaction mix incorporating 0.125 U Taq polymerase (Fisher BioReagents) and a 1X PCR buffer system was added to 5  $\mu$ l of PCR products prior to incubation for a further 45 min at 72°C.

## **7. Detection and analysis of PCR products**

Diluted A-tailed PCR products were mixed to 8.8  $\mu$ l Hi-Di™ formamide (Applied Biosystems) and 0.2  $\mu$ l GS600LIZ Size Standard (Applied Biosystems). After incubation at 95°C for 5 min, the samples were loaded onto an ABI Prism 3730 and a 3500 DNA Analyzer 48-capillary array system (Applied Biosystems). The G5 matrix filter DS-33 was used to detect the five dyes 6-FAM™ (blue), VIC™ (green), NED™ (yellow), PET™ (red) and LIZ™ (orange). The samples were injected for 15 sec at 1,600 V. Separations were performed at 15,000 V for 30 min with a run temperature of 63°C using the POP™-7 Polymer for 3730 (Applied Biosystems), run through a 50 cm capillary array (Applied Biosystems). Following data collection, samples were analysed with GeneMapper v4.0 (Applied Biosystems).

## 8. SNP genotyping methods

The populations Fon, Bariba, Yoruba, Ahizi and Yacouba were genotyped using 96 Y-SNPs on a BioMark™ HD system (Fluidigm, USA) as described in (Fortes-Lima et al., 2015). Y-SNP haplogroups were assigned according to ISOGG Y-DNA Haplogroup Tree 2015 (<http://www.isogg.org/tree/>) updated in April 2015. All other populations were genotyped using classic SNaPshot technics using a hierarchical approach. In total, 14 haplogroup lineages were detected and grouped in 7-Main and 12-Detailed classes for ML-analyses (Supp Data 6); they were used in combination with 4 sets Y-STR profiles (full CombYplex, Y-filer, CombYplex\_M1 and CombYplex\_M2) in PredYmale program to calculate how accurately a haplogroup lineage could be assigned.

## 9. Sequencing

Each locus was sequenced for the Male 1 control sample. Primers for sequencing are reported in Supplementary data 3. Each PCR product was sequenced in two reactions using forward and reverse PCR primers. The sequence reaction was performed with the BigDye Terminator v3.1 Cycle Sequencing Kit (Applied Biosystems). Sequence products were run on an ABI 3730 DNA Analyzer (Applied Biosystems). Sequences were analysed using Sequence Scanner Software v1.0 (Applied Biosystems) and BioEdit Sequence Alignment Editor version 7.2.5.

### SUPPLEMENTARY DATA 3

## 10. Forensic parameters and discrimination indexes: population grouping and comparative analyses

For GEO and HAPLO grouped samples, the following diversity parameters were calculated: haplotype diversity (HD) was calculated using  $HD = \frac{n}{n-1} (1 - \sum xi^2)$ , where  $n$  = the number of haplotypes in the dataset and  $xi$  = the frequency of the  $i$ th haplotype (Nei et al., 1981), gene diversity (GD) was calculated analogously to HD where  $n$  and  $xi$  denote the total number of samples and the relative frequency of the  $i$ th allele (Nei, 1973), discrimination capacity (DC) was defined as the ratio between the number of different haplotypes and the total number of haplotypes:  $DC = \frac{N_{diff}}{N}$  where  $N_{diff}$  was the number of different haplotypes,  $N$  was the sample size, and match probability (MP) was calculated as the sum of squared haplotype frequencies  $MP = \sum pi^2$  where  $pi$  was the frequency of the  $i$ th haplotype. Haplotype number ( $n$ ) and haplotype frequencies were estimated using Arlequin v 3.5.2.2 (Excoffier and Lischer, 2010). We represented the distribution of Y-STR haplotypes according to their haplogroup class by PCA: analyses were carried out using R software v 2.15.3 (R CoreTeam, 2017) and ade4 packages (Dray and Dufour, 2007). In addition, we performed Linear Discriminant Analyse (LDA) using the MASS package (Venables and Ripley, 2002) to estimate the proportion of haplogroups that were classed to a satisfactory precision. For LDA analysis, about 75% of individuals per haplogroup class taken randomly are used to train the model, while the remaining 25% is used to validate the trained classifier by testing its efficiency. This procedure was run 100 times. Given that the ML training and the split between training and validation datasets are heuristic, all the scores are averaged over 100 trials. We tested haplogroup prediction on the most represented haplogroup classes in our sample: E1a, E1a1a, E1a1b, G, I, J, and R1a1a1 (and on the collapsed root E group, including E1a, E1a1a and E1a1b).

## 11. Predicting haplogroups using machine-learning approaches: PredYMaLe

Haplogroups are usually defined by a given set of SNPs, but here, we explore whether they could also be recovered from the phylogenetic information contained in the Y-STR haplotype profiles alone. Different methods have been developed to predict haplogroups based on STRs, such as the Bayesian-based haplogroup predictor (<http://www.hprg.com/hapest5/index.html>) or Nevgen (<https://www.nevgen.org>), but neither of these is based on generalized ML models such as those proposed here. Here, similarly to the work of Schlecht et al., (2008), albeit with a higher resolution, we developed a generalist ML-based approach to the problem of haplogroup assignment from Y-STR profiles, then applied it to the particular case of the CombYplex profiles. We also assessed whether it performs better than the more common linear discriminant analysis.

We ran a pre-pilot study to test the efficiency of seven ML models (detailed in Bishop, 2006) so the fittest ML models could be implemented in PredYmale (details of pre-pilot study in Supplementary Data 4). Three models were eventually selected: Support Vector Machines (SVM), Random Forest Classifiers and k-Nearest Neighbors (kNN). These models follow the same concept: they build a classifier (a function) that maps a point in the problem space (here, a sample defined by its repeat counts for a given set of STRs) to a given class (here, a haplogroup). It should be noted that naive Bayes classifiers, a common method to address the problem of linking a set of STR markers to a haplogroup, and tested in a pilot run, have been constantly outperformed by SVMs and Random Forest Classifiers.

#### SUPPLEMENTARY DATA 4

**Support Vector Machines** (SVM) are classifiers that linearly partition the problem space by determining the frontier of the hyperplane maximizing its distance to the training samples (Cortes and Vapnik, 1995). Although SVMs were originally designed to discriminate between only two classes, they can be used in multi-class classification problematics (Chih-Wei Hsu

and Chih-Jen Lin, 2002), the problem being then divided in as many one-versus-all sub-problems as there are classes, which are solved independently. These partial classifiers are then merged to define the final classifier. Concretely, each sample in the training set is represented in the problem space by a point whose coordinates are the number of repetitions for each STR. Samples with close characteristics will cluster together. The SVM will determine a set of hyperplanes maximizing the margin between the classes. New points (i.e. unlabelled samples) are classified in either class depending on where they find themselves with regards to these hyperplanes.

**Random Forest Classifier** decision trees (Breiman, 1993) are linear classifiers that partition the problem space by defining a tree of binary conditions based on the features of a sample. Each new sample is then run through this tree of questions until it reaches a leaf, containing its predicted haplotype. Since a decision tree tends to over-fit the dataset it has been trained with, it might encounter difficulties generalizing when confronted with new samples. The random forest model (Ho, 1995) was developed to alleviate this limitation. At first, it trains multiple independent trees on several distinct subsets of the training data. Then, their outputs are averaged to define the final classifier. To improve the efficiency of random forests, we trained them with the AdaBoost boosting algorithm (Freund and Schapire, 1997). AdaBoost successively trains several copies of a base classifier (here a random forest) on the same dataset, and the training is adapted over generations to force the classifier to focus on hard to classify samples. Finally, all the generated classifiers over the generations are weighted according to their performances and combined to produce the final classifier. In our case, the learning process generates a decision tree defining questions on the number of repetitions of each STR. Depending on the answer, the sample to be classified will fall in one of the haplogroups. A notable advantage of this method is that its architecture (a sequence of questions) is easy for a human to understand, making the classification process transparent.

The **k-nearest neighbour** algorithm (also known as k-NN) is a non-parametric classification method. To produce a prediction for an unlabelled point, the algorithm combines the labels of the  $k$  closest points from the learning dataset according to a voting system. There are many ways to adapt the algorithm to the problem at hand, for instance by choosing the distance used, by applying a preliminary dimension reduction, by weighting the votes and so on. An advantage of the k-NN is that its error rate in a multi-class classification problem is proved to be bounded as an expression of the Bayes error rate, giving it a solid theoretical ground.

### **Implementing PredYMaLe**

We developed PredYMaLe (Predicting Y-lineages using machine learning models), a graphical interface to our automatic labelling solution, available at <https://gitlab.com/delehef/predymale/>. It is implemented in Python using the scikit-learn machine learning library and the Qt5 GUI library, and is available for GNU/Linux, macOS and Windows. PredYMaLe can be used on any Y-STR dataset where every sample is represented as a set of numerical repeat values (*e.g.* CombYplex, PPY23, etc.). Empty or null values are deliberately not supported in PredYMaLe: to avoid biases stemming from an imperfect dataset, we advise users to remove or insightfully fix erroneous profiles. The predicted labels can be exported to a CSV file for easy interoperability with other programs.

### **Procedure**

We tested whether haplogroups could be predicted using the three selected ML models implemented in the PredYMaLe program, and the three different Y-STR profiles (CombYplex full, CombYplex\_M2, and Y-filer). Each model was trained and evaluated using the HAPLO dataset (503 individuals, 7 Main and 12 Detailed Hg classes considered, 19 populations) and

according to the same protocol. The dataset was normalized in the [0; 1] range to avoid numerical discrepancies influencing the final result. Similar to the LDA analyses, 75% of the samples were used to train the model, while the remaining 25% were used to evaluate the trained classifier by testing its efficiency. Given that ML training, as well as the split between training and validation datasets are heuristic, all the scores are averaged over 100 trials. This also alleviates score outliers and offers a better interpretation of the performances of multiple models on real datasets. For that purpose, we performed two runs of analyses: for the first run, individuals were considered to belong to one of seven major haplogroup classes (E1a-M33, E1b1a-M2\*, E1b1b-M215, G-M201, I-M170, J-M304, and R1a1a1-M269 called *MainHg*), and for the second run, to one of twelve more detailed haplogroup classes (E1a-M33, E1b1a1-M2\*, E1b1a7-M191, E1b1a7a-U174, E1b1a8a-U209, E1b1a8a1-U290, E1b1b1-M35\*, G-M201, I-M170, J-M304, R1b1a1a2-M269, and R1b1a1a2a1a2a1b1a1-M167 called *DetailedHg*). The poorly represented haplogroup classes (e.g. F-189, and E1b1b1b1a-M81) could not be included in the procedure.

**Validation:** The evaluation process gives a score to a model, reflecting the efficiency of its predictions. We used the standard success score defined as  $s = n_C / n_T$ , where  $n_C$  is the number of successfully labelled validation samples and  $n_T$  the total count of validation samples. One success rating noted ‘score’ considers prediction as correct only if the predicted label of the validation sample matches the expected one.

## RESULTS

### 1. CombYplex: from polymorphisms to discrimination power

The CombYplex polymorphism was assessed based on 996 samples. All CombYplex profiles are available in Supplementary data 5. As expected, we observed an increasing level of polymorphisms from the less discriminative set of M1 markers (mean allele number: 6; Table 1) to the most discriminative M2 set (mean allele number: 9; Table 1). Forensic parameters were calculated for the GEO and HAPLO sample groups defined above (Table 2). GD and HD were greater than 0.999 for all GEO and HAPLO sub-groups using full CombYplex profiles. As expected, when M1 and M2 were analysed independently, M2 was always more discriminant than M1, with MP values oscillating from 0.001 (all populations) to 0.003 (Europe) using the GEO sample, and from 0.007 (Hg R) to 0.14 (Hg F) using the HAPLO sample. Indexes of discrimination capacity and match probability were observed in line with these values.

#### SUPPLEMENTARY DATA 5

#### TABLE 1, TABLE 2

### *Inter-haplogroup comparative analyses: PCA and LDA*

We tested whether CombYplex and Y-filer profiles could easily discriminate between haplogroup classes using the HAPLO sample (Supplementary 6). For this aim, we performed a PCA with seven haplogroup classes (MainHg) and a LDA (Table 3). PCA results based on CombYplex showed that haplogroup classes are well-discriminated along the two first axes (Figure 2a, especially R1b1a1 and E1a1a), but also along the second and third axes (Figure 2b, G, and I). LDA scores reach 94% in average, and oscillate from excellent (100 for E1a-

M33, E1b1a-M2, G-201, J-M304, R1b1a1a2-M269), to very good (95 for I-M170), correct for one of the less represented classes (62% for E1b1b-M35\*).

#### SUPPLEMENTARY DATA 6

##### FIGURE 2a,b

In comparison, discrimination of haplogroup classes appears less efficient using Y-filer profiles, both on F1xF2 and the F2xF3 axes (Figure 3a, b) but also using LDA (81% on average).

##### FIGURE 3a,b

These results provide evidence of the high resolving power of the 32 Y-STR CombYplex profile, not only for investigating paternal lineages but also for discriminating among haplogroups. Based on these encouraging results, we assessed whether haplogroup classes can be predicted using an ML approach based on CombYplex, Y-filer, CombYplex\_M1 Y-STR and CombYplex\_M2 Y-STR profiles.

## 2. Haplogroup prediction (HP) using Y-STR profiles and PredYMaLe program

We tested whether haplogroup classes can be predicted using an ML-based approach on CombYplex, Y-filer, CombYplex\_M1 Y-STR and CombYplex\_M2 Y-STR profiles. Results from the first run (seven major haplogroup classes (E1a-M33, E1b1a-M2\*, E1b1b-M215, G-M201, I-M170, J-M304, and R1a1a1-M269 called *MainHg*) were very informative on the three methods and the four datasets tested. Although HP scores using SVM and Random Forest are similar, SVM performed slightly better than Random Forest (Table 3); on average,

these two methods gave much better results than kNN: Random Forest/SVM HP average 3 methods 90-97%; kNN HP average 3 methods: 52-73%; Table 3).

TABLE 3

Compared to classic LDA (73 - 94%), SVM and Random Forest models perform systematically better, whatever the STR dataset, and especially using CombYplex. This results illustrates the combined impact of the marker number and the mutation rate range chosen on assignment accuracy. However, LDA performs better than kNN also for the three methods tested here. From the four STR datasets tested, we noted a noticeable performance of CombYplex (SVM: 97%) compared with M1 (SVM: 96%) and Y-filer (SVM: 95%), the M2 subset being systematically declassified (SVM: 92%, RF: 90%, kNN: 52%); when all E classes are collapsed, HP scores are very high (SVM et RF: 96-100%). A strong heterogeneity in HP scores is observed between haplogroup classes, even when the best method (SVM) is considered with the best STR combination (CombYplex): the G (67%) and E1b1b (67%) branches give the lowest HP scores compare to all others branches (100%). These two haplogroup classes represent the least represented ones (respectively N=9, 12), thus, suggesting the strong influence of sample size on the efficiency of HP. By analyzing confusion matrices for the best combination SVM/CombYplex and the worst combination kNN/M2, we observe clear differences in misclassification profiles (Figure 4): for the best combination, only two misclassifications are observed: E1b1a for E1b1b, and R for G. In contrast, 5 miss-targeted classifications are observed for kNN/M2, illustrating the incapability of this model/STR panel to associate an STR profile to a defined haplogroup class (especially for G: 0% HP).

FIGURE 4

No classifier exhibits a particularly skewed behavior regarding either of the metrics; all of them, on both datasets, follow the same pattern: F1-score and markedness stay close, while the informedness tends to score lower, denoting conservative classifiers. Therefore, defining the best classifier as the one with the best overall scores is straightforward. For a more detailed insight, Supplementary Data 7 (Supp tables 7a-7c) contain the per-class, per-dataset and per-classifier precisions, recalls, F1-scores, informednesses and markednesses.

Supplementary Data 7 (Supp table 7a à 7c)

The second run aimed to test the impact of sub-branch on haplogroup assignation accuracy score. We used a maximum resolution by considering the 12 most represented haplogroup branches (E1a-M33, E1b1a1-M2\*, E1b1a7-M191, E1b1a7a-U174, E1b1a8a\*-U209, E1b1a8a1-U290, E1b1b1-M35\*, G-M201, E1b1b1b1a, I J, R1b1a1a2-M269, R1b1a1a2a1a2a1b1a1-M167) and the two best models selected from the first run: SVM and Random Forest (Table 4). Per-class, per-dataset and per-classifier precisions, recalls, F1-scores, informednesses and markednesses are given in Supplementary Data 7 (Supp tables 7d-7f).

TABLE 4

Supplementary Data 7 (Supp tables 7d à 7f)

The average HP scores are high for both models and the four datasets, but they are lower than those from the first run, probably due to the smaller sample sizes and the close genetic affinity

of the different classes. Better prediction performances are observed for Random Forest, all STR datasets considered, with the highest average HP score obtained for CombYplex. The lowest scores are observed for M2 with an average HP score of 71% for Random Forest; this Y-STR dataset also has higher heterogeneity in HP scores between classes (from 27% for E1b1a1 to 100% E1a-M33; Table 4). By analyzing the confusion matrices for the best combination (Random Forest/CombYplex) and the worse (SVM/M2), we noticed that misclassification profiles are different (Figure 5). For Random Forest/CombYplex, misclassifications occur mainly across phylogenetically neighbors E1b1a and R1b1a1a2 branches. In contrast, for SVM/M2, misclassifications are associated with very diverse branches on the whole Y-chromosome phylogenetic tree (e.g. hg G), reflecting the impact either of highly mutating markers, the lower number of STR loci in this panel or the lack of association between STR profile and Y-haplogroup due to the impact of additional molecular mechanism as gene conversion.

FIGURE 5

## DISCUSSION

In this paper, we assess whether a panel of well-balanced Y-STR mutations, built around two sub-STR panels (from  $3.85 \times 10^{-04}$  to  $1.45 \times 10^{-02}$  mutation/locus/generation), associated with machine learning (ML) approaches can efficiently predict haplogroups. We developed the 32 Y-STR panel "CombYplex" and genotyped it on 996 male individuals from three continents (West and South Africa, West Europe, South America) to explore and confirm the discrimination capacity of the full, M1 and M2 panels, using classing forensic and statistics parameters. Then, we developed the ML approach PredYMaLe (Predicting Y-lineages using

ML models) and tested it on an assembled panel of 503 individuals, for which Hg and Y-filer information were also available in our database allowing a direct comparison of Y-STR assemblies.

### **STR panels and ML classifiers: an ideal association?**

We have demonstrated noticeable differences in prediction scores between STR panels and ML methods. Among all ML classifiers, SVM and Random Forests give better and more homogeneous prediction scores (90-97%) compared with kNN (52-97%) for this dataset, independently of the panels analysed.

When performing basal branch analyses (7-classes), mutationally well-balanced panels (CombYplex, Y-filer) and mutationally average panels (M1) performed better than the M2 panel, which was systematically outperformed. This result suggests that mutationally well-balanced or average STR panels should be preferred when analysing basal branches. The lower performance of M2 could imply either that assignment accuracy is affected by homoplasia using M2, due to the high mutation rate of the panel, or by the low number of STRs analysed (14 STRs). The latter argument is less probable since the 15 selected STRs of the Y-filer profiles gave better results.

When moving toward terminal branches (12-classes), mutationally well-balanced STR panels (CombYplex, Y-filer) performed better than M1 and M2 panels. M1 composed solely of average mutating STRs (18 STRs) were less performant due to its lack of discrimination power, giving equivalent results to M2 with four additional STR loci. Assignment accuracies for M1 and M2 decrease for the less represented classes, reflecting the need for the largest training set possible, and also a well-balanced STR panel with a sufficient number of STR loci when exploring closely related phylogenetic branches.

### **Variation in performance accuracies across Hg classes**

We showed that some haplogroups (e.g. E1a, I, J) have very distinct and unambiguous Y-STR profiles leading to 100% assignment accuracy scores, while others haplogroups (e.g. G, E1b1b) are more prone to misclassification within the STR panels and datasets analysed here. The impact of complexifying molecular mechanisms, such as gene conversion (Rozen et al., 2003), which potentially affect these profiles cannot be excluded (Balaresque et al., 2014) and could be further investigated. However the consistently worst scores of misclassification for the G and E1b1b haplogroups is likely to be the simple consequence of their small sample size. If the low accuracy of less well represented classes is problematic, empirical trends suggest that results are instantly improved when more training data are available. By running PredYmale with 10 additional G profiles collected recently, we observed that the prediction accuracy score reaches 83%, illustrating that prediction accuracy is significantly improved when more training data are available. We encourage users to train and use PredYmale on their own datasets, to learn about the prediction scores expected for the part of the tree explored. Given that PredYmale computations are rather fast, users should not hesitate to use larger datasets, or to adapt their STR panels to attain the best prediction scores.

### **Using PredYMaLe with other STR panels**

Our results demonstrate the need to find a good equilibrium between the number of markers, their mutability and the sample size of the training set according to the tree structure

considered. When analysing basal branches, well-balanced STR panels or average mutating STR panels can be selected preferably with SVM or Random Forest classifiers to ensure higher prediction scores. The M1 panel, an average mutating STR panel, gives very good results. Since these STRs have generally simpler motifs or low repeat counts, they can be extracted from whole-genome sequencing data using pre-existing tools (STRait Razor, Warshauer et al., 2013) and used to predict basal branches.

When moving toward terminal branches, mutationally well-balanced STR panels associated with SVM or Random Forest classifiers can be selected. In both cases, a minimal number of markers ( $> 20$ -30 STRs) is required to guarantee the best prediction scores possible. In forensic genetics, two commercial kits are commonly used, PPY23 (Purps et al., 2014) and Y filer<sup>®</sup> Plus (Gopinath et al., 2016). We have briefly tested whether our program could be confidently used with these panels by running PredYmale on published data. Based on our previous conclusions, we have only included the most represented classes ( $N > 20$ ). We analysed 451 individuals from five basal branches (E1b1b, G, I, J, R) for PPY23 (Pamjav et al., 2017, Heraclides et al., 2017), and 282 individuals from four basal branches (G, I, J, R1) for Y filer<sup>®</sup> Plus (Lacerenza et al., 2017). The average prediction scores obtained with SVM and Random Forest reached 98.5% for PPY23 and 97% for Y-filer plus (equiv. sample for ComBYplex reaches 98.5%). These results confirm the high prediction scores obtained with the SVM and Random Forest classifiers, for the three mutationally well balanced panels, for basal branches and sufficiently large training sets.

### **Predicting Hg using ML approaches: SVM, random forest and nearest neighbours classifiers**

By developing an ML program (PredYMaLe), designed to predict haplogroups using any Y-STR profiles, we show that ML models, especially SVM and Random Forest, give much

better HP results compared to alternative ML methods, including Bayesian, or Neural Network-based models. Interestingly these two classifiers have been reported to perform quite well for many other biological data (Fernández-Delgado et al., 2014). An interesting observation resides in the large variance of scores depending on the algorithm used: naive Bayes methods giving the worst results, while SVMs reach excellent precisions. The low accuracy of naive Bayes-based methods, in this case, can be explained by the fact that these algorithms consider features independently, and so cannot capture the information contained in their covariance patterns. SVMs, on the other hand, by maximizing the margin between the training classes, typically give excellent results as long as first, the problem is linearly well separable, which seems to be the case in this study, and second, that there is no consequent overlap between the different classes. Were it not the case, one can apply the “kernel trick” (Aizerman et al., 1964), which uses Mercer’s theorem to computationally cheaply immerse the dataset in a much larger space, where classes that are not linearly separable in the original space might become linearly separable.

In conclusion, support vector machines, random forests and nearest-neighbors classifiers are interesting alternatives to Bayesian or Neural networks classifiers to predict Y-haplogroups. Future users should note that although we developed and mostly used PredYmale with datasets featuring Y-STR profiles sampled with the CombYplex kit, the underlying ML concepts in our tool can be used on any STR panel (using STR repetition counts). We encourage users to train and use PredYmale on their own datasets regardless of the typing method.

## ACKNOWLEDGMENTS

We thank all DNA donors and volunteers associated with the sampling sessions. We also warmly thanks Prof. Maria Cátira Bortolini for giving us access to Brazilian samples, to Prof. Antoine Gessain for the Guyanese Noir Marrons. This work was supported by a Maturation research grant (CB's post-doctoral position), Research and Post-graduate Teaching Pole (PRES), the University Toulouse III (11.007), the LABEX DRIIHM (Investing in a future programme, ANR-11-LABX-0010), the OHM Haut Vicdessos, the Spanish Ministry of Economy and Competitiveness's grants (CGL2010-15191/BOS and CGL2014-53985-R) and the National Research Foundation Grant IFR160623173836 (MED). This work was performed using HPC resources from CALMIP (grant P1434). FD was supported by a PhD studentship (INSA, France), AM by a PhD studentship (Ministry of research, French government), NS by La Estancia de Otoño HOCR Cia. Ltda. (grant number 201509), CLH by a Spanish's research contract. CFL was supported by the EUROTAST Marie Curie Initial Training Network (EU FP7/2007-2013, grant no. 290344) and the Sven and Lilly Lawski's Foundation (N2019-0040). Ethics approvals were obtained: from the Senate Research Committee of the University of the Western Cape for South African samples under (ethic number 15-4-97, DC-2011-1436), from the Ethics Committee of the Faculty of Health Sciences, University d'Abomey-Cavali, Benin for the Beninese samples (ethic number 07/T4/2015/CE/FSS/UAC, 30th October 2015), from the University Bioethics committee (Sede de Investigación Universitaria, SIU) for the Colombian samples (ethic number 09-12-225 form), from the research ethics committee of the Universidade Federal do Rio Grande do Sul (Resolution no. 98002/1998) for the Brazilian samples Brazilian Ethics Commission, CONEP ethic number 1333/2002). Other samples from Africa were collected in the 80s or before, and ethics approval were not requested at that time; however, all participants were

volunteers with the purpose of collaborating with scientific studies, gave oral consent for the collection, and the confidentiality of their personal information has been preserved, following Helsinki Declaration.

## REFERENCES

Aizerman, M. et al. (1964) 'Theoretical foundations of the potential function method in pattern recognition learning', *Automation and Remote Control*, 25, pp. 821-837.

Athey, T.W. (2006) 'Haplogroup Prediction from Y-STR Values Using a Bayesian-Allele-Frequency Approach', *Journal of Genetic Genealogy*, 2, pp. 34-39.

Austerlitz, F. and Heyer, E. (1998) 'Social transmission of reproductive behavior increases frequency of inherited disorders in a young-expanding population.', *Proceedings of the National Academy of Sciences of the United States of America*. The National Academy of Sciences, 95(25), pp. 15140–15144.

Balaresque, P. et al. (2008) 'Dynamic nature of the proximal AZFc region of the human Y chromosome: multiple independent deletion and duplication events revealed by microsatellite analysis.', *Human mutation*, 29(10), pp. 1171–80, doi: 10.1002/humu.20757.

Balaresque, P. et al. (2014) 'Gene conversion violates the stepwise mutation model for microsatellites in y-chromosomal palindromic repeats.', *Human mutation*, 35(5), pp. 609–17, doi: 10.1002/humu.22542.

Bishop, C. M. (2006) *Pattern recognition and machine learning*. Springer.

Bowden, G. R. et al. (2008) 'Excavating past population structures by surname-based sampling: the genetic legacy of the Vikings in northwest England.', *Molecular biology and evolution*, 25(2), pp. 301–9, doi: 10.1093/molbev/msm255.

Breiman, L. et al. (1984) *Classification and regression trees*. Chapman & Hall/CRC, p368.

Cadamuro, V. C. et al. (2015) 'Determined about sex: sex-testing in 45 primate species using a 2Y/1X sex-typing assay.', *Forensic science international: Genetics*, 14, pp. 96–107, doi: 10.1016/j.fsigen.2014.09.010.

Calafell, F. and Larmuseau, M. H. D. (2017) 'The Y chromosome as the most popular marker in genetic genealogy benefits interdisciplinary research', *Human Genetics*, 136(5), pp. 559–573, doi: 10.1007/s00439-016-1740-0.

Chaix, R. et al. (2008) 'Genetic traces of east-to-west human expansion waves in Eurasia', *American Journal of Physical Anthropology*, 136(3), pp. 309-317, doi: 10.1002/ajpa.20813.

Chih-Wei Hsu and Chih-Jen Lin (2002) 'A comparison of methods for multiclass support vector machines', *IEEE Transactions on Neural Networks*, 13(2), pp. 415–425, doi: 10.1109/72.991427.

Cortes, C. and Vapnik, V. (1995) 'Support-Vector Networks', *Machine Learning*. Kluwer Academic Publishers-Plenum Publishers, 20(3), pp. 273–297, doi: 10.1023/A:1022627411411.

Dray, S. and Dufour, A.-B. (2007) 'The ade4 Package: Implementing the Duality Diagram for Ecologists', *Journal of Statistical Software*, 22(4), pp. 1–20, doi: 10.18637/jss.v022.i04.

Excoffier, L. and Lischer, H. E. L. (2010) 'Arlequin suite ver 3.5: a new series of programs to perform population genetics analyses under Linux and Windows', *Molecular Ecology Resources*. John Wiley & Sons, Ltd (10.1111), 10(3), pp. 564–567, doi: 10.1111/j.1755-0998.2010.02847.

Fernández-Delgado, M. et al. (2014) 'Do we Need Hundreds of Classifiers to Solve Real World Classification Problems?', *Journal of Machine Learning Research*, 15, pp. 3133–3181.

Fortes-Lima, C. et al. (2015) 'Genetic population study of Y-chromosome markers in Benin and Ivory Coast ethnic groups.', *Forensic science international: Genetics*, 19, pp. 232–237, doi: 10.1016/j.fsigen.2015.07.021.

Freund, Y. and Schapire, R. E. (1997) 'A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting', *Journal of Computer and System Sciences*, Academic Press, 55(1), pp. 119–139, doi: 10.1006/JCSS.1997.1504.

Gill, P. et al. (1994) 'Identification of the remains of the romanov family by DNA analysis', *Nature Genetics*, 6(2), pp.130-135, doi: 10.1038/ng0294-130.

Gopinath, S. et al. (2016) 'Developmental validation of the Yfiler ® Plus PCR Amplification Kit: An enhanced Y-STR multiplex for casework and database applications', *Forensic Science International: Genetics*, 24, pp. 164–175, doi: 10.1016/j.fsigen.2016.07.006.

Hanson, E. et al. (2012) 'Performance evaluation and optimization of multiplex PCRs for the highly discriminating OSU 10-locus set Y-STRs.', *Journal of forensic sciences*, 57(1), pp. 52–9, doi: 10.1111/j.1556-4029.2011.01910.

Heraclides, A. et al. (2017) 'Y-chromosomal Analysis of Greek Cypriots Reveals a Primarily Common pre-Ottoman Paternal Ancestry With Turkish Cypriots', *PLoS One*, 12(6):e0179474, doi: 10.1371/journal.pone.0179474.

Heyer, E. et al. (2009) 'Genetic diversity and the emergence of ethnic groups in Central Asia.', *BMC Genetics*, 10 (49), pp. 1-8, doi: 10.1186/1471-2156-10-49.

Heyer, E. et al. (2015) 'Patrilineal populations show more male transmission of reproductive success than cognatic populations in Central Asia, which reduces their genetic diversity.', *American Journal of Physical Anthropology*, 157(4), pp. 537-543, doi: 10.1002/ajpa.22739.

Ho, T. K. (1995) 'Random Decision Forest', in *Proceedings of the 3rd International Conference on Document Analysis and Recognition*. Montréal, pp. 278–282.

Jannuzzi, J. et al. (2020) 'Male lineages in Brazilian populations and performance of haplogroup prediction tools', *Forensic Science International: Genetics.*, 44, pp. 1-7, doi: 10.1016/j.fsigen.2019.102163.

Jobling, M. A. and Tyler-Smith, C. (2003) 'The human Y chromosome: An evolutionary marker comes of age', *Nature Reviews Genetics*, 4, 598–612, doi: 10.1038/nrg1124.

Kayser, M. et al. (2004) 'A Comprehensive Survey of Human Y-Chromosomal Microsatellites', *The American Journal of Human Genetics*, 74(6), pp. 1183–1197, doi: 10.1086/421531.

Kayser, M. (2017) 'Forensic use of Y-chromosome DNA: a general overview', *Human Genetics*, 136(5), pp. 621–635, doi: 10.1007/s00439-017-1776-9.

King, T. E. et al. (2007) 'Thomas Jefferson's Y chromosome belongs to a rare European lineage', *American Journal of Physical Anthropology*, 132(4), pp. 584–589, doi: 10.1002/ajpa.20557.

King, T. E. et al. (2014) 'Identification of the remains of King Richard III', *Nature Communications*. Nature Publishing Group, 5: 5631, pp. 1-8, doi: 10.1038/ncomms6631.

King, T. E. and Jobling, M. A. (2009a) 'Founders, drift, and infidelity: The relationship between y chromosome diversity and patrilineal surnames', *Molecular Biology and Evolution*, 26(5), pp. 1093-1102, doi: 10.1093/molbev/msp022.

King, T. E. and Jobling, M. A. (2009b) 'What's in a name? Y chromosomes, surnames and the genetic genealogy revolution', *Trends in Genetics*, 25(8), pp. 351-360, doi: 10.1016/j.tig.2009.06.003.

Kivisild T. (2017) 'The study of human Y chromosome variation through ancient DNA', *Human Genetics*, 136, pp. 529–546, doi: 10.1007/s00439-017-1773-z.

Lacerenza, D.S. et al. (2017) 'Investigation of extended Y chromosome STR haplotypes in Sardinia.' *Forensic Science International: Genetics*, 27, pp. 172-174, doi: 10.1016/j.fsigen.2016.12.009.

Martinez-Cadenas, C. et al. (2016) 'The relationship between surname frequency and Y chromosome variation in Spain', *European Journal of Human Genetics*, 24(1), pp. 120–128, doi: 10.1038/ejhg.2015.75.

Mirabal et al. (2010) 'Human Y-Chromosome Short Tandem Repeats: A Tale of Acculturation and Migrations as Mechanisms for the Diffusion of Agriculture in the

Balkan Peninsula', *American Journal of Physical Anthropology*, 142, pp 380-390, doi:10.1002/ajpa.21235.

Nei, M. (1973) 'Analysis of Gene Diversity in Subdivided Populations', *Proceedings of the National Academy of Sciences*, 70(12), pp. 3321–3323, doi: 10.1073/pnas.70.12.3321.

Nei, M. et al. (1981) 'Polymorphism and evolution of the Rh blood groups', *The Japanese Journal of Human Genetics*, 26, 263–278, doi: 10.1007/BF01876357.

Pamjav, H., et al. (2017) 'A Study of the Bodrogeköz Population in North-Eastern Hungary by Y Chromosomal Haplotypes and Haplogroups.', 292(4), pp. 883-894, doi: 10.1007/s00438-017-1319-z.

Pardo-Seco, J. et al. (2019) 'Biogeographical informativeness of Y-STR haplotypes', *Science Bulletin. Elsevier*, 64(19), pp. 1381–1384, doi: 10.1016/J.SCIB.2019.07.025.

Parson, W. et al. (2016) 'Massively parallel sequencing of forensic STRs: Considerations of the DNA commission of the International Society for Forensic Genetics (ISFG) on minimal nomenclature requirements', *Forensic Science International: Genetics*, 22, pp. 54–63, doi: 10.1016/j.fsigen.2016.01.009.

Purps, J. et al. (2014) 'A global analysis of Y-chromosomal haplotype diversity for 23 STR loci', *Forensic Science International: Genetics*, 12, pp.12-23, doi: 10.1016/j.fsigen.2014.04.008.

R CoreTeam (2017) 'R: A language and environment for statistical computing'. Vienna, Austria: R Foundation for Statistical Computing.

Ralf, A. et al. (2019) 'Forensic Y-SNP analysis beyond SNaPshot: High-resolution Y-chromosomal haplogrouping from low quality and quantity DNA using Ion

AmpliSeq and targeted massively parallel sequencing', *Forensic Science International: Genetics*, 41, pp. 93-106, doi: 10.1016/j.fsigen.2019.04.001.

Rozen, S. et al. (2003) 'Abundant gene conversion between arms of palindromes in human and ape Y chromosomes', *Nature*, 423(6942), pp.873-876, doi: 10.1038/nature01723.

Schlecht, J. et al. (2008) 'Machine-learning approaches for classifying haplogroup from Y chromosome STR data', *PLoS Computational Biology*, 4(6): e1000093, doi:10.1371/journal.pcbi.1000093.

Sehović et al. (2017) 'Network analysis on the in silico assigned Y chromosome haplogroups in Western Balkan populations', *Genetics & Applications*, 1(2), pp. 36-43, doi: 10.31383/ga.vol1iss2pp36-43.

Sobrino, B., Brión, M. and Carracedo, A. (2005) 'SNPs in forensic genetics: A review on SNP typing methodologies', *Forensic Science International*, 154 (2-3), pp. 181-194, doi: 10.1016/j.forsciint.2004.10.020.

Venables, W. N. and Ripley, B. D. (2002) *Modern Applied Statistics with S*. New York, NY: Springer New York (Statistics and Computing), doi: 10.1007/978-0-387-21706-2.

Verdu, P. et al. (2010) 'Limited dispersal in mobile hunter-gatherer Baka Pygmies', *Biology Letters*, 6, pp. 858–861, doi: 10.1098/rsbl.2010.0192.

Warshauer, D. H. et al. (2013) 'STRait Razor: A length-based forensic STR allele-calling tool for use with second generation sequencing data', *Forensic Science International: Genetics*, 7(4):409-17, doi: 10.1016/j.fsigen.2013.04.005.

Young, K.L. et al. (2011) 'Paternal Genetic History of the Basque Population of Spain',  
Human Biology, 83(4), pp. 455-475.

Journal Pre-proof

**FIGURES Legends**

**Figure 1a.** CombYplex M1 profile of male control; two artifacts can occasionally be observed on the M1 electropherogram: in the polymorphism zone of the DYS588 locus (blue dye) and in the polymorphism zone of the DY508 locus (yellow dye), as shown here.

**Figure 1b.** CombYplex M2 profile

**Figure 2a.** PCA for CombYplex F1xF2

**Figure 2b.** PCA for CombYplex F2xF3

**Figure 3a.** PCA for Y-filer F1xF2

**Figure 3b.** PCA for Y-filer F2xF3

**Figure 4.** Confusion matrices for the first run on MainHg (7 haplogroup classes) for CombYplex/SVM and M2/k-Nearest Neighbors.

**Figure 5.** Confusion matrices for the second run on DetailedHg (12 haplogroup classes) for CombYplex/Random Forest/ and M2/SVM.





Fig 3

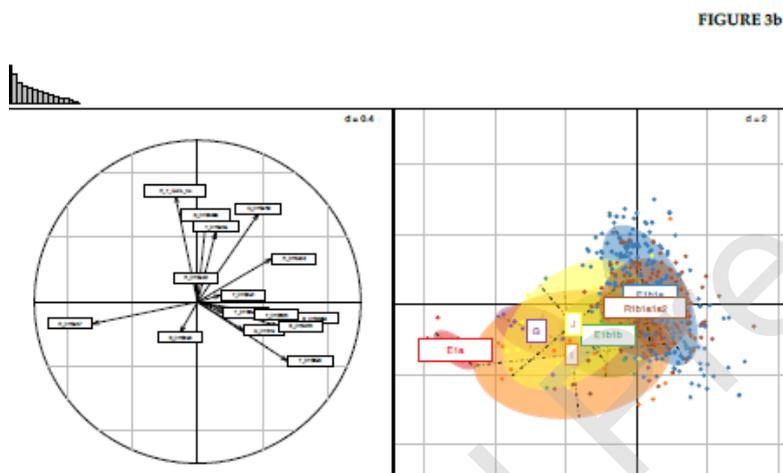
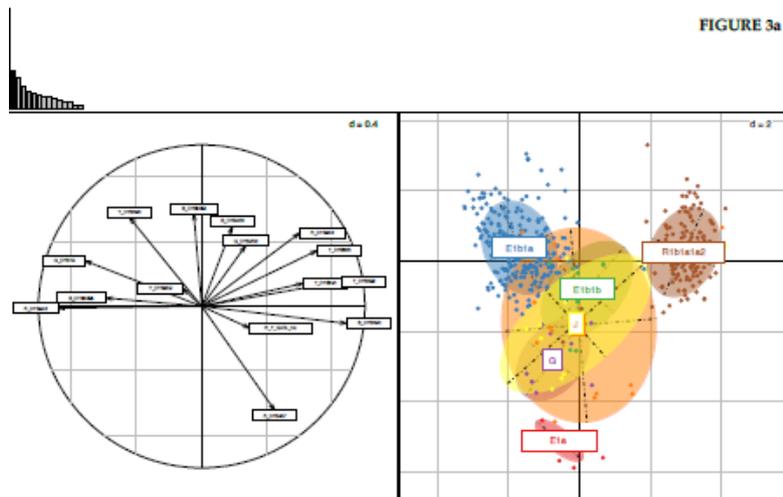


Fig 4

FIGURE 4

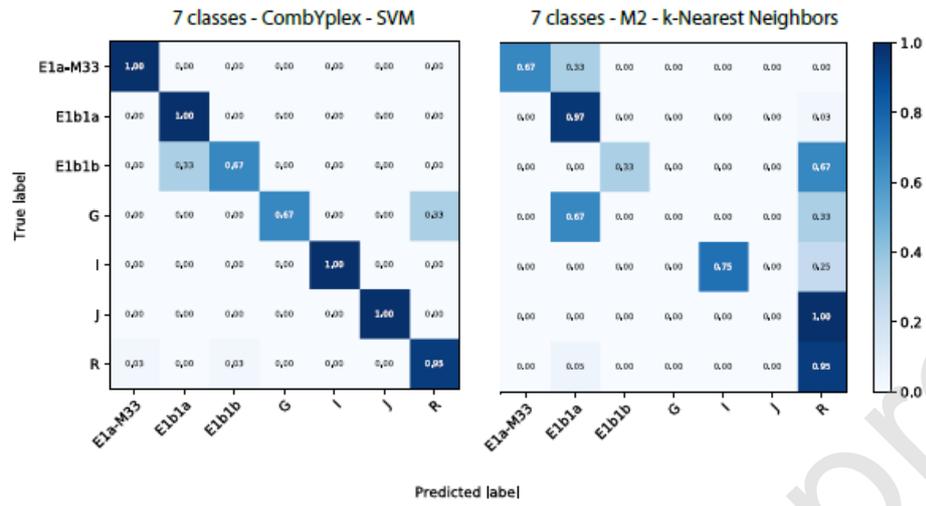
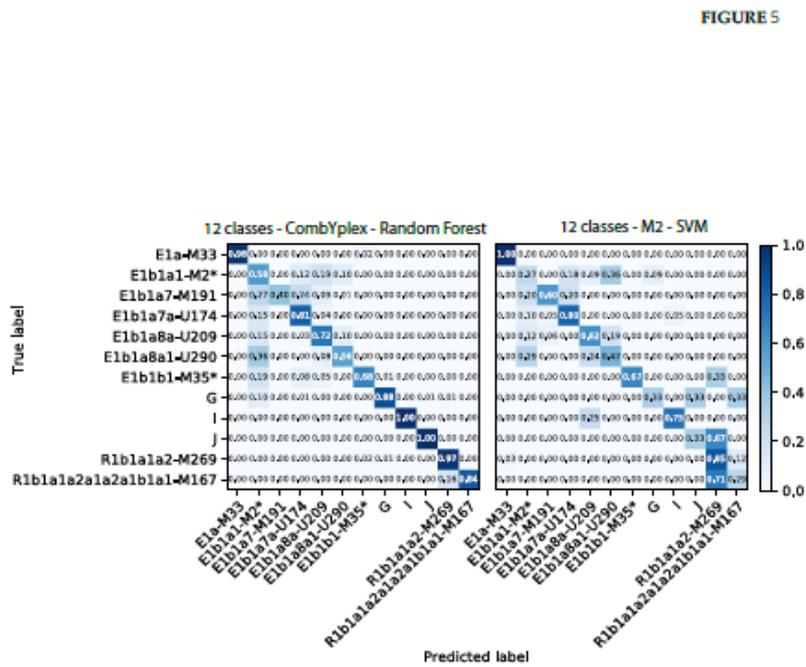


Fig 5



## Tables

**Table 1.** CombYplex M1 (a) et M2 (b): markers, molecular structures, primers and amplification conditions. \*Dyes: Blue: FAM; Green: VIC; Yellow: NED; Red: PET. Nomenclature is given according to the following papers: (Kayser et al., 2004; Gusmão et al., 2006; Parson et al., 2016 and the STRidER Reference database: <https://strider.online/>).

a) Com bYple x_M 1																			
M1 Markers	M uta tio n ra te	Repeat structure	Nom encl atur e used	D ye *	Primer F orward				Primer Reverse				According to litterature			Human ref. NC_0002 4.9 (UCSC)		Obs erve d	Male control 1
					Name	Sequence (5'-3')	T m ( ° C)	[ μ M ' ]	Name	Sequence (5'-3')	T m ( ° C)	[ μ M ' ]	Alle le R an ge	A lle le c ou nt	Ex pe cte d siz e ra nge (b	Alle le	A mp lic on (bp )	Alle le Ran ge	Alle le



			dERS TRs ENFS I Refer ence data base																
<b>DYS587</b>	2,6 2E- 03	(CAATA) <sub>n</sub> [(CA GTA)(CAATA)] 3	Gus mao et al. 2006	V I C	VIC_D YS587 _F2	cttctttgga aagtagcatt tcat	5 8 , 0, 8 1	DYS58 7_R2 queu e	aaagtctgacaa tgagaagggttt ctaagttcagg	6 8 , 0, 8 9	8- 1 6	9	19 1- 22 2	1 1 1	19 1	8-16	11	191	
<b>DYS575</b>	3,9 1E- 04	(AAAT) <sub>n</sub>	Gus mao et al. 2006	V I C	VIC_D YS575 _F2	cagaggttcc agtaagctta gatca	6 0 , 0, 3 3	DYS57 5_R2	cattgatgggctt taggtga	5 9 , 0, 3 9	8- 1 2	5	26 0- 27 6	1 0	26 8	8-12	10	268	
<b>DYS578</b>	9,9 5E- 04	(AAAT) <sub>n</sub>	Kays er et al. 2004	V I C	DYS57 8_F	gaggcggaa ctttcagtga g	6 0 , 0, 5 0	VIC_D YS578 _R2	cagaagtcacct gtgttcaa	6 0 , 0, 5 1	7- 1 0	4	30 5- 31 7	9	31 3	6-11	9	313	
<b>DYS632</b>	3,9 7E- 04	(CATT) <sub>n</sub>	Gus mao et al. 2006	N E D	NED_ DYS63 2_F	cacagtttca gtcttgcatt g	5 9 , 0, 0 3	DYS63 2_R	tctgggcaacag aaggagac	6 0 , 0, 1 4	8- 1 0	3	10 6- 11 4	9	11 0	7-11	9	110	
<b>DYS508</b>	3,0 3E- 03	(TATC) <sub>n</sub>	Gus mao et al. 2006	N E D	NED_ DYS50 8_F	acaatggca atcccaaatt c	5 9 , 0, 4 6	DYS50 8_R	gaacaaataag gtgggatggat	5 9 , 0, 4 1	8- 1 5	1 1	16 5- 19 3	1 1	17 7	8-15	11	167	
<b>DYS640</b>	3,9 8E-	(AAAT) <sub>n</sub>	Kays er et	N E	NED_ DYS64	ggaaaaacc atgagatcct	5 9 , 0, 2 9	DYS64 0_R	aagcccgttcat attttaaagac	5 7 , 0, 2 1	9- 1 5	5	25 2-	1 1	26 0	7-13	11	260	

	04		al. 2004	D	0_F	gtc	, 8			, 9	3		26 8							
<b>DYS511</b>	1,5 2E- 03	(GATA) <sub>n</sub>	Kays er et al. 2004	N E D	NED_ DYS51 1_F	tggggtgga tgtgtaggta ga	6 0, , 2	0, 3	DYS51 1_R	tctggttgtgcct tagatttga	5 9, , 7	0, 3	9- 1 4	6	30 7- 32 7	9	30 7	7-14	10	311
<b>DYS577</b>	4,1 1E- 04	(ATTC) <sub>n</sub>	Kays er et al. 2004	P E T	PET_D YS577 _F	tttttctacgt gtgtatccac taacc	5 9, , 8	0, 1 5	DYS57 7_R	gtgtccccagcc ctgtta	5 9, , 5	0, 1 5	8- 1 1	4	10 0- 11 2	9	10 4	6-12	9	106
<b>DYS556</b>	1,5 9E- 03	(AATA) <sub>n</sub>	Gus mao et al. 2006	P E T	PET_D YS556 _F	tcaccaatg acattttaca gca	5 9, , 1	0, 6	DYS55 6_R	ttggttagtgtaa tgcatccag	5 7, , 7	0, 6	8- 1 2	5	15 6- 17 2	1 1	16 8	8-13	11	156
<b>DYS517</b>	3,2 1E- 03	(AAAG) <sub>n</sub> N <sub>13</sub> (A AAG) <sub>3</sub>	Kays er et al. 2004	P E T	PET_D YS517 _F2	aactgacca gcaaaaatg ttaa	5 7, , 9	0, 5	DYS51 7_R2	tgtctgagacct acaagattgc	5 7, , 1	0, 5	1 0- 1 8	9	21 3- 24 5	1 5	23 3	9-18	13	154
<b>DYS565</b>	2,0 9E- 03	(ATAA) <sub>n</sub>	Kays er et al. 2004	P E T	PET_D YS565 _F2	ccaggaagc agtgttgcac	5 9, , 8	0, 3	DYS56 5_R2	gcagttctctgcc tgtatgg	5 8, , 5	0, 3	9- 1 4	6	28 0- 30 0	1 2	29 2	9-14	12	292
<b>DYS538</b>	3,9 4E- 04	(GATA) <sub>n</sub>	Kays er et al. 2004	P E T	PET_D YS538 _F	ttggggaaa acagatggt gt	6 0, , 2	1, 7	DYS53 8_R	ccaaataccat cataggaagaa	5 9, , 2	1, 7	9- 1 3	5	33 9- 35 5	1 0	34 3	8-13	10	343

M  
e  
a  
n  
6

b) Complex x_M2	M2 Markers	Mutation rate	Repeat structure	Nomenclature used	Dye*	Primer Forward			Primer Reverse			Acc. To literature			Human ref. NC_0002 4.9 (UCSC)		Observed	Male control 1
						Name	Sequence (5'-3')	Tm ( $^{\circ}$ C)	Name	Sequence (5'-3')	Tm ( $^{\circ}$ C)	Allele Range	Allele count	Expected size range (bp)	Allele	Amplification (bp)		
						SRY				FAM	FAM_SRY_F2	gcgaaactc agagatcag caag	60,081	SRY_R1	tgtgcctcctgg aagaatgg	61,089		
UTY				FAM	FAM_UTXU	cagtgttacc agccttaac	53,027	UTY_R	ggcaggcttact ttttagag	52,021					91			
UTX				FAM	FAM_UTXU_F1	ag		UTX_R	tctgtggactag gtttgtggt	55,011					120			



60	2E-03	(ATAG) <sub>n</sub>	mao et al. 2006, Pars on 2016 FSI Genetics	I C	YS460_F	tatcatttatt atgtat	7 , 1	4	0_R	gaatctgacacc	9 , 0	4	1 3		11 9	0	7			
DYS442	9,78E-03	(TATC) <sub>2</sub> (TGTC) <sub>3</sub> (TATC) <sub>n</sub>	Gusmao et al. 2006	V I C	VIC_D YS442_F2	tgcaaaatc acggaacca a	6 1	0, 1 5	DYS442_R2	caagccactgca aatgtca	5 9 , 4	0, 1 5	9- 1 6	8	17 3- 20 1	1 2	18 5	8-16	12	118
DYS510	5,99E-03	(GATA) <sub>3</sub> N <sub>12</sub> (GATA) <sub>n</sub> N <sub>13</sub> (GGAT) <sub>4</sub> N <sub>9</sub> (GATA) <sub>3</sub>	Kays et al. 2004; here only the last (GATA) <sub>n</sub> polymorphism is mea	V I C	VIC_D YS510_F	ttttcctccc ttaccacag a	5 8 , 7	0, 4 8	DYS510_R	tctggagaagac agaactgtca	5 9 , 1	0, 4 8	9- 1 5	7	24 5- 26 9	1 1	25 3	8-15	11	253

			sure d assu ming that the rest of the moti f is mon omo rphic .																	
<b>DYS541</b>	3,9 2E- 03	$(TATC)_n(TTC)_1(TATC)_3$	Kays er et al. 2004	V I C	VIC_D YS541 _F	catcattaat tctatctgttc atccat	5 8 , 8	0, 4 5	DYS54 1_R	tggataaagaac acctttaagaag c	5 9 , 3	0, 4 5	1 0- 1 5	6	31 0- 33 0	1 2	31 8	6-15	12	272
<b>DYS576</b>	1,4 3E- 02	$(AAAG)_n$	Gus mao et al. 2006 , STRi dERS TRs ENFS I Refer	N E D	NED_ DYS57 6_F	ccaagcaac atagcaaga cct	5 9 , 4	0, 1 3	DYS57 6_R	aagcgtattgtc ttggctttt	5 9 , 4	0, 1 3	1 3- 2 2	1 0	10 8- 14 4	1 7	12 4	13- 25	19	132

			ence data base																	
<b>DYS5 13</b>	6,0 9E- 03	(TATC) <sub>n</sub>	Gus mao et al. 2006	<b>N E D</b>	NED_ DYS51 3_F2	tgttgtaaaa atgactactg tggtatg	5 8 , 6	0, 2 2	DYS51 3_R2	ccacatcagcac tattacttaact a	5 8 , 9	0, 2 2	9- 1 5	7	29 4- 31 8	1 2	30 6	9-15	12	310
<b>DYS4 58</b>	8,3 6E- 03	(GAAA) <sub>n</sub>	Gus mao et al. 2006  STRi dERS TRs ENFS I Refer ence data base	<b>N E D</b>	NED_ DYS45 8_F	tgggtggtg gaggttactg t	6 0 , 3	0, 1 2	DYS45 8_R	cccaaagttctg gcattacaa	6 0 , 0	0, 1 2	1 1- 2 4	1 4	18 3- 23 5	1 6	20 3	11- 24	18	211
<b>DYS4 81</b>	4,9 7E- 03	(CTT) <sub>n</sub>	STRi dERS TRs ENFS I Refer ence data base	<b>P E T</b>	PET_D YS481 _F	aggaatgtg gctaacgct gt	5 9 , 8	0, 2	DYS48 1_R	accagaaggttg caagactca	5 9 , 9	0, 2	1 8- 3 2	1 5	10 9- 15 1	2 2	12 1	18- 32	22	121

DYS6 12	1,4 5E- 02	(CCT)n(CTT) <sub>1</sub> ( TCT) <sub>n</sub> (CCT) <sub>1</sub> (T CT) <sub>n</sub>	Adap ted from STRi dERS TRs ENFS I Refer ence data base ; here [CCT ] <sub>n</sub> =5 CTT [TCT] ] <sub>n</sub> =4 CCT [TCT] ] <sub>n</sub> : only the last (TCT) ] <sub>n</sub> poly mor	P E T	PET_D YS612 _F	ccccatgcc agtaagaat a	5 9 , 8	1, 2 5	DYS61 2_R	tgaggaaggc aaaagaaaa	5 9 , 8	1, 2 5	1 9- 3 1	1 3	18 6- 22 2	2 5	20 4	16- 31	26	207
------------	------------------	---	--	-------------	----------------------	-----------------------------	------------------	--------------	--------------	-------------------------	------------------	--------------	-------------------	--------	---------------------	--------	---------	-----------	----	-----

			phism is measured assuming that the rest of the motif is monomorphic.																
<b>DYS44</b>	5,4 5E-03	(ATAG) <sub>n</sub>	Gusmao et al. 2006	PET	PET_DYS444_F	catagaatg aaaggtgtg aacca	59,0 ,45 ,0	DYS44 4_R	tgccattcaaac tcacgttg	60,9 ,45 ,7	9-16	8	26 4-28 2	14	27 4	8-16	12	266	
<b>DYS533</b>	5,0 1E-03	(ATCT) <sub>n</sub>	Gusmao et al. 2006 STRIDERS TRs	PET	PET_DYS533_F	attcatctaa catctttgtc atctacc	58,0 ,95 ,2	DYS533 3_R	ttaacttgcccttt tgcaccc	59,0 ,95 ,2	9-14	6	33 4-35 4	12	34 6	8-14	12	346	



**Table 2.** Forensic parameter estimates for GEO and HAPLO samples for the full CombYplex, M1 and M2 and Y-filer. Parameters calculated: Genetic Diversity or Haplotype Diversity (GD/HD), Discrimination Capacity (DC), and Match Probability (MP).

**Table 2.**

Population (Geo sample)	CombYplex total					CombYplex M1				CombYplex M2							
	N	n	GD/HD	DC	MP	n	GD/HD	DC	MP	n	GD/HD	DC	MP				
All pop	996	916	0,9998	0,9196	0,0012	607	0,9964	0,6094	0,0053	889	0,9998	0,8926	0,0013				
South America : native (Palikur)	6	6	0,9999	1	0,1666	4	0,9630	0,6667	0,2778	6	0,9999	1	0,1667				
South America : admixed	107	96	0,9986	0,8972	0,0118	84	0,9921	0,7850	0,0197	92	0,9982	0,8598	0,0127				
Africa native	444	391	0,9995	0,8806	0,0029	242	0,9917	0,5450	0,0124	374	0,9994	0,8423	0,0033				
Africa admixed	56	52	0,9982	0,9286	0,0210	45	0,9953	0,8036	0,0268	52	0,9981	0,9286	0,0210				
Europe	383	368	0,9998	0,9608	0,0030	253	0,9916	0,6606	0,0123	364	0,9998	0,9504	0,0029				
Haplogroup (Haplo sample)	CombYplex total					CombYplex M1				CombYplex M2				Y-filer			
Total Hg	503	n	GD/HD	DC	MP	n	GD/HD	DC	MP	n	GD/HD	DC	MP	n	GD/HD	DC	MP
E1a	15	14	0,9956	0,9333	0,0756	12	0,9891	0,8000	0,0933	13	0,9919	0,8667	0,0844	10	0,8889	0,6667	0,2000
E1b1a	275	244	0,9992	0,8873	0,0049	192	0,9958	0,6982	0,0093	238	0,9989	0,8655	0,0053	228	0,9988	0,8291	0,0056
E1b1b	12	12	1	1	0,0833	11	0,9931	0,9166	0,0972	11	0,9931	0,9167	0,0972	10	0,9877	0,8333	0,1111
F	7	7	1	1	0,1429	7	1	1	0,1429	7	1,0000	1	0,1428	7	1	1	0,1429
G	9	9	1	1	0,0987	8	0,9843	0,8750	0,1562	9	1,0000	1	0,0987	9	1	1	0,1250
I	14	13	0,9949	0,9286	0,0816	13	0,9949	0,9285	0,0816	13	0,9949	0,9286	0,0816	14	1	1	0,0714
J	12	12	1	1	0,0833	11	0,9931	0,9167	0,0972	12	1,0000	1	0,0833	11	0,9931	0,9167	0,0972
R1b1a1a2	159	152	0,9997	0,9560	0,0070	97	0,9810	0,6100	0,0291	151	0,9996	0,9497	0,0070	142	0,9989	0,8931	0,0081

N = Number of samples; n = number of distinct haplotypes; HD: haplotype diversity (gene diversity); DC: discrimination capacity; MP, match probability

**Table 3.** Prediction scores (%) for seven haplogroup classes using three machine learning methods (SVM, Random Forest, k Nearest Neighbors) and LDA on four Y-STR datasets (CombYplex, M1, M2, Y-filer kit). For LDA, 10 individuals have been removed for Y-filer kit due to missing data; DYS502 has been removed from M1 analyses due to the lack of polymorphism.

**Table 3.**

Haplogroup	N	Method	Prediction score (in %)			
			Full CombYplex	M1	M2	Y-filer
<b>E1a-M33</b>	15	SVM	100	100	100	100
		Random Forest	97	99	83	99
		k Nearest Neighbors (kNN)	67	100	67	67
		LDA	100	100	100	97
<b>E1b1a</b>	275	SVM	100	99	97	99
		Random Forest	100	100	97	100
		k Nearest Neighbors (kNN)	99	100	97	100
		LDA	99	97	98	100
<b>E1b1b</b>	12	SVM	67	33	67	67

		Random Forest	28	28	28	54
		k Nearest Neighbors (kNN)	33	33	33	33
		LDA	62	61	<b>55</b>	75
<i>All E collapsed</i>	302	SVM	100	100	96	96
		<i>Random Forest</i>	100	100	97	100
		<i>k Nearest Neighbors (kNN)</i>	99	100	93	100
<b>G</b>	9	SVM	<b>67</b>	67	<b>0</b>	<b>67</b>
		Random Forest	71	75	5	69
		k Nearest Neighbors (kNN)	67	67	0	33
		LDA	100	88	67	88
<b>I</b>	14	SVM	100	100	100	75
		Random Forest	99	98	79	74
		k Nearest Neighbors (kNN)	75	100	75	75
		LDA	95	94	81	44
<b>J</b>	12	SVM	100	100	67	67
		Random Forest	98	100	13	39

		k Nearest Neighbors (kNN)	67	100	0	67
		LDA	100	100	14	67
<b>R1b1a1a2-M269</b>	159	SVM	95	98	93	98
		Random Forest	97	95	97	98
		k Nearest Neighbors (kNN)	100	98	95	98
		LDA	100	100	99	96
<b>Average</b>	<b>496</b>	<b>SVM</b>	<b>97</b>	<b>96</b>	<b>92</b>	<b>95</b>
		<b>Random Forest</b>	<b>97</b>	<b>96</b>	<b>90</b>	<b>95</b>
		k Nearest Neighbors (kNN)	<b>73</b>	97	<b>52</b>	<b>68</b>
		LDA	94	91	73	81

**Table 4.** Prediction scores (%) for twelve haplogroup classes using the two best machine learning methods (SVM and Random Forest) on four Y-STR datasets (CombYplex, M1, M2, Y-filer kit).

**Table 4.**

Haplogroup	N	Method	Prediction score (in %)		
			CombYplex	M2 only	Y-filer
E1a-M33	15	SVM	100	100	100
		Random Forest	98	90	99
E1b1a1-M2*	44	SVM	45	27	27
		Random Forest	58	46	37
E1b1a7-M191	17	SVM	40	60	80
		Random Forest	40	40	60
E1b1a7a-U174	79	SVM	75	80	90
		Random Forest	81	70	87
E1b1a8a-U209	66	SVM	75	62	56
		Random Forest	72	74	70

<b>E1b1a8a1-U290</b>	69	SVM	35	47	47
		Random Forest	56	59	63
<b>E1b1b1-M35*</b>	10	SVM	100	67	67
		Random Forest	68	32	48
<b>G-M201</b>	9	SVM	67	33	67
		Random Forest	88	28	92
<b>I</b>	14	SVM	100	75	75
		Random Forest	100	83	72
<b>J</b>	12	SVM	100	33	33
		Random Forest	100	32	43
<b>R1b1a1a2-M269</b>	134	SVM	85	85	94
		Random Forest	97	99	91
<b>R1b1a1a2a1a2a1b1a1-M167</b>	25	SVM	86	29	0
		Random Forest	84	60	58
<b>Average</b>	<b>494</b>	SVM	71	64	67
		<b>Random Forest</b>	<b>79</b>	<b>71</b>	<b>74</b>