

Forensic Science International: Genetics

An empirical investigation into the effectiveness of genetic genealogy to identify individuals in the UK --Manuscript Draft--

Manuscript Number:	FSIGEN_2019_449R1
Article Type:	Research Paper
Section/Category:	Regular Paper
Keywords:	Forensic Genealogy; Genetic Genealogy; SNP; GEDmatch
Corresponding Author:	Jim Thomson Eurofins Forensic Services Teddington, United Kingdom
First Author:	Jim Thomson
Order of Authors:	Jim Thomson Tim Clayton John Cleary Maurice Gleeson Debbie Kennett Michelle Leonard Donna Rutherford
Abstract:	<p>The use of genetic genealogy techniques to identify Joseph James DeAngelo as the prime suspect in the Golden State Killer case in 2018 has opened up a new approach to investigation of cold cases. Since that breakthrough, genetic genealogy methods have been reported to be applied to around 100 cases. To date, all of these reports relate to investigations in the US, where the high uptake of “direct-to-consumer” (DTC) genetic testing by individuals conducting private ancestral research has provided the necessary publicly available data for successful forensic investigations. We have conducted a study to assess the likely effectiveness of genetic genealogy techniques if applied to investigations in the UK. Ten volunteers provided their own SNP array data, downloaded from a DTC provider of their choice. These data sets were anonymised and uploaded to the GEDmatch Genesis genealogy website, mimicking data sets from unsourced crime samples or unidentified human remains. A team of experienced genealogists then attempted to identify the donors of the anonymised data sets by working with matches on the database and identifying points where the matches’ trees intersect to determine their shared family lineages which were further investigated using traditional resources (such as birth, marriage, death and census records, social media and online family trees). Through these methods, four of the ten donors were identified, at least to the level of one of a set of siblings. This confirms that, despite the over-representation of US citizens on publicly accessible genealogy databases, there is still potential for effective use in investigations outside the US where legislation permits. One of our four identified individuals was of Indian heritage (via St Vincent and the Grenadines) highlighting that in the right circumstances individuals of non-European origin can be identified.</p>
Suggested Reviewers:	
Opposed Reviewers:	
Response to Reviewers:	

Author responses to Reviewer Comments for FSIGEN_2019_449 An empirical investigation into the effectiveness of genetic genealogy to identify individuals in the UK

Reviewer #1: This paper details a useful investigation of the relative usefulness of genetic genealogy to identify unknown individuals from high quality samples. A similar study was also undertaken and published in Buzzfeed News which the authors may or may not have been aware of, but this has not been referenced. Although the focus of this study was about the technique's usefulness in the UK, the outcome was similar.

Response: Description and reference for this (and reference to another journalistic study by Kristen Brown article) added at lines 76-81

While the authors clearly recognised potential ethical and governance issues that would need to be addressed before such a process could be widely accepted, they have not mentioned other factors that can influence the use of the technique, such as the quality and quantity of available material for analysis, potentially limited to sexual assaults with ejaculation or unidentified deceased from which an appropriate amount of material can be recovered, along with individuals with the relevant genealogical skills

Response: Paragraph to comment on these practicalities added to Introduction at lines 44-55

The authors have utilised the GEDmatch Genesis tools and it would be helpful if there was a short summary of these features in the text and, given the recent website changes to require an opt-in, as well as the purchase of the site by Verogen, it would be helpful to add something about the current size of the database accessible by law enforcement and any potential changes that may have been revealed

Response: New paragraph describing relevant GEDmatch tools added to Materials and Methods section at lines 142-154

Response: Discussion updated to include updated figures for GEDmatch opt-in (line 464), details of GEDmatch acquisition (lines 468-474) and reference to Othram DNA Solves database for law enforcement investigations (lines 481-484)

More importantly, can the authors quantify the number of genealogy hours employed for each of the cases, or at least a median and range for the set?

Response: A paragraph detailing estimated time spent by the genealogy team on the investigations has been added (lines 229-240)

Reviewer #2: The article provides a persuasive rationale for why research into the efficacy of genetic genealogy outside the US is needed, both in the introduction and concluding paragraph. The limitations of the study are generally well stated, and the results are both novel and useful. A few minor points to consider are detailed below:

- More detail related to the DTC tests needs to be given in the methods (I know it is mentioned at the bottom of Table 1) – from which companies were the tests purchased and how big were the SNP arrays,

Response: Additional section on Genetic Testing by DTC providers added to M&M. Lines 108-117.

... and a comment on whether there is any suggestion that the different types of DTC uploads affected the power of the GEDmatch search?

Response: A sentence has been added to comment that the study was too small for any differences between outcomes based on DTC provider to be considered significant (lines 448-450).

In this respect, we do not think the different chips made any difference to the outcomes of the searches but we do not have any empirical evidence and it is not something we considered when doing the analyses. The critical factor in each case was the nature of the cousin matches at GEDmatch. GEDmatch do not reveal much about the way they do their matching but the thresholds are very permissive with matches potentially going right down to 7 cM. There might be some discrepancies with the low-level matches but these would not have been useful anyway for our purposes and we do not generally work with small segments under about 15 cM anyway because there is such a high false-positive rate.

There are far more Ancestry testers at GEDmatch than from any other company. That is mainly because Ancestry has the largest database (16 million testers). 23andMe have tested over 10 million people but these people are largely testing for health rather than ancestry and are far less likely to opt in to cousin matching. The three people in our study who tested at 23andMe all had lots of matches with people who had tested at Ancestry so we do not think they were missing out on any matches. We believe the fact that we did not solve any of the 23andMe uploads was just a chance finding in a very small dataset. Another factor might be that the volunteer subjects who used 23andMe were potentially less likely to be interested in genealogy, and hence less likely to have closer relatives in the databases.

· Table 1 – The legend here has to include what the ethnicity codes mean since it will not be intuitive to most people.

Response: Legend revised with additional information

· Since all volunteers had previously purchased a DTC test themselves, do you worry that this is self-selecting a segment of the population along, for example, socio-economic lines, where relatives are also more likely to have had a test? Similarly, we know close relatives often get tests together (either because they see the results from a close relative and want to do it, or because they are given as gifts from a satisfied ‘customer’), so this again may bias the results in relation to subject 1. Neither of which invalidate the study, but it is something to consider

Response: Additional paragraph on the potential for bias in the sample set added to the discussion (lines 455-466)

· Line 203 – you need to make it clearer that AD is the name you have given to the individual who was a match since to start with I assumed it was some sort of acronym, especially as SA is always written in quotation marks

Response: Clarified by comment in the results section (lines 248-252) and modification of sentence at line 267

· It would be much easier to understand the familial links for subject 3 if a family tree was included.

Response: Family tree added (figure 1)

Highlights

- Ten volunteer UK residents provided high density SNP data downloaded from a direct-to-consumer testing company.
- Data were uploaded anonymously to GEDmatch and genetic genealogy techniques used to attempt to identify the volunteers by name
- Four out of a sample of ten UK residents were identified by name or as one of a group of named siblings.
- Demonstration that genetic genealogy techniques can be effective in investigation of individual identity for cases in the UK

1 **An empirical investigation into the effectiveness of genetic genealogy to identify**
2 **individuals in the UK.**

3

4 Jim Thomson^a, Tim Clayton^a, John Cleary^{b,c}, Maurice Gleeson^b, Debbie Kennett^{b,d},
5 Michelle Leonard^b, Donna Rutherford^b.

6

7 ^a Eurofins Forensic Services, Teddington, TW11 0LY UK

8 ^b ISOGG UK and Ireland

9 ^c Heriot-Watt University, Boundary Road North, Edinburgh, Scotland, EH14 4AS, UK

10 ^d University College London, Gower Street, London WC1E 6BT, UK

11

12 Corresponding author email: jimthomson@eurofins.co.uk

1 **An empirical investigation into the effectiveness of genetic genealogy to identify** 2 **individuals in the UK.**

3

4 **Abstract**

5 The use of genetic genealogy techniques to identify Joseph James DeAngelo as the
6 prime suspect in the Golden State Killer case in 2018 has opened up a new approach to
7 investigation of cold cases. Since that breakthrough, genetic genealogy methods have
8 been reported to be applied to around 100 cases. To date, all of these reports relate to
9 investigations in the US, where the high uptake of “direct-to-consumer” (DTC) genetic
10 testing by individuals conducting private ancestral research has provided the necessary
11 publicly available data for successful forensic investigations. We have conducted a study
12 to assess the likely effectiveness of genetic genealogy techniques if applied to
13 investigations in the UK. Ten volunteers provided their own SNP array data, downloaded
14 from a DTC provider of their choice. These data sets were anonymised and uploaded to
15 the GEDmatch Genesis genealogy website, mimicking data sets from unsourced crime
16 samples or unidentified human remains. A team of experienced genealogists then
17 attempted to identify the donors of the anonymised data sets by working with matches
18 on the database and identifying points where the matches’ trees intersect to determine
19 their shared family lineages which were further investigated using traditional resources
20 (such as birth, marriage, death and census records, social media and online family
21 trees). Through these methods, four of the ten donors were identified, at least to the
22 level of one of a set of siblings. This confirms that, despite the over-representation of US
23 citizens on publicly accessible genealogy databases, there is still potential for effective
24 use in investigations outside the US where legislation permits. One of our four identified
25 individuals was of Indian heritage (via St Vincent and the Grenadines) highlighting that in
26 the right circumstances individuals of non-European origin can be identified.

27

28 **Keywords**

29

30 Genetic Genealogy, Forensic Genealogy, GEDmatch, SNP

31

32 **Introduction**

33

34 The rapid rise in prominence of genetic genealogy methods applied to forensic
35 investigations in 2018 and 2019 has had a transformative impact on investigation of cold
36 cases for which unsourced DNA evidence is available [1]. This investigative approach,
37 known in various circles as investigative genetic genealogy or forensic genealogy, takes
38 advantage of the rapidly expanding resources available to amateur genealogical
39 researchers investigating their bio-ancestral origins and family histories. Direct-to-
40 consumer (DTC) genetic testing using high density SNP arrays has become cheap and
41 widely available, fuelling a rapid increase in the size of searchable SNP databases held
42 by the DTC companies and by other independent organisations, such as GEDmatch.

43

44 To utilise these resources for forensic investigations, sufficient good quality DNA
45 assumed to originate from the perpetrator of a crime (or from unidentified remains) is
46 needed. This requirement may limit this approach to cases where semen or relatively
47 large bloodstains attributable to the perpetrator are available, or to bodies or body parts
48 from which sufficient DNA can be recovered. From such samples, SNP genotype data
49 compatible with those generated by the DTC providers must be generated. Typically, the

50 high density SNP arrays used for these analyses include approximately 600,000 –
51 730,000 SNPs, dependent on the supplier, most of which are autosomal but which often
52 also include X and Y chromosome and mitochondrial SNPs [2]. These data can then be
53 uploaded to a permissible genealogy database such as GEDmatch, whereupon
54 individuals with the relevant genealogical skills are needed to conduct the required
55 searches and pursue the genealogical investigations.
56

57 Although publicly available genetic data had previously been used to assist in cold case
58 identifications [3–5] it was the identification of Joseph James DeAngelo as the prime
59 suspect in the Golden State Killer case which alerted many to the potential of such
60 investigations. DeAngelo was identified following SNP genotyping of a semen sample
61 taken from a victim of the killer in the 1980s. The data were uploaded to the GEDmatch
62 genetic genealogy database, and investigative work by Barbara Rae-Venter, a genetic
63 genealogist working with the FBI, finally led to the identification of DeAngelo who was
64 arrested and charged with first degree murder in April 2018. [6–8]
65

66 Between April 2018 and January 2019 a total of 28 cases were publicised wherein the
67 law enforcement agencies announced identification of DNA from a suspected
68 perpetrator with the aid of genetic genealogy [9].

69 These cases and a number of subsequent reports all related to crimes committed in the
70 USA where the uptake of genetic genealogy tests by the public has been highest. Erlich
71 *et al* [10] determined that nearly 60% of genealogy searches in a large database of
72 1.28m individuals would return a match to a relative with at least 100 cM of shared DNA,
73 a level usually indicative of a 3rd cousin (3C) relationship or closer. However, this
74 empirical study utilised a test set comprised of subjects already on the database and so
75 is subject to any bias inherent in the database constitution (such as over-representation
76 of US residents of European ancestry). In two other journalistic studies, also focussed
77 on identification of US residents, Aldhous [11] reported that 6 out of 10 BuzzFeed
78 employees who provided DNA results from an unspecified DTC genetic testing provider
79 were identified following genetic genealogy investigations, and Brown [12] reported the
80 successful identification of herself by an independent genealogist. Despite these useful
81 investigations, the question of whether such investigative methods would yield
82 successful results if applied to searches outside the USA (relating to (a) unidentified
83 human remains or (b) the victims or perpetrators of crimes) remains untested. This study
84 seeks to provide preliminary information to address this question based on a small
85 convenience sample of ten individuals to assess the potential effectiveness of genetic
86 genealogy methods in identifying UK residents.
87

88 **Materials and Methods**

89 *Selection of subjects*

90
91
92 Ten volunteer subjects were recruited from *<information removed>* staff. The primary
93 criterion was that the individuals had privately purchased and completed a DNA
94 ancestry/genealogy test from one of the major direct-to-consumer (DTC) providers. All
95 volunteers were fully briefed on the project aims and potential implications of third-party
96 scrutiny of their DNA records and associated genealogy and provided full consent. The
97 project was also approved by our in-house ethical approval procedures. The only other
98 criterion applied was that all subjects were UK residents with at least three years of
99 continuous residency. The recruitment process was not intended to provide a

100 representative sample of the UK population, but represented a convenience sample
101 suitable for an indicative feasibility assessment.

102
103 For each subject, sex and date of birth were recorded. Subjects also provided
104 information on their country of birth; their own self-defined ethnicity (using United
105 Kingdom Home Office self-defined ethnicity codes) [13]; and the country of birth and
106 ethnicity (if known) of both of their parents and all of their grandparents.

107
108 *Genetic testing by DTC providers*

109
110 Genetic tests were procured from either Ancestry (Lehi, UT, USA) or 23andMe
111 (Sunnyvale, CA, USA) by the individual volunteers using the standard UK web-based
112 order processes. Tests were conducted between April 2016 and January 2019 using the
113 then current SNP panels. For all seven Ancestry users, this was the Ancestry v2 chip
114 based on the Illumina OmniExpress Plus chip with ~669,000 SNPs [14]; for two
115 23andMe users, this was the 23andMe v4 chip, a bespoke customised Illumina chip with
116 ~602,000 SNPs[15] ; and for one 23andMe user, this was the 23andMe v5 chip, the
117 Illumina GSA chip with~640,000 SNPs [16].

118
119
120 *Upload of data to GEDmatch Genesis*

121
122 All subjects downloaded their own raw DNA data files from their DTC provider. These
123 data files were transferred to the project manager. For each subject, a new account was
124 created on the GEDmatch website using a random user name and a generic email
125 address unattributable to the subject. DNA SNP data from each subject was uploaded to
126 the GEDmatch Genesis database and all uploads were changed to a privacy status of
127 “Research” which is an option “*provided primarily for artificially created research kits*”
128 [17] and which ensures that the subject’s DNA will not be included in match results of
129 other users. Following successful upload, the kit numbers generated by GEDmatch for
130 each subject were provided to the genealogy team, enabling them to conduct the
131 required searches. All searches were conducted using the GEDmatch Genesis beta
132 website (which was launched on 20 December 2018). GEDmatch Genesis and
133 GEDmatch were merged on 1st June 2019 and the website has reverted to the original
134 GEDmatch domain name but the new Genesis features are now integrated into the
135 website [18].

136
137 *GEDmatch searching*

138
139 All GEDmatch searches and subsequent genealogical investigations were conducted on
140 a *pro bono* basis by a team of five highly-experienced genetic genealogists.

141
142 GEDmatch is a platform for testers who have used one DTC provider to compare their
143 data with testers from all others, offering a number of online tools for its users, some of
144 which are free to use while others require a small monthly subscription. The base
145 function is a search of the kit of interest for matches across the whole available database
146 (“One-to-many” comparison), which can be followed by detailed comparisons between a
147 kit and each of its matches (“One-to-one” comparison, for autosomal DNA plus a specific
148 tool for the X chromosome). Another free to use tool allows the user to see whether
149 some of a kit’s matches also match each other (“People who match both or 1 of 2 kits”),
150 which enables networks of shared relationships to be built around the subject, and a

151 related subscription tool maps “Clusters” of mutually related kits. These network and
152 cluster building tools were the most important tools for the analysis presented here,
153 along with the shared cM data, as they indicated where intersections between the family
154 trees of the subject and matches may be found.

155
156 Following upload, “one-to-many” searches were carried out for each subject to provide
157 the list of 3000 closest matches in the GEDmatch database (estimated size at time of
158 study: 1.2 million) [19]. Population admixture proportions were estimated using the
159 Eurogenes K13 model [20], one of a number of admixture calculators included in the
160 GEDmatch toolkit. Each kit was assessed for runs of homozygosity by using the “Are
161 your parents related?” tool.

162 163 *Genetic genealogy investigation strategy*

164
165 Investigation of the identity of each subject progressed in several distinct stages:

- 166
167 1. Assessment of the amount of shared DNA with the top matches and the degrees
168 of relationship predicted from them using the Shared cM Project tool that supplies
169 ranges of known relationships observed within crowd-sourced data. The
170 probabilities of predicted relationships quoted below are based on this tool. [21]
- 171 2. Identification of clusters of shared matches that might indicate a common
172 ancestor (cluster analysis).
- 173 3. Examination of family trees of the top matches to identify points of intersection,
174 likely to indicate common ancestors. Some matches were linked to online family
175 trees but these varied in completeness and accuracy. Genealogical work using
176 standard genealogical records (e.g. birth, marriage and death records, census
177 records, etc.) or public domain information (e.g. online obituaries, electoral
178 registers, social media etc.) followed to verify, correct and extend existing trees,
179 and to build new trees from scratch where none existed.
- 180 4. Following identification of possible common ancestors, trees were built forward in
181 time from them to identify candidates for the person of interest fitting the profiles
182 of gender, age, circumstantial evidence, and other information.

183
184 The pilot study ended at the beginning of June 2019 and all the profiles were then
185 deleted from GEDmatch.

186 187 **Results**

188 189 *Subjects*

190
191 The ten volunteer subjects included six males and four females aged between 22 and
192 56. Seven were born in England or Wales with parents and grandparents also declared
193 as White British from the UK. One was born in England to a Chilean father and English
194 mother, although declared her father’s and paternal grandfather’s ethnicity as White
195 British (reported as members of a British expatriate community in Chile). One was born
196 in England to parents and grandparents of Indian ethnicity who had been resident in St
197 Vincent and the Grenadines. One was born in Romania to Romanian parents and
198 grandparents.

199
200
201

	Subject 1	Subject 2	Subject 3	Subject 4	Subject 5
Assigned Name Initials	"AW"	"JS"	"SA"	"SE"	"MM"
Sex	M	F	M	F	M
Date of birth	1977	1997	1964	1991	1977
Country of birth	England	England	England	Romania	Wales
Self-declared ethnicity	W1	W1	W1	W9	W1
Father	England, W1	England, W1	England, W1	Romania, W9	Wales, W1
Mother	England, W1	England, W1	England, W1	Romania, W9	Scotland, W1
Paternal GF	England, W1	England, W1	Wales, W1	Romania, W9	Wales, W1
Paternal GM	England, W1	England, W1	England, W1	Romania, W9	Wales, W1
Maternal GF	England, W1	England, W1	England, W1	Romania, W9	Scotland, W1
Maternal GM	England, W1	England, W1	England, W1	Romania, W9	Scotland, W1
DTC Provider	Ancestry	Ancestry	Ancestry	23andMe	Ancestry

	Subject 6	Subject 7	Subject 8	Subject 9	Subject 10
Assigned Name Initials	"GP"	"AB"	"BM"	"AJ"	"JB"
Sex	M	M	F	F	M
Date of birth	1983	1982	1992	1984	1963
Country of birth	England	England	England	England	England
Self-declared ethnicity	W1	W1	W1	A1	W1
Father	England, W1	England, W1	Chile, W1	St.Vincent, A1	England, W1
Mother	Guernsey, W1	England, W1	England, W1	St.Vincent, A1	England, W1
Paternal GF	Unknown	England, W1	Chile, W1	St.Vincent, A1	Shetland, W1
Paternal GM	Unknown	England, W1	Chile, W1	St.Vincent, A1	Shetland, W1
Maternal GF	Guernsey, W1	England, W1	Ireland, W2	St.Vincent, A1	England, W1
Maternal GM	Guernsey, W1	England, W1	Ireland, W2	St.Vincent, A1	England, W1
DTC Provider	23andMe	Ancestry	23andMe	Ancestry	Ancestry

202
203

204 Table 1. Details of the ten volunteer subjects. Ethnicity is defined using the United
205 Kingdom Home Office self-defined ethnicity codes [13]; W1=British; W2=Irish; W9=Any
206 other white background (other than British, Irish or Gypsy/Irish traveller); A1=Indian

207
208

209 *One-to-Many matches*

210

211 One-to-many comparison searches were conducted for each subject. Outcomes are
212 summarised in table two, with matches categorised by the total shared cM of DNA
213 between the subject and the nearest matches. A greater amount of shared DNA is
214 indicative of a closer relationship between the subject and the matched sample. One
215 subject ("AW") was identified as having two close relatives (likely to be a parent or child,
216 and a sibling) on the GEDmatch database. Two subjects had closest matches in the 80-
217 500 cM range, suggestive of second or third cousin (2C-3C) relationships and a further
218 three had closest matches with 50-80 cM shared, suggestive of 3C-4C. The remaining
219 four subjects had no matches sharing more than 50 cM although, along with all other
220 subjects, had several or many (up to 859) more distant matches in the 20-50 cM range
221 suggestive of 4C or more distant relationships.

222
223

Range cM	Subject 1	Subject 2	Subject 3	Subject 4	Subject 5	Subject 6	Subject 7	Subject 8	Subject 9	Subject 10	Likely relationship
~3570	1	0	0	0	0	0	0	0	0	0	Parent or child
2200-3300	1	0	0	0	0	0	0	0	0	0	Sibling (full)
1200-2200	0	0	0	0	0	0	0	0	0	0	Half siblings, aunts, uncles, grand-parents/children
550-1200	0	0	0	0	0	0	0	0	0	0	1C ~ 1C removed
80-500	0	0	1	0	0	0	0	0	2	0	approx 2C~3C or 1C removed
50-80	2	0	0	0	3	0	0	2	4	1	approx 3C~4C
30-50	6	1	7	8	7	16	6	36	2	5	approx 4C or more distant (could be closer)
20-30	108	99	81	74	71	199	67	823	2	41	approx 5C or more distant (could be closer)

Table 2. Summary of one-to-many searches. For each subject, the numbers of close matches in each cM range is shown.

Time spent on identification investigations

The genealogical team spent 2-3 hours on a triage analysis of each kit to judge its likely chances of success. 5 kits with low matches were judged as requiring potentially over 100 hours of work to identify distant cousins of the matches and these were not pursued, although a start was made in some cases. In total some 200-300 hours were put in by the genealogical team on pursuing the 5 cases judged most promising, including reporting and team meetings. The quickest case was solved in around 3 hours, while more complex ones needed between 50-100 hours to solve. The genealogical team were aware that the targets were associated with <information removed>, but had no staff list to work from so that the identifications of the targets' names had to arise out of the genealogical tracing.

Outcomes of identification investigations

Following the initial one-to-many matching, a degree of triage was applied with further investigations focussed mainly on those individuals with closer nearest matches. The results for all subjects are discussed below starting with the four subjects who were identified by the investigation either by name, or as one of a set of named siblings. In the following descriptions, initials in quotation marks (e.g. "AW") refer to the false identities assigned to the ten subjects under investigation, and other initials (e.g. AD) refer to real individuals identified as matches or likely relatives of those subjects during the investigations.

Subject 1 (assigned initials "AW")

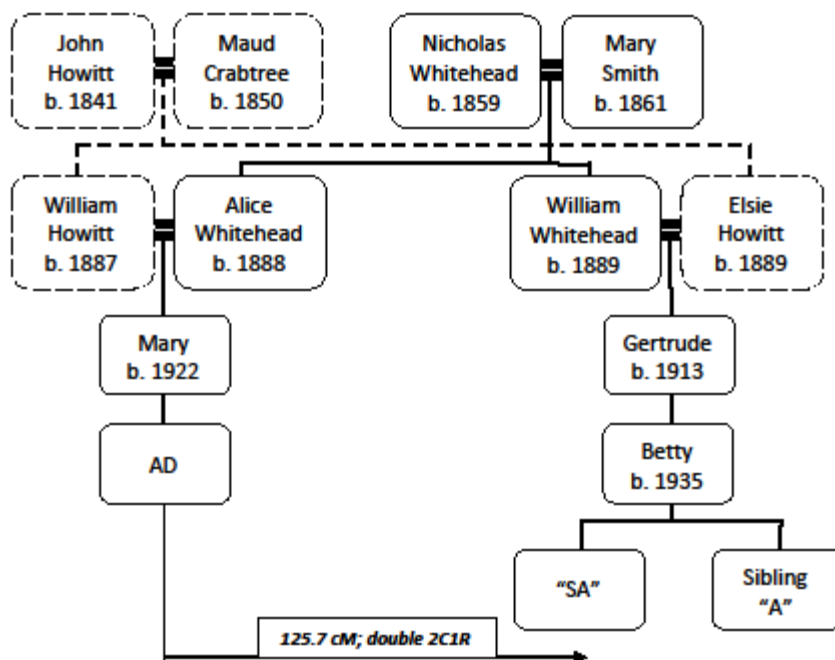
One parent/child match and one sibling match were identified on GEDmatch making this the most straightforward identification case. Cross-referencing of recent electoral registers and birth, marriage and death indexes enabled a candidate family group to be identified. The person of interest "AW" was identified as one of two twin brothers. The investigation did not identify whether these twins were monozygotic or dizygotic which would have implications for any further DNA identification testing were this a real investigative case. There was a third brother who could be eliminated as being "AW", as

263 he also had a genetic profile in GEDmatch matching in the expected sibling range to
264 "AW".

265
266 *Subject 3 (assigned initials "SA")*

267 There was a good top match at 125 cM to an individual called AD, predicting
268 approximately 85% probability of being 3C or closer to "SA", allowing the descendants of
269 common great-great-grandparents to be targeted in the search. AD's paternal side was
270 eliminated by an identified close match not shared with "SA", supported by a small but
271 likely significant 13 cM X chromosome match between AD and "SA", suggesting that the
272 connection may be found on both persons' maternal trees (since, both being male, they
273 had received X chromosomes only from their mothers). The shared matches for the two
274 gave no further clear leads: for example, one promising shared match, JB, appeared
275 likely to connect to AD through an unknown parent, meaning this parental connection
276 would need to be solved before the connection to AD and "SA" could be solved in turn.

277 A genealogical search strategy was adopted constructing a tree forwards in time to
278 identify all descendants of AD's maternal great-great-grandparent couples. The search
279 was limited to living male descendants whose immediate ancestor in the tree was their
280 mother, and so could have shared X chromosome segments with AD (ruling out any
281 male to male descents) to identify candidates to be "SA" or candidates for targeted
282 testing to find closer relatives. This led to candidate siblings, one of whom was Subject
283 3, whose tree intersected with AD's tree at two points, sharing two of AD's great-great-
284 grandparent couples. The candidate siblings were younger than AD, making them 2C1R
285 (second cousins one removed) to AD along two routes, therefore double 2C1R to AD
286 (Figure 1).
287



288
289

290 Figure 1: Tree of recent ancestry of "SA". "SA" and his closest match AD are second
291 cousins once removed (2C1R) twice, because each descends from two ancestral
292 couples due to marriages between two pairs of siblings. Names have been changed.

293
294

295 *Subject 9 (assigned initials "AJ")*

296

297 The top match for "AJ" (364 cM – likely between 2C-1C1R) was identifiable from
298 GEDmatch as KW, resident in High Wycombe, UK and related to members of a
299 community of Indo-Vincentians (people from St Vincent and the Grenadines of Indian
300 origin). "AJ" was correctly identified as Subject 9, and was closely related to a number of
301 people in the St Vincent and Grenadines Indian community in the UK, who are
302 descended from 19th century Indian indentured settlers in St Vincent, named B and K,
303 and whose network of family and social relationships and online family trees gave the
304 clues that led to identifying the female subject. Although the identification was
305 successful, the genealogists had cautioned that they could not guarantee they had
306 identified all the descendants of B and K who might fit the profile because the relevant
307 birth, marriage and death records in St Vincent are not accessible from the UK. They
308 suggested that identifying Subject 9's mother, who is also likely to be descended from
309 the founding couple, would help to corroborate the identification. Subject 9 was more
310 distantly related to her matches through her identified paternal line than predicted from
311 the amount of shared DNA, likely due to pedigree collapse through endogamy,
312 increasing observed levels of shared DNA across the family network, causing the search
313 to be broadened to include more distant relatives of the matches who were also
314 descended from the founder couple.

315

316 *Subject 10 (assigned initials "JB")*

317

318 The top match at 53 cM was not especially close, with an approximately 50-60%
319 possibility of being related to "JB" as 4C or closer, and the other matches were likely to
320 be more distant. The closest match, AR, was identifiable as a resident of Aberdeenshire,
321 Scotland, with her paternal ancestry from the Shetland Islands and maternal from north-
322 east Scotland.

323

324 Several other close matches with "JB" were also Shetlanders or had recent ancestors
325 (up to grandparents) from Shetland. A 5000-person tree was constructed for the
326 Shetland cluster, drawing upon some very well constructed and referenced trees made
327 available publicly by family historians with a deep knowledge of Shetland communities.

328

329 "JB" had a large cluster of matches indicating an ancestor from the Shetland Islands,
330 who could be a parent, grandparent or great-grandparent. This community is
331 endogamous so the matches within the cluster share several sets of ancestors in
332 common. Mapping the intersections of the shared relationships between clusters of
333 matches allowed modelling of the places on the trees under construction where "JB"
334 could be. "JB" had more distant matches with non-Shetlanders, which were assumed to
335 represent one side of his ancestry (paternal or maternal), and marriages on the tree
336 between Shetlanders and non-Shetlanders were identified and profiled. The
337 identification was made when one of the Shetland match's trees intersected with a non-
338 Shetland match's tree, which turned out to be on Subject 10's mother's side and from
339 Yorkshire, the father's side being the Shetland side of the family. Through this route, "JB"
340 was identified as probably being one of three brothers. One of these siblings was

341 Subject 10. There was also a sister who can be eliminated on the assumption that the
342 target person is male.

343
344

345 Two further subjects (5 and 8) had promising matches likely to be at the 3C-4C level but
346 in neither case could informative trees be developed and neither subject was identified.

347
348

Subject 5 (assigned initials "MM")

349

350 Three top matches shared 50-80 cM which suggested that identification might be
351 possible. However, "MM" was not identified, but the 69 cM shared with the top match,
352 WB, indicated an approximately 75% chance that the ancestors in common with that top
353 match were 3C or closer, and just a 5% chance that the relationship was 4C or more
354 distant. Therefore, the ancestors in common were likely to be within WB's 8 pairs of
355 great-great-grandparents. This considerably limited the pool of potential persons to be
356 "MM". There seemed a strong likelihood that "MM" was Welsh, or might be the child of
357 Welsh parents. However, there was also a possibility of endogamy in these Welsh family
358 trees which would make the relationship more distant than predicted.

359

Subject 8 (assigned initials "BM")

360

361 "BM" had two matches with just over 50 cM shared DNA and who could be identified.
362 These matches were likely to be 3C-4C and there were good chances of capturing 100%
363 of the potential ancestors in common. However, neither match had available trees and
364 neither match shared X-DNA with "BM". Most of the other matches were too distant to
365 make successful searching likely.

366
367

368 The remaining four subjects (Subject 2, 4, 6 and 7) were the least promising with no
369 matches sharing more than 50 cM. For all four of these individuals, cluster analysis and,
370 where appropriate, some preliminary investigations of ancestral trees were conducted
371 but in no cases could substantial headway be made to indicate the identity of the
372 subjects. Relationships between 4th cousins (4C) imply common ancestors among 16
373 great-great-great-grandparental couples and a typical person may not be able to identify
374 all of their ancestors at this level, leaving blocked lines on their trees and the possibility
375 that the common ancestors of interest can be on those blocked lines. Consequently,
376 match lists with the closest matches at likely 4C distance present poor chances for
377 successful identification. Summaries of the reported results for these four cases are
378 shown below.

379

Subject 2 (assigned initials "JS")

380
381

382 The top genetic matches to "JS" were predicted to be no closer than 4C, and were likely
383 in many cases to be more distant than that. None of the top matches were matches to
384 each other so they were likely to converge with "JS" on different lines several
385 generations back, making identifying a small pool of candidates to profile challenging.

386

Subject 4 (assigned initials "SE")

387
388

389 The top matches were all more than 80% likely to be no closer than 4C, with no
390 clustering or trees to go on for further analysis.

391

392 *Subject 6 (assigned initials "GP")*

393

394 The top genetic matches to "GP" were predicted to be no closer than 4C, and were likely
395 in many cases to be more distant than that. "GP" may have a recent ancestor with
396 Ukrainian Jewish ancestry, but it was not possible to quantify how recently that ancestor
397 lived.

398

399 *Subject 7 (assigned initials "AB")*

400

401 The top matches to AB were particularly low, with only a 12% or less chance of any of
402 them being at the 3C level, and most were likely to be more distant than 4C. This
403 rendered the possibility of identifying ancestors in common between AB and his matches
404 remote. Some clusters of matches offered prospects for research as they have closely
405 related family clusters of their own with trees, but as their relationships to AB were likely
406 very distant, any reconstructed trees were likely to be partial up to the necessary degree
407 of ancestry and may not have included the ancestors in common.

408

409 **Discussion**

410

411 The success of genetic genealogy methods in identifying victims and perpetrators of
412 crime in the United States has been widely reported and is of significant interest to other
413 jurisdictions. Whilst there are undoubtedly significant ethical and legislative hurdles and
414 considerations to address, which in some jurisdictions may prevent such approaches
415 being considered, there are also practical questions regarding the applicability of the
416 method outside of the US.

417

418 We have provided some preliminary data to demonstrate the efficacy of the method in
419 identifying UK residents, The demonstration that four out of a convenience sample of ten
420 volunteers could be identified either by name, or as one of a set of siblings, is a useful
421 indicator that the predominance of US individuals on GEDmatch and other genealogical
422 databases does not preclude successful identification of individuals from outside the US.

423

424 Although GEDmatch does not publish the composition of their database by country of
425 origin, co-founder Curtis Rogers has confirmed that "our number one users are in the
426 United States, then Canada, then England, and then Australia" [22]. Although it is
427 difficult to source reliable data on the proportion of US, Canadian and Australian
428 residents with British ancestral origins it is evident that significant numbers of all of these
429 nationals will have British ancestry and are likely to have living relatives in the UK. Of all
430 countries other than the US, it seems likely that residents of the UK will be amongst
431 those with the highest chance of having relatives present on GEDmatch or other
432 permissible databases whether those relatives are currently resident in the UK or
433 elsewhere. In this context, it is notable, although not statistically significant, that the
434 subject in our study with Romanian ancestry generated no matches sharing more than
435 50 cM and no clustering or existing trees to give grounds for a successful identification.
436 Without more detailed knowledge of the demographic composition of GEDmatch or other
437 databases it is hard to predict the effectiveness of such searching methods if applied to
438 residents of other countries although it seems likely that residents of European countries
439 with higher representation in the US population, such as Germany and Ireland, may be
440 most likely to lead to successful identifications.

441

442 The identification of one individual of Indian origin in our sample set is noteworthy in that
443 it demonstrates that, in the right circumstances, searches for individuals of non-
444 European origin can be successful. In this example (subject 9, from a small Indo-
445 Vincentian population), and in the case of subject 10, whose family originated from
446 Shetland (an island group 170km north of the Scottish mainland), the investigations were
447 significantly assisted by well-documented genealogical records and family trees
448 developed for these localised populations by interested historians and genealogists.
449 Both of these individuals are now resident in the south of England, far from the
450 geographical origins of the population groups which led to their identification, again
451 highlighting that in a forensic investigation it is hard to predict any likelihood of
452 successful identification without at least an initial triage and review of matches and
453 cluster information.

454
455 The small size of this sample, and the possibility of bias introduced by selecting subjects
456 from employees of a scientific services company who may fall into socio-economic
457 groups more likely to have relatives who have also carried out such tests, precludes any
458 extrapolation of the success rates obtained here to predict likely success in the UK
459 population at large. In this respect we also acknowledge that subject 1 was aware that
460 his parent and brother had already carried out such testing prior to this study, but his
461 inclusion in the volunteer sample set was not influenced by this prior knowledge as we
462 accepted all volunteers meeting our three year residency criteria.

463
464 The small sample set also precludes any further interpretation of the observation that all
465 four of the identified subjects were tested by Ancestry whereas the three subjects tested
466 by 23andMe remained unidentified.

467
468 The study reported here was completed prior to the change to GEDmatch's Terms of
469 Service and Privacy Policy which has significantly altered the landscape for genetic
470 genealogy searching, at least in the near future. In May 2018, following the initial reports
471 of successful searches in the Golden State Killer case, GEDmatch amended their
472 previously non-specific Terms of Service and Privacy Policy conditions to specifically
473 allow law enforcement usage but only to "(1) identify a perpetrator of a violent crime
474 against another individual; or (2) identify remains of a deceased individual" [23]. "Violent
475 crime" is defined here as homicide or sexual assault. However, in May 2019 the policy
476 was again changed, this time requiring that users "opt-in" if they consent to their DNA
477 data to be used for Law Enforcement purposes [24]. All current users at that time were
478 therefore automatically set as "opted-out" of law enforcement searches effectively
479 reducing the sample set available for such searching from ~1.2 million to zero at that
480 point. Since then, it is reported that about 200,000 users have opted in to the LE option
481 (by December 2019) [25], but clearly the overall efficacy of investigative genetic
482 genealogy is highly dependent on the available data set and at present the results
483 described in this study will not be representative of LE searches of the much smaller
484 dataset now available through GEDmatch. In December 2019 GEDmatch was
485 purchased by the forensic genomics company Verogen, who have announced plans to
486 add new features to improve the functionality of the database and to make the database
487 more secure. Following the purchase, all EU users were opted out of LE matching and
488 had to re-consent to comply with GDPR (General Data Protection Regulation). In the
489 short term this will further reduce the potential utility of the database for law enforcement
490 agencies in the EU [26].

491

492 Although to date the majority of reported successful searches have utilised GEDmatch,
493 other options are available to counterbalance this reduced access to GEDmatch data. In
494 January 2019 another DTC company FamilyTreeDNA (FTDNA) announced that it was
495 now testing samples for the FBI and allowing them to upload profiles to its database [27].
496 FTDNA currently allows US law enforcement agencies to register for uploads and
497 considers work with non-US agencies on a case-by-case basis [28]. Another company,
498 Othram Inc, which specialises in whole-genome sequencing for forensic purposes,
499 launched a new database in December called DNA Solves that allows users to upload
500 their data to help solve crimes and identify victims [26,29].

501

502 In summary, this demonstration of the effectiveness of investigative genetic genealogy
503 provides the initial evidence needed by regulatory bodies outside the US to prompt
504 further consideration of the ethical and regulatory frameworks to enable safe and
505 socially-acceptable introduction of investigative genetic genealogy with appropriate
506 regard for privacy and data security issues. This can facilitate the introduction of this
507 valuable investigative approach to assist in the identification of unidentified human
508 remains as well as the victims and perpetrators of crime in the UK and elsewhere.

509

510

511

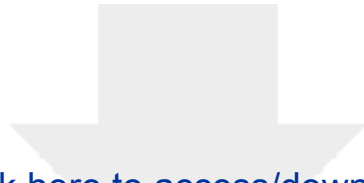
512 **References**

513

- 514 [1] D. Kennett, Using genetic genealogy databases in missing persons cases and to
515 develop suspect leads in violent crimes, *Forensic Sci. Int.* 301 (2019) 107–117.
516 <https://doi.org/10.1016/j.forsciint.2019.05.016>.
- 517 [2] International Society of Genetic Genealogy Wiki, Autosomal DNA Testing
518 Comparison Charts, *Int. Soc. Genet. Geneal. Wiki.* (2016).
519 https://isogg.org/wiki/Autosomal_DNA_testing_comparison_chart (accessed
520 January 3, 2020).
- 521 [3] T. Arango, The Cold Case That Inspired the ‘Golden State Killer’ Detective to Try
522 Genealogy, *New York Times.* (2018).
523 <https://www.nytimes.com/2018/05/03/us/golden-state-killer-genealogy.html>.
- 524 [4] S. Augenstein, The tale of the abandoned girl’s DNA that led to a notorious cold
525 case, *Forensic Mag.* (2017) 1–9.
526 [https://web.archive.org/web/20191020024455/http://www.forensicmag.com/article/
527 2017/02/tale-abandoned-girls-dna-led-notorious-cold-case](https://web.archive.org/web/20191020024455/http://www.forensicmag.com/article/2017/02/tale-abandoned-girls-dna-led-notorious-cold-case).
- 528 [5] M. Kayser, Forensic use of Y-chromosome DNA: a general overview, *Hum.*
529 *Genet.* 136 (2017) 621–635. <https://doi.org/10.1007/s00439-017-1776-9>.
- 530 [6] C. Phillips, The Golden State Killer investigation and the nascent field of forensic
531 genealogy, *Forensic Sci. Int. Genet.* 36 (2018) 186–188.
532 <https://doi.org/10.1016/j.fsigen.2018.07.010>.
- 533 [7] T. Arango, A. Goldman, T. Fuller, To Catch a Killer: A Fake Profile on a DNA Site
534 and a Pristine Sample - The New York Times, *New York Times.* (2018).
535 [https://www.nytimes.com/2018/04/27/us/golden-state-killer-case-joseph-
536 deangelo.html](https://www.nytimes.com/2018/04/27/us/golden-state-killer-case-joseph-deangelo.html).
- 537 [8] C.J. Guerrini, J.O. Robinson, D. Petersen, A.L. McGuire, Should police have
538 access to genetic genealogy databases? Capturing the Golden State Killer and
539 other criminals using a controversial new forensic technique, *PLoS Biol.* 16 (2018)
540 e2006906. <https://doi.org/10.1371/journal.pbio.2006906>.
- 541 [9] E.M. Greytak, C.C. Moore, S.L. Armentrout, Genetic genealogy for cold case and
542 active investigations, *Forensic Sci. Int.* 299 (2019) 103–113.

- 543 <https://doi.org/10.1016/j.forsciint.2019.03.039>.
- 544 [10] Y. Erlich, T. Shor, I. Pe'er, S. Carmi, Identity inference of genomic data using
545 long-range familial searches, *Science*. 362 (2018) 690–694.
546 <https://doi.org/10.1126/science.aau4832>.
- 547 [11] P. Aldhous, We tried to find 10 buzzfeed employees just like cops did for The
548 Golden State Killer, *BuzzFeed News*. (2019) 1–15.
549 [https://www.buzzfeednews.com/article/peteraldhous/golden-state-killer-dna-](https://www.buzzfeednews.com/article/peteraldhous/golden-state-killer-dna-experiment-genetic-genealogy)
550 [experiment-genetic-genealogy](https://www.buzzfeednews.com/article/peteraldhous/golden-state-killer-dna-experiment-genetic-genealogy) (accessed January 3, 2020).
- 551 [12] K. V. Brown, A Researcher Needed Three Hours to Identify Me From My DNA,
552 *Bloomberg*. (2019) 1–5. [https://www.bloomberg.com/news/articles/2019-04-12/a-](https://www.bloomberg.com/news/articles/2019-04-12/a-researcher-needed-three-hours-to-identify-me-from-my-dna)
553 [researcher-needed-three-hours-to-identify-me-from-my-dna](https://www.bloomberg.com/news/articles/2019-04-12/a-researcher-needed-three-hours-to-identify-me-from-my-dna).
- 554 [13] HM Government, Criminal Justice System Exchange Data Standards Catalogue,
555 2014.
556 [https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attac-](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/691544/self-defined-ethnicity-18plus1.pdf)
557 [hment_data/file/691544/self-defined-ethnicity-18plus1.pdf](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/691544/self-defined-ethnicity-18plus1.pdf).
- 558 [14] R.A. Canada, Ancestry V2 Chip, Haplogroup. (n.d.).
559 <https://haplogroup.org/ancestry-v2-chip/>.
- 560 [15] R A Canada, 23andMe V4 Chip, Haplogroup. (n.d.).
561 <https://haplogroup.org/23andme-v4-chip/>.
- 562 [16] 23andMe, Int. Soc. Genet. Geneal. Wiki. (n.d.). <https://isogg.org/wiki/23andMe>.
- 563 [17] GEDmatch raw data upload utility, (n.d.).
- 564 [18] Kitty Cooper, GENESIS Basics: GEDmatch reinvented part 1, *Kitty Cooper's Blog*.
565 (2019). [https://blog.kittycooper.com/2019/02/genesis-basics-gedmatch-](https://blog.kittycooper.com/2019/02/genesis-basics-gedmatch-reinvented-part-1/)
566 [reinvented-part-1/](https://blog.kittycooper.com/2019/02/genesis-basics-gedmatch-reinvented-part-1/).
- 567 [19] Genealogical Database Growth Slows, *DNA Geek*. (2019).
568 <https://thednageek.com/genealogical-database-growth-slows/>.
- 569 [20] D.F. Serviski, Updated Eurogenes K13 now at GEDmatch, *Eurogenes Genet.*
570 *Ancestry Proj.* (2013) 2–6. [http://bga101.blogspot.com/2013/11/updated-](http://bga101.blogspot.com/2013/11/updated-eurogenes-k13-at-gedmatch.html)
571 [eurogenes-k13-at-gedmatch.html](http://bga101.blogspot.com/2013/11/updated-eurogenes-k13-at-gedmatch.html).
- 572 [21] DNA Painter, (n.d.). <https://dnainter.com>.
- 573 [22] C. Curtis, 10094220 @ www.abc.net.au, (2018).
574 [https://www.abc.net.au/news/science/2018-08-13/online-genealogy-police-dna-](https://www.abc.net.au/news/science/2018-08-13/online-genealogy-police-dna-databases-golden-state-killer/10094220)
575 [databases-golden-state-killer/10094220](https://www.abc.net.au/news/science/2018-08-13/online-genealogy-police-dna-databases-golden-state-killer/10094220).
- 576 [23] D. Kennett, Updates to the Terms of Service and Privacy Policy at GEDmatch,
577 *Cruwys News*. 101 (2018). [https://cruwys.blogspot.com/2018/05/updates-to-](https://cruwys.blogspot.com/2018/05/updates-to-terms-of-service-and-privacy.html)
578 [terms-of-service-and-privacy.html](https://cruwys.blogspot.com/2018/05/updates-to-terms-of-service-and-privacy.html).
- 579 [24] A. Vaughan, DNA database opts a million people out from police searches, *New*
580 *Sci.* (n.d.). [https://www.newscientist.com/article/2203857-dna-database-opts-a-](https://www.newscientist.com/article/2203857-dna-database-opts-a-million-people-out-from-police-searches/)
581 [million-people-out-from-police-searches/](https://www.newscientist.com/article/2203857-dna-database-opts-a-million-people-out-from-police-searches/).
- 582 [25] M. Taylor, Verogen CEO: 'GEDmatch Will Be Improved, Not Changed,' *Forensic*.
583 (2019). [https://www.forensicmag.com/559058-Verogen-CEO-GEDmatch-Will-Be-](https://www.forensicmag.com/559058-Verogen-CEO-GEDmatch-Will-Be-Improved-Not-Changed/)
584 [Improved-Not-Changed/](https://www.forensicmag.com/559058-Verogen-CEO-GEDmatch-Will-Be-Improved-Not-Changed/).
- 585 [26] Justin Petrone, Forensic Genomics Market Advances Due to Consumer
586 Databases, *Technology Innovation, Genomeweb*. (2020).
587 [https://www.genomeweb.com/sequencing/forensic-genomics-market-advances-](https://www.genomeweb.com/sequencing/forensic-genomics-market-advances-due-consumer-databases-technology-innovation#.XiV96DP7S71)
588 [due-consumer-databases-technology-innovation#.XiV96DP7S71](https://www.genomeweb.com/sequencing/forensic-genomics-market-advances-due-consumer-databases-technology-innovation#.XiV96DP7S71).
- 589 [27] S. Hernandez, family-tree-dna-fbi-investigative-genealogy-privacy @
590 www.buzzfeednews.com, (2019).
591 [https://www.buzzfeednews.com/article/salvadorhernandez/family-tree-dna-fbi-](https://www.buzzfeednews.com/article/salvadorhernandez/family-tree-dna-fbi-investigative-genealogy-privacy)
592 [investigative-genealogy-privacy](https://www.buzzfeednews.com/article/salvadorhernandez/family-tree-dna-fbi-investigative-genealogy-privacy).
- 593 [28] FamilyTreeDNA, Law Enforcement Matching – Frequently Asked Questions, *Fam.*

594 Learn. Cent. (2019). [https://www.familytreedna.com/learn/ftdna/law-enforcement-](https://www.familytreedna.com/learn/ftdna/law-enforcement-faq/)
595 [faq/](https://www.familytreedna.com/learn/ftdna/law-enforcement-faq/).
596 [29] Andrea Leinfelder, The Woodlands-based Othram applies DNA sequencing to aid
597 investigations, *Houst. Chron.* (2019).
598 [https://www.houstonchronicle.com/business/article/The-Woodlands-based-](https://www.houstonchronicle.com/business/article/The-Woodlands-based-Othram-applies-DNA-sequencing-14934741.php)
599 [Othram-applies-DNA-sequencing-14934741.php](https://www.houstonchronicle.com/business/article/The-Woodlands-based-Othram-applies-DNA-sequencing-14934741.php).
600
601
602
603
604
605
606
607
608



Click here to access/download
RDM Data Profile XML
DataProfile_4534714.xml

