

Serveur Académique Lausannois SERVAL serval.unil.ch

Author Manuscript

Faculty of Biology and Medicine Publication

This paper has been peer-reviewed but does not include the final publisher proof-corrections or journal pagination.

Published in final edited form as:

Title: The use of Bayesian Networks and simulation methods to identify the variables impacting the value of evidence assessed under activity level propositions in stabbing cases.

Authors: Samie L, Champod C, Taylor D, Taroni F

Journal: Forensic science international. Genetics

Year: 2020 Jun 11

Issue: 48

Pages: 102334

DOI: [10.1016/j.fsigen.2020.102334](https://doi.org/10.1016/j.fsigen.2020.102334)

In the absence of a copyright statement, users should assume that standard copyright protection applies, unless the article contains an explicit statement to the contrary. In case of doubt, contact the journal publisher to verify the copyright status of an article.

The use of Bayesian Networks and simulation methods to identify the variables impacting the value of evidence assessed under activity level propositions in stabbing cases

1. Introduction

Technical analytical developments have made it possible to analyse very low amounts of DNA. The drawback of this progress is that evaluation of the results has become increasingly complex when activity level propositions are taken into account. Indeed, low quantity DNA traces can be the result of a secondary transfer or even a tertiary transfer. Variables related to transfer, persistence and recovery are becoming increasingly important considerations in evidence evaluation, as the relevant question arising in Court is no longer “who is the source of this DNA?” but rather “how did this DNA get there?” [1-2]. These considerations were at the basis of the ENFSI guideline for evidence reporting [3] and the advice given to DNA reporting officers to systematically report trace DNA evidence considering activity level propositions adopting a likelihood ratio (LR) approach. This article is focused on how to interpret traces with small amounts of DNA, typically found following the touching of an object. This is not a simple task as it requires the probabilistic consideration of multiples variables such as transfer, persistence, recovery and background level of DNA. Many experts face practical difficulties when assessing these variables due to the perception that every case has a unique set of circumstances and numerical assignments are critically dependent on the specificities of the case. Biedermann *et al.* [4] explained further why forensic scientists may struggle with these assessments. A common argument put forward is: “because each case has its own feature, the use of numerical values from experimental studies performed under controlled (laboratory) conditions cannot be used for evaluation in real life case”. Besides, experiments can only cover a limited number of options for any particular variable, so it is difficult to envisage experiments taking into account all possible variations. However, although each case may have different data for these variables, this does not mean that the LR would be affected by all possible variations of these variables. Identifying the variables that impact on the LR will help forensic scientists to focus on a limited number of variables of interest in order to limit time and cost of the required data acquisition. The key task that will be explored in this paper is to identify the variables that have a significant impact on the weight to be assigned to the DNA findings.

The objective of this paper is to present a methodology to support forensic scientists in the evaluation of their results given activity level propositions. This study will show how to identify

32 the variables impacting on the LR taking advantage of simulation methods. The purpose is to show
33 how to reduce the number of variables that require consideration. Once this is done then this
34 provides a focus for further data collections which can be used to inform distributions that can be
35 used for evaluative court-going purposes.

36 Taylor *et al.* [5] built a Bayesian Network (BN) [6] allowing the transfer mechanisms to be
37 considered. They offer a way to compute the LR associated with the DNA findings considering
38 activity level propositions. However, in order to be used in casework, the BN nodes need to be
39 quantified by probabilities, also informed with adequate data. The number of variables in the BN
40 from [5] is large and a case-specific data acquisition to inform the parameters on all of them would
41 be out of reach. This paper provides a method on how to identify the key variables that truly impact
42 on the LR by performing simulations based on the BN. It is an extension of the simulation approach
43 adopted by Taylor *et al.* [7] on body fluids attribution. Because the BN construction in [5] is the
44 basis of the present contribution, the reader is advised to refer to it as we will not fully describe
45 here its construction. In this paper, we will limit the description of the BN to the few modifications
46 that we have introduced, the BN parametrisation and the simulations techniques that we will use.

47 To illustrate this method, the scenario of a stabbing attack with a knife is used. It can be
48 summarized as follows: A victim is found dead in his flat following a stabbing incident carried out
49 by an offender using a knife. The exact day of the stabbing is uncertain. At the crime scene, a knife
50 is recovered and believed to be the attacker's weapon. The knife had been properly secured. It is
51 believed by the investigators, in line with pathologist's assessment, that the stabbing occurred
52 about 2 days (± 0.5 day) before the discovery of the crime scene.

53 Based on the elements of the investigation, not related to DNA evidence, a person of interest (POI)
54 is arrested and suspected to be the offender who stabbed the victim. Two days after the reporting
55 of the incident, a DNA swab is taken from the unstained smooth plastic handle of the knife with a
56 view to detect potential trace DNA left by the offender.

57 The prosecution's proposition (denoted H_p) is that the POI was the person who used the knife to
58 stab the victim.

59 In this paper, to gain some generalisation, two options are studied to reflect upon the defence point
60 of view (denoted H_{d_1} and H_{d_2} respectively):

- 61 – In the first defence proposition, H_{d_1} , the possibility of a secondary transfer is explicitly considered.
62 The POI claims that, he shook hands with an unknown person, probably few hours before the

63 stabbing, and it is that unknown person who is the real offender (who will be called alternative
64 offender AO).

65 – In the second defence proposition, Hd₂, it is alleged that the POI didn't stab the victim, but
66 someone else unrelated to him did it (AO). In addition, the POI denies any prior encounter with
67 either the victim, the knife or the AO. This represents a situation in which the POI is claiming no
68 direct link with the offence.

69 Exploring these two defence's propositions will allow to show that, depending on the propositions
70 of interest, the variables that have the strongest impact on the LR might be different.

71
72
73 The above-described circumstances provide the elements of what we will call the "initial case".
74 Again, with a view to explore further that this specific set of circumstances, we will develop two
75 additional cases. The first will adopt circumstances favouring the transfer of the POI's DNA (a
76 "high transfer case"). The second adopts circumstances that are less favourable to the transfer of
77 POI's DNA (an "low transfer case"). These sets of circumstances are further described in section
78 2.4.

80 2. Methodology

81 The methodology adopted for this research is decomposed into several stages. First, we
82 constructed a Bayesian network (BN) allowing a scientist to assess DNA findings associated with
83 the stabbing attack scenario (including the possibility to change the defence proposition). The
84 construction is mainly based on Taylor *et al.* [5] but has been adapted in order to carry out
85 simulations. Then the conditional probability tables (CPTs) for each node of the network have
86 been informed based on the data available from the literature. The parameters have been modelled
87 in a way to (a) allow a Bayesian update in the light of new data and (b) to reflect when applicable
88 the fact that the amount of data may be sparse or limited.

89 As such the constructed BN can be used to compute a likelihood ratio for a given case as done in
90 [5]. However, we would like to go further by exploring impact and limitations of data on the LR
91 values. A similar approach was adopted by [7]. This will be done using simulations resampling
92 from the underpinning distributions used to inform the CPTs. The sampling will be carefully
93 chosen by exploring one variable after the other. That is done in order to identify which variable
94 (or node in the BN) has the most significant impact on the variations observed on the LRs. The

95 isolation of the variables that have the most bearing on the variation of the LR from one simulation
96 to the other is important to inform a data acquisition strategy. As it will transpire later, the
97 developed BN has multiple variables and forensic laboratories will certainly not be able to
98 systematically investigate the underpinning data for all of them. We will show that this simulation
99 methodology allows us to successfully identify the key variables that should be the focus of future
100 data acquisition. Such methods have already been applied successfully in areas such as DNA and
101 fibres [8, 9]. Here, it is proposed to adopt such techniques to establish a baseline to inform future
102 data acquisition campaigns.

103 In the next section, we will firstly present the development of the Bayesian network on the basis
104 of what we call the “initial case”. It will be used as the starting point from which all simulations
105 will be carried out. Then, the method used to perform the simulations, from this “initial case”, will
106 be presented.

107 2.1. *The Bayesian network, the underpinning data to inform the CPTs and variable instantiations for* 108 *the “initial” case*

109 In this section, we present the BN, its variables, their states and the data used to inform the CPTs.
110 In addition, we will use this BN with specific variable instantiations representing the circumstances
111 of a case. We have called it the “initial case”. Initially, it is given this set of circumstances that we
112 will explore how the variables of the BN impact on the LR values. The other two cases (high
113 transfer and low transfer) will be dealt with in a second step of the study.

114 The Object-Oriented Bayesian Network (OOBN), illustrated in Figure 8 of Taylor *et al.* [5] is
115 reproduced here in Figure 1 and will be used in our study. It was adapted to be easily used in
116 simulations.

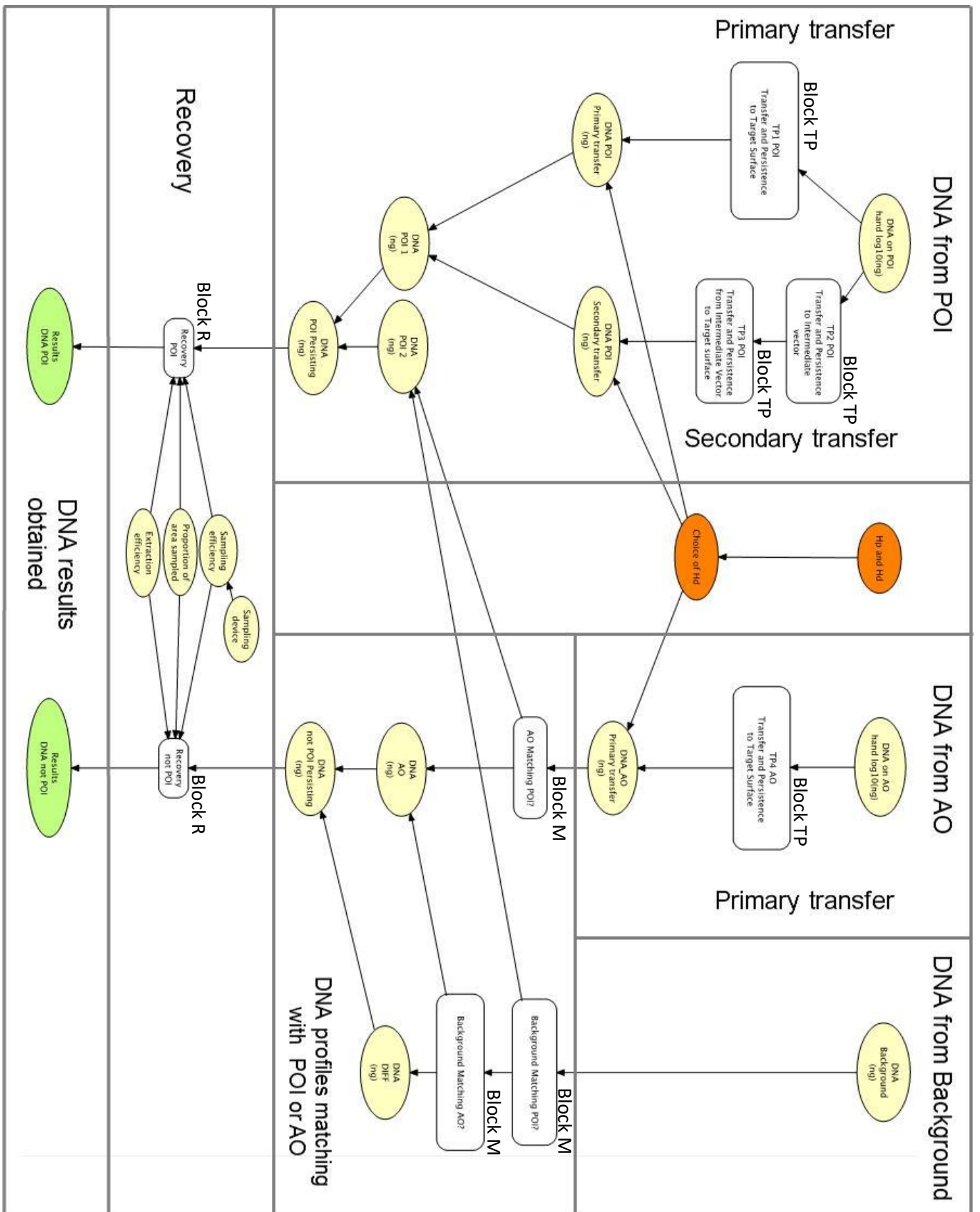


Figure 1: Bayesian network used to evaluate the findings under activity level propositions involving primary vs secondary transfer events adapted from [5].

120 Each TP block contains variables involved in transfer and persistence event, more specifically:

- 121 • Block TP1 contains variables involved in the transfer and persistence of POI’s DNA
122 from POI’s hand to the knife handle that occurs during the stabbing.
- 123 • Block TP2 contains variables involved in the transfer and persistence of POI’s DNA
124 from POI’s hand to AO’s hand that occurs during the handshake.
- 125 • Block TP3 contains variables involved in the transfer and persistence of POI’s DNA
126 from AO’s hand to the knife handle that occurs during the stabbing.
- 127 • Block TP4 contains variables involved in the transfer and persistence of AO’s DNA
128 from AO’s hand to the knife handle that occurs during the stabbing.

129 Variables included in each block called TP (for Transfer Persistence), M (for matching) and R
130 (for recovery) are described in Table 1. The BN has been developed using Hugin Researcher
131 (version 8.6, www.hugin.com).

132

Block	Variable
Block TP	Proportion of Contact
	Transfer proportion
	Nature of surface (for both target and shedding surface)
	Type of Contact
	Environment
	Days
	Decay factor of the loss of DNA (node named ‘alpha2’), function of Environment.
Block M	Match Probability
Block R	Sampling efficiency
	The proportion of area sampled
	Extraction efficiency

133 *Table 1: Variables taken into account in the blocks of the BN described in Figure 1.*

134

135 Regarding the states of the variables, they can be instantiated depending on the case information
136 available.

137 Table 2 details the states for each variable according to the case circumstances of the “initial
138 case” and the type of sampling used by the laboratory.

139

Variable	Possible states	Instantiated state for “initial case”	Explanation for the choice of the instantiated state
Nature of target surface	Hand Smooth Rough	Smooth	TP1-TP3-TP4: The surface of the knife handle is smooth.
		Hand	TP2: It’s a handshake
Nature of shedding surface	Hand Smooth Rough	Hand	TP1-TP2-TP3-TP4: The hand is the primary source of DNA
Vigour of contact	Passive Pressure Friction	Friction	TP1-TP3-TP4: It is assumed that type of contact is friction when you stab a person.
		Pressure	TP2: It is assumed that the type of contact is pressure when you shake hands.
Environment	Favourable Poor	Favourable	TP1-TP3-TP4: The target surface is the knife. It is assumed that the knife being kept in a paper bag by the crime scene investigators shortly after it was seized. The environment (the paper bag) is considered favourable in the sense that DNA will be preserved.
		Poor	TP2: The intermediate surface is a hand. It is assumed that this surface can be considered as an unfavourable environment because of the high risk of contact with other surfaces, resulting in a loss of DNA.
Sampling device	Tapelift Swab	Swab	A swab was used to take samples from the knife

140 *Table 2: Variables with their associated states and the instantiated state corresponding to the circumstances of*
141 *the “initial case”.*

142

143 Table 3 presents the non-instantiated variables and the data that will be used to inform their
144 CPTs.

Variable	Data informing the variable
Days	Data 1
Proportion of area sampled	Data 2
Transfer proportion	Data 3
Sampling efficiency	
Extraction efficiency	
DNA quantity on hands	
Background	

145 *Table 3: Non-instantiated variables and corresponding data used to inform their CPTs.*

146 The data that will be used to inform their CPTs are referred to as Data 1 to Data 3 are detailed
147 hereinafter.

148

149 The outcomes of the DNA analysis are set in two variables that gives the amount of DNA in
 150 ng, respectively for the POI and for not POI (Table 4).

Variable	Possible states (ng)
Results DNA POI	interval node; 0 to 0.1 in steps of 0.01 0.1 to 1 in steps of 0.1 1 to 5 in steps of 0.5 5 to 10 in steps of 1 10 to 25 in steps of 5 25 to 1000 1000 to inf
Results DNA not POI	interval node; 0 to 0.1 in steps of 0.01 0.1 to 1 in steps of 0.1 1 to 5 in steps of 0.5 5 to 10 in steps of 1 10 to 25 in steps of 5 25 to 1000 1000 to inf

151 *Table 4: Possible results for the amount of DNA in ng respectively for POI and not POI.*

152

153 *2.2. Representing our lack of knowledge in the BN*

154 Within the nodes in the BN there is uncertainty as to what the state of nature was in any
 155 particular instance. This can be thought of as comprising two distinct parts, the uncertainty
 156 surrounding the state that applies to a specific case, and the uncertainty surrounding the
 157 suitability of our data to model the world. The node ‘DNA on hands’ is a good example to
 158 further explain this idea.

159 In the stabbing scenario being considered there is uncertainty surrounding the amount of DNA
 160 on the hands of the POI (or AO) and we deal with this in the BN by treating the ‘DNA on hands’
 161 node as a distribution that is meant to reflect our uncertainty through the amount of DNA that
 162 the general population will possess on their hands. This distribution reflects our prior belief on
 163 the amount of DNA that the suspect (or AO) had on their hands at the time of the offence.
 164 Within the BN architecture the POI and AO have separate ‘DNA on hands’ nodes to reflect the
 165 fact that they are different people and can have different amounts of DNA on their hands (the
 166 posterior distribution of which would be obtained after instantiation of case information). In
 167 order to model the distribution of DNA on hands in the population, we must take a sample from
 168 the population and measure the amount of DNA on peoples’ hands. The second component of

169 our uncertainty is then how representative our sample is of the general population. It is only
 170 this second component of uncertainty that will benefit from additional research and sampling
 171 i.e. will potentially reduce the sensitivity of the LR to the data underlying these nodes. Consider
 172 a scenario where the results have very little power to distinguish between the propositions (and
 173 hence a LR close to one would be obtained). If this is due to the shape of the distributions used
 174 to model the population (and not how representative the sample is of the general population),
 175 then no amount of additional sampling will help to increase discrimination.

176 In our sampling schemes (which we list below) there are different aspects of case information
 177 and experimental data that we focus on during our simulations that represent both aspects of
 178 uncertainty.

179

180 Data 1: Days

181 To account for the uncertainty on the number of days between the stabbing and the crime scene
 182 attendance, and the number of days (or hours) between the stabbing and the alleged handshake,
 183 the variable “Days” is modelled by a Gamma distribution $Ga(\alpha, \beta)$ that has been discretized on
 184 a range of possibilities (from 0 to 31 days). The choice for a Gamma distribution allowed to
 185 easily model events from 0 to infinity. The parameters α and β (shape and rate) are calculated
 186 using [10] based on a mean and a variance (set by the circumstances) (Table 5). The variance
 187 is set to 0.5 in order to account for an uncertainty of about two days maximum. Note that in [5]
 188 we modelled the persistence of DNA (or the reduction of the amount of DNA over time) using
 189 an exponential decay curve set by a variable called ‘alpha2’. The parameter of the decay curve
 190 depends on the state of the variable “Environment” (either ‘poor’ or ‘favourable’) and is based
 191 on [Raymond et al, 2009]. The parameter ‘alpha2’ is set to 0.022 if the environment is
 192 favourable or to 0.052 if the environment is poor.

193

Variable	Discretized states	Mean and variance	Explanation	$Ga(\alpha, \beta)$ modelling the variable
Days	interval node; 0, 0.5 then 1 to 31 in steps of 1	$\mu=2$ $\sigma^2=0.5$	TP1-TP3-TP4: The item was examined around two days after the offence.	$Ga(8, 4)$
		$\mu=0.5$ $\sigma^2=0.5$	TP2: It is assumed that the handshake was made less than 12 hours but more than 2 hours before the stabbing.	$Ga(0.5, 1)$

194 *Table 5: Variable “Days” with its states and parameters (mean and variance and associated Gamma*
 195 *distributions) used to inform the CPTs.*

196

197 Data 2: Proportion of area sampled

198 This variable is modelled using a Beta distribution $Be(\alpha, \beta)$, whose parameters are estimated
 199 from a mean and variance for the proportion. It represents the proportion the touch surface that
 200 is sampled (using a swab or a tapelift). It allows to account for the fact that the whole (100%)
 201 of the touched surface may not be sampled. Note that it does account for the uptake efficiency
 202 of the swab or tapelift. The variability on the latter is accounted directly in the variable “Transfer
 203 proportion”. The parameters α, β are computed from mean and variance according to [11]. Table
 204 6 summarizes these parameters. As for “Days”, this distribution is used here to reflect an
 205 uncertainty regarding the circumstances of the case. If we had no doubt regarding them, it would
 206 be unnecessary to carry out such simulations.

207

Variable	Discretized states	Mean and variance of data	Explanation	Be(α, β) based on mean and variance
Proportion of area sampled	interval node; values from 0 to 1 in steps of 0.1	$\mu=0.95$ $\sigma^2=0.001$	We have assumed that almost the entire knife handle is sampled, representing about 95% of its surface.	Be(44.17, 2.32)

208 *Table 6: Variable “Proportion of area sampled” with the explained mean and variance associated to the data*
 209 *used to inform the CPTs, with the parameters of the beta distribution based on the mean and variance.*

210

211 Data 3: Transfer proportion, Sampling efficiency, Extraction efficiency, DNA quantity on
 212 hands and Background

213 We have adopted a full Bayesian strategy to inform the parameters associated with these
 214 variables. It means that we have initially set a prior probability distribution for the variables.
 215 Then, based on data from the scientific literature, we have updated these distributions, leading
 216 to posterior distributions that will be used to inform the CPTs associated with them.

217 The prior distributions set for each variable are presented in Table 7.

218 For the variables “Transfer proportion”, “Sampling efficiency” and “Extraction efficiency”, a
 219 so-called flat prior $Be(1,1)$ has been chosen to start from a uninformed situation between 0 and
 220 1.

221 The variable “DNA quantity on hands” is transformed in \log_{10} (from $-\infty$ to $+\infty$) with a prior
222 distribution based on a mean quantity of 2ng of DNA and a large variance (variance of 1). The
223 mean quantity of 2ng is what appears to be a reasonable amount of DNA. A large variance was
224 chosen to account for the paucity of data available at this stage. Note that the above method
225 with these parameters led to only positive values.

226 The variable “Background” is modelled using a Gamma distribution whose parameters have
227 been estimated [10] from a mean quantity of background DNA (0.5 ng) and an associated
228 variance (variance of 5). The mean of 0.5 ng is viewed as a reasonable upper bound quantity.
229 The large variance has been chosen to be large reflecting the paucity of data available. Ten prior
230 observations have been drawn from that Gamma distribution to act as our prior data counts. We
231 will then update these counts with the data obtained from the literature. 10 was chosen to reflect
232 the paucity of the sample sizes available in the literature. Indeed, we couldn’t claim more in the
233 prior counts than what is actually published.

234 The data used to update the above prior distribution for the variable “Transfer proportion” are
235 based on [5] where the authors estimated the parameters of the Beta distribution modelling
236 “Transfer proportion” based on simulations from the original data from Daly [12], Bontadelli
237 [13] and Goray [14]. To reflect the fact that a laboratory will conduct only a limited number of
238 experiments to inform that variable, we randomly selected 51 data from the beta distribution
239 obtained in [5] to act as the dataset that we will use to update our prior distribution. 51 is the
240 number of experiments done by Daly [12].

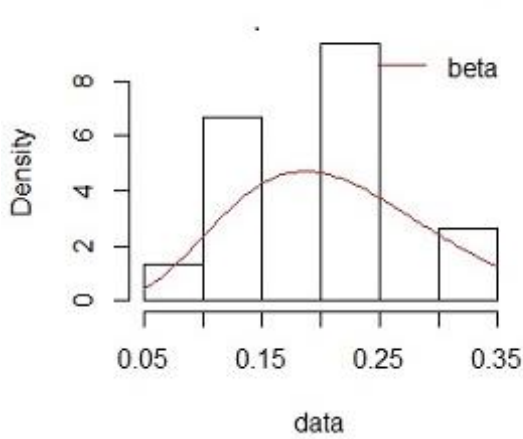
241 For the variable for the “Sampling Efficiency”, we will distinguish the “swab” from the
242 “Tapelift”. In [5], the parameters $Be(25, 20)$ of the beta distribution representing the data for
243 “Swab” condition were based on [15]. As before, we will randomly draw a limited sample (21)
244 from that distribution to carry out the Bayesian update of our prior distribution. Indeed, 21
245 experiments were done in [15]. For the “Tapelift” condition, we have used directly the data
246 from [15].

247 For the variables “Extraction efficiency”, “DNA quantity on hands” and “Background”, we
248 have used directly the data from the literature (Table 7). We have used these data to fit,
249 respectively, a Beta, a Normal and a Gamma distribution (Figure 2). Note that only 15 data
250 points from [17] have been used to inform the variable “Extraction efficiency”, hence the poor

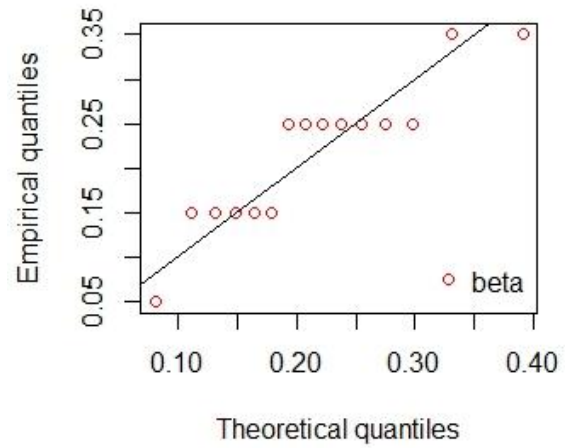
251 quality of the fit (Figure 2a). It will serve as a starting point, keeping in mind that this variable
252 will probably be flagged up following the simulations as a variable requiring more data to
253 inform it.

254 (a) Extraction efficiency

Histogram and theoretical densities



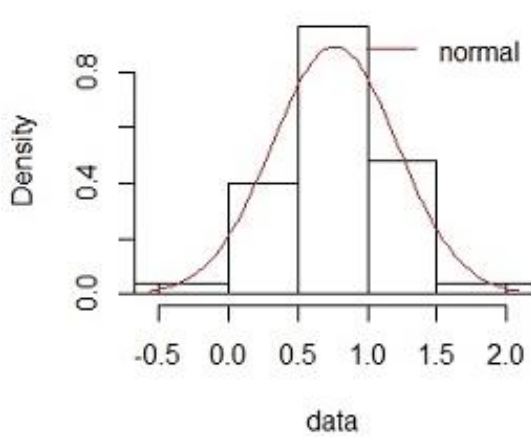
Q-Q plot



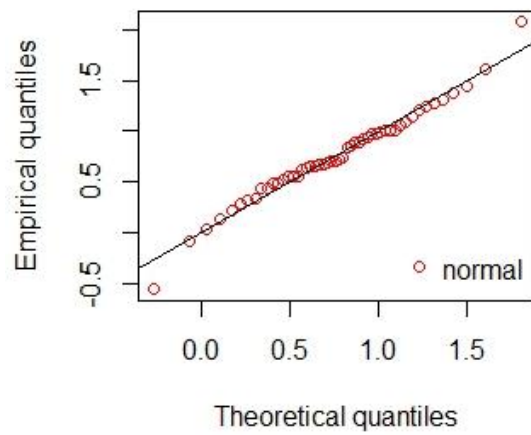
255

256 (b) Log10 of Quantity of DNA on hands

Histogram and theoretical densities



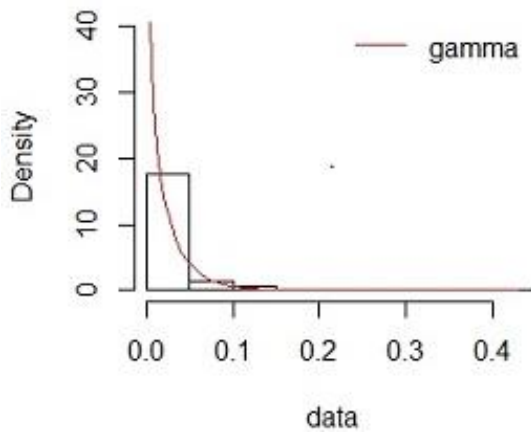
Q-Q plot



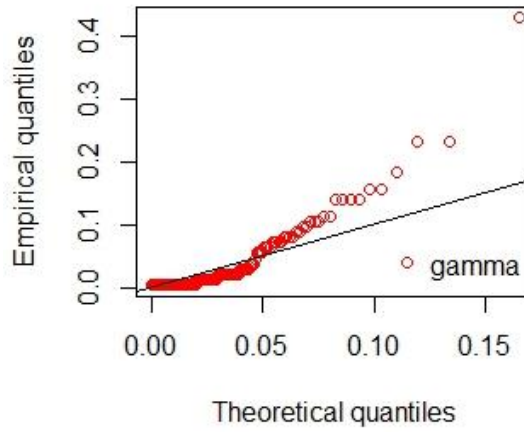
257

258 (c) Background

Histogram and theoretical densities



Q-Q plot



259

260

261

Figure 2: Histogram and theoretical densities and Q-Q plots for data of the variables “Extraction efficiency” (a), “DNA quantity on hands” (b) and “Background” (c).

262

Beta, Normal and Gamma distributions were updated using standard Bayesian methods [16]

Variable	Discretized states	Prior distribution	Data Source	Number of data points	Posterior distribution
Transfer proportion	interval node; from 0 to 1 in steps of 0.05	Be(1,1)	51 random sample from the Beta distribution defined in [5]	51	Hand(rough)/passive: Be(1.89,3.83) Hand(rough)/pressure: Be(1.99,2.78) Hand(rough)/friction: Be(2.02, 2.98) Smooth/passive: Be(1.55,32.65) Smooth/pressure: Be(1.73,33.43) Smooth/friction: Be(1.43, 2.10)
Sampling efficiency	interval node; from 0 to 1 in steps of 0.05	Be(1,1)	State “Tapelift”:[15] State “Swab”: 21 random sample from beta distribution from [5]	21	Swab: Be(14.36,11.79) Tapelift: Be(3.15,17.11)
Extraction efficiency	interval node; from 0 to 1 in steps of 0.05	Be(1,1)	[17]	15	Be(5.79,18.43)
Log10(DNA quantity on hands)	interval node; -inf to -1.5 -1.5 to 3.5 in steps of 0.1 3.5 to inf	N(0.3,1)	[13]	50	N(0.764,0.004)
Background	interval node; 0 to 0.1 in steps of 0.01 0.1 to 1 in steps of 0.1 1 to 5 in steps of 0.5 5 to 10 in steps of 1	10 data randomly selected from the distribution Ga(0.08,0.16)	[18]	301	Ga(0.6, 30.15)

263 *Table 7: Variables, discretized states, prior distribution, data source to update them and posterior distribution*
264 *used to inform their CPTs.*

265 2.3. *Methods used to perform simulations from the “initial case”*

266 In the previous section, we have presented the BN that captures the “initial case”. As such it
267 can be used to compute a LR for any DNA outcome (Table 4) in terms of quantity of DNA
268 corresponding to the POI and non-corresponding to the POI. For a given outcome, we will
269 obtain one LR that encapsulates the knowledge that has been used to inform the CPTs.

270 Using the Bayesian update mechanism presented before (under Data 3), we have used specific
271 data sets from published studies. The numbers of data points used in these studies are rather
272 sparse and can be seen as small subsets from larger and unknown populations. Through
273 simulations, we would like to show the impact (if any) on the LR of such limited samples. To
274 do that, we have resampled with replacement the data points that have been used to carry out
275 the Bayesian update. The simulation method is detailed later in this section. At each simulation
276 then, the above-described posterior distributions are recomputed, the BN CPTs are updated and
277 LRs obtained for any DNA outcome. We carried out that task 100 times. Hence for a given
278 DNA outcome, we have now 100 LRs. These LRs will have a range that can be characterized
279 [e.g. IQR, min to max]. Note that these 100 LRs represents 100 slightly different scenarios.

280 A first set of results out of these simulations will show these ranges (one per possible DNA
281 outcome and for each of the defence propositions). They reflect how the LRs vary based on the
282 state of knowledge and understanding for sets of data points used to inform the CPTs.

283 The second question we will address through simulation is to identify which variables impact
284 the most on the LR. This identification of impacting variables will allow prioritisation of further
285 data acquisitions. That approach stems from the realisation that, given the complexity of the
286 problem, we cannot expect systematic acquisition of large datasets for all the variables
287 identified in the BN. To make that selection, we will carry out resampling on a variable by
288 variable basis (keeping all the other variables constant). For each set of simulations (100
289 simulations per variable), we will measure the ranges (for each DNA outcome and considering
290 each of the defence propositions). These simulations on a variable per variable basis allow
291 pinpointing of the variables that have the most impact on the LRs. These shall then constitute
292 the focus for further data acquisition, because additional datasets have the potential to reduce
293 the observed ranges of LRs and improve on their robustness.

294 The simulations were performed in Rstudio (Version 1.1.463) [19] with R (version 3.5) [20]
 295 combined with RHugin (version 8.4) [21]. The library RHugin specifically allows to liaise
 296 directly with the Hugin inference engine. It allows then to load the BN as in Hugin and to
 297 interact directly with it without resorting to the Hugin GUI. It means that all the captured
 298 dependencies between the variables in the BN are duly maintained and are part of the
 299 computation.

300 Table 8 presents the variables that will be used for the simulation exercise (either jointly or
 301 separately) with an indication of the method of simulation that will be applied. These three
 302 simulation methods are described below.

303

304

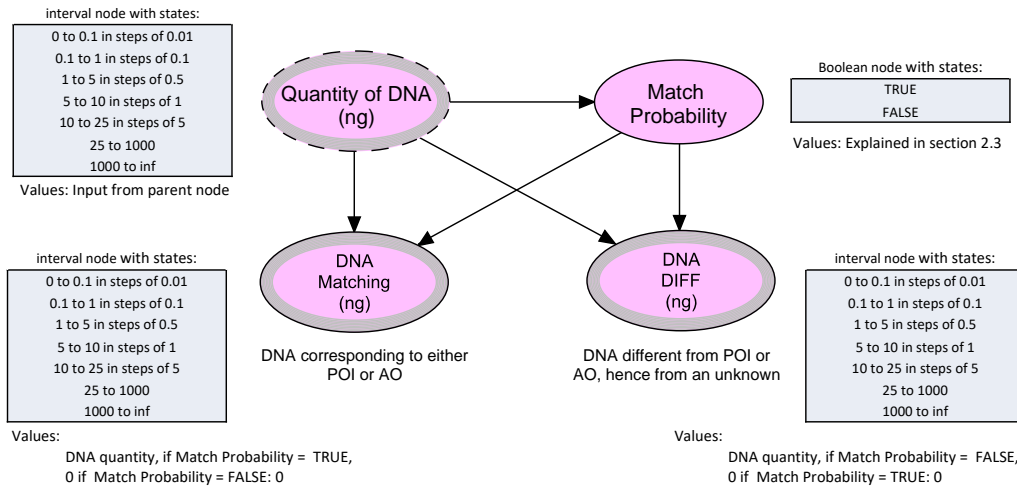
Variable	Simulation method
Days	1
Proportion of area sampled	
Transfer proportion	2
Sampling efficiency	
Extraction efficiency	
DNA quantity on hands	
Background	3
Match probability	

305 *Table 8: Each variable associated with a method of simulation.*

306 *Simulation method 1:* For each simulation, the mean for each variable (“Days” and “Proportion
 307 of area sampled”) is resampled from a Normal distribution with a mean set as per the initial
 308 case but allowing a variance, respectively of 0.5 and 0.01 around it. For instance, the mean for
 309 the number of days in the TP1 block is 2. For each simulation then, the mean will be obtained
 310 by randomly selecting from a sample from a $N(2,0.5)$, until the mean is positive (For
 311 “Proportion of area sampled”, the mean is resampled until a value between 0 and 1 is obtained).
 312 From that mean, and keeping the variance constant, we estimate, as before, the parameters of
 313 the Gamma distribution (for “Days”). The values for the variance were chosen because, in the
 314 authors’ opinion, they adequately reflect the amount of uncertainty surrounding the timeframes.

315 *Simulation method 2:* For each of these variables that were subjected to the Bayesian update,
 316 the simulations are carried out by resampling (with replacement) from original data presented
 317 in the previous section.

318 *Simulation method 3*: The block M of the BN, presented in Figure 3, shows a dependence
 319 between the node “Match probability” and the node “Quantity of DNA” as in [5].



320

321 *Figure 3: Block M with its four nodes and associated states and values.*

322 In that block, only the node “Match Probability” requires input probabilities (the node “Quantity
 323 of DNA” being set by its parent). Its state is TRUE when the DNA profile is matching, FALSE
 324 otherwise. The probability of being TRUE is set by the match probability which in turn depends
 325 on the quantity of DNA. Linked to the quantity of DNA is the number alleles that the profile
 326 will show. To compute the match probability, we have accounted for that relationship between
 327 the quantity of DNA and the number of alleles.

328 Hence, we will first establish the relation between quantity of DNA and the number alleles and
 329 then propose a way to compute match probability as a function of the number of alleles. The
 330 number of alleles corresponding to a given quantity of DNA is modelled based on the empirical
 331 observations [22] made between the quantity of DNA and the minimum and maximum numbers
 332 of detected alleles (Table 9). A Gamma distribution is used to represent the numbers of detected
 333 alleles as a flexible modelling distribution for counts between 0 and infinity. The parameters of
 334 the Gamma distribution are informed on the mean and variance obtained from [22]. These
 335 Gamma distributions will be used at each simulation to generate a given number of alleles
 336 corresponding to a given quantity of DNA (sampled from its own distribution).

337

Quantity (ng)	Minimum to maximum numbers of detected alleles	Mean and variance	Ga(κ, θ) modelling the Number of alleles
0-0.01	0 to 1	$\mu=1; \sigma^2=0.1$	Ga(10, 10)

0.01-0.02	0 to 2	$\mu=1; \sigma^2=0.4$	Ga(2.5, 2.5)
0.02-0.03	1 to 3	$\mu=2; \sigma^2=0.4$	Ga(10, 5)
0.03-0.04	2 to 4	$\mu=3; \sigma^2=0.4$	Ga(22.5, 7.5)
0.04-0.05	3 to 5	$\mu=4; \sigma^2=0.4$	Ga(40, 10)
0.05-0.06	4 to 6	$\mu=5; \sigma^2=0.4$	Ga(62.5, 12.5)
0.06-0.08	5 to 7	$\mu=6; \sigma^2=0.4$	Ga(90, 15)
0.08-0.1	6 to 8	$\mu=7; \sigma^2=0.4$	Ga(122, 17.5)
0.1-0.2	9 to 22	$\mu=15; \sigma^2=1.5$	Ga(150, 10)
0.2-0.3	20 to 30	$\mu=25; \sigma^2=1.4$	Ga(446, 17.9)
0.3-0.4	29 to 32	$\mu=31; \sigma^2=0.5$	Ga(1922, 62)
>0.4	32	the number of alleles is 32	the number of alleles is 32

338 *Table 9: Range of quantities of DNA (ng) with the associated minimum and maximum numbers of alleles observed*
339 *empirically. For each number of alleles (min-max), a mean and variance is set to reflect these ranges and the*
340 *Gamma distributions parameters are obtained.*

341 To obtain the match probability for a given number of alleles (associated with a given quantity
342 of DNA), we proceeded as follows:

- 343 (1) 5 million full DNA profiles (32 alleles in total) are randomly generated based on allelic
344 frequencies of the NGMSelect kit [23] for the Swiss Caucasian population.
- 345 (2) From the above profiles, 27046 partial DNA profiles are created using the allelic
346 degradation model from Hicks et al [24]. It means that for each number of alleles (1 to
347 31), we have a collection of partial DNA profiles (from 6 to 1000) depending on the
348 number of alleles.
- 349 (3) For each of these profiles, their match probability (MP) is computed with a θ of 0.02
350 using the allele frequencies from [23].
- 351 (4) For a given number of alleles, at each simulation, we draw the match probability from
352 the collection of match probabilities associated with the drawn number of alleles
353 (corresponding to a given quantity of DNA).

354
355 During the simulation process, the number of alleles and the associated MPs are resampled only
356 if this node is taken into consideration. Otherwise, its CPTs is set once for all simulations.

357 2.4. Method used to perform simulations beyond the “initial case”

358 The “initial case” represents the circumstances of the case, typically set by the chosen states in
359 the nodes that have been instantiated. The BN also allows us to explore other sets of
360 circumstances. Hence, we can repeat the simulation process to explore how each node impacts

361 on the LRs under any set of circumstances. We have chosen to investigate two additional
 362 scenarios deviating from the conditions of the “initial case”. The first will adopt circumstances
 363 that will favour the transfer of the POI DNA (a rough surface to increase the deposition, a short
 364 time delay between the stabbing and the collection of the swabs, and favourable environmental
 365 conditions). The second adopts circumstances that are less favourable to the transfer of DNA.
 366 The states of the nodes for each scenario are presented in Table 10. The other nodes related to
 367 the case circumstances were kept constant with the same states as for the “initial case”.

368

Node	“Initial case” Instantiated state or mean	“High transfer case” Instantiated state or mean	“Low transfer case” Instantiated state or mean
Nature of target surface	TP1-TP3-TP4: Smooth TP2: Hand	TP1-TP3-TP4: Rough TP2: Hand	TP1-TP3-TP4: Smooth TP2: Hand
Days between both transfers or transfer and recovery	TP1-TP3-TP4: 2 TP2: 0.5	TP1-TP3-TP4: 0.5 TP2: 0.5	TP1-TP3-TP4: 20 TP2: 0.5
Environmental conditions	TP1-TP3-TP4: Favourable TP2: Unfavourable	TP1-TP3-TP4: Favourable TP2: Unfavourable	TP1-TP3-TP4: Unfavourable TP2: Unfavourable

369 *Table 10: Choice for node instantiations and mean for the “initial case”, for the “high transfer case” and for the*
 370 *“low transfer case”.*

371 2.5. Method used for the analysis of the simulation results

372 Following a set of simulations (one for each of the initial, the high transfer and the low transfer
 373 case), we obtain 100 LRs for each combination of results (quantity of POI and not POI DNA –
 374 36 x 36 possibilities), each proposition retained for the defence (Hd₁ and Hd₂) and for each
 375 node considered (10 nodes and the case with all simulated nodes considered jointly). It
 376 represents a total of 28,512 combinations and 2,851,200 LRs.

377 A dedicated Shiny application (https://lydie-samie.shinyapps.io/DNA_Activity/) has been
 378 designed to allow the visualisation of these results for each possible combinations of variables.

379 To explore which node (or variable) has impact on the LRs, we ordered them by range. That
 380 will be done by conditioning on the defence proposition. It is important to stress that with this
 381 approach, we aim at identifying the variables that have the most impact across all possible
 382 outcomes. The variable by variable analysis allows us to identify which variables contribute
 383 significantly to the whole. By significant, we mean that the range of log₁₀ (LRs) (meaning the
 384 difference between the maximum log₁₀ (LR) and the minimum log₁₀ (LR)) produced by

385 sampling a variable exceeds 1 order of magnitude of \log_{10} (LR). The range was chosen to cover
386 up to extreme situations even if their probability of occurrence is low. At this stage, the issue
387 to screen among variables with a rather wide net, instead of being too strict in their selection by
388 adopting a more stringent criterion such as the interquartile range.

389

390 2.6. *Performing simulations adapting the number of data points informing the conditional* 391 *probability tables*

392 As proposed in [7] we will mimic an increase of knowledge, to do so we have increased the
393 counts informing the CPTs by a constant factor (retaining the observed proportions).

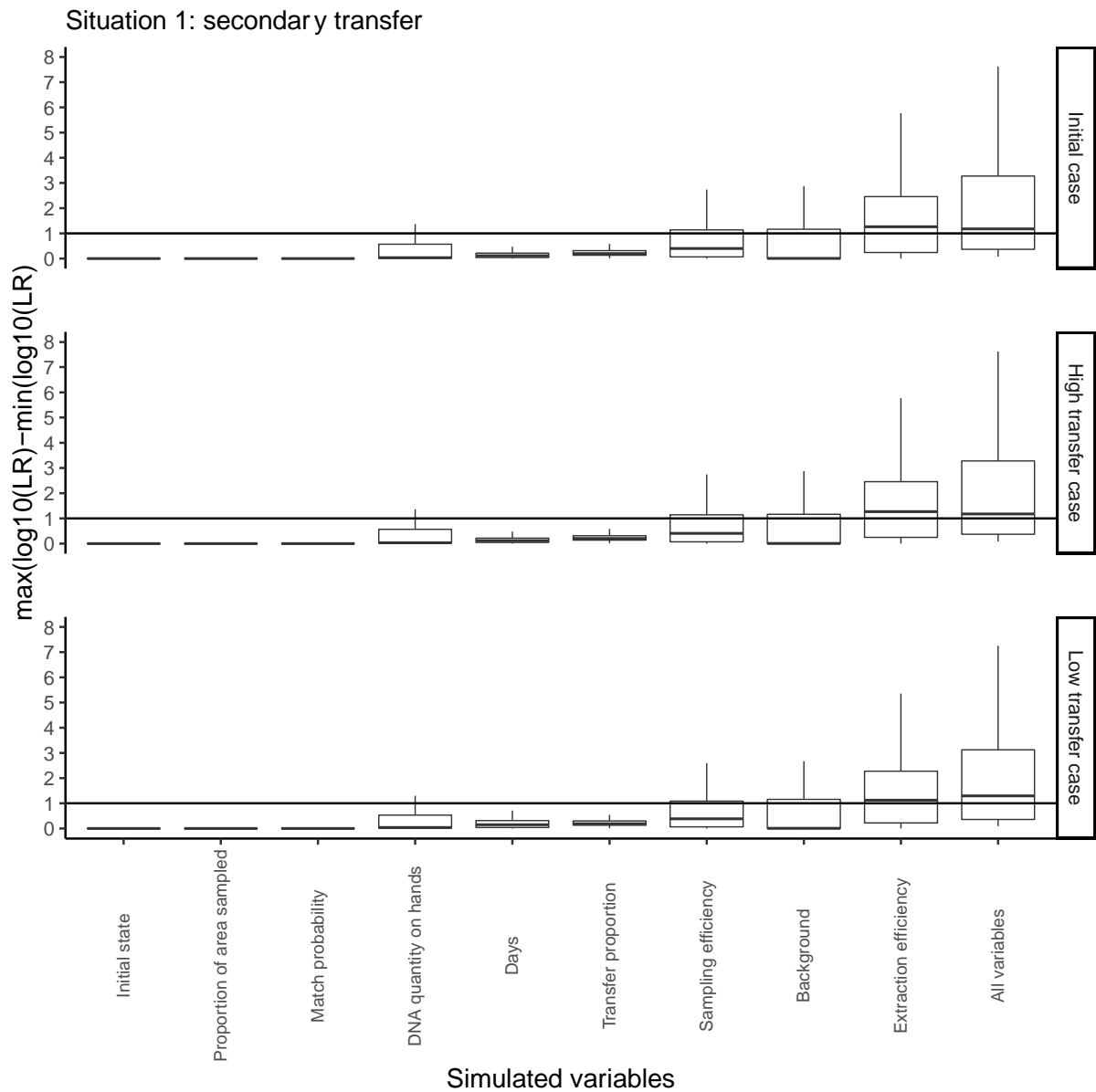
394 The variables (“Background”, “DNA quantity on hands”, “Sampling efficiency”, “Extraction
395 efficiency” and “Transfer Proportion”) that proved to be significant based on the method
396 described in section 2.4, have been studied using simulations. To simulate an increase in the
397 size of datasets, we have adapted the simulation process (*Simulation method 2*) by tripling the
398 number of data points resampled and then used to inform the corresponding CPTs. The original
399 numbers of datapoints used for each variable are given in Table 7. It means that the relative
400 proportions associated with each state of each variable remain the same, but the counts are
401 multiplied by 3 as if the study had been conducted on a larger sample. The factor of 3 is what
402 we considered reasonable for an operational laboratory.

403 **3. Results and Discussion**

404 3.1. *Ranges of simulated LR_s obtained under Hd₁*

405 Hd₁ stipulates that the defence alleges that the DNA corresponding to the POI’s profile is the
406 consequence of a secondary transfer. The DNA corresponding to POI’s profile is the
407 contribution of the POI following the secondary transfer and the DNA different from POI’s
408 profile is to the potential joint contribution of the background and the DNA of the alternative
409 offender (AO). Figure 4 illustrates the ranges of \log_{10} (LRs) for all outcomes (i.e. all considered
410 amounts of POI and not POI’s DNA, each possibility gives a range of LR_s after 100
411 simulations) considering respectively all variables jointly and then each variable simulated in
412 turn. The results are shown for the three cases considered: “initial”, “high transfer” and “low
413 transfer”). Regardless of the case, the significant variables (with an impact of more than an
414 order of magnitude) are: “DNA quantity on hands”, “Extraction efficiency”, “Background”,

415 and “Sampling efficiency”. Two of these variables “Extraction efficiency” and “Sampling efficiency”
 416 “efficiency” relates to the laboratory choices regarding their sampling devices and extraction
 417 techniques. The results obtained for the three cases are similar, particularly the results when
 418 “Initial case” and “High transfer case” are considered.



419

420 *Figure 4: Boxplots presenting the ranges (min-max) expressed in log₁₀ of the LRs obtained following 100*
 421 *simulations under Hd₁. The panels present the results for the three cases considered. The horizontal line drawn at*
 422 *the difference of 1 (in log₁₀) set the limit above the variable considered will be declared as having a significant*
 423 *effect on the global variability shown when “all variables” are resampled jointly.*

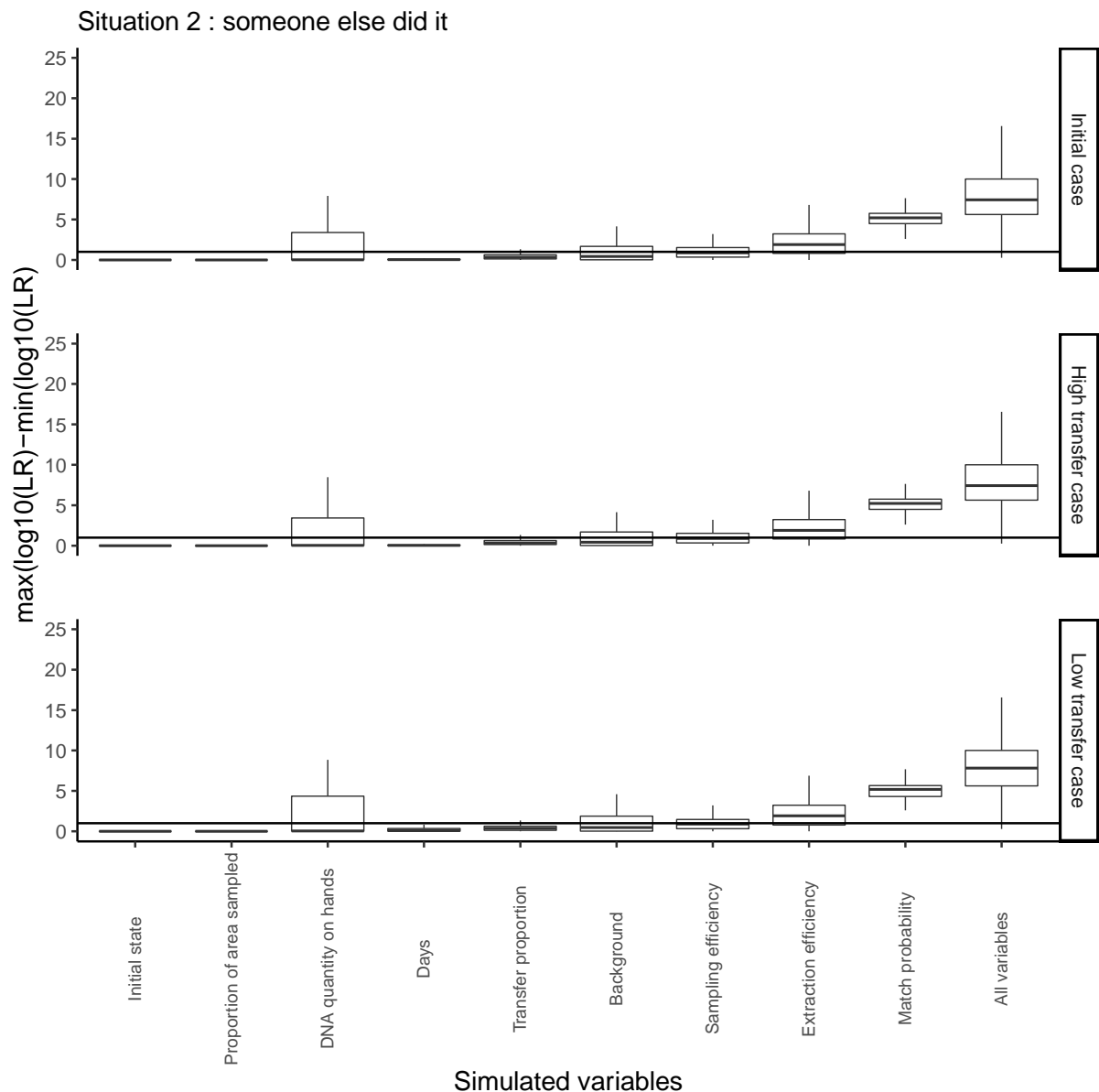
424 Hence, when secondary transfer is alleged, the current state of knowledge regarding the transfer
 425 and persistence of DNA is sparse and induces a large variability on the simulated LRs. One way
 426 to overcome this problem would be to increase the underpinning data (see section 3.4 below).
 427 It is rather easy for a laboratory to increase its knowledge base associated with some of the

428 variables involved (see [25] for example regarding the acquisition of data in relation to sampling
429 and extraction efficiency).

430 3.2. *Ranges of simulated LR_s obtained under situation 2 (H_{d2})*

431 In situation 2 the defence alleges that the DNA corresponding to POI's profile is due to the
432 contribution of an unknown alternative offender (AO) (H_{d2}). In this situation, the possibility
433 for a secondary transfer for the POI is not retained. Figure 5 illustrates the ranges of log₁₀ (LR_s)
434 plotted considering the same variables as before. The variable by variable analysis allows to
435 identify the variables contributing to more than one order of magnitude. They are the following:
436 "DNA quantity on hands", "Match probability", "Extraction efficiency", "Proportion of
437 transfer", "Sampling efficiency" and "Background", regardless of the scenario considered.

438



439

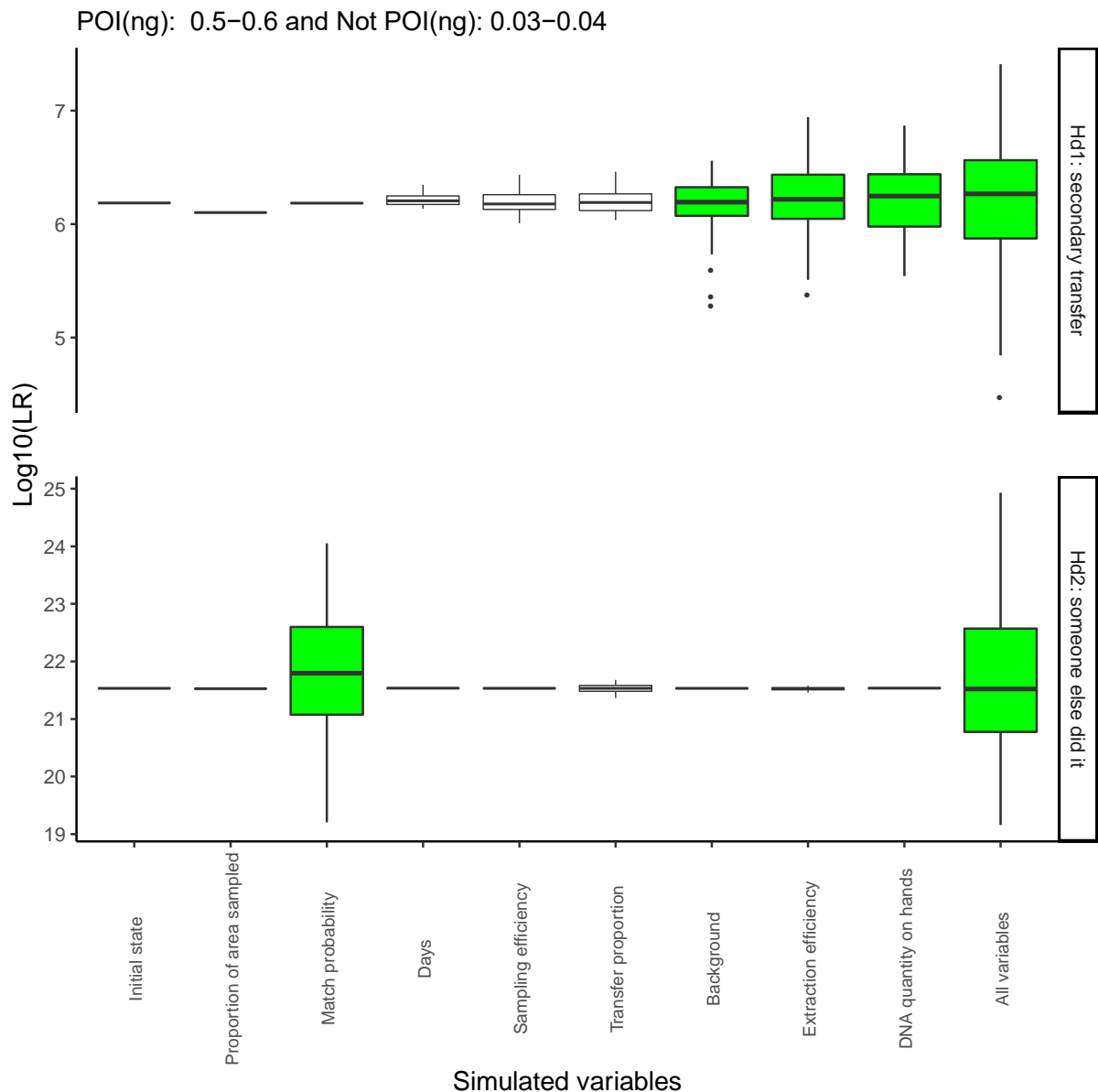
440 *Figure 5: Boxplots presenting the ranges (between min-max) expressed in log₁₀ of the LRs obtained following*
 441 *100 simulations under Hd₂. The panels present the results for the three cases considered. The horizontal line drawn*
 442 *at the difference of 1 (in log₁₀) set the limit above the variable considered will be declared as having a significant*
 443 *effect on the global variability shown when “all variables” are resampled jointly.*

444 In addition to the variables identified for Hd₁, the variables “Match probability” and “Transfer
 445 proportion” comes into play under Hd₂. This is due to the fact that under that defence
 446 proposition, the DNA corresponding potentially to the POI may arise from the background level
 447 of DNA on the surface. That was much less critical under Hd₁ as the secondary transfer was
 448 dominating the considerations compared to the background.

449 The above results are showing the most impacting variables in a global sense, regardless of the
 450 outcome observed in a given case (i.e. the respective amounts of DNA corresponding to POI
 451 and to the not POI). It helps us to identify the variables that ought to receive attention if we

452 want to treat the problem globally considering all possible outcomes. Also the case (initial, high
453 transfer and low transfer) considered as no impact on the impacting variables.

454 However, in a given case with observed outcomes, the variables that will have the most impact
455 will vary and could be different from the above globally selected variables. To illustrate this
456 point, we present one example (with 0.5-0.6 ng of POI's DNA with 0.03-0.04 ng of not POI's
457 DNA) leading to a different selection of the most contributing variables (Figure 6). The reader
458 can refer to the Shiny application to select any set of outcomes and explore the most impacting
459 variables. In this chosen particular case, under Hd₁, the results show that the most significant
460 variables are "DNA quantity on hands", "Background" and "Extraction Efficiency", whereas
461 under Hd₂, the most impacting variable is the "match probability". As before, a variable is
462 declared to be significant (boxplots coloured in green in Figure 6) if the range (the whole height
463 of the boxplots shown in Figure 6) is above 1 order of magnitude of log₁₀ (LR).



464

465 *Figure 6: Boxplots of the LR_s (expressed in Log₁₀) obtained after 100 simulations of each variable in both*
 466 *situations (Hd₁ and Hd₂) with a quantity of POI's DNA between 0.5 and 0.6 ng with between 0.03 and 0.04 of not*
 467 *POI's DNA. The boxplots corresponding to the significant variables are coloured in green.*

468 The purpose of setting up a rather complex simulation regime was twofold: (1) facing a large
 469 networks of connected variables, the simulations allow us to identify where there is a
 470 knowledge gap by identifying the most impacting variables and (2) to appraise the range of
 471 variations the computed LR_s may take given the limited data points that have informed each
 472 conditional probability table in the Bayesian network.

473 For the example in Figure 6, we note that when considering Hd₁ the variables significantly
 474 impacting are only 2 and relates to transfer mechanisms. Under Hd₂, the match probability
 475 variable is dominating on all the other variables, making variables in relation to transfer less
 476 impacting in favour of the variable in relation to the source of the DNA. This is not surprising

477 as Hd₂ stipulated that the POI has no previous encounter with the knife and should his/her DNA
478 profile be present, it must be from the AO whose profile, by chance, matches the POI or some
479 background DNA which, by chance, matches the POI. Given that we are considering a
480 relatively large amounts of DNA for the POI (0.5-0.6 ng), the chance of an adventitious
481 correspondence is very small. The range of LRs observed stems from the fact that the simulation
482 process accounts for the distribution of the number of alleles that could be derived from that
483 quantity of DNA and for each of them the match probabilities that could potentially be
484 associated. In other words, it accounts for a large range of potential matching DNA profiles
485 each of them having different match probabilities depending on the number of alleles and their
486 designations. In this paper the match probabilities are used as a plug-in to explore the relative
487 importance of the variables, they do not guide precisely the sub-source level likelihood ratio.
488 They are reasonable in terms of order of magnitude for single profiles, but they fail to account
489 for the impact of mixtures on the sub-source level likelihood ratio. In a given case, with a
490 specified DNA allelic profile, that variability will be extremely reduced as the match probability
491 node will be directly informed with the sub-source level likelihood ratio computed for the case
492 at hand. It will lead to different overall likelihood ratio compared to the values reported in
493 Figure 6. When mixtures are involved, the node “Match probability” will be informed by 1 over
494 the LR obtained for that mixture considering sub-source propositions.

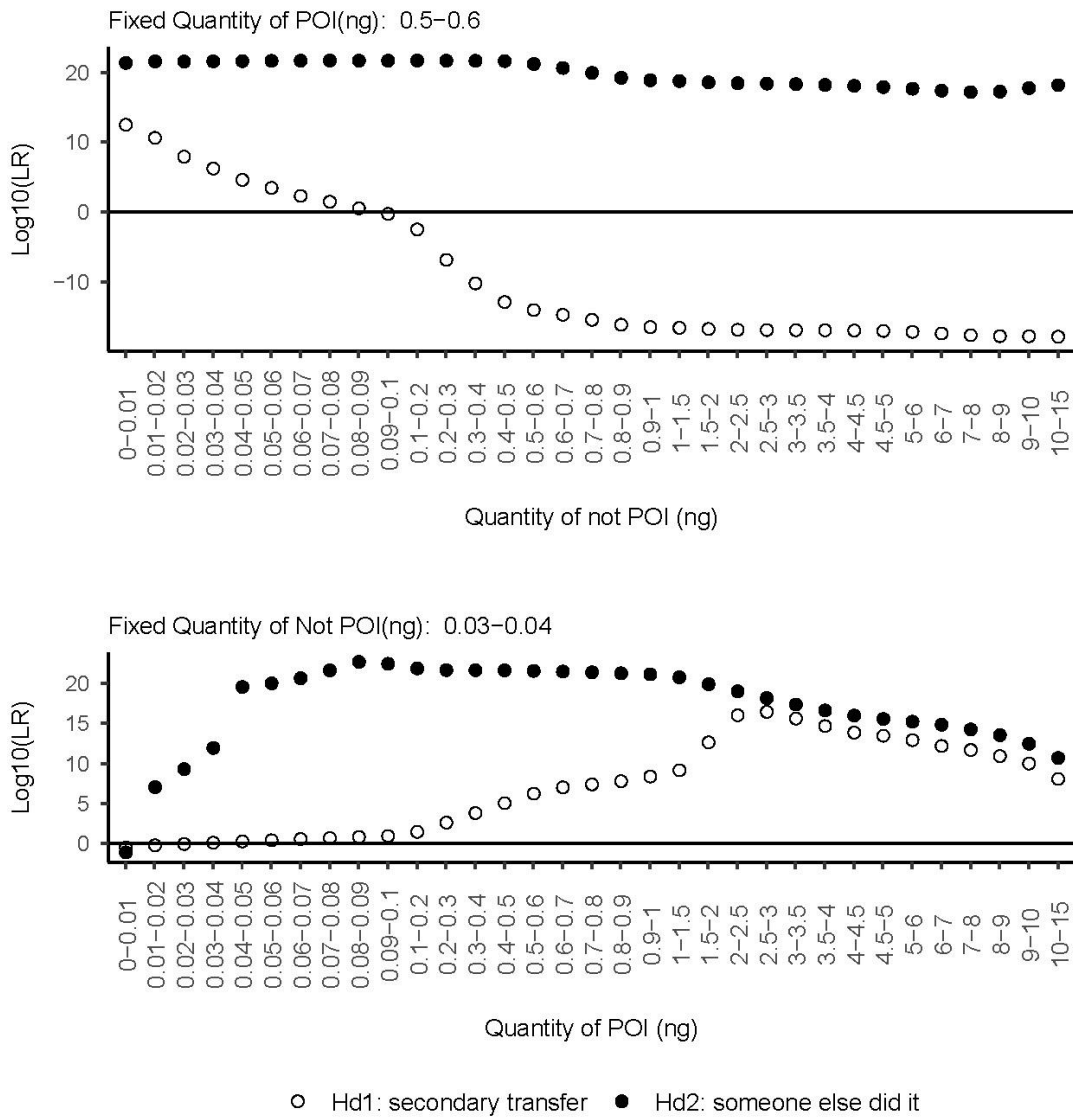
495 In this example, when considering Hd₁ and Hd₂, the overall range of the LRs is of 2 orders of
496 magnitude but for a median log₁₀(LR) of about 6 under Hd₁ and about 22 under Hd₂. The 2
497 orders of magnitude have then a greater impact under Hd₁ rather than Hd₂ with regards to how
498 it might influence the decision making of the recipient. We shall see below (in section 3.4) that
499 one way to reduce the range is to increase the number of data points informing the impacting
500 variables, that would be worth the investment for a case involving Hd₁ and not for a case
501 involving Hd₂.

502 Exploring the results in the Shiny application allows one to observe that, for some specific cases
503 (a given amount of POI’s DNA and not-POI’s DNA), other variables may have a higher impact,
504 hence be more critical to the case. That observation calls for a case specific approach if needed
505 once a specific outcome has been observed in a given case. The general trend though is given
506 by the variables identified globally.

507

508 3.3. Impact on the LR of the quantities of DNA, POI and not POI respectively

509 To show the respective impact of the quantity of POI's DNA and not POI's DNA on the LR
 510 obtained (summarized by their median), we can show two situations linked to the above
 511 example. In the first, we will maintain the quantity of not POI's DNA (0.03-0.04 ng) and vary
 512 the amount of POI's DNA. In the second, we will keep fixed the amount of POI's DNA (0.5-
 513 0.6 ng) and vary the not POI's DNA quantity. Both are shown in Figure 7.



514

515 *Figure 7: LR's obtained when the amount of DNA (POI and non-POI) are respectively varied for a given amount*
 516 *of DNA that remain fixed.*

517

518 When the quantity of DNA of the POI is kept fixed (top graph of in Figure 7), under Hd₂, the
519 presence of non-matching DNA has little impact on the LR. It is the fact that Hd₂ stipulates that
520 the POI has no link whatsoever of with the knife, hence the probability that some background
521 DNA or AO DNA, would correspond to the POI, is the key consideration. The amounts have
522 almost no effect. Under Hd₁, however, the LR will be above 1 (the findings will provide support
523 for Hp) with a maximum above one billion, but lower than 1 over the match probability (value
524 under Hd₁)¹. When the quantity of not POI's DNA is increased, the LR drops gradually. The
525 clinching point (when the LR is equal to 1 or 0 on a log₁₀ scale) is with 0.07–0.08 ng of non-
526 matching DNA. Above that amount, the non-matching DNA becomes more compatible with
527 the stabbing activity, hence the findings overall would lend support for the defence.

528 When the quantity of DNA of the not POI is kept fixed and we increase the amount the quantity
529 of DNA corresponding to the POI (bottom graph of in Figure 7), under Hd₂, the likelihood ratio
530 gradually increases with the increased amount DNA corresponding to the POI. The maximum
531 LR is obtained when the quantity is the most expected when handling a knife (0.07–0.08 ng).
532 Above that quantity, the LR will gradually reduce. Under Hd₁, the LR will increase only when
533 the quantity of DNA corresponding to POI reaches the point that it is more compatible with the
534 POI stabbing scenario than under the POI handshaking scenario. That clinching point is at the
535 same quantity as the quantity of not POI's DNA (0.03–0.04 ng). That is logical because we
536 have modelled the transfer probabilities in the same way for both POI and AO. Then, adding
537 quantity of DNA corresponding to the POI will gradually increase the LR up until the quantity
538 that is best expected under the primary transfer and not under the scenario of a secondary
539 transfer. Then it will reduce again. The LRs lending support for the defence spans over a large
540 range of POI corresponding DNA quantities despite the presence of non-matching DNA. It
541 stems from the fact that the non-matching DNA is a quantity that is more compatible with
542 background level than with a quantity you would expect following primary transfer.

543 The above considerations will change as a function of the choice of quantities, hence in the
544 Shiny application, the one can find the same representations as in Figure 7 but for any given
545 choice of POI or not POI quantities of DNA.

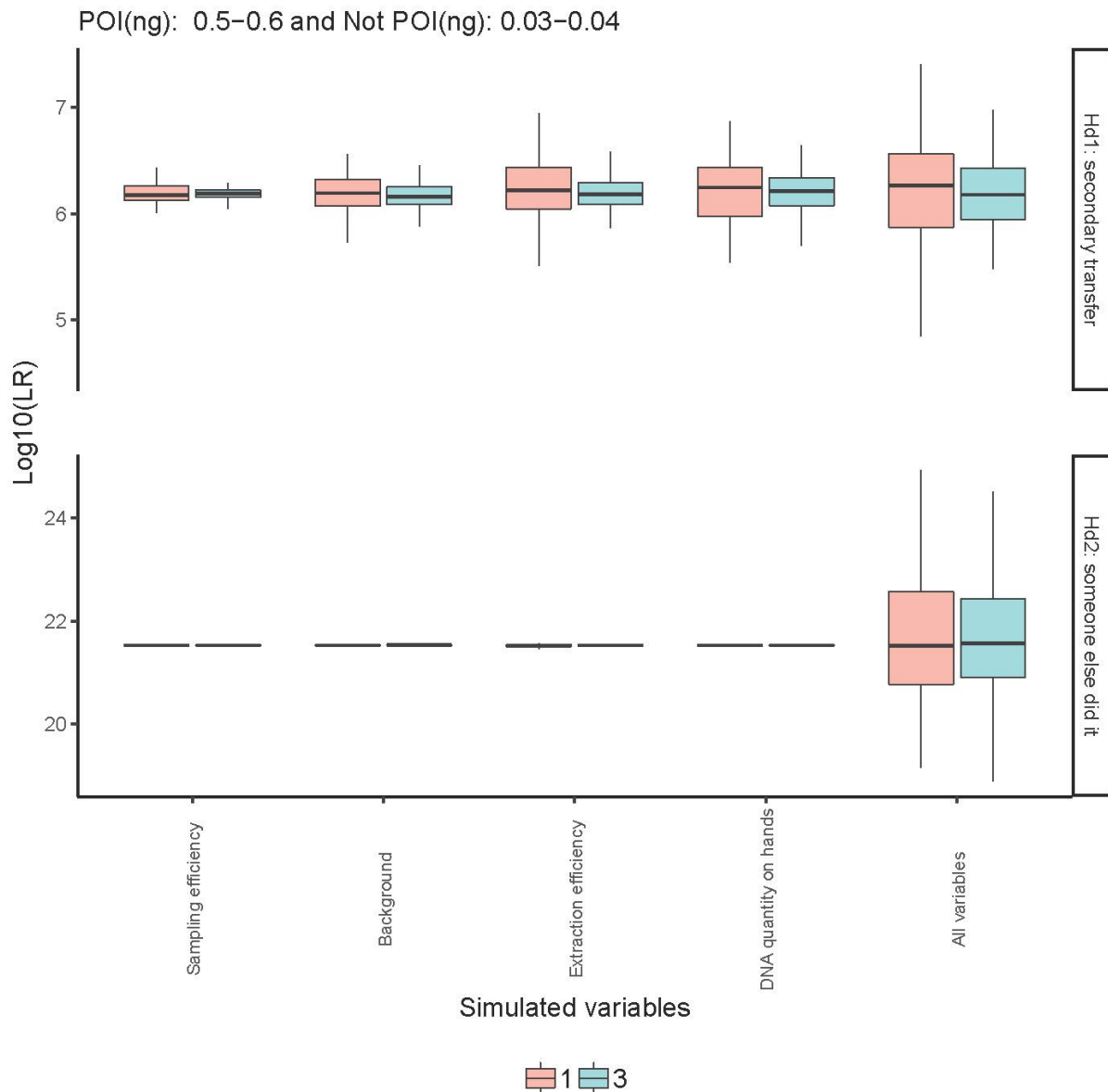
¹ Interestingly, while this is true of the values we found, there is nothing stopping the LR from being greater than 1 over the match probability (MP). If the probabilities of the transfers are extremely low, they can quite legitimately (in theory) lead to an LR that exceeds 1/MP.

546 3.4. *Increasing the datasets informing impacting variables*

547 Should we judge that the overall variation on the ranges of LRs obtained is too high (such in
548 the example shown in Figure 6 under Hd_1), one way to reduce it is to increase the number of
549 data points informing the most impacting variables. Indeed, large variations in LRs translate
550 limitations on the size of the datasets constituting the knowledge used to inform the CPTs. But
551 before rushing into conducting specially designed experiments to acquire additional data to
552 inform the relevant conditional probability tables, we can again take advantage of the simulation
553 strategy developed in this research. It will allow to assess if the analytical investment is worth
554 the effort, hence lead to a reduction of the ranges of LRs that will be beneficial for the case at
555 hand.

556 We could proceed considering any combination of results (POI and not POI's DNA) or jointly
557 for all of them. In Figure 8 we show the results for the case presented before (POI's DNA of
558 0.5-0.6 ng and not POI's DNA of 0.03-0.04 ng).

559



560

561 *Figure 8: Boxplots of the LRs (expressed in Log₁₀) obtained after 100 simulations of each variable in both*
 562 *situations (Hd₁ and Hd₂) with a quantity of POI's DNA between 0.5 and 0.6 ng with between 0.03 and 0.04 of not*
 563 *POI's DNA. The boxplots in red are with the actual data, those in blue are obtained after increasing the*
 564 *underpinning data counts by a factor of 3 on the significant variables only.*

565

566 As expected, when probabilities are informed by an increased number of experiments, the
 567 ranges of likelihood ratios decrease. Under Hd₁, the reduction is greater than under Hd₂ due to
 568 the fact that the increase of data affected the significant variables (identified in Figure 4). For
 569 Hd₂ however, the increase of data does not have a strong impact on the range of LR observed
 570 for “All variables”. This is due to the dominant impact on the LR of the match probability that
 571 is not affected by this increase in data. If Hd₂ is the defence proposition, there is no benefit in
 572 acquiring more data on the other variables that were significant under Hd₁.

573 These simulations allow an assessment, in advance of investing in a large number of
 574 experiments to inform the CPTs, if the benefits, in terms of reduction of LR range, would be
 575 meaningful for the case at hand.

576 *3.5. The discretisation of continuous variables*

577 During the analysis setting up of nodes within the BN and the assignment of probabilities we
 578 have discretised the continuous variables, such as those dealing with DNA amounts, proportions
 579 of transfer, etc. This is a small step away from a fully Bayesian analyses of the data (and the
 580 idea of having parameters for the continuous distributions of variables or their priors), which
 581 would keep all continuous variables as continuous distributions. This is largely due to the
 582 limitations of the BN software. It would be possible to maintain continuous variables, and
 583 perhaps utilise stochastic sampling techniques, in a more customisable software, at the cost of
 584 a loss of comprehensibility. The discretisation itself could have an effect on the sensitivity of
 585 the LR to the data. As an initial investigation into the potential effects of different discretizations
 586 we performed simulations of the initial case using a BN whose states of the variables associated
 587 with the quantity of DNA (in ng) were discretized in three different ways as described in Table
 588 11, called “New discretization 1” and “New discretization 2”. Moving from the initial
 589 discretization, to New discretization 1, then to New discretization 2 represents increasing
 590 coarseness in describing the data.

Variable	Possible states for the new discretization 1	Possible states for the new discretization 2
Results DNA POI	interval node; 0 to 0.01 0.01 to 0.09 in steps of 0.02 0.09 to 0.1	interval node; 0 to 0.01 0.01 to 0.1 in steps of 0.03 0.1 to 0.2
Results DNA not POI	0.1 to 0.2	0.2 to 0.8 in steps of 0.3
DNA Matching	0.2 to 1 in steps of 0.2 1 to 1.5	0.8 to 1 1 to 1.5
DNA DIFF	1.5 to 4.5 in steps of 1 4.5 to 5	1.5 to 4.5 in steps of 1.5 4.5 to 5
Background	5 to 6 6 to 9 in steps of 1.5	5 to 6 6 to 10 in steps of 2
Quantity on hands (ng)	9 to 10 10 to 15 15 to 25 25 to 1000 1000 to inf	10 to 15 15 to 25 25 to 1000 1000 to inf

591

592 *Table 11: New discretizations of the states for the non-instantiated interval variables.*

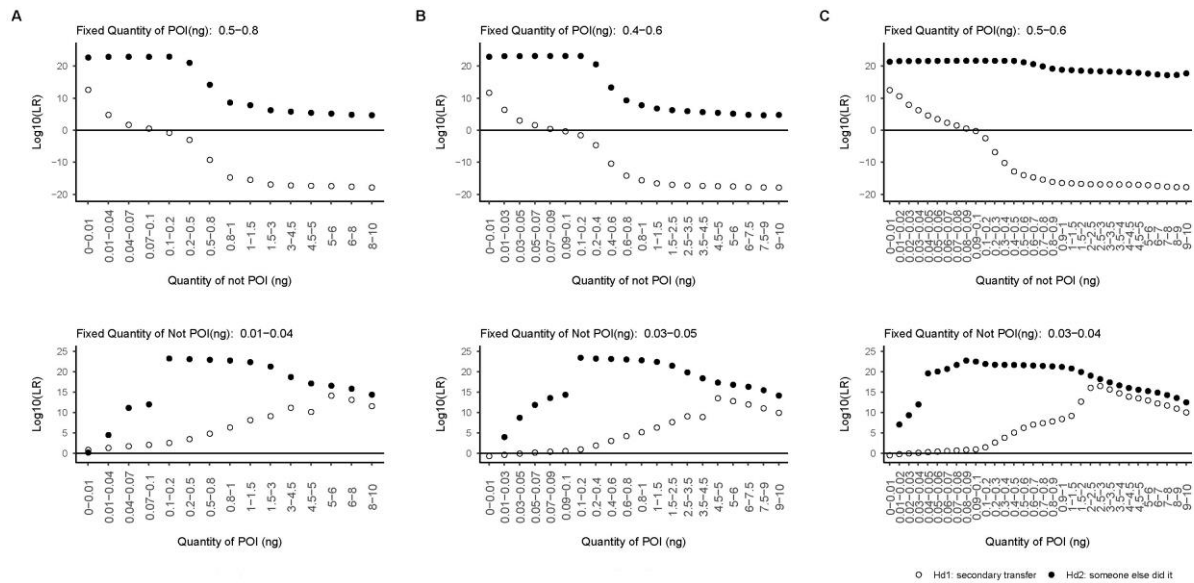
593

594 The results show the following trends regarding the LRs obtained (Figure 9):

- 595
- 596 • In situation 1 (Hd1: secondary transfer), for a fixed DNA quantity of POI (Figure 9 top
597 plots), or not POI (Figure 9, bottom plots), the trends taken by the LRs are similar,
regardless the chosen discretization as shown by the white dots.
 - 598 • For situation 2 (Hd2: someone else did it, black dots in in Figure 9), we have similar
599 observations when the quantity of POI rises for a fixed quantity of not POI (bottom
600 plots). However, for a fixed quantity of POI (top plots), we observed that the presence
601 of non-matching DNA has little impact on the LR for the initial discretization (plot C).
602 This is due to the fact that Hd₂ stipulates that the POI has no link whatsoever of with
603 the knife, hence the probability that some background DNA or AO DNA, would
604 correspond by chance to the POI, is the key consideration that is not impacted by the
605 quantity of not POI. For coarser discretization (plot A and B), the log₁₀ (LR) decreases
606 more drastically when an increased quantity of not POI is obtained.
 - 607 • We note that at the junction between bins of different sizes the LRs may drop or increase
608 abruptly. For example, bottom plots A under Hd₂ (black dots), when moving from 0.07-
609 0.1 ng with the next bin at 0.1-0.2 ng, we observe an increase of the LR. The same is
610 seen on plots under B (bin 0.09-0.1 to bin 0.1-0.2).

611

612



613

614 *Figure 9: LRs obtained when the amount of DNA (POI and non-POI) are respectively varied for a given amount*
 615 *of DNA that remain fixed with the BN whose states of the variables were discretized with the new discretization 2*
 616 *(A) or the new discretization 1 (B) or the initial discretization (C). The top plots refer to the situation where the*
 617 *quantity that is varied is the DNA from not POI and the quantity of POI is fixed (the states varies only because of*
 618 *the discretization). The bottom plots give the LRs for a fixed quantity of not POI and a varying quantity of DNA*
 619 *corresponding to the POI.*

620

621 During our work on discretization, we noted that care must be exercised when changing bin
 622 sizes. The relative sizes of adjacent bins have an impact. This is due to the fact the probabilities
 623 for the nodes associated with the quantity of transferred DNA and the quantity of recovered
 624 DNA are obtained by multiplying the previous quantities by a discount factor linked to the loss
 625 of DNA. These probabilities, obtained through multiplication, will be affected by the sizes of
 626 adjacent bins. In some cases (not reported in this paper), our discretization choices led to LRs
 627 that were misleading (cases under Scenario 2 with LRs below 1). These results highlight the
 628 need for a careful choice of the states. There is a need in the future to investigate how BN can
 629 be constructed in a way that is less affected by multiplication factors.

630

631 To further illustrate the impact of discretization on the obtained LRs, we compared, in Table
 632 12, the $\log_{10}(\text{LR})$ for four contrasting outcomes, in both scenarios:

- 633
- 634 • High quantity of POI (2 ng) with high quantity of not-POI (2 ng)
 - 635 • High quantity of POI (2 ng) with low quantity of not-POI (0.03 ng)
 - 636 • Low quantity of POI (0.03 ng) with high quantity of not-POI (2 ng)
 - Low quantity of POI (0.03 ng) with low quantity of not-POI (0.03 ng).

Situation	Outcomes	Log10 (LR) New discretization 2	Log10 (LR) New discretization 1	Log10 (LR) Initial discretization
Situation 1 (secondary transfer)	High quantity of POI (2 ng) High quantity of not-POI (2 ng)	-15	-16	-16
	High quantity of POI (2 ng) Low quantity of not-POI (0.03 ng)	9	12	18
	Low quantity of POI (0.03 ng) High quantity of not-POI (2 ng)	-18	-18	-18
	Low quantity of POI (0.03 ng) Low quantity of not-POI (0.03 ng)	1	1	0.8
Situation 2 (someone else did it)	High quantity of POI (2 ng) High quantity of not-POI (2 ng)	21	20	18
	High quantity of POI (2 ng) Low quantity of not-POI (0.03 ng)	21	21	19
	Low quantity of POI (0.03 ng) High quantity of not-POI (2 ng)	-18	-17	-13
	Low quantity of POI (0.03 ng) Low quantity of not-POI (0.03 ng)	4	8	12

638

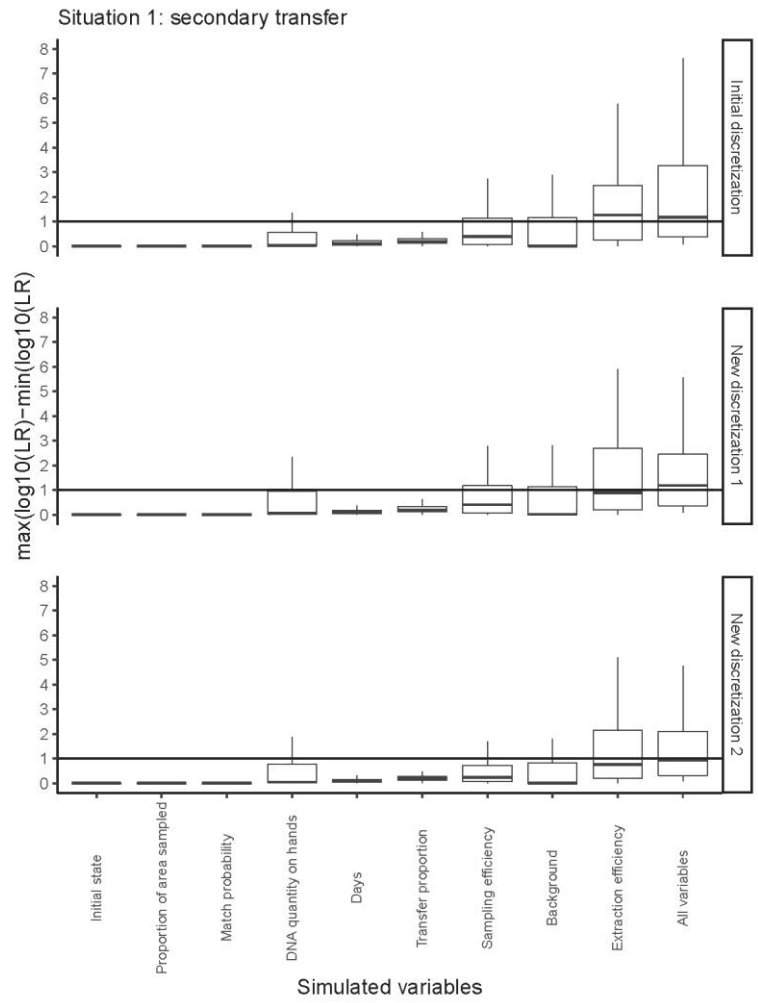
639 *Table 12: LRs (in log10) obtained for four outcomes, depending on the situation and the chosen discretization of*
640 *the states.*

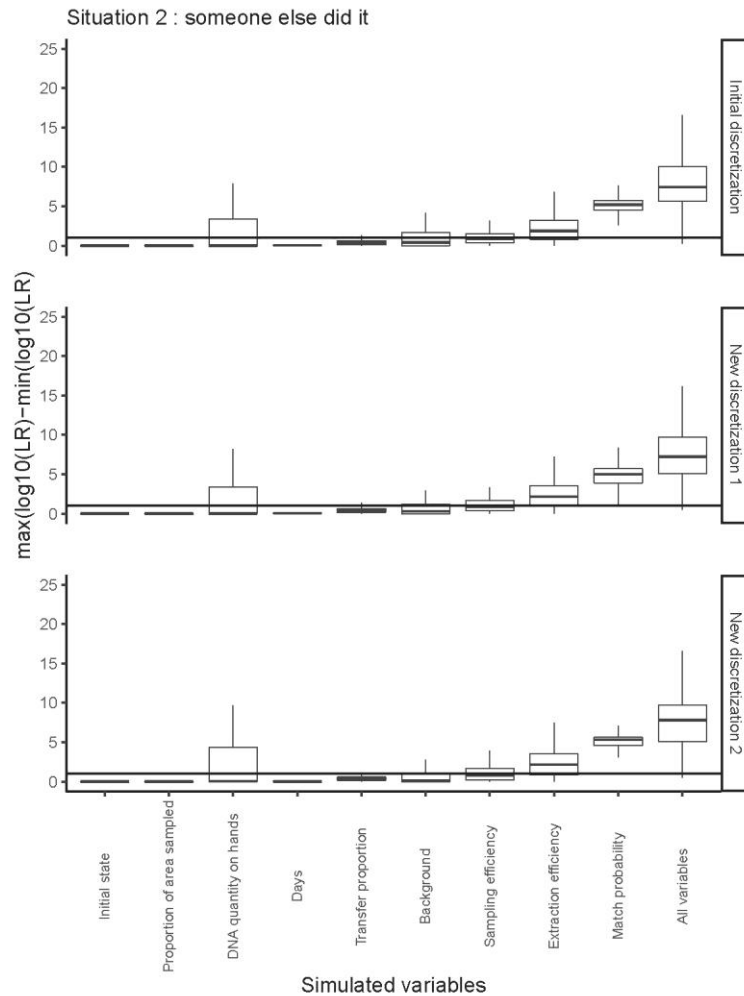
641

642 We observe that the choice of discretization can have an impact on the value of the LRs,
643 depending on the choice of DNA quantities as findings. The fewer states there are, the less
644 discrimination of the propositions is achieved, which means that smaller LRs can be obtained.

645 For example, when a low quantity of not-POI DNA is considered, coarser discretization
646 (moving from initial to new discretization 2), leads to lower LR because of the loss of
647 discrimination capability between the propositions.

648 Regarding the selection of the impacting variables, we observe an impact of discretization that
649 is linked to the loss of discrimination capability when coarser discretization is adopted (Figure
650 10). This is especially the case for the new discretization 2. For all discretizations, the same
651 variables are passing the threshold set to be declared significant, but we can expect that if an
652 even coarser discretization would be chosen, some significant variables would drop below the
653 threshold because of the loss of discrimination.





656

657 *Figure 10: Boxplots presenting the ranges (min-max) expressed in log10 of the LR's obtained following 100*
 658 *simulations under Situation 1 and Situation 2. The panels present the results for the initial case using either the*
 659 *initial discretization of the states or the two new discretizations presented in this part. The horizontal line drawn*
 660 *at the difference of 1 (in log10) set the limit above the variable considered will be declared as having a significant*
 661 *effect on the global variability shown when "all variables" are resampled jointly.*

662

663 4. Conclusion

664 The ENFSI guideline [3] is advising forensic scientist to evaluate biological traces with a low
 665 level of DNA considering activity level propositions. For that task, specific variables such as
 666 transfer, persistence, recovery and background need to be considered. Many experts face
 667 practical difficulties when considering these variables for various perceived reasons. Typically,
 668 expert will indicate that:

- 669 – the number of variables at play is overwhelming and unmanageable;

- 670 – every case represents a unique set of circumstances and any numerical assignment
671 of probabilities is critically dependant on the specificities of the case;
- 672 – the paucity of current published studies to inform the parameters associated with
673 these variables cannot reasonably be compensated by reasonable data acquisition
674 campaigns.

675 In this study, we show that Bayesian networks can handle this complexity efficiently, when
676 coupled with simulation techniques, they can be used to identify the most impacting variables,
677 hence reducing the data acquisition burden by directing the laboratory to the key issue.

678 The method has been applied to a scenario involving trace DNA recovered from knife handles
679 where the prosecution alleges that the person of interest (POI) stabbed a victim. The findings
680 considered take the form of given quantities of DNA (in ng) corresponding or not the POI. As
681 a general tendency, regardless of the findings, we showed that when the defence claims that the
682 POI has nothing to do with the incident, the match probability associated with the POI will
683 dictate most of the weight to be assigned to the findings. If the POI defence is invoking the
684 possibility of secondary transfer, the key variables are associated with the sampling, the
685 extraction efficiency, the background and the quantity of DNA on the hands.

686 Simulation techniques can also be used to assess the merit of increasing the knowledge base (in
687 terms of size of studies carried out) when the significant variables had been identified. We
688 presented preliminary results on the impact of the choice of discretization of the variables.
689 Discretization can have an impact on the LR_s and potentially on the choice of impacting
690 variables, mainly due to the loss of discriminability between the propositions when a too coarse
691 discretisation is adopted. In our view, the number states and their ranges should be chosen
692 carefully in a way that avoids losing information (e.g. merging states to a point where
693 discrimination is lost).

694 Finally, we noted that the identification of significant variables depends on the obtained DNA
695 results and this selection may be refined on a case by case basis. To allow one exploring all
696 possibilities, a dedicated Shiny application has been designed ([https://lydie-
697 samie.shinyapps.io/DNA_Activity/](https://lydie-samie.shinyapps.io/DNA_Activity/)).

698 **5. Acknowledgements**

699 We would like to express our gratitude to Marco De Donno for all his assistance in performing
700 the simulation.

701 **6. Bibliography**

- 702 [1] F. Taroni, A. Biedermann, J. Vuille, N. Morling, Whose DNA is this? How relevant a question? (a note
703 for forensic scientists), *Forensic Sci. Int. Genet.* 7 (2013) 467–470.
- 704 [2] C. Champod, DNA transfer: informed judgment or mere guesswork?, *Frontiers in Genetics.* 4 (2013) 300
705 in A. Biedermann, J. Vuille, F. Taroni, DNA, Statistics and the Law: A Cross-Disciplinary Approach to
706 Forensic Inference, *Frontiers Media.* 4 (2014) 22-24.
- 707 [3] S. Willis et al., ENFSI guideline for evaluative reporting in forensic science, European Network of
708 Forensic Science Institutes, Dublin (2015). http://enfsi.eu/wp-content/uploads/2016/09/m1_guideline.pdf
709 (accessed May 6, 2019).
- 710 [4] A. Biedermann, C. Champod, G. Jackson, P. Gill, D. Taylor, J. Butler, N. Morling, T. Hicks, J. Vuille, F.
711 Taroni, Evaluation of Forensic DNA Traces When Propositions of Interest Relate to Activities: Analysis
712 and Discussion of Recurrent Concerns, *Frontiers in Genetics.* 7 (2016) 1–12.
- 713 [5] D. Taylor, A. Biedermann, L. Samie, K.-M. Pun, T. Hicks, C. Champod, Helping to distinguish primary
714 from secondary transfer events for trace DNA, *Forensic Sci. Int. Genet.* 28 (2017) 155-177.
- 715 [6] F. Taroni, A. Biedermann, S. Bozza, P. Garbolino, C. Aitken, Bayesian Networks for Probabilistic
716 Inference and Decision Analysis in Forensic Science, 2nd edition, *Statistics in Practice*, Chichester: John
717 Wiley & Sons, Ltd (2014).
- 718 [7] D. Taylor, T. Hicks, C. Champod. Using Sensitivity Analyses in Bayesian Networks to Highlight the
719 Impact of Data Paucity and Direct Future Analyses: A Contribution to the Debate on Measuring and
720 Reporting the Precision of Likelihood Ratios. *Science & Justice*, 56 (2017) 402-410.
- 721 [8] P. Ka-Man, Interprétation des profils génétiques obtenus à partir des traces de contact, PhD thesis, School
722 of Criminal Justice, University of Lausanne, 2016.
- 723 [9] R. Palmer, The Evaluation of Fibre Evidence in the Investigation of Serious Crime, PhD thesis, School
724 of Criminal Justice, University of Lausanne, 2016.
- 725 [10] M. Khodabina, A. Ahmadabadi, Some properties of generalized gamma distribution, *Mathematical*
726 *Sciences* 4 (2010) 9-28
- 727 [11] D. Dufresne, The Beta Product Distribution with Complex Parameters, *Commun. Stat. —Theory*

- 728 Methods. 39:5 (2010) 837–854
- 729 [12] D.J. Daly, C. Murphy, S.D. McDermott, The transfer of touch DNA from hands to glass, fabric and wood,
730 Forensic Sci. Int. Genet. 6 (2013) 41–46.
- 731 [13] L. Bontadelli, Study of DNA Shedder Quality, MSc thesis, School of Criminal Justice, University of
732 Lausanne, Lausanne, 2009.
- 733 [14] M. Goray, R.J. Mitchell, R.A.H.V. Oorschot, Investigation of secondary transfer of skin cells under
734 controlled conditions, Leg. Med. 12 (2010) 117–120.
- 735 [15] T.J. Verdon, R.J. Mitchell, R.A.H.V. Oorschot, Evaluation of tapelifting as a collection method for
736 contact DNA, Forensic Sci. Int. Genet. 8 (2014) 179–186.
- 737 [16] J.M. Bernardo and A.F.M Smith, Bayesian Theory, John Wiley & Sons, Inc (2000)
- 738 [17] E. Butts, Exploring DNA Extraction Efficiency Forensics@NIST 2012 Meeting, Gaithersburg
739 https://www.nist.gov/sites/default/files/documents/oles/3_Butts-DNA-extraction-2.pdf (2017).
- 740 [18] G.E. Meakin, E. V. Butcher, R.A.H. van Oorschot, R.M. Morgan, The deposition and persistence of
741 indirectly-transferred DNA on regularly-used knives, Forensic Sci. Int. Genet. Suppl. Ser. 5 (2015) e498-
742 e500.
- 743 [19] RStudio Team (2015). RStudio: Integrated Development for R. RStudio, Inc., Boston, MA URL
744 <http://www.rstudio.com/>.
- 745 [20] R Core Team (2019). R: A language and environment for statistical computing. R Foundation for
746 Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- 747 [21] Kjell Konis. (2017). RHugin: RHugin. R package version 8.4. <http://rhugin.r-forge.r-project.org>] to link
748 with the developed Hugin network.
- 749 [22] L. Samie, T. Hicks, V. Castella, F. Taroni, Stabbing simulations and DNA transfer, Forensic Science
750 International: Genetics 22 (2016) 73–80.
- 751 [23] C. Gehrig, B. Balitzki, A. Kratzer, C. Cossu, N. Malik, V. Castella, Allelic proportions of 16 STR loci –
752 including the new European Standard Set (ESS) loci – in a Swiss population sample, Int. J. Legal Med.
753 128 (2013) 461–465.
- 754 [24] T. Hicks, F. Taroni, J. Curran, J. Buckleton, O. Ribaux, V. Castella, Forensic Science International :
755 Genetics Use of DNA profiles for investigation using a simulated national DNA database : Part I . Partial
756 SGM Plus 1 profiles, Forensic Sci. Int. Genet. 4 (2010) 232–238.

757 [25] Samie L., Champod C., Glutz V., Garcia M., Castella, V. & F. Taroni, The efficiency of DNA extraction
758 kit and the efficiency of recovery techniques to release DNA using flow cytometry, *Science & Justice*, 59
759 (2019) 405-410