

An introductory guide to Evaluative Reporting



National Institute of
Forensic Science
Australia New Zealand

About NIFS

The National Institute of Forensic Science (NIFS) is a directorate within the Australia New Zealand Policing Advisory Agency (ANZPAA). Our strategic intent is to Promote and Facilitate Excellence in Forensic Science. Our role is to deliver support to the forensic science community in the areas of co-ordination, innovation, information management, education and training, and quality. Our program of work is underpinned by a Strategic Plan approved by the ANZPAA Board of Australia and New Zealand Police Commissioners. The Australia New Zealand Forensic Executive Committee (ANZFEC), comprising representatives from the forensic service providers of our funding agencies, has oversight of the delivery of the Strategic Plan via our annual Business Plan.

About this document

Title: An introductory guide to evaluative reporting
Date: ©JUNE 2017
Available from: www.nifs.org.au
Enquiries: secretariat.nifs@anzpaa.org.au

Acknowledgements

This document was prepared by a working group convened under the NIFS Operational Effectiveness Program 2016/2017. The working group members are:

Dr Kaye Ballantyne (Victoria Police, Australia)
Dr Joanna Bunford (New South Wales Health Pathology, Australia)
Dr Bryan Found (Victoria Police, Australia)
Inspector David Neville (Queensland Police Service, Australia)
Dr Duncan Taylor (Forensic Science South Australia)
Mr Gerhard Wevers (Institute of Environmental Science and Research, New Zealand)
Mr Dean Catoggio (National Institute of Forensic Science, Australia New Zealand)

The working group would like to thank the following people who contributed material for the preparation of this document:

Ms Eva Bruenisholz (National Institute of Forensic Science, Australia New Zealand)
Mr Robert Morgan (National Institute of Forensic Science, Australia New Zealand)
Mr Matthew Carmody (Australia New Zealand Policing Advisory Agency)
Ms Caroline Gibb (Victoria Police, Australia)

Cover image

Forensic scientists examining gun 2012
Cultura Creative (RF)/ Alamy Stock Photo

Contents

01	Purpose
01	Introduction to evaluative reporting
02	The basics
04	Propositions and expectations
05	Analysis and findings
07	Evaluation and reporting
10	Transitioning to evaluative reporting
13	Conclusion
13	Resources
14	Appendix A – Firearms case example
17	Appendix B – Paint case example
20	Appendix C – Fingerprint case example
22	Appendix D – Glass case example



Purpose

The purpose of this guide is to provide an introduction to evaluative reporting in forensic science. It covers the basics of evaluative reporting, including terminology, and uses examples in the appendices to illustrate the concepts.

The guide also provides managers with advice on implementing evaluative reporting where they consider it appropriate. It provides a useful list of resources, including where to find relevant books and journal articles, comprehensive guides, and training courses on the subject.

Introduction to evaluative reporting

Evaluative reporting is a formalised thought process that enables the evaluation of scientific findings given two opposing (or competing) propositions. It is a way of providing a strength of the findings of an examination given those alternative propositions.

It can be used by comparative forensic science disciplines where you are forming an opinion based on your observations, or where a decision has to be made. It is generally not used for factual reporting, such as drug identification.

Evaluative reporting is a means of dealing with uncertainty and provides a balanced approach to evidence interpretation. Properly applied, cognitive bias can be minimised and opinions can be updated in a logical way on receipt of new information. As such, the use of evaluative reporting in forensic science could assist in addressing some of the issues highlighted by the President's Council of Advisors on Science and Technology (PCAST) and the National Academy of Science reports on foren-

sic science in the United States.^{1,2} Organisations such as the European Network of Forensic Science Institutes (ENFSI), the Royal Statistical Society (UK) and the Association of Forensic Service Providers have issued position statements and guidelines around its use.^{3,4}

In some laboratories, such as the Netherlands Forensic Institute (NFI) and The Institute of Environmental Science and Research (ESR) in New Zealand, evaluative reporting is used for all disciplines and cases. Other laboratories apply the framework in specific disciplines, such as trace evidence analysis at Forensic Science South Australia and the Forensic and Analytical Science Service (New South Wales), or handwriting and signature examination (Victoria Police Forensic Services Department). The traditional comparative disciplines such as fingerprint examination and firearms and toolmark analysis are also being reported evaluatively in some agencies, with methods of reporting moving from conventional identification statements to more probabilistic expressions.

¹ *Report on Forensic Science in Criminal Courts: Ensuring Scientific Validity of Feature-Comparison Methods*, Executive Office of the President of the United States Council of Advisors on Science and Technology (PCAST) September 2016

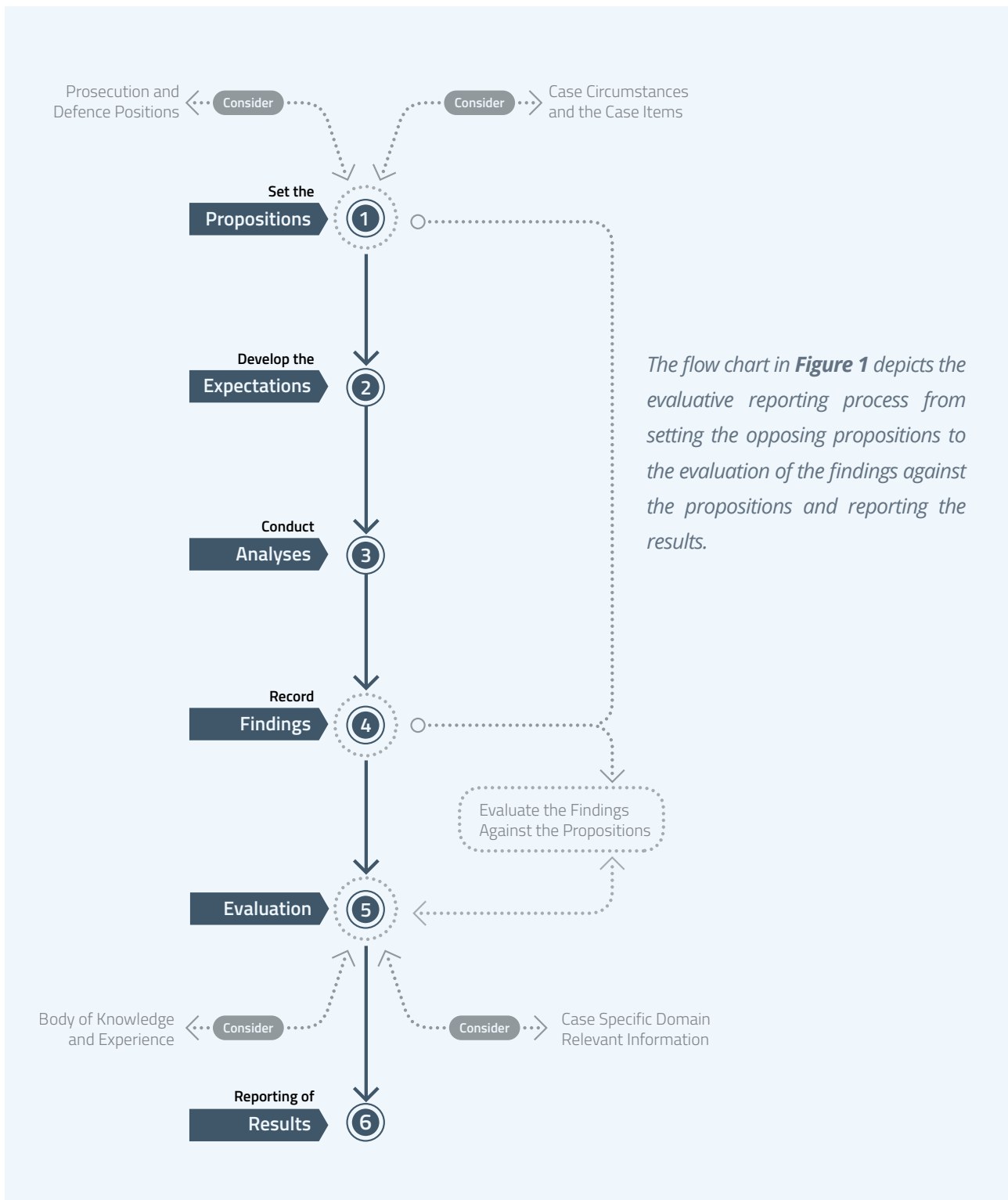
² *Strengthening Forensic Science in the US: A Path Forward*, Committee on Identifying the Needs of the Forensic Sciences Community, National Research Council, National Academies Press, 2009.

³ *ENFSI Guideline for Evaluative Reporting in Forensic Science*, European Network of Forensic Science Institutes 2015 v3.0.

⁴ Jackson, G., Aitken, C., Roberts, P., 2014, 'Case Assessment and Interpretation of Expert Evidence', Practitioner Guide No. 4, Working Group on Statistics and Law of the Royal Statistical Society

The basics

Figure 1:
The Evaluative Reporting Process



A summary of the principles and terminology associated with each step in the evaluative reporting process is provided below. The concepts behind the process are discussed in more detail in subsequent sections. Several case examples are provided in the appendices to help the reader understand how the process can be applied in different comparative disciplines.

1. Propositions

The propositions are based on the case circumstances, the defence/prosecution scenarios and the availability of case items, and it is against these propositions that the findings will be evaluated. The propositions should not change as a result of observations, although are free to change as more information regarding case circumstances, or information from defence and prosecution is obtained.

Prosecution position

The explanation being put forward by the prosecution.

Defence position

The explanation being put forward by the defence. Note that there may be no scenario offered by defence, in which case, where possible, the responsibility of assigning a reasonable defence proposition falls to the individual carrying out the evaluation.

Case circumstances

The framework of circumstances surrounding the alleged crime and the items relating to individuals who may be involved. Case circumstances are typically considered to be relevant details required to evaluate the observations.

Case items (exhibits)

The items collected as part of the investigation of the alleged crime that are relevant to the propositions.

2. Expectations

The list of expected outcomes based upon the propositions.

3. Analyses

The tests carried out on case items, and any samples taken from them. Typically, the type of tests will be dictated by the propositions and expected outcomes being considered and the state/nature of the samples.

4. Findings

The observations or results from the laboratory analyses.

5. Evaluation

The consideration of all observations in light of propositions, relevant information, limits of the testing procedure used, knowledge and experience. This may involve the assignment of a numerical value for the probability of the findings given the competing propositions or may be non-numerical and instead a statement of relative support for one proposition over the other.

Body of knowledge and experience

The base of information that the analyst will draw on to assess components of their overall evaluations. This may include published literature, databases of characteristic frequencies, knowledge of analysis limitations, professional experience or personal knowledge, beliefs and assumptions.

Case Specific domain relevant information

Includes information that is not directly a case circumstance, but is still important to the evaluation of the observations. These include aspects such as the time between case items being obtained and examined, or the manner of item collection and handling.

6. Reporting

The explanation of the whole process for the fact-finder, including the information and method used to form propositions, the available case items, the analyses carried out on them, the resulting observations, the method of evaluation and the conclusions drawn from it. Assumptions made during the evaluative processes should be made clear to the fact-finder.

Propositions and expectations

Propositions (hypotheses) are formulated from the case information provided and are mutually exclusive and exhaustive in the context of the case. Explanations can be considered as a less formal exploration of alternatives for the findings and, according to Evett et al (2000), can be a statement of the obvious, prescriptive, speculative or fanciful.

The propositions typically align with the prosecution and the defence positions and are normally denoted Hp and Hd, respectively. Alternatively, they may align to different case scenarios that are not attributed to either a prosecution or defence position and are denoted H1 and H2.

Propositions must be mutually exclusive meaning that they cannot occur at the same time, like tossing a coin that can only ever land on heads or tails. They can be amended or refined if new information becomes available as the case progresses.

The formulated propositions will lead to the development of a list of possible (and most expected) outcomes, at a minimum two, covering the most relevant situations based on the case circumstances. This may assist with case assessment as part of the submission process, or will direct what features of the exhibits will be examined. Expectations are testable outcomes, while propositions refer to the higher level scenarios.

The firearms case example at figure 2 shows how propositions and expectations can be derived from the case scenarios, case circumstances and the case items. In this example, it is alleged that the suspect shot the deceased. The forensic practitioner is presented with a bullet recovered from the deceased and a firearm that was seized from the suspect.

One proposition would infer a relationship between the suspect's firearm and the bullet recovered from the deceased:

H1 *The recovered bullet was fired using the suspect's firearm.*

and the alternate proposition would infer no relationship

H2 *The recovered bullet was fired using another firearm.*

The expectations based upon the propositions infer a relationship between two observable data sets. In the firearm example, the two data sets are the marks observed on the bullet recovered from the deceased and the marks observed on the test fired bullets from the suspect's firearm.

The expected outcomes given H1 might be ***"The recovered bullet and bullets fired from the suspect's firearm show the same class characteristics and microscopic detail"*** and the expected outcomes given H2 might be ***"The recovered bullet and bullets fired from the suspect's firearm show different class characteristics and/or microscopic detail"***.

Hierarchy of propositions

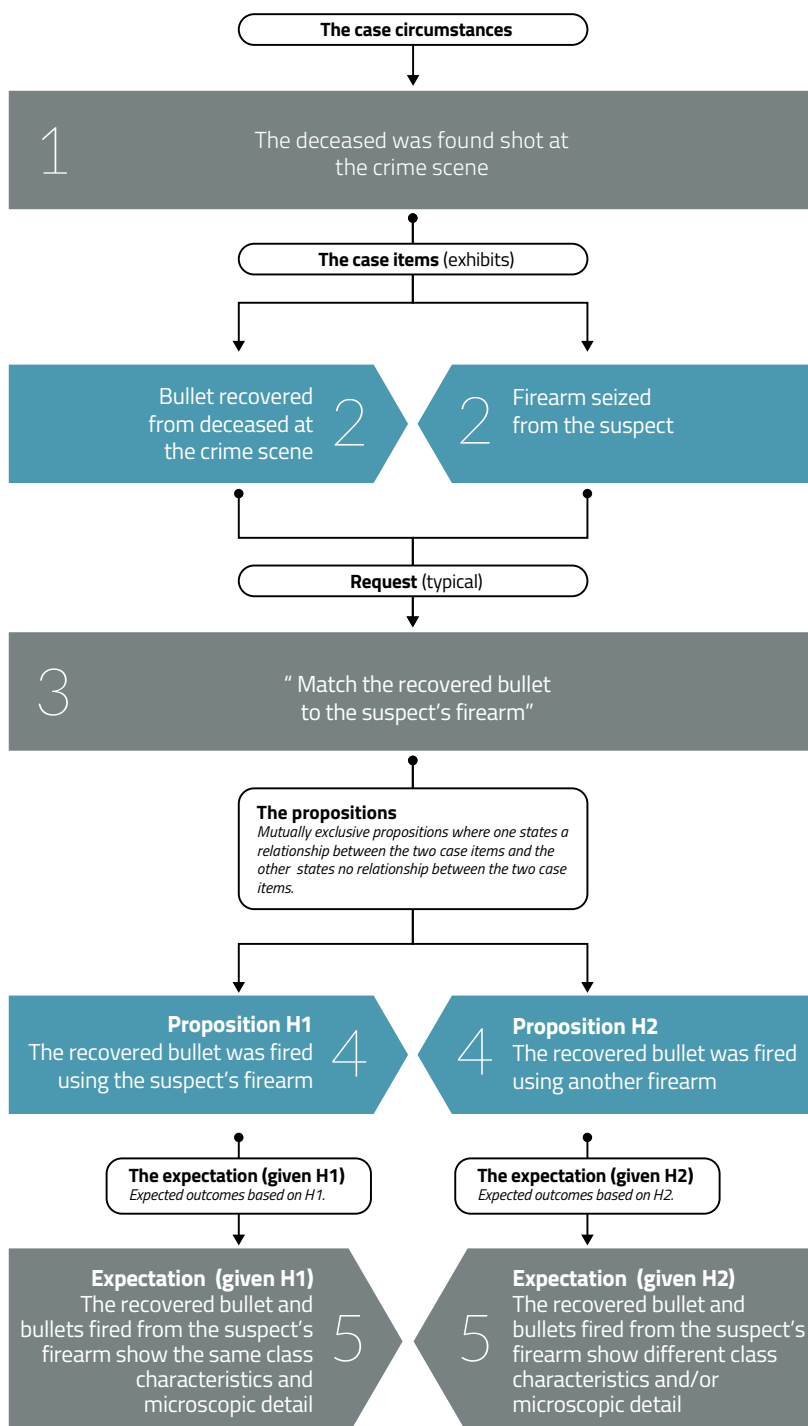
The hierarchy of propositions refers to the level at which the propositions are aimed. The highest level in this hierarchy is the offence level, which addresses the guilt or innocence of the suspect. Below this sits activity level, where propositions speak to how the evidence arose. Source level propositions instead refer to where the evidence came from.

Forensic scientists most commonly report on source level propositions, but may find that activity level propositions are more relevant to the issue at hand in some cases. As the propositions will necessarily change between the two levels, the assessment of the probabilities of the findings will change, potentially resulting in different strengths of the findings between activity and source level.

Analyses and findings

Figure 2:

A firearms case example demonstrating how the propositions and expectations can be derived from the case circumstances and the items that are available for testing.



The specific analyses conducted depend upon the propositions and expectations being considered and the state and nature of the case items and samples.

For example, the firearms case in figure 2 would involve the use of an optical microscope to observe and compare the marks appearing on the bullet recovered from the deceased with those appearing on the bullets fired from the suspect's firearm. Marks, also referred to as striations or grooves, are left on a bullet as it is fired down the barrel of a firearm. Typically, multiple test firings of the suspect's firearm would be made to observe the amount of variability in marks produced by that firearm.

The findings are the recorded observations that may include photographs, measurements and subjective descriptions of the features appearing on the recovered bullet and those appearing on the test fired bullets.

A discussion about

Traditional identification model and evaluative reporting

Traditional reporting

In many traditional comparative forensic disciplines, most laboratories interpret evidence so that conclusions will typically fall into the following three categories; elimination, inconclusive and identification. This is a classic “falling off the cliff” approach to evidence interpretation, where the findings are inconclusive until enough evidence exists that it passes a threshold and the findings become conclusive (either elimination or identification). Using this method, a result that almost reaches the threshold for identification is reported in the same way as a result that almost reaches the threshold for elimination (i.e. inconclusive).

Effectively, any findings that fail to be conclusive will be inconclusive. The inconclusive range is large and covers findings that are almost an elimination to almost an identification. If you imagine a scale from 0 to 100, where 0 is an elimination and 100 is an identification, then the inconclusive range falls between 0 and 100. This is a very large range, containing a lot of information and important evidence that might be lost with this approach to evidence interpretation. This method has most likely originated as a means of dealing with the uncertainty of comparisons that don't meet the required “match” threshold.

Evaluation using a logical framework

When our belief in the truth of an event is not absolute, we say that uncertainty exists. We deal with uncertainty every day, from decisions on what to wear, what to buy, how to get to work or when to accept an invitation. A commonly used example is the weather, and our uncertainty around whether or not it will rain. When we compare our belief in some event to another real-world event for which we are familiar (such as the outcomes of rolling dice, tossing a coin, or making a gamble) then we are said to be assigning a probability to that event occurring.

We can update our belief in an event occurring as we obtain more information. In the example of assigning a probability of rain, we may look outside and see that it is cloudy, and then watch a weather forecast which predicts rain for that day. Each of these additional pieces of information can be considered and we can assign a probability for rain, conditioned on what we now know. This same thought process is used when applying evaluative reporting.

The fundamental principles of evaluative reporting or interpretation are that (i) the crime must be considered to have occurred within a framework of circumstances, (ii) that the findings must be considered in light of at least two competing propositions that will be guided by the case circumstances and (iii) that the role of the expert is to comment on the probability of their findings, given these propositions and not on the propositions themselves. It is this last principle that allows the fact-finders to combine aspects of evidence they hear during the course of the trial with their judgement in their deliberations. This framework of evidence evaluation is commonly referred to as evaluative reporting, but may also be referred to as the likelihood ratio approach, logical thinking or Bayesian inference.

Those that undertake evaluative reporting find that they are able to provide informative opinions for findings that have traditionally been considered ‘inconclusive’ and in ways that can be standardised between analysts within, and between, organisations and forensic disciplines.

Experience based assignment of likelihood ratios

The key to evaluative reporting is the consideration of the probability of the findings given two competing propositions. Assignment and reporting of a likelihood ratio is the most common usage of this framework at present. It is a common misapprehension that the use of the evaluative reporting framework requires the existence of large population databases to assign empirically derived, numeric likelihood ratios. Whilst this is clearly the preferred situation, as it provides the most robust and transparent manner of estimating evidentiary weight, it is possible to estimate the weight of the findings based on experience, published estimates or combined estimates of probability from multiple practitioners.

The assignment of numeric probabilities to generate explicit likelihoods can be performed in two ways. The first is to utilise population data, collected from samples of known origin, to estimate the frequency of the observed characteristics within the population of interest. This may involve estimating the frequency of particular paint types on cars, the amount of glass in the area with a defined refractive index, or the number of fingerprints with particular numbers and types of minutiae. In such cases, the likelihood ratio can be assigned and precision estimated.

The likelihood may be reported numerically, such as:

“It is 500 times more likely to observe the correspondence between the questioned and known prints if proposition H1 is true than if proposition H2 is true”

Alternatively, the numeric estimate may be translated to a verbal scale. There are a variety of scales in use, although many share similar wording and levels (see table 1 for an example of a commonly used scale). When translated to a verbal equivalent, the opinion may be reported as:

“the findings provide strong support for proposition H1, as opposed to proposition H2”
Or
“the results are far more probable given proposition H1 than given proposition H2”.

Table 1: An example of a reporting scale which includes numerical values and verbal equivalents.

Verbal Conclusion <i>(support for or against proposition 1)</i>	Likelihood Ratio <i>(LR)</i>
Extremely strong support against	< 0.000001
Very strong support against	0.000001 – 0.001
Strong support against	0.001 – 0.01
Moderate Support against	0.01 – 0.1
Slight support against	0.1 – 1
Neutral	1
Slight support for	1-10
Moderate Support for	10-100
Strong support for	100-1,000
Very strong support for	1,000-1,000,000
Extremely strong support for	>1,000,000

Some disciplines or evidentiary types may not have robust population databases, or are unable to estimate feature frequencies in a reliable manner. In such instances, practitioners can utilise their experience and knowledge of the relevant population to subjectively estimate the likelihood. The firearm example provided in Appendix A demonstrates this subjective assessment. This numeric estimate may be reported as described or may be translated across to a verbal scale for reporting purposes. Some laboratories modify the verbal scale to include the traditional “identification” and “exclusion” opinions for the findings when it is subjectively determined that the probability of the findings given a proposition is zero. An example of this would see the reporting of an exclusion when a firearm of interest has a different calibre or rifling characteristics to the questioned bullet.

Alternatively, practitioners may avoid formal expression of a numeric probability, and subjectively select a verbal qualifier to express the degree of support provided by the findings. When assigning subjective probabilities, it may not be possible to consistently and accurately separate between levels on the scale for some forms of pattern comparison, and so a reduced scale can assist with consistency between practitioners. For example, some forensic handwriting comparisons are reported using only the “very strong support” and “qualified support” levels, where qualified support infers a higher level of uncertainty in the opinion.

If the subjective assignment of probabilities based on a practitioner’s judgement and experience is utilised, it is critical that within- and between- practitioner consistency in evidence⁵ evaluation is checked. As the choice of support can vary depending on the quality of the findings, the frequency of the features within the population, or a combination of the two, it is important that all examiners are provided with training and guidance regarding the assignment of probabilities and selection of level of support. Examiners should be able to articulate how they differentiate between the levels, and why they have chosen a particular level for a particular finding. Failure to clearly define how to use the scale may result in inconsistency between examiners, where some use the scale as a measure of evidential quality and others in their level of certainty in their conclusion.

Reporting the results

Forensic practitioners normally prepare a report to assist the triers of fact in their adjudications. The report is an explanation of the whole process, including the information and method used to form propositions, the available case items, the analyses carried out on them, the resulting observations, the method of evaluation and the conclusions drawn from it.

When documenting the analyses and reporting the results it is important that the practitioner convey the nature of the information used to evaluate the finding. This information may include the results of empirical research, databases, surveys or practitioner experience and knowledge. Assumptions made during the evaluative processes should be made clear to the fact-finder along with the limitations of the testing. The appendices provide examples of common phrases for reporting of the results for different forensic disciplines.

⁵ The term evidence has a specific legal definition in some jurisdictions. However, it is common for scientific literature to also refer to the findings of forensic analysis as evidence. Therefore, both the terms are sometimes used interchangeably in this text.

Transitioning to evaluative reporting

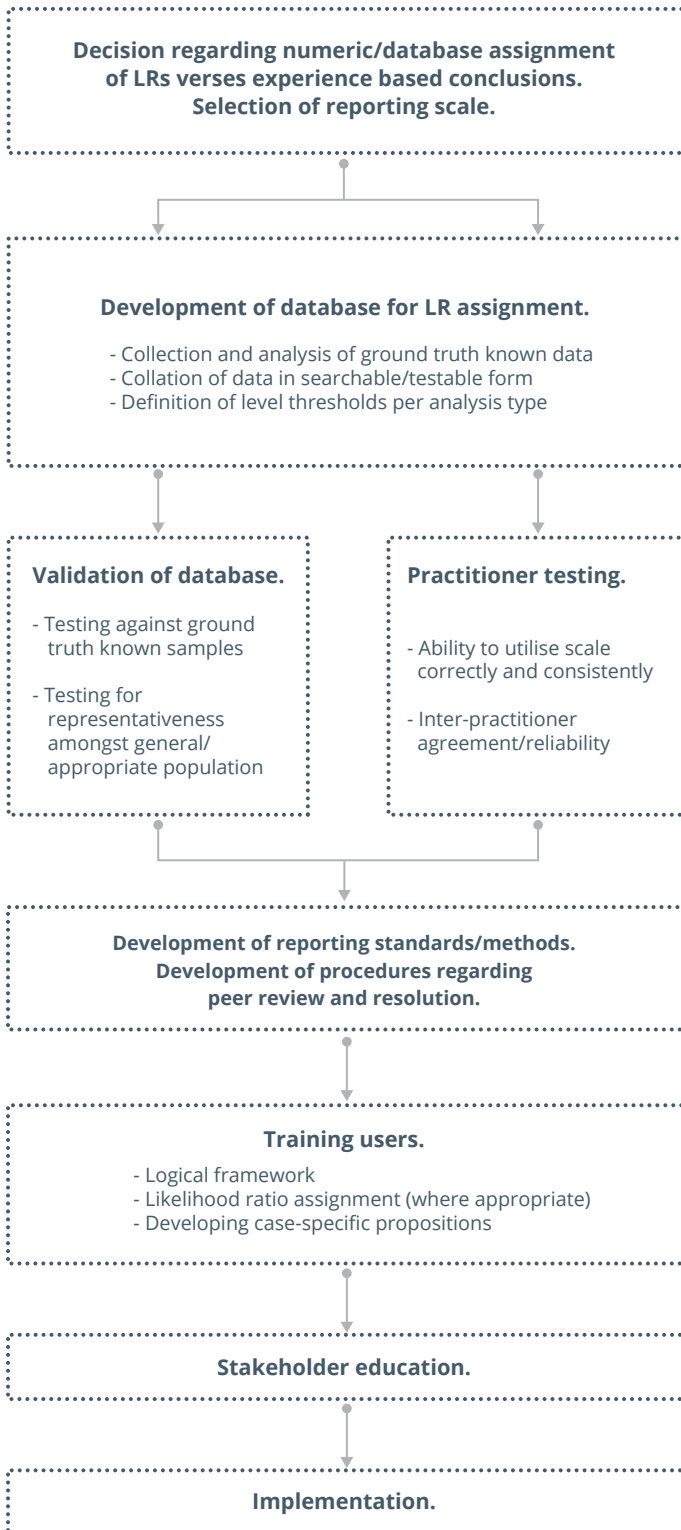


Figure 4:
Flow chart of considerations for implementing evaluative reporting.

Achieving a smooth transition from traditional reporting frameworks to evaluative reporting requires an investment in practitioner education and training, stakeholder engagement and in the case of a shift to numeric LR use, the development of frequency databases from known materials.

Whilst the level of preparation and investment may seem daunting for some disciplines, the process may be accomplished in stages, with key implementation points staggered to allow training, education and feedback to practitioners and stakeholders, including the judiciary.

Good project management will assist with a successful transition to evaluative reporting. The application of the Plan-Do-Check-Adjust (PDCA) cycle will support continuous improvement and refinement of evaluative reporting training and processes.⁶

Some considerations for implementing evaluative reporting are discussed below and summarised in the flowchart shown at figure 4.

⁶ The Plan-Do-Check-Adjust (PDCA) cycle, otherwise known as Plan-Do-Check-Act, is a model for business improvement based upon the scientific method.

Selection of a reporting scale

The selection of an appropriate reporting scale is an important step in transitioning to evaluative reporting. The scale should be suitable for use across all disciplines. Research may demonstrate that the use of different scales is appropriate for some types of evidence. The scale may include numerical and/or verbal equivalents.

Consideration may be given to the ranges of weights of the findings that will be encountered, the number of categories required, and ability to discriminate between the different scale levels within the different evidence types. Laboratories which already use a particular scale for some disciplines such as DNA or trace evidence may find it easiest to expand the use of these scales across additional areas (as shown in the firearms example in Appendix A and Table 1).

Development of a database

Individual disciplines must then assess their own methods and evidence types to determine how the scale will be applied to the opinions and conclusions drawn. Disciplines using numeric estimates of the population frequency will need to develop databases for their relevant populations, or determine if published frequency estimates are applicable to their population. Collaboration and data sharing between laboratories and jurisdictions can greatly assist in the collation of sufficient data for reliable frequency estimates. Databases should be validated by testing against separate population samples, and checked for random sampling assumptions.

Definition of level thresholds

Disciplines without numeric estimates, or where the independence of characteristics and assignment of likelihood ratios is not yet possible, will need to assess their analytical and interpretational methods to determine how the weight of the findings can be reliably determined between cases and practitioners. Methods relying on defined features, such as fingerprint minutiae, microscopic stria within toolmarks, or randomly acquired characteristics in footwear impressions, may be able to provide guidance regarding the number and type of characteristics that might be required for each level on the verbal scale. It should, however, be noted that the use of defined thresholds is not generally recommended due to the effect of a threshold on cognitive search strategies, and the variance in frequencies between particular characteristics. Evidence types without clearly definable features, those without knowledge of population frequencies or characteristic variances, and those with varying forms and combinations of features may require more consideration to the

assignment of levels of support. For example, within handwriting and signature analysis the weight of the findings may vary depending on the quality and quantity of unknown (questioned) writings, the variation observed within the known writing exemplars, and the exact features and letter forms involved. Thus, providing lists or numbers of line features may not be possible – instead, examples of writing comparisons that were determined to show different levels of support for varying propositions may be collated to use as training materials. It should also be noted that the separation between numeric and purely verbal disciplines or techniques may not be clear cut – in some areas, numeric estimates may be available for sub-source and source level propositions, but verbal descriptors may be required for activity level propositions due to a lack of data or models at the present time.

Practitioner testing

The internal use of the framework, whether the assignment of LRs or the use of the verbal scale, should be verified for use within each area. It is critical to confirm that there is consistency and correctness in the use of the framework amongst all staff. Laboratories currently reporting within the evaluative reporting paradigm have noted that consistency amongst trained staff is high – any variation in assignment of approximate LR or verbal descriptor is generally between adjacent levels on the scale. If, when tested, staff vary more than this in their assignment of the weight for the same comparison with the same propositions, additional training or guidance in the use of the framework may be required. The collection of variance rates for opinions may also be an important consideration for the fact-finder – high variance between practitioners may indicate lower reliability, which may indicate lower probative value. Policies should be developed to address variation between propositions, interpretations or conclusions, in the same way that differences in opinion under peer review are addressed. There should also be a mechanism to disclose any potential differences to courts under a fully transparent reporting model. Preparing reporting standards, policies and procedures

Policies, procedures and reporting standards should be prepared to support consistent application of evaluative reporting methodology among practitioners. This should include procedures for peer review and resolution. This documentation should be continually updated as reporting standards and processes are refined.

Education and training

Following establishment of set criteria for either the assignment of the numeric likelihood ratio, or of the verbal descriptor, extensive training in evaluative reporting must occur for all practitioners that will be utilising the framework. The educational program should cover the Bayesian paradigm, the development of case-specific propositions, the evaluation of evidentiary weight under the framework, the reporting and explanation of the framework and its application to the case in question. Laboratories may find useful guidance material for staff in the practitioner guides produced by the Royal Statistical Society (UK) and ENFSI, covering the fundamentals of probability and statistical evidence, the use and logic of inferential reasoning and the interpretation of expert evidence. A series of articles within *Science and Justice*, the *Journal of The Chartered Society of Forensic Sciences*, also provide helpful information regarding the theory and application of the evaluative reporting framework within an operational casework laboratory. References to these documents can be found in the subsequent “Resources” section.

It has previously been noted that the development of case-specific propositions is perhaps the most difficult aspect of learning evaluative reporting framework. The propositions chosen will vary depending on the information known within the case, the ability for a particular technique to address aspects of the evidence and the questions being posed to the forensic practitioner (if any). Emphasising the difference between propositions and explanations has proven helpful (Evetts et al 2000), as has emphasising the difference between source and activity level propositions. Given the difficulty that may be encountered in formulating propositions, value may be obtained in introducing propositional logic education prior to moving to full evaluative reporting, to enable practitioners to gain experience and insight into the various propositions that may be encountered and assessed within their own casework. It will also provide an appreciation of the changing nature of propositions as a case develops – initial assessments may only allow source level propositions to be addressed due to a lack of case information. Subsequent discussions with informants, prosecutors and defence may allow a refinement in the propositions considered, with a shift to activity level if required.

Within disciplines not used to evaluative reporting, the identification and training of an evaluative reporting champion, who is also a discipline expert, can be critical to the uptake of the approach and engagement from staff. This champion, with a deep understanding of the logical framework, the development of propositions and the assignment of LR or application of the scale to their discipline, can be invaluable in mentoring staff and providing on-going training and assistance. In the same manner, the preparation of model cases reported under the evaluative framework can be a useful training aid for both reporting staff and stakeholders. Such cases may aid staff in visualising the new paradigm, and how they may apply it to their own cases.

Stakeholder engagement

Stakeholder education is a critical part of introducing the evaluative framework. Courts and informants may be accustomed to receiving reports with identification statements, and may find adjusting to a probabilistic framework difficult if the transition is not accompanied by education programs around the framework and reporting models. If reports and statements are accompanied by descriptive explanations of the terminology and scale, then the adjustment may be lower, but pre-education regarding the philosophy and operation of the framework enables a more comprehensive understanding of the logic and reasoning behind the new reporting style.

Implementation

While implementing a full evaluative reporting model across multiple disciplines within an operational casework laboratory may seem a daunting task at the outset, there are laboratories worldwide who have successfully made the transition. Specialist education from practitioners accustomed to reporting probabilistically will aid those new to this paradigm, and enable potential pitfalls and difficulties to be identified early, mitigating risk associated with the change. A staged approach, guided by highly trained users proficient in the logical framework, can assist in easing the resourcing burden upon reporting disciplines. When successfully implemented, evaluative reporting provides a logical framework for conducting forensic examinations and communicating to the courts and other criminal justice system stakeholders the weight of the findings given the propositions.

Conclusion

Evaluative reporting combines and applies practitioner's knowledge and experience with available, relevant information to evaluate the probability of the findings in light of two competing propositions that are formed to assist the Court. It shows the clear path of reasoning undertaken by the practitioner in reaching their conclusion, and this transparency assists not only the practitioner (in formulating their opinion), but also the court in assessing the expert testimony. Further explanation of evaluative reporting can be found in the examples located in the appendices and in the resources listed below.

Resources

COMPREHENSIVE GUIDES - the following references provide a comprehensive overview on the use of evaluative reporting in forensic science. Both are available online:

S.M. Willis, L. McKenna, S. McDermott, G. O'Donell, A. Barrett, B. Rasmusson, A. Nordgaard, C.E.H. Berger, M.J. Sjerps, J.-J. Lucena-Molina, G. Zadora, C. Aitken, T. Lovelock, L. Lunt, C. Champod, A. Biedermann, T.N. Hicks, F. Taroni, ENFSI Guideline for Evaluative Reporting in Forensic Science, European Network of Forensic Science Institutes, 2015

'Case Assessment and Interpretation of Expert Evidence', Practitioner Guide No. 4, Working Group on Statistics and Law of the Royal Statistical Society

JOURNAL ARTICLES AND BOOKS - the following references provide background information and an overview of the theory behind evaluative reporting:

Buckleton, J, Nichols, R, Triggs, C, and Wevers, G, An Exploratory Bayesian model for firearm and tool mark interpretation. *AFTE Journal*, 37 4 (2005) 352-361

Bunch, S G, Consecutive Matching Striation Criteria: A General Critique, *Journal of Forensic Sciences*, 45 5 (2000) 955-962.

Bunch S & Wevers G, Application of Likelihood Ratios for Firearm and Tool Marks Analysis, *Science and Justice*, 2013 Vol 53.2, pp 223-229

Cook, R., Evett, I.W., Jackson, G., Jones, P.J. & Lambert, J.A., 1998, 'A model for case assessment and interpretation', *Science & Justice*, vol. 38, pp. 151-156

Cook, R., Evett, I.W., Jackson, G., Jones, P.J. & Lambert, J.A., 1998, 'A hierarchy of propositions: deciding which level to address in casework', *Science & Justice*, vol. 38, pp. 231-239

Champod, C., Lennard, C., Margot, P., Stoilovic, M., 2016 *Fingerprints and Other Ridge Skin Impressions* 2nd Ed, CRC Press Boca Raton p70-116

Curran, Hicks, Buckleton. *Forensic Interpretation of Glass Evidence*. CRC Press. 2000.

Evett, I.W., Jackson, G. & Lambert, J.A., 2000, 'More on the hierarchy of propositions: exploring the distinction between explanations and propositions', *Science & Justice*, vol. 40, pp. 3-10

Evett, I.W., Jackson, G., Lambert, J.A. & McCrossan, S., 2000, 'The impact of the principles of evidence interpretation on the structure and content of statements', *Science & Justice*, vol. 40, pp. 233-239

Jackson, G., Jones, S., Champod, C. & Evett, I.W. 2006, 'The nature of forensic science opinion – a possible framework to guide thinking in investigations and in court proceedings', *Science & Justice*, vol. 46, no. 1, pp. 33-44

Kerkhoff, W et al, The likelihood Ratio Approach in Cartridge Case and Billet Comparison, *AFTE Journal* 2013, Vol 45:3, pp284-289.

Lucy, D., 2013. *Introduction to statistics for forensic scientists*. John Wiley & Sons.

Newton, A, The Association between a Paint Flake and a Wheelbarrow on the Basis of Toolmarks, *AFTE Journal*, Vol 45:3, pp245-251.

Stoney, D A, What Made Us Ever Think We Could Individualize Using Statistics, *Journal of The Forensic Science Society*, 31 2 (1991) 197-199.

Association of Forensic Science Providers, 2009, 'Standards for the Formulation of Evaluative Forensic Science Expert Opinion', *Science & Justice*, vol. 49, pp.161

Berger, C.E.H., 2010, 'Criminalistics is reasoning backwards – Logically correct reasoning in forensic reports and in the courtroom', *Nederlands Juristenblad*, vol. 85, pp784-789

Marquis et al., 2016, 'Discussion on how to implement a verbal scale in a forensic laboratory: Benefits, pitfalls and suggestions to avoid misunderstandings', *Science & Justice*, col. 56, issue 5, 364-370

Robertson, B., Vignaux, G.A. & Berger, C.E.H., (2016 *'Interpreting Evidence – Evaluating Forensic Science in the Courtroom 2nd Edition'*, Wiley (ISBN: 978-1-118-49243-7)

Lindley, D.V., 2013, *'Understanding Uncertainty, Revised Edition'*, Wiley (ISBN: 978-1-118-65012-7)

Wevers, G., Neel, M. and Buckleton, J. A comprehensive statistical analysis of striated tool mark examinations part 2: Comparing known matches and known non-matches using likelihood ratios, *AFTE Journal*, 43 2 (2011) 137-145.

FORMAL TRAINING COURSES - formal training courses in the evaluation of forensic evidence are available. The University of Lausanne, Switzerland, offers the following (note: other institutions may also offer courses):

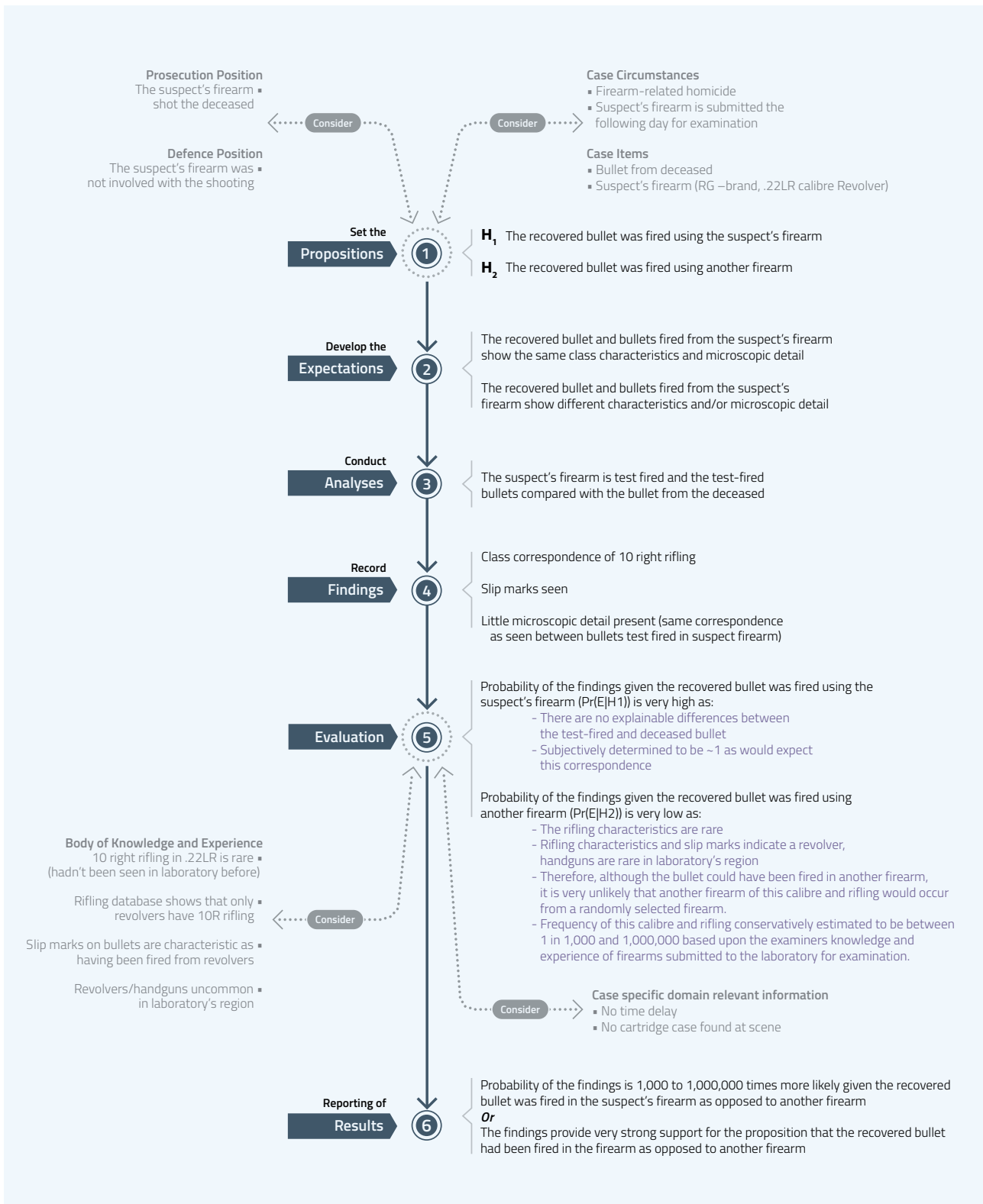
Certificate of Advanced Studies in Statistics and the Evaluation of Forensic Evidence.

Essentials of Bayesian Networks in Forensic Science.

ONLINE FREE ACCESS COURSES - short courses in subjects including statistics and logical thinking are available online. Websites sites such as 'Coursera' and 'edX' are two such providers.

Appendix A – Firearms case example

Figure 5:
Flow diagram depicting the steps in evaluative reporting for a firearms case example.



Let us imagine that the forensic practitioner is presented with a bullet that is recovered from a deceased and a firearm seized from the suspect. The examination request is to determine if the recovered bullet could have been fired by the suspect's firearm. This would lead to the establishment of two competing propositions, one that perhaps aligns to the prosecution position **"the recovered bullet was fired in the suspect's firearm"** and one that perhaps aligns to the defence proposition **"the recovered bullet was fired in another firearm"**. Although, the practitioner cannot determine who actually fired the weapon, matching the ammunition to the weapon provides vital facts for the investigation or the trier of fact.

Typically, when a firearm is fired, its surface (e.g. firing pin and barrel) leaves characteristic marks or striations on the bullet. These striations can display features that may be characteristic of the class of firearm (i.e. features shared between brands and models of firearms such as calibre, riffling, groove width, number and direction) and other microscopic detail that may be left by a particular firearm.

Using this knowledge, the expected outcome aligned to the prosecution proposition would be **"the recovered bullet and bullets fired from the suspect's firearm show the same class characteristics and microscopic detail"** and the expected outcome aligned to the defence proposition would be **"the recovered bullet and bullets fired from the suspect's firearm show different characteristics and/or microscopic detail"**.

When conducting the analysis, the examiner needs to first determine if the absence of shared class characteristics excludes the crime scene bullet as having been fired from the suspect's firearm (e.g. different calibre bullet). In this example, the bullet found is a .22LR calibre with 10 right riffling. The firearm submitted is a .22LR calibre revolver with a barrel having 10 right riffling. As the class characteristics cannot exclude that this bullet was shot from this firearm, it is time to look at microscopic detail. This is done by firing 10 bullets in the submitted firearm to determine the intra-variability of microscopic detail left by the barrel of the submitted firearm and then comparing them with the microscopic detail on the recovered bullet and the test-fired bullets.

Minimal microscopic detail was found, so if the practitioner was to use traditional identification reporting, the result would be inconclusive. However in a probabilistic approach, it is possible to integrate external information such as the frequency of appearance of the type of firearm to add weight to these observations. The complete reasoning is explained in figure 5.

The strength of the findings comes from the size of the "alternative source" population that would share the same pertinent features seen in the crime scene sample. If this population is large (as the features are common), then the likelihood ratio will be low as there is a relatively high probability of a random match. The more features that are present, the smaller the alternative source population becomes as there is less probability of obtaining the same pertinent features from a sample taken at random. Survey data will assist in determining the size of the "alternative source" population. Subjectively assessing likelihood ratios basically comes down to assessing the size of this "alternative source" population as will be shown in this example, and other examples contained in this guide.

In this firearm example, let's assume we have a very unusual calibre bullet found at a crime scene with unusual riffling. A firearm was submitted to the laboratory of the same unusual calibre and riffling. A comparison of the bullets test fired in the firearm with the bullet recovered from the deceased showed a correspondence of class characteristics only with little microscopic detail present (the test-fired bullets also showed the same amount of correspondence).

How strong are the findings?

Pr(E|H1) “the probability of the findings (E) given the bullet was fired using the suspect's firearm”.

In Figure 5, a correspondence was observed of class characteristics and both recovered bullet and test-fired bullets were relatively featureless microscopically. There were also no significant differences seen.

The probability of observing a correspondence of class characteristics given the bullet was fired using the suspect's firearm will be very high as we would expect to find this correspondence. We can assume the probability is 1 (or very close to). The strength of the findings will largely depend on **Pr(E|H2)**.

Pr(E|H2) “the probability of the findings (E) given the bullet was fired using another firearm”.

We know from experience and training that firearms of this calibre and rifling are very rare. Therefore, the probability of finding a correspondence of class characteristics given the bullet was not fired using the suspect's firearm will be very low, (this would be based on the frequency of firearms of this calibre and rifling in the firearm population).

The ratio of these two subjective probabilities would result in a large likelihood ratio as we are dividing a relatively large probability by a small probability. As the LR will be greater than 1, the findings would provide support for an association between the bullet and the firearm. The next step is to determine the strength of the findings or the level of support that would be provided.

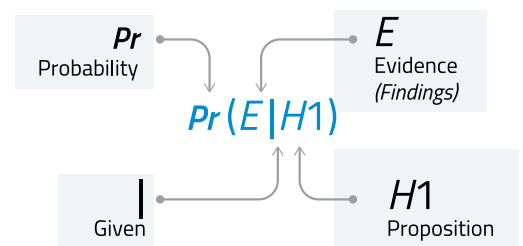
This final conclusion will be up to the examiner based on their experience and training. However, it could be conservatively estimated that between 1 in 1,000 and 1 in 1,000,000 firearms would share these characteristics.

$$Pr(E|H2) = \frac{1}{1000} \text{ to } \frac{1}{1,000,000}$$

$$Pr(E|H1) = \sim 1$$

$$LR = \frac{Pr(E|H1)}{Pr(E|H2)} = \frac{1}{\frac{1}{1000}} \text{ to } \frac{1}{\frac{1}{1,000,000}} = 1000 \text{ to } 1,000,000$$

What does the formula mean?



Thus the probability of the findings is 1000 to 1,000,000 times more likely given the bullet was fired using the suspect's firearm than using another firearm. An example of the wording used in a statement could be (depending on the verbal scale):

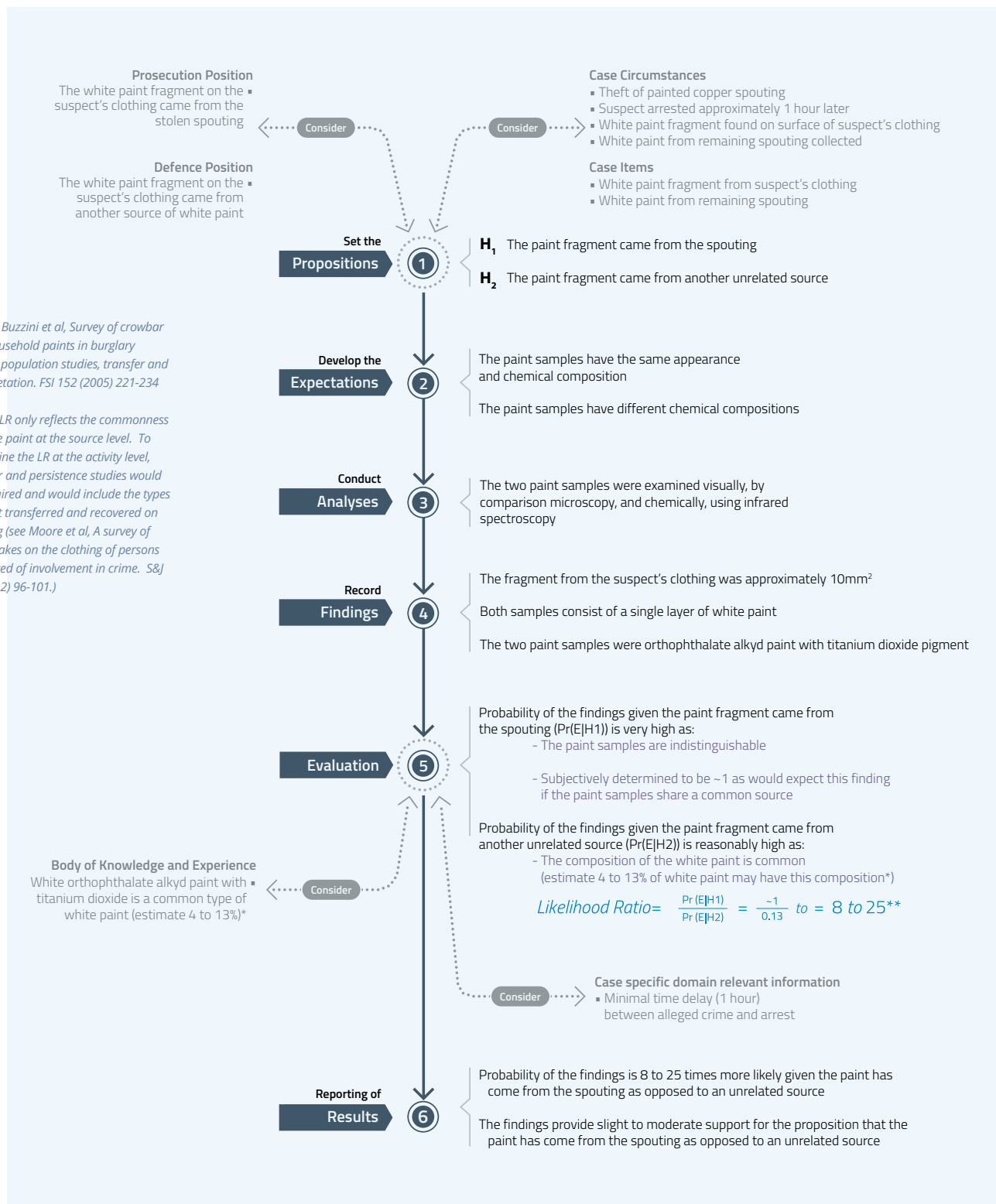
“In interpreting these findings, I have given consideration to the probability of observing this correspondence given the bullet was fired using the firearm, as opposed to observing this correspondence given the bullet was fired in another firearm.

In my opinion, this calibre of bullet and rifling is very uncommon, and as such, the probability of observing this correspondence given the bullet was fired using another firearm is extremely low. Conversely, the probability of these findings are high given the bullet was fired using this firearm.

Therefore in my opinion, the findings provide very strong support for the proposition that the bullet had been fired using the firearm as opposed to another firearm.”

Appendix B – Paint case example

Figure 6: Flow diagram depicting the steps in evaluative reporting for a paint analysis case example.



Paint can come from many sources, such as vehicles, architectural sources, tools etc. The source of paint, the number of layers, the paint layer colour and the paint chemical composition are all usually considered when interpreting paint trace evidence. The more corresponding layers and colours there are between two paint samples, the larger the likelihood ratio becomes. This is because the probability of randomly obtaining these findings from another source becomes smaller.

Consider a comparison of two single-layered white paint samples. White paint typically lacks the various tints, tones and shades that coloured paints have. White paint will also typically contain only one of two white pigments. Therefore, a single-layered white paint fragment has weaker evidential value compared to that of single-layered coloured paint fragment. One could consider conservatively that 4 to 13% of orthophthalate alkyd white paints chosen at random from the population could share the same appearance and chemical composition as the paint in question. (Buzzini et al, Survey of crowbar and household paints in burglary cases – population studies, transfer and interpretation. Forensic Science International, vol 152 (2005) 221-234).

So how strong are the findings?

If the two paint samples have come from the same source, we would expect a correspondence of appearance and chemical composition. Therefore, the probability (Pr) of the findings (E) given the paint has come from the same source (prosecution proposition H1) will be high and can be given a probability of approximately one. This can be expressed in the following formula:

$$Pr(E|H1) = \sim 1$$

If the paint samples have come from different sources, then we have to consider the potential population of sources of this type of paint. We previously estimated that 4 to 13% of white paint taken at random could also have the same appearance and chemical composition as the paint in question. Therefore, the probability of the findings given the paint has come from an unrelated source of white paint (Pr(E | H2)) is approximately 0.04 to 0.13.

$$Pr(E|H2) = 0.04 \text{ to } 0.13$$

The likelihood ratio becomes:

$$LR = \frac{Pr(E|H1)}{Pr(E|H2)} = \frac{\sim 1}{0.13} \text{ to } \frac{\sim 1}{0.04} = \sim 8 \text{ to } 25$$

In other words, the probability of the findings is approximately eight to twenty-five times greater given the two white paint samples have come from the same source as opposed to another source of white paint taken at random. It is not necessary for the numerical value of the likelihood ratio to be reported and in its place a corresponding verbal scale could be used, for example:

"In my opinion, this correspondence provides slight to moderate support for the proposition that these two samples of white paint have come from the same source, compared to coming from different sources".

An example of reporting this evidence could be:

“In interpreting the findings, I have given consideration to the probability of observing this correspondence given the two paint samples have come from the same source of white paint, as opposed to observing this correspondence given the two paint samples have come from two unrelated sources of white paint.

In my opinion, this type of white paint is relatively common, and as such, the probability of observing this correspondence given the two samples of paint have come from different sources of white paint is reasonably high. Conversely, the probability of the findings is high given the two samples of white paint have come from the same source.

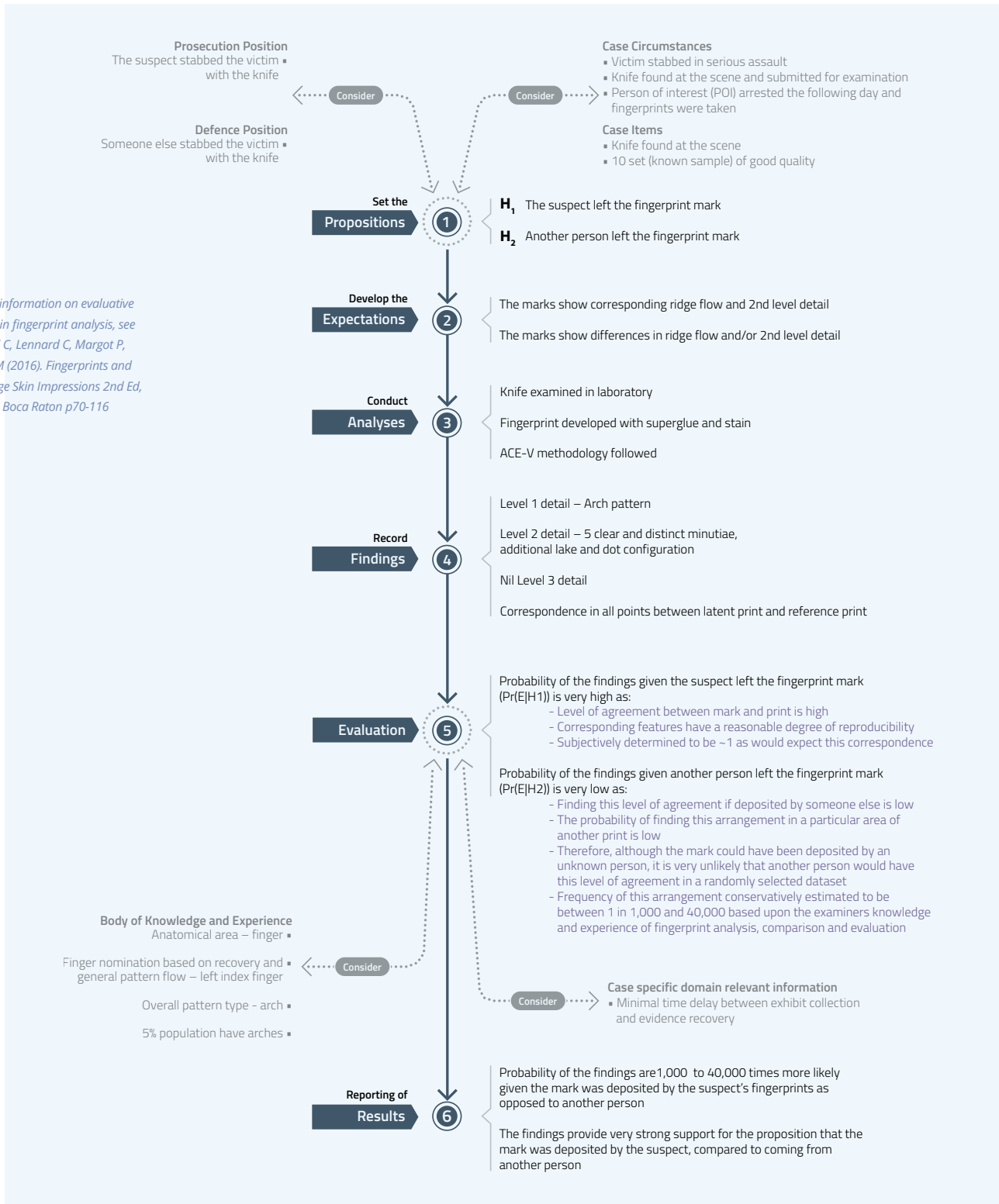
Therefore, in my opinion, this correspondence provides slight to moderate support for the proposition that these two samples of white paint have come from the same source, compared to coming from different sources”.

The above example has only considered the probability of the findings at the source level, that is could these samples of paint come from the same source? To determine the probability of the evidence at the activity level, that is the suspect having come into contact with the spouting, resulting in paint being transferred and retained on his clothing, requires transfer and persistence studies and would include the types of paint recovered from clothing. A good example of an appropriate survey was carried out by Moore et al (A survey of paint flakes on the clothing of persons suspected of involvement in crime. Science and Justice, vol 52 (2012) 96-101.)

As can be seen, the strength of the findings comes from the commonality or rarity of the paint in question in the relevant population. For example, one can easily subjectively assess that the probability of obtaining two unrelated twenty-layered, multi-coloured paint samples is extremely unlikely and would result in an extremely high likelihood ratio, or more commonly subjectively assessed as conclusive.

Appendix C – Fingerprint case example

Figure 7: Flow diagram depicting the steps in evaluative reporting for a fingerprint case example.



For more information on evaluative reporting in fingerprint analysis, see Champod C, Lennard C, Margot P, Stoilovic M (2016). *Fingerprints and Other Ridge Skin Impressions 2nd Ed*, CRC Press Boca Raton p70-116

In this fingerprint example, a person is stabbed during an assault. A suspect is arrested and a knife is recovered from him. His fingerprints are taken upon arrest and the knife is collected and submitted to the laboratory. The examination request is to determine if fingermarks found on the knife could have been left by the suspect.

Generally, when a fingertip touches a surface, it leaves a fingermark comprising of friction ridges. A comparison between the friction ridges of the fingermark with friction ridges of a fingerprint collected from a suspect can establish whether a fingermark was left by a specific finger of a person.

Three levels of details can be observed:

- Level 1:** **General pattern**
(loops, whorls, arches - similar to class characteristics)
- Level 2:** **Minutiae or ridge characteristics**
(lake, dots, fork - individual characteristics)
- Level 3:** **Ridge pores**
(individual characteristics)

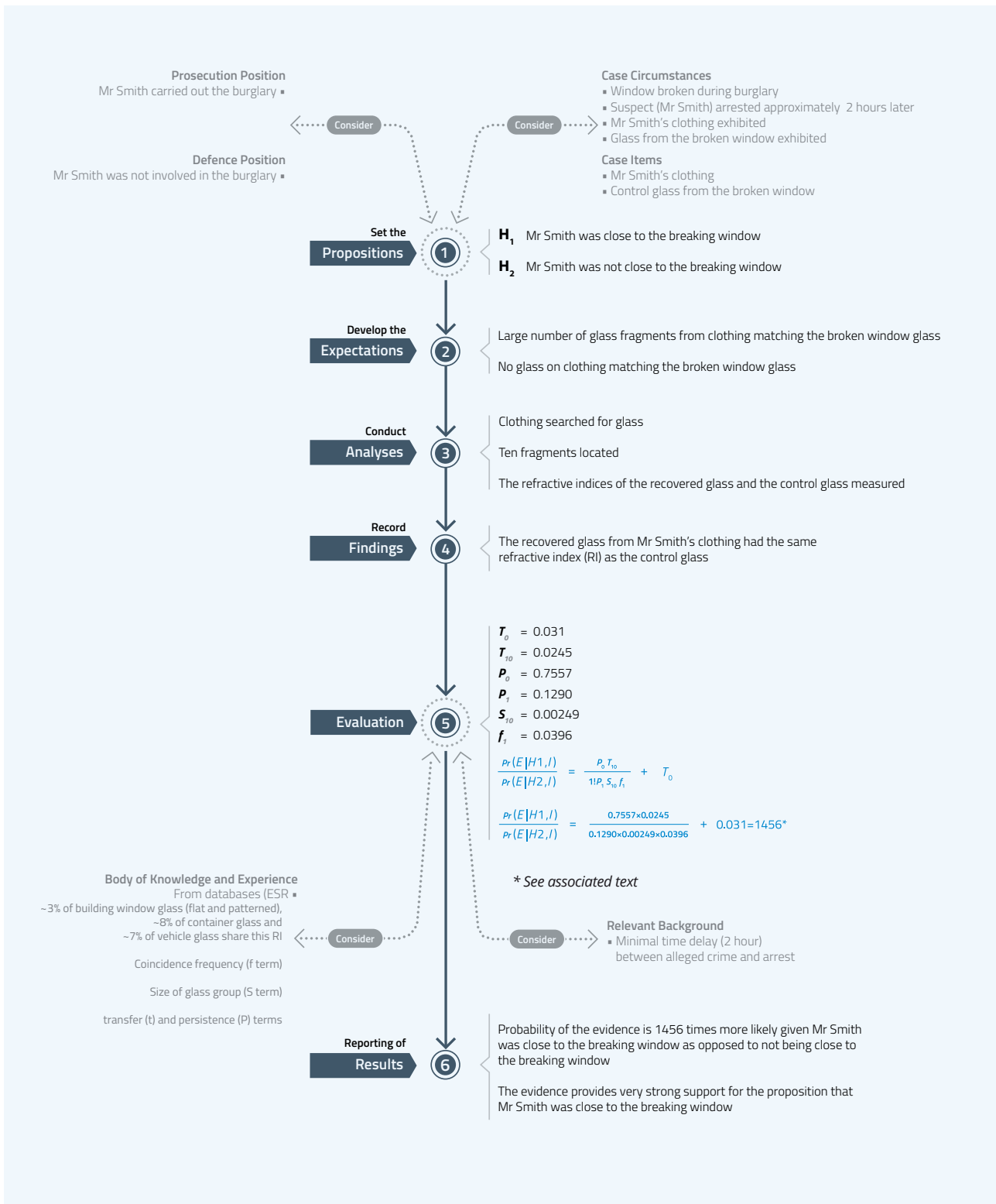
The comparison examination is done by first determining if the level 1 findings on the mark are similar to level 1 findings of the suspect's fingerprints. In this case, one fingermark is found and its level 1 is an arch. Four of the suspect's fingerprints are arches.

As the class characteristic cannot exclude the suspect's fingerprints, the next step is to look at level 2 characteristics. Five minutiae are observed on the fingermark. Five same minutiae are observed in the same arrangement on one of the fingerprints of the suspect and no differences between the mark and the print are observed. Level 3 examination did not reveal sufficient detail.

At this point, if an expert were to report the strength of the findings in the traditional identification model, it would be inconclusive. However, in a probabilistic approach, it is possible to consider external information such as the frequency of appearance of these minutiae in this arrangement in the population to add weight to these observations. Please see the flowchart for the complete reasoning.

Appendix D – Glass case example

Figure 8:
Flow diagram depicting the steps in evaluative reporting for a glass analysis case example.



Likelihood ratios from glass evidence can be assigned if the relevant databases have been compiled. Consider a case where clothing from Mr Smith, who is suspected of breaking a window, is submitted to the laboratory along with a sample of glass from the broken window. We are informed that his clothing was collected 2 hours after the window was alleged to have been broken. A search of the clothing reveals the presence of ten fragments of glass. Using refractive index measurements, it is found that all ten recovered fragments have the same refractive index as the broken window.

Unlike the previous firearm and paint examples where we only considered the rarity of certain pertinent features, for activity level glass evidence (i.e. was Mr Smith close to the breaking window?) we also need to consider other factors, such as;

- *How many glass fragments would typically be transferred to clothing when glass is broken.*
- *The period of time from when the glass was broken to when the clothing was exhibited.*
- *How many of the transferred glass fragments would persist (be retained) on the clothing during this time period (T term).*
- *The number of pre-existing groups of glass from different sources that are found on clothing i.e. background glass (P term).*
- *How common the matching glass type is amongst other sources of broken glass (f term).*

These factors all require databases and studies which are outside the scope of this document.

How are these factors used?

The prosecution position is that Mr Smith broke the window. If he did break the window, then the glass on his clothing could have got there as a result of him being close to the breaking window, resulting in ten fragments (T_{10}) persisting on his clothing after the glass was transferred. We assume that all ten fragments of glass have come from the broken window and that there were no pre-existing groups of glass already present (P_0) on his clothing.

Mr Smith could also have been close to the breaking window, however, no glass was transferred (T_0) and there was already one pre-existing group (P_1) of ten fragments (S_{10}) on his clothing, which just happened to have the same refractive index (f_1) as the broken window.

In summary,

- T_0 *the probability of no fragments being transferred*
- T_{10} *the probability of ten fragments being transferred*
- P_0 *the probability of no existing (persisting) group of fragments*
- P_1 *the probability of one existing (persisting) group of fragments*
- S_{10} *the probability of an existing group of ten fragments (group size 10)*
- f_1 *coincidence frequency of the pre-existing "matching" glass*

This gives the following formula

$$Pr(E|H1,I) = P_0 T_{10} + 1!P_1 S_{10} f_1$$

The various scenarios are added as they are mutually exclusive events (i.e. the ten fragments of glass can't be transferred as well as already being present on the clothing). The defence position would be that the group of ten fragments did not come from the broken window and there was already a single group of glass (P_1) consisting of ten fragments on Mr Smith's clothing (S_{10}) that just happened to match the broken window (f_1). Therefore,

$$Pr(E|H2,I) = 1!P_1 S_{10} f_1$$

Thus the likelihood ratio becomes

$$\frac{Pr(E|H1,I)}{Pr(E|H2,I)} = \frac{P_0 T_{10} + P_1 S_{10} f_1}{P_1 S_{10} f_1} = \frac{P_0 T_{10}}{P_1 S_{10} f_1} + T_0$$

Using the following values, the likelihood ratio is assigned as follows (these have been calculated from ESR's New Zealand glass databases and glass programs).

T_0	= 0.031	$\frac{Pr(E H1,I)}{Pr(E H2,I)} = \frac{P_0 T_{10}}{1!P_1 S_{10} f_1} + T_0$
T_{10}	= 0.0245	
P_0	= 0.7557	$\frac{Pr(E H1,I)}{Pr(E H2,I)} = \frac{0.7557 \times 0.0245}{0.1290 \times 0.00249 \times 0.0396} + 0.031 = 1456$
P_1	= 0.1290	
S_{10}	= 0.00249	
f_1	= 0.0396	

i.e. the evidence is 1456 times more likely given the clothing was close to the breaking window as opposed to the clothing not being close to the breaking window.

Reducing the potential population of other sources.

In the above glass example, the frequency of the glass was 0.0396. This actually relates to how common that refractive index is in ESR's glass database and how close it matches the refractive index of the control glass. However, there are other characteristics that may be present on recovered glass that can also reduce the frequency and therefore increase the likelihood ratio.

From ESR's databases, the refractive index of the first example shares the same refractive index as ~3% of building window glass (flat and patterned), ~8% of container glass and ~7% of vehicle glass. If we were fortunate to have some fragments with a flat original surface, then this would reduce the potential source of other glasses to flat building glass and vehicle glass as patterned and container glass do not have flat surfaces.

If we anneal the recovered glass and find it is toughened, then all non-toughened glass sources are removed (annealing the glass removes the stresses within a glass fragment and can be used to distinguish between toughened and non-toughened glass).

If the recovered glass had an original thickness that was the same as the control (say 2.9mm), then if the glass on the clothing did not come from the broken window, it had to have come from another source of glass that was flat, toughened glass with a thickness of 2.9mm and shares the same refractive index. The following table (Table 2) shows the effect that reducing the potential sources of other glass has on the likelihood ratio.

Table 2:
Effect of Glass Source in the Likelihood Ratio (LR)

	BUILDING GLASS (%)	CONTAINER GLASS(%)	VEHICLE GLASS (%)	FREQUENCY (f) (%)	LR
RI ONLY	3	7.7	7.3	0.396	1456
RI, FLAT	2.4	0	7.3	0.0083	6945
RI, FLAT, TOUGHENED	0.1	0	4.2	0.0037	15578
RI, FLAT, TOUGHENED, 2.9MM THICK	0	0	0.21	0.0002	320224

ANZPAA
Australia New Zealand
Policing Advisory Agency



Level 6 Tower 3, World Trade Centre,
637 Flinders Street,
Docklands VIC 3008
Australia
secretariat.nifs@anzpaa.org.au
www.nifs.org.au