

Explaining the Likelihood Ratio in DNA Mixture Interpretation

Mark W. Perlin
Cybergenetics, Pittsburgh, PA

December 29, 2010

*In the Proceedings of Promega's
Twenty First International Symposium on Human Identification*

Cybergenetics © 2010



Contact information:

Mark W. Perlin, PhD, MD, PhD
Cybergenetics
160 North Craig Street
Suite 210
Pittsburgh, PA 15213
USA
(412) 683-3004
(412) 683-3005 FAX
perlin@cybgen.com

Abstract

In DNA identification science, the likelihood ratio (LR) assesses the evidential support for the identification hypothesis that a suspect contributed their DNA to the biological evidence. The LR summarizes the sensitivity and specificity of a statistical test. The LR logarithm is a standard information measure for stating the support for a simple hypothesis (i.e., a single assertion relative to its logical alternative).

After Alan Turing's LR methods cracked the German Enigma code during World War II, LR usage became widespread. The LR is ubiquitous in the physical, biological, social, economic, computer and forensic sciences. First introduced into biological identification through paternity testing, the LR enjoys unparalleled international usage as the most informative DNA mixture statistic.

Yet American crime labs avoid the LR, and prefer to report DNA inclusion statistics that they find easier to explain in court. Such "inclusion" methods (variously termed PI, CPI, CPE or RMNE) use less of the DNA data, typically discarding a million-fold factor of identification information. Thus highly informative DNA mixture evidence can be reported as "inconclusive" or assigned an unrealistically low match score. Unfortunately, minimizing DNA evidence leads to a failure to identify criminals, with an adverse effect on public safety.

To make the LR more acceptable to American analysts and their juries, we need more intuitive ways to explain the LR. Fortunately, the LR can be expressed (by Bayes theorem) in several equivalent ways. Stated in plain English, these alternative formulations include:

1. the *information gain* in the identification hypothesis from the DNA data,
2. how well the identification hypothesis *explains the data*, relative to its alternative, and
3. our *increased belief* in a match to a suspect, based on the inferred evidence genotype.

The second LR formulation prevails in forensic DNA. While natural for computers and statisticians, non-mathematicians often find its formulas opaque. In this paper, we describe the other two formulations as intuitive ways to explain the LR simply and accurately. Moreover, these other approaches avoid the dread "transposed conditional." Using DNA case examples, we show how to easily understand the LR, present it in court, and deflect superficial challenges.

For the American public to benefit from the full protective power of DNA identification information, analysts must be able to confidently explain the LR. This paper shows them how.

Table of Contents

Abstract	2
Introduction.....	4
History	4
Presenting the LR in four different ways	5
Hypothesis form	5
Likelihood form.....	5
Genotype form	6
Match form	7
DNA mixtures	7
Quantitative mixture interpretation.....	8
TrueAllele computer inference	9
Case and validation examples	10
A 7% minor contributor from a victim's fingernail	10
Multiple amplifications of low-level three contributor DNA.....	11
Large variation in human interpretation.....	12
Computers preserve DNA identification information.....	12
Computers are a million times more informative	13
Human review discards most DNA information	13
LR-based databases for DNA investigation	14
SWGAM mixture guidelines	15
The LR and forensic science	15
Conclusion.....	16
Resources	16
Appendix.....	17
Genotype probability notation	17
Equivalent LR formulations	17
LR equivalence proofs	19
References	23
Figure Legends.....	25
Figures.....	26

Introduction

The likelihood ratio (LR) appears in many fields of biological, information, physical and social science. The LR is a standard measure of information that summarizes in a single number the data support for a hypothesis. It is a way of accounting for all the evidence in favor of or against a particular hypothesis (or proposition) (1). The LR is also the match statistic that is used in DNA reporting (2-4). The LR's good legal and scientific standing underlies forensic science's credibility in court. Importantly, the LR quantifies how our belief in a hypothesis changes after observing experimental data.

Yet the likelihood ratio has not yet gained traction in the United States for reporting on DNA evidence. There are several reasons for its relative unpopularity. First, outside of the DNA area, most forensic science disciplines do not yet have a LR available. Moreover, within DNA identification, forensic analysts sometimes find the LR hard to explain. However, all DNA match statistics (including inclusion) are likelihood ratios (5). The stronger likelihood ratio methods can preserve DNA match information, while the weaker ones discard match information. Since our scientific goal is to preserve information, it can be very helpful to know how these LR approaches differ.

Without a likelihood ratio, highly informative DNA can be misreported. For example, one might incorrectly state that such evidence is inconclusive (which lets criminals go free). This paper introduces an easier way to explain the likelihood ratio. By providing a "plain language" wording, this less forbidding approach may help forensic scientists become more comfortable with the LR, so that they can regularly use LRs in DNA reporting and testimony.

History

Let us briefly review the (somewhat British) history of the likelihood ratio. It was in the 18th century that the Reverend Thomas Bayes first came up with the idea of updating one's beliefs (or probability) based on data (6). He showed how to use evidence to revise our belief in a hypothesis. "Bayes theorem," popularized in the 19th century by Pierre-Simon Laplace (7), is the cornerstone of statistical inference based on mathematical probability.

In the 1940s, Alan Turing (the father of computer science) used likelihood ratios for the Enigma code breaking project (8). Jack Good, a statistician who worked with Alan Turing, ushered the LR into mainstream scientific thought with his classic book "Probability and the Weighing of Evidence" (1). His beautifully written Chapter 6 of that book describes the modern LR, and is informative reading for scientists.

In the 1970s, Dennis Lindley, a Bayesian statistician in England, introduced likelihood ratios into forensic science in a rigorous way, starting with glass evidence (9). Lindley's "Understanding Uncertainty" is an outstanding book written for lawyers, judges

and nonscientists (10). Many nonspecialists have found this book to be an intuitive introduction to probability concepts. Over the last two decades, John Buckleton (11), Ian Evett (12), Bruce Weir (13) and others (14) have brought the likelihood ratio into the interpretation of DNA and mixtures.

Presenting the LR in four different ways

We examine four different forms of the likelihood ratio. Each form has its own mathematical formula and scientific interpretation. Even though they appear to be very different, these four LR forms are actually equivalent to one another. We prove this equivalence in the Appendix.

Hypothesis form

Here is the original LR form that appeared sixty years ago in Chapter 6 of Jack Good's classic book (1). This is the *hypothesis* form of the LR. It focuses on the identification hypothesis, which in our case for DNA is "the suspect contributed to the evidence." We start off with our prior belief about that hypothesis, based on randomly selected people in a population *before* we examine any data in a case. *After* we have seen the data, we update our belief in the hypothesis.

Now we look at the odds of the identification hypothesis, given that we have seen the data (numerator), relative to what we knew before (denominator). The LR is the ratio of these two numbers. In other words, what was the information gain based on the data?

$$\text{information gain in hypothesis} = \frac{\text{Odds}(\textit{hypothesis}|\textit{data})}{\text{Odds}(\textit{hypothesis})}$$

We can state this odds ratio of information gain in plain English. Suppose that the factor was a billion (we use this "billion" LR number in our examples throughout the paper). We could say "the evidence increased our belief that the suspect contributed to the DNA by a factor of a billion." This LR sentence resembles ordinary language.

Likelihood form

The *likelihood* form below is an equivalent LR (Appendix), but is the crux of our problem. This LR expression can discomfort many DNA analysts, who question whether the technical math can be explained to a jury in a way that is readily understood. The form is based on the "likelihood" concept, which is the conditional probability of the observed data, assuming some hypothesis (15). The likelihood is a mathematical way of saying how much a particular hypothesis explains the data.

$$\text{information gain in likelihood} = \frac{\text{Prob}(\text{data}|\text{identification hypothesis})}{\text{Prob}(\text{data}|\text{alternative hypothesis})}$$

This LR form supposes that there is an *alternative* hypothesis that someone else (other than the suspect) contributed to the evidence. These two LR hypotheses (either the suspect contributed or he did not) are exhaustive and mutually exclusive. Statisticians and computers often contrast two hypotheses, using (what is for them) straightforward mathematics. The LR here is a ratio of two likelihoods: the probability of the data given the identification hypothesis, divided by the probability of the data given the alternative hypothesis.

In plain language, we might state the LR as "the probability of observing the evidence assuming that the suspect contributed to the DNA is a billion times greater than the probability of observing the evidence assuming that someone else was the contributor."

Since the time of Chaucer and Shakespeare, England has long celebrated prolix descriptive phrasing. Americans, however, often prefer to express themselves more succinctly. So while British juries and analysts may savor the exactitude of a long LR sentence imparting conditional probabilities, that approach might not travel well across the Atlantic. George Bernard Shaw once said "England and America are two countries separated by a common language." Perhaps that linguistic distinction helps explain why LRs are currently not as widely used in this country.

There is a real risk of someone transposing a conditional probability in court, and thus invalidating the DNA evidence. Defense attorneys may intentionally try to put the wrong words into the mouth of an expert witness, while prosecutors may achieve this by accident. The likelihood concept is subtle, since the "probability of the data" under different hypothesis does not even form a probability distribution – the likelihood numbers do not add up to one.

Genotype form

Here is a third form of the likelihood ratio that is mathematically equivalent to the others (Appendix). If we know what a genotype is (we'll see this visually later on), we can state the genotype information gain at the suspect's genotype. Now this LR *genotype* form is starting to look a bit more understandable. Since there is no conditional probability in this form, there is no way that we can inadvertently transpose a conditional.

$$\text{information gain in genotype} = \frac{\text{Prob}(\text{evidence genotype})}{\text{Prob}(\text{coincidental genotype})}$$

The LR genotype form simply compares the probability of the evidence genotype, relative to a coincidental genotype, as evaluated at the suspect's genotype. In other words, before we saw the case data, there was a random population genotype based

on the product rule (e.g., $2pq, p^2$). After we've seen the data, that genotype has been updated (by Bayes theorem). We might read the math in natural language as "at the suspect's genotype, the evidence genotype is a billion times more probable than a coincidental genotype."

This genotype form works well for people who understand genotypes, or perhaps can see a helpful picture (as we shall see in Figure 5). But maybe the jury does not know about genotypes, or there is no picture to show them. So let us introduce another LR form.

Match form

The *match* form is the simplest way to write the LR in straightforward English. Again, this form is mathematically equivalent to the other three LR expressions (Appendix).

$$\text{information gain in match} = \frac{\text{Prob}(\textit{evidence match})}{\text{Prob}(\textit{coincidental match})}$$

The match LR form addresses the question, "How much more does the suspect match the evidence than some random person?" That is, "What is the match information gain?" The formulation emphasizes DNA match, without explicitly mentioning hypotheses, data or genotypes. Most people have an intuitive idea about what it means for things to match. This LR form compares the probability of an evidence match to a coincidental one.

In regular language, the LR match form says "a match between the suspect and the evidence is a billion times more probable than a coincidental match." That statement is certainly straightforward language. It may even be more comprehensible than the usual random match probability (RMP) phrasing for single source DNA. We will continue to use this match form of the LR, illustrating its use with some examples in the remainder of this paper.

DNA mixtures

DNA mixtures occur when more than one person contributes to a biological specimen. Mixtures are common in forensic DNA practice. Their interpretation is interesting because the data can suggest more than one possible allele pair value for a contributor genotype at a locus.

A likelihood ratio compares an evidence match relative to coincidence. However, each mixture interpretation method uses its own likelihood function to explain the data. Therefore, different interpretation methods extract varying DNA information, producing different LR match scores (16). Let us review three representative methods.

(a) Random match probability is done on single source DNA evidence, or when there is a clear major contributor to the mixture. We write the LR here as *one* over the probability of a coincidental match (as computed using the product rule). RMP describes the chance of seeing a random match in the population.

$$\text{random match LR} = \frac{1}{\text{Prob}(\textit{coincidental match})}$$

(b) The "inclusion" mixture interpretation method is popular because analysts find it easy to understand and explain. However, as we shall soon see, inclusion diffuses genotype probability over many allele pair possibilities, most of which have no support in the data. The result is a *small* matching genotype probability relative to a coincidental match, producing a small LR.

$$\text{inclusion LR} = \frac{\textit{small match probability}}{\text{Prob}(\textit{coincidental match})}$$

(c) Quantitative interpretation methods can preserve more DNA identification information by inferring a *large* matching genotype probability for those allele pairs having data support. Relative to a coincidental match, the LR is therefore greater.

$$\text{quantitative LR} = \frac{\textit{large match probability}}{\text{Prob}(\textit{coincidental match})}$$

Note that the denominator of coincidental match is the same in all three mixture interpretation methods. The LR strength differences occur in the numerator, with the probability of an inferred genotype based on the evidence. The DNA data indicates how much weight to assign each allele pair. More informative methods use more of the data to place greater probability on the correct allele pair solution.

Quantitative mixture interpretation

Figure 1 shows quantitative mixture data at a short tandem repeat (STR) genetic locus. The x-axis is DNA fragment length (in base pair), while the y-axis shows relative fluorescent units (rfu). We see two taller alleles (28, 30) that might come from a major contributor, and two shorter alleles (29, 32.2) perhaps from a minor contributor. The question we ask here is "what are the underlying genotypes?" That is, "how can we infer the major and minor contributor genotypes at this locus?"

Most DNA analysts in the US use *qualitative* thresholds (17), shown in Figure 2. The threshold level of 50 rfu shown here is lower than a typical stochastic threshold, which would raise this threshold three times higher to 150 rfu and make the STR data disappear entirely (18). Applying a threshold to quantitative data slices away information. The quantitative data is lost, in this case leaving four all-or-none "allele"

events¹. Forming all possible allele pairs would produce ten candidate pairings of those four "allele" events. This genotype listing diffuses the probability across ten allele pairs (most of which are not feasible), and thus reduces the likelihood ratio (5).

Quantitative mixture interpretation does not use thresholds or "allele" events. Rather, a computer system proposes all possible combinations of explanatory variables (allele pairs, mixture weight, stutter, relative amplification, peak uncertainty, degraded DNA, etc.) in order to generate peak patterns that can be compared with the experimental data (16). In Figure 3, we see the same STR data (green), now with a quantitative superimposed pattern (gray) that fits well.

The computer constructs a quantitative pattern by hypothesizing particular values for explanatory parameters. When this pattern explains the data very well (as seen in the figure), the likelihood function returns a high value. A high likelihood confers a higher probability to the pattern's hypothesized parameters, such as the genotype allele pair values.

When we propose genotype allele pairs or other parameters that do not explain the data well, we observe incorrect patterns that do not resemble the data. These ill-fitting patterns produce a low likelihood, which give very low (or no) probability. Intermediate patterns, that only loosely fit the data, yield likelihood values in between.

This complete search across all parameter values, following the laws of probability, is how computers infer genotypes (and other parameters) (19). The result is a probability distribution over all possible allele pairs, with probability weights that are not equal (20). A quantitative DNA interpretation method considers all possibilities, and describes experimental uncertainty through scientific probability. A valid statistical inference is not permitted to change the observed quantitative data (e.g., threshold operations are not allowed), but rather tries to explain all of it mathematically (21).

TrueAllele computer inference

We will use Cybergene TrueAllele[®] Casework system in our examples. TrueAllele is a quantitative computer interpretation method (16). The system conducts statistical search using a high dimensional probability model with thousands of explanatory variables. The genotype random variable is of primary interest, because its probability distribution is the only evidence component that will enter into a likelihood ratio.

TrueAllele preserves all the identification information in the DNA evidence. The computer objectively infers genotypes without ever seeing a suspect. Only afterwards does it make any comparison with a suspect, many suspects, or an entire country's convicted offender database of possible suspects.

¹ An observed data peak, regardless of threshold, need not actually be a true allele.

TrueAllele can work with any number of mixture contributors (e.g., 2, 3, 4, 5). The system has probability models for PCR stutter (22), allele imbalance, degraded DNA, etc. – all the STR data factors that forensic analysts routinely examine. Most importantly, TrueAllele calculates the uncertainty of every peak. This is critical, because by knowing the uncertainty around each peak, the system then knows to what extent the genotype patterns are accounting for the data.

I gave a live demo seven years ago at a Promega symposium (23). The Macintosh laptop solved a two person mixture, finding the minor contributor genotype in 30 seconds. Much of Cybergenetics research for five years after that presentation centered on how to calculate the uncertainty at every STR peak. Think of data uncertainty as a bell curve around each peak that describes exactly how confident we are in its height. That peak uncertainty modeling is what lets us confidently proceed with reliable genotype inference² (24). These standard data uncertainty methods first appeared in computational statistics about twenty years ago (25).

The TrueAllele system was created over ten years ago. The software evolved into its 25th version two years ago. TrueAllele has been used on over a 100,000 evidence samples. The technology is offered as a product, a service, or as a combined product and service, depending on end-user needs.

Case and validation examples

A 7% minor contributor from a victim's fingernail

The landmark case Commonwealth v. Foley was the first time that a mathematically rigorous computer interpretation of DNA mixtures was admitted into evidence after a pretrial hearing and used in court (26-28). The fingernails of the victim contained a 6.7% minor component of an unknown contributor. Pennsylvania state trooper Kevin Foley, boyfriend of the deceased's estranged wife, was charged with the homicide.

The original inclusion statistic from a national laboratory was a LR of 13,000. An independent expert's obligate allele method gave a LR of 23 million. The TrueAllele computer applied quantitative inference to report a LR of 189 billion. Three different match statistics, all leading to one verdict (29).

Let us state the computer's result using straightforward match LR language. First, we list our assumptions, such as having two contributors to the DNA mixture, including the known victim. Then, in plain language, we say: "A match between Mr. Foley and the fingernails is 189 billion times more probable than a coincidental match to an unrelated Caucasian." That is the likelihood ratio, stated in a match form that everyone can understand.

² What if mixture software didn't work out the data uncertainty? Such simple software would only be guessing about data confidence, and so its inferred genotypes would often be wrong. An invalid inference method invites court challenge.

Multiple amplifications of low-level three contributor DNA

I testified this past summer at Oxford Crown Court in Regina v. Broughton – the Queen of England against an arsonist. The biological evidence was a low template mixture of three DNA contributors, taken from a fuse. Orchid Cellmark had amplified the sample in triplicate. Accounting for the post-PCR enhancement, the pre-enhanced peak heights would have been well under threshold. As we see at locus vWA in Figure 4, each amplification has a highly dissimilar peak pattern because of considerable stochastic PCR variation.

No match score was found by human review, but one was needed for court. When this happens, groups often call on Cybergene to process the mixture data and determine an informative LR value. We applied TrueAllele computation to the data, examining all three amplifications with a joint likelihood function (16). The computer spent considerable time working out the peak uncertainty, modeling the variance distribution by trying out all possibilities.

Figure 5 shows a powerful way of visualizing the LR, using the genotype form. We see the locus vWA genotypes (probability distributions over allele pairs) before and after TrueAllele inference.

- The population distribution of the allele pairs based on the product rule ($2pq$, p^2) is shown in brown. This small amount of probability at each allele pair is what we believe *before* observing the data.
- *After* looking at all the quantitative data, the computer updates its genotype belief, changing its probability distribution (from the population) to whatever the data has indicated, as shown in blue. There was a probability gain at some allele pairs, and a loss at others.

The computer did not know the suspect genotype when it solved the problem, so its inference was entirely objective. It produced an objectively inferred genotype, i.e., allele pair probability distribution of the unknown contributor. We now identify the suspect's allele pair [14, 18] at the vWA locus, slide over a window (shown in red), and look only at this particular allele pair. That is what the likelihood ratio tells us to do – focus solely on the suspect's genotype, since other allele pairs are not relevant to the LR. We see that the posterior genotype probability (blue) is six times higher than the prior population probability (brown). So the LR at vWA is about 6.

I showed this picture in a TrueAllele visual interface to the prosecutor. He was then able to explain the LR by himself (quite successfully, and without my intervention) to his fellow prosecutors and police, using visualizations like this at the other loci. He enjoyed having a solid grasp of the fundamental concepts. He was not an expert in the underlying science or mathematics, but this visual form of the LR was entirely obvious to him and his colleagues.

The picture visualizes the genotype form of the LR. For the match LR statement, we list our assumptions, such as a co-ancestry theta value of 1%, and the presence of three contributors. In understandable words, we can now state the LR as "a match between Mr. Broughton and the fuse is 3 million times more probable than a coincidental match to an unrelated Caucasian." This LR is described in ordinary English, and is mathematically correct.

Large variation in human interpretation

LR methods vary, as Dr. John Butler of the National Institute of Standards and Technology (NIST) has shown (30). His classic slide from five years ago presented LR results from independent human interpretations done by over 50 laboratories on a single two-contributor mixture sample. There was a range from an inclusion LR of 31 thousand (10^4) to a more quantitative LR method of 213 trillion (10^{14}). These match scores represent very different likelihood ratios spanning *ten orders of magnitude* (10^{10} , or $10^{14}/10^4$) produced from the identical DNA data.

Computers preserve DNA identification information

We show in Figure 6 LR comparison results from a recent mixture interpretation study (28). Dr. Margaret Kline of NIST prepared the DNA samples. The data were a series of mixture combinations (90:10, 70:30, 50:50, 30:70 and 10:90) of known genotypes. There were two different pairs of individuals, serially diluted at 1 nanogram, 1/2 ng, 1/4 ng and 1/8 ng for a total of 40 prepared mixtures.

Let us first focus on the blue scatter plot. For each point, the x-axis shows the amount of unknown culprit DNA on a logarithmic scale: 10 picograms, 100 pg up to 1000 pg. We can determine this quantity by multiplying the total DNA amount times the mixture weight. The y-axis, also on a log scale, shows the likelihood ratio: thousand, million, billion, trillion, quadrillion and so on.

The blue scatter plot presents the TrueAllele quantitative interpretation results. For each two-person mixture, the computer assumes the known victim profile and solves for an unknown genotype (probability distribution). As we move leftward from 1000 pg, we see that down to about 100 pg, all the DNA match information is preserved. Then, leftward from 100 pg down to 10 pg, there is a predictable linear decrease in LR match information. At about a million-to-one, the jury "convincing" likelihood ratio level (31), the regression line crosses at 15 pg, a measure of TrueAllele's genotyping sensitivity.

The red scatter plot shows the LRs for an inclusion mixture interpretation method. As expected, below 150 pg inclusion no longer reaches a "convincing" LR of a million-to-one; the DNA identification information has gone. That relative paucity of derived information is why crime labs (using "threshold" interpretation methods) tend to not interpret evidence much below 100 pg. Computer search with a highly explanatory

quantitative probability model does not have this qualitative review limitation. It can therefore reliably achieve ten times greater sensitivity, reaching down to 15 pg of DNA.

Computers are a million times more informative

A TrueAllele mixture validation study, done collaboratively with Dr. Barry Duceman of the New York State Police, will soon be published in the Journal of Forensic Sciences (16). The paper shows how probabilistic computer interpretation that uses all of the quantitative data preserves likelihood ratio information. Comparison was made with mixture review methods that use "thresholds", finding that such qualitative approaches often discard considerable identification information.

In Figure 7, the x-axis lists eight adjudicated mixture items in cases without a known victim genotype. The y-axis gives the likelihood ratio on a logarithmic scale, 10^5 , 10^{10} , 10^{15} , etc. The TrueAllele computer LR values are shown (blue) for each case. Also shown are the human review inclusion LR scores (orange). These scores were the LR values reported in the case folder from the calculated CODIS combined probability of inclusion (CPI) match statistics³. Comparison was made using the same population databases, without co-ancestry theta correction.

We see that, on average, the computer (blue) LRs of about 10 trillion (10^{13}) preserve identification information. But human inclusion review of the same case mixture data (orange) averages only 10 million (10^7). Relative to quantitative TrueAllele interpretation, using thresholds typically discards a factor of a million (10^6 , or $10^{13}/10^7$) of DNA identification information.

All match statistics are likelihood ratios and can be explained within the same scientific framework. Therefore, the relative efficacy of mixture interpretation methods can be compared in studies (such as these) using the LR logarithm as a universal information measure.

Human review discards most DNA information

More dramatically, Dr. Duceman and I then looked at what happened when we didn't assume that human review produced any match score (32, 33). We simply examined all the results (both computer and human) for the 86 mixture items. Figure 8 lists these items on the x-axis, ordered by decreasing match information. The y-axis again shows DNA match information, as measured by log likelihood ratio.

The TrueAllele computer inferred match information for each item is shown as the large (blue) background. In the different foreground colors, RMP (gray), CLR (green) and CPI (orange), we see the analyst's LR result for each case. Threshold-

³ The LR comparisons for another eight case items, each having a known victim genotype, comparing TrueAllele with the combined likelihood ratio (CLR) method, are not shown here. They are reported in the JFS paper.

based review lost about two thirds of the information when there was an answer. But, most importantly, fewer than 30% of the items were even assigned a match score.

Thus, over 70% of qualitatively reviewed mixture items did not have a match score. This fact has productivity implications. Suppose that a lab needed to find at least one LR match statistic in a mixtures-only case in order to introduce evidence at court. The lab would have to keep processing items until they were lucky enough to find one yielding a match score, sort of like gambling with DNA dice.

With a 70% LR failure rate, probability tells us that a lab must process 3.33 items for every one that gets an LR value. However, in this study, the computer successfully determined an LR for every mixture item. So sequential item testing with "threshold" mixture interpretation would impose a 233% burden of unnecessary work (effort, cost, time, etc.). Moreover, the resulting match statistic would be far less informative than the computer's LR on the same data. Indeed, labs that limit evidence submission to three items would produce no match score at all in a third of the cases.

LR-based databases for DNA investigation

The likelihood ratio also applies to investigative DNA databases. Any highly informative DNA investigative database should use quantitative LR matching.

The allele database approach discards information (34). Its information handling is just like CPI. In fact, it simply lists "included" alleles to give a very weak representation of the genotype probability distribution. Some identification exists there, but it's not very informative. Like inclusion, an allele-based DNA database makes very poor use of the data. The CODIS-like approach can only store and match information-poor mixture genotypes.

In contrast, a probabilistic genotype database preserves more of the DNA evidence information. Such databases can store and match those genotypes (as probability distributions). When a new convicted offender or evidence genotype is uploaded to the database, a likelihood ratio is then incrementally computed (27). This information-rich DNA database approach exploits the sensitivity and specificity that likelihood ratios are known for throughout science.

The TrueAllele system provides this LR-based evidence versus suspect (e.g., convicted offender) genotype match. When Cybergenetics reanalyzed the World Trade Center data (35), we showed how the TrueAllele LR database could be used for disaster victim identification.

The same probabilistic genotyping and LR methods can be used with kinship data to find missing people. In the same way, TrueAllele completely automates familial search, with no human involvement or additional costs, finding DNA matches in the

background. Cybergenetics can customize LR-based DNA database matching for any state or country, in accordance with the laws and regulations of their jurisdiction.

SWGDAM mixture guidelines

The 2010 SWGDAM DNA mixture interpretation guidelines (18) provide for reliable scientific computing in paragraph 3.2.2. The paragraph essentially says that a stochastic threshold is not necessary when using a validated probabilistic genotype method. That provision can help labs make better use of their DNA mixture evidence data. In jurisdictions that require a match statistic, "threshold" methods often end up rendering more than half of the mixture items unusable.

The new SWGDAM guidelines also describe an alternative "stochastic threshold" approach. But stochastic thresholds raise the qualitative peak cutoff, as many labs have recently observed. A higher threshold discards more peak data, and so fewer evidence items can be reported with a match statistic. Removing more peak data leads to a higher false negative error rate (24); with mixtures, the error rate can exceed 100% (falsely excluded alleles per locus). With low-level mixtures, such as property crimes, the information yield can be greatly diminished.

Fortunately, via paragraph 3.2.2, the SWGDAM mixture interpretation guidelines let labs use probability modeling to preserve DNA identification information. Moreover, we can measure the LR efficacy improvement, because all match statistics are likelihood ratios and can therefore be compared numerically.

The LR and forensic science

What is the point of forensic science? Why do taxpayers fund it in the first place? Certainly the public and the police believe that labs work hard to preserve all the identification information that is present in the DNA evidence.

Forensic scientists want to provide accurate results. If the true DNA match number is a trillion to one, we want to report a trillion to one, and not a million to one. If the true number is a million to one, we don't want to say the data are inconclusive. The point is to serve the criminal justice system for law enforcement and the courts in an objective way to help protect the public from crime.

Forensic DNA science is about DNA labs and analysts continually progressing to employ the most informative methods available, so that they can bestow the benefits of science on the public. Science professionals care deeply about the safety of society, and want to do the most accurate job possible.

The likelihood ratio is an essential tool for preserving and accurately presenting DNA match evidence. American forensic scientists need to communicate with their

judges and juries in the English language. Testifying experts can be wary of conditional probability (risking transposed conditionals) or arcane sentences that extend on for many words. A simple solution is to state the LR by saying "A match between the suspect and the evidence is a billion times more probable than a coincidental match."

Conclusion

We can rearrange the likelihood ratio into a variety of different mathematical forms (Appendix). Many analysts would prefer to not talk about the "probability of data" ratio that might work well for computers and statisticians, but can seem so opaque to a non-specialist. Fortunately, the likelihood ratio can be made easy to understand and explain in court. I have done this, and have taught scientists and lawyers how to explain it to others. The approach described here is to state the LR in an appropriate form.

The match form of the LR is especially understandable, and is therefore useful in reports and testimony. This LR form uses the usual denominator for coincidental match, just like the RMP statistic. In the numerator, we find the strength of match between the evidence and suspect. This match strength decreases with weaker interpretation methods (e.g., inclusion), but stays high and is preserved with more informative likelihood methods that better explain the data (e.g., TrueAllele). When we take the ratio of these two probabilities to form a LR, the DNA identification information is preserved. That preservation of evidence is the primary purpose of forensic science.

Resources

Handouts for this (and other cited) presentations can be downloaded from our website on the Presentations page (<http://www.cybgen.com/information/presentations.shtml>). For all our recent scientific presentations and posters, Cybergenetics creates a web page that provides the abstract, handout, transcript and narrated movie of the slides. This dissemination enables a review after the talk, access for those who were unable to attend the meeting, and continuing education possibilities.

If you are interested in reading our scientific papers, you can download the manuscripts from our website (<http://www.cybgen.com/information/publications.shtml>). We also provide course lectures and supporting materials for scientists and lawyers who are interested in learning more about quantitative mixture interpretation and the likelihood ratio (<http://www.cybgen.com/information/courses.shtml>).

You may want to see quantitative TrueAllele interpretation and its LR reporting for challenging DNA cases that people cannot solve. If so, you can send Cybergenetics some interesting case data (at no charge) for TrueAllele processing and a follow up customized webinar. If you have further questions about DNA identification science, please contact me by email (perlin@cybgen.com).

Appendix

Genotype probability notation

Let us write down the genotypes and their probabilities that we shall use here. Each probability mass function (pmf) $q(x)$, $r(x)$ and $s(x)$ describes the uncertainty of its respective genotype evidence Q , population R and suspect S . The evidence genotype likelihood function $\lambda_Q(x)$ is not a pmf, but is used in the construction of pmf $q(x)$.

genotype	notation	source	role	probability	data
Q	$\lambda_Q(x)$	evidence	likelihood	$\Pr\{d_Q Q=x\}$	d_Q
Q	$q(x)$	evidence	posterior	$\Pr\{Q=x d_Q\}$	d_Q
R	$r(x)$	population	prior	$\Pr\{R=x d_R\}$ or $\Pr\{Q=x\}$	d_R
S	$s(x)$	suspect	comparison	$\Pr\{S=x d_S\}$	d_S

The domain of each function is the set G of possible genotype values at a locus, so that each $x \in G$ is an allele pair. The range of each function is the set of nonnegative real numbers.

Equivalent LR formulations

The *hypothesis form* of the likelihood ratio (LR) expresses the information gained in the hypothesis H odds by having observed data (1)

$$[1] \quad LR = \frac{O(H|d_Q, d_R, d_S)}{O(H)}$$

Here, hypothesis H is that the suspect contributed to the DNA evidence, and the DNA data comprises the questioned evidence d_Q , the reference population allele frequencies d_R and suspect profile d_S .

For the *likelihood form*, standard Bayesian rearrangements (36) tell us that the LR can also be written as the ratio of conditional probabilities (Proof A)

$$[2] \quad LR = \frac{\Pr\{d_Q|H, d_R, d_S\}}{\Pr\{d_Q|\bar{H}, d_R, d_S\}}$$

where \bar{H} is the alternative hypothesis that someone else contributed to the evidence.

Suppose that there is uncertainty in the evidence genotype Q having pmf $q(x)$ or in suspect genotype S with pmf $s(x)$. Then this genotype uncertainty can be expressed in the LR as

$$[3] \quad LR = \frac{\sum_{x \in G} \lambda_Q(x) \cdot s(x)}{\sum_{x \in G} \lambda_Q(x) \cdot r(x)}$$

where $\lambda_Q(x)$ is the likelihood function of the evidence genotype Q and $r(x)$ is the pmf of reference population genotype R (Proof B). Although this LR shares many useful features of the match LR approximation (27), this exact LR equation uses likelihood function $\lambda_Q(x)$ instead of posterior probability $q(x)$. When genotype Q is inferred using a population prior R , likelihood λ_Q and posterior q are easily inter-converted by renormalizing with prior r , since $q(x) \propto \lambda_Q(x) \cdot r(x)$.

A general genotype form of the LR is

$$[4] \quad LR = \sum_{x \in G} \frac{q(x) \cdot s(x)}{r(x)}$$

This genotype formulation describes the LR solely in terms of the posterior pmfs of genotypes Q , S , and R (27). This genotype form is exact when the reference population genotype R is used as the prior for inferring genotype Q and coancestry is not considered (Proof C).

The suspect genotype S is often definite, with all its probability mass placed at a single allele pair x_s . This happens when genotype S is from a reference sample or suspect database.

In this case, the *genotype form* of the LR then reduces to the single term

$$[5] \quad LR = \frac{q(x_s)}{r(x_s)} = \frac{\Pr\{Q = x_s | d_Q\}}{\Pr\{Q = x_s\}}$$

at the suspect genotype value x_s (Proof D). Note that $r(x_s)$ is the prior probability of the genotype at value x_s , while $q(x_s)$ is the posterior probability of x_s after having examined the DNA data. Thus, in this "definite suspect" situation, the LR becomes a genotype probability ratio that expresses the gain in identification information provided by the data (37).

The *match form* of the LR arises with a definite suspect genotype, written as

$$[6] \quad LR = \frac{q(x_s) \cdot s(x_s)}{r(x_s)}$$

at suspect allele pair x_s (Proof E). The numerator gives the probability of a match between the evidence and suspect genotypes, while the denominator is the probability of a coincidental match.

LR equivalence proofs

Proof A. We expand the LR odds definition [1] in the standard way into probability ratios

$$\begin{aligned} LR &= \frac{O(H|d_Q, d_R, d_S)}{O(H)} \\ &= \frac{\Pr\{H|d_Q, d_R, d_S\} / \Pr\{\bar{H}|d_Q, d_R, d_S\}}{\Pr\{H\} / \Pr\{\bar{H}\}} \end{aligned}$$

Rearranging denominators we have

$$= \frac{\Pr\{H|d_Q, d_R, d_S\} / \Pr\{H\}}{\Pr\{\bar{H}|d_Q, d_R, d_S\} / \Pr\{\bar{H}\}}$$

By Bayes theorem, the posterior probability of H can be interchanged with its likelihood, renormalizing appropriately. Doing this separately for numerator and denominator, we obtain

$$= \frac{\Pr\{d_Q|H, d_R, d_S\} / \Pr\{d_Q\}}{\Pr\{d_Q|\bar{H}, d_R, d_S\} / \Pr\{d_Q\}}$$

Canceling out the total probability factors $\Pr\{d_Q\}$ yields the desired equation [2].

Proof B. We start from equation [2] with the conditional probability form of the LR

$$LR = \frac{\Pr\{d_Q|H, d_R, d_S\}}{\Pr\{d_Q|\bar{H}, d_R, d_S\}}$$

Using the law of total probability (or, "extending the conversation"), we consider every possible allele pair $x \in G$ for genotype Q .

$$\begin{aligned} &= \frac{\sum_{x \in G} \Pr\{d_Q|H, d_R, d_S, Q=x\} \cdot \Pr\{Q=x|H, d_R, d_S\}}{\sum_{x \in G} \Pr\{d_Q|\bar{H}, d_R, d_S, Q=x\} \cdot \Pr\{Q=x|\bar{H}, d_R, d_S\}} \end{aligned}$$

The likelihood's probability of the data $\Pr\{d_Q|\dots\}$ is unaffected by hypothesis H or \bar{H} . In the numerator, the evidence genotype Q under hypothesis H that the suspect contributed to the evidence, $\Pr\{Q = x|H, \dots\}$ becomes the suspect's genotype S . Similarly, in the denominator the genotype Q under hypothesis \bar{H} that someone else contributed $\Pr\{Q = x|\bar{H}, \dots\}$ becomes the population genotype R . We therefore derive the ratio

$$\begin{aligned} & \sum_{x \in G} \Pr\{d_Q|d_R, d_S, Q = x\} \cdot \Pr\{S = x|d_R, d_S\} \\ &= \frac{\sum_{x \in G} \Pr\{d_Q|d_R, d_S, Q = x\} \cdot \Pr\{S = x|d_R, d_S\}}{\sum_{x \in G} \Pr\{d_Q|d_R, d_S, Q = x\} \cdot \Pr\{R = x|d_R, d_S\}} \end{aligned}$$

Eliminating noninfluential conditioning variables, we then have that

$$\begin{aligned} & \sum_{x \in G} \Pr\{d_Q|Q = x\} \cdot \Pr\{S = x|d_S\} \\ &= \frac{\sum_{x \in G} \Pr\{d_Q|Q = x\} \cdot \Pr\{R = x|d_R\}}{\sum_{x \in G} \Pr\{d_Q|Q = x\} \cdot \Pr\{R = x|d_R\}} \end{aligned}$$

Substituting in our notation for the evidence likelihood $\lambda_Q(x)$, and posterior pmfs for the suspect $s(x)$ and population $r(x)$ genotypes, we obtain the desired equation [3] result

$$LR = \frac{\sum_{x \in G} \lambda_Q(x) \cdot s(x)}{\sum_{x \in G} \lambda_Q(x) \cdot r(x)}$$

Proof C. The evidence genotype prior probability is the same as a population genotype pmf $r(x)$, since that is our belief before we see any evidence data. The denominator of equation [3]

$$\sum_{x \in G} \lambda_Q(x) \cdot r(x)$$

is a Bayes normalization constant. Bayes theorem then tells us that

$$\frac{\lambda_Q(x)}{\sum_{y \in G} \lambda_Q(y) \cdot r(y)} = \frac{q(x)}{r(x)}$$

We can therefore algebraically rewrite equation [3] by folding the denominator into the numerator's summation as:

$$\begin{aligned}
 LR &= \frac{\sum_{x \in G} \lambda_Q(x) \cdot s(x)}{\sum_{x \in G} \lambda_Q(x) \cdot r(x)} \\
 &= \sum_{x \in G} \left[\frac{\lambda_Q(x)}{\sum_{y \in G} \lambda_Q(y) \cdot r(y)} \right] \cdot s(x) \\
 &= \sum_{x \in G} \left[\frac{q(x)}{r(x)} \right] \cdot s(x)
 \end{aligned}$$

This rearrangement proves equation [4], that

$$LR = \sum_{x \in G} \frac{q(x) \cdot s(x)}{r(x)}$$

Proof D. When S is a definite genotype, the pmf $s(x)$ places all of its probability mass at one allele pair x_s , so that

$$s(x) = \begin{cases} 1 & x = x_s \\ 0 & \text{otherwise} \end{cases}$$

Therefore, the summation of equation [4] reduces to a single nonzero term at x_s

$$\begin{aligned}
 LR &= \sum_{x \in G} \frac{q(x) \cdot s(x)}{r(x)} \\
 &= \sum_{x \in G} \frac{q(x)}{r(x)} \cdot s(x) \\
 &= \frac{q(x_s)}{r(x_s)} \cdot 1 + \sum_{\substack{x \in G \\ x \neq x_s}} \frac{q(x)}{r(x)} \cdot 0 \\
 &= \frac{q(x_s)}{r(x_s)} + 0
 \end{aligned}$$

So, with definite genotype S having a unique allele pair x_s , we have equation [5]

$$LR = \frac{q(x_s)}{r(x_s)}$$

Note: Considering our original probability definitions, we see that in this case the LR is simply the posterior-to-prior genotype probability ratio

$$LR = \frac{\Pr\{Q = x_s | d_Q\}}{\Pr\{Q = x_s\}}$$

The other allele pair possibilities are not relevant to the LR, just the genotype probability gain at the matching allele pair x_s . This result can be also derived directly from Bayes

theorem via $LR = \frac{\Pr\{d_Q | Q = x_s\}}{\Pr\{d_Q\}} = \frac{\Pr\{Q = x_s | d_Q\}}{\Pr\{Q = x_s\}} = \frac{q(x_s)}{r(x_s)}$.

Proof E. The general form of LR equation [4] is a summation of $q(x) \cdot s(x) / r(x)$ genotype pmf values. With a definite suspect genotype S , $s(x_s) = 1$, and so (after eliminating the zero terms) we arrive at equation [6]

$$LR = \frac{q(x_s) \cdot s(x_s)}{r(x_s)}$$

The numerator is the probability of a match between the (independent) evidence and suspect genotypes. The denominator is the probability of a coincidental match that finds the suspect in the population genotype. So, with a definite suspect allele pair x_s , we obtain the match form of the LR.

References

1. Good IJ. Probability and the Weighing of Evidence. London: Griffin, 1950.
2. Balding DJ, Nichols RA. DNA profile match probability calculation: how to allow for population stratification, relatedness, database selection and single bands. *Forensic Sci Int.* 1994;64(2-3):125-40.
3. Evett IW, Buffery C, Willott G, Stoney D. A guide to interpreting single locus profiles of DNA mixtures in forensic cases. *J Forensic Sci Soc.* 1991 Jan-Mar;31(1):41-7.
4. Harbison S, Buckleton J. Applications and extensions of subpopulation theory: a caseworkers guide. *Science & Justice.* 1998;38:249-54.
5. Perlin MW. Inclusion probability is a likelihood ratio: implications for DNA mixtures (poster #85). *Promega's Twenty First International Symposium on Human Identification, 2010; San Antonio, TX.* 2010.
6. Bayes T, Price R. An essay towards solving a problem in the doctrine of chances. *Phil Trans.* 1763;53:370-418.
7. Laplace PS. *Theorie analytique des probabilites.* Paris: Ve. Courcier, 1812.
8. Good IJ. Studies in the history of probability and statistics. XXXVII. A.M. Turing's statistical work in World War II. *Biometrika.* 1979;66(2):393-6.
9. Lindley DV. A problem in forensic science. *Biometrika.* 1977;64(2):207-13.
10. Lindley DV. *Understanding Uncertainty.* Hoboken, NJ: John Wiley & Sons, 2006.
11. Buckleton JS, Triggs CM, Walsh SJ, editors. *Forensic DNA Evidence Interpretation.* Boca Raton, FL: CRC Press, 2004.
12. Evett IW, Weir BS. *Interpreting DNA Evidence: Statistical Genetics for Forensic Scientists.* Sunderland, MA: Sinauer Assoc, 1998.
13. Curran JM, Triggs C, Buckleton J, Weir BS. Interpreting DNA mixtures in structured populations. *J Forensic Sci.* 1999;44(5):987-95.
14. Gill P, Brenner CH, Buckleton JS, Carracedo A, Krawczak M, Mayr WR, Morling N, Prinz M, Schneider PM, Weir BS. DNA commission of the International Society of Forensic Genetics: Recommendations on the interpretation of mixtures. *Forensic Sci Int.* 2006;160:90-101.
15. Edwards AWF. *Likelihood.* Expanded ed. Baltimore: Johns Hopkins University, 1992.
16. Perlin MW, Legler MM, Spencer CE, Smith JL, Allan WP, Belrose JL, Duceman BW. Validating TrueAllele® DNA mixture interpretation. *Journal of Forensic Sciences.* 2011;56(November):in press.
17. SWGDAM. Short Tandem Repeat (STR) interpretation guidelines (Scientific Working Group on DNA Analysis Methods). *Forensic Sci Commun (FBI).* 2000 July;2(3).
18. SWGDAM. Interpretation guidelines for autosomal STR typing by forensic DNA testing laboratories. 2010.
19. Gilks WR, Richardson S, Spiegelhalter DJ. *Markov Chain Monte Carlo in Practice:* Chapman and Hall, 1996.
20. Gelman A, Carlin JB, Stern HS, Rubin D. *Bayesian Data Analysis.* Boca Raton, FL: Chapman & Hall/CRC, 1995.

21. O'Hagan A, Forster J. Bayesian Inference. Second ed. New York: John Wiley & Sons, 2004.
22. Perlin MW, Szabady B. Linear mixture analysis: a mathematical approach to resolving mixed DNA samples. *J Forensic Sci.* 2001;46(6):1372-7.
23. Perlin MW. Simple reporting of complex DNA evidence: automated computer interpretation. *Promega's Fourteenth International Symposium on Human Identification*, 2003; Phoenix, AZ. 2003.
24. Perlin MW. Reliable interpretation of stochastic DNA evidence. *Canadian Society of Forensic Sciences 57th Annual Meeting*, 2010; Toronto, ON. 2010.
25. Gelfand AE, Smith AFM. Sampling-based approaches to calculating marginal densities. *J Amer Statist Assoc.* 1990;85:398-409.
26. Martin W. *Commonwealth of Pennsylvania v. Kevin Foley*. Indiana County; 2009.
27. Perlin MW, Kadane JB, Cotton RW. Match likelihood ratio for uncertain genotypes. *Law, Probability and Risk.* 2009;8(3):289-302.
28. Perlin MW, Sineelnikov A. An information gap in DNA evidence interpretation. *PLoS ONE.* 2009;4(12):e8327.
29. Perlin MW, Cotton RW. Three match statistics, one verdict (A78). *AAFS 62nd Annual Scientific Meeting*, 2010 February 22-27; Seattle, WA. *American Academy of Forensic Sciences*; 2010. p. 63.
30. Butler JM, Kline MC. NIST Mixture Interpretation Interlaboratory Study 2005 (MIX05), Poster #56. *Promega's Sixteenth International Symposium on Human Identification*, 2005; Grapevine, TX. 2005.
31. Koehler JJ. When are people persuaded by DNA match statistics? *Law and Human Behavior.* 2001;25(5):493-513.
32. Belrose JL, Duceman BW. New York State Police validation of a statistical tool for genotype inference and match that solves casework mixture problems (A79). *AAFS 62nd Annual Scientific Meeting*, 2010; Seattle, WA. *American Academy of Forensic Sciences*; 2010. p. 64.
33. Perlin MW, Duceman BW. Profiles in productivity: Greater yield at lower cost with computer DNA interpretation (Abstract). *Twentieth International Symposium on the Forensic Sciences of the Australian and New Zealand Forensic Science Society*, 2010 September; Sydney, Australia. 2010.
34. Niezgodna SJ, Brown B. The FBI Laboratory's COmbined DNA Index System Program. *Sixth International Symposium on Human Identification*, 1995; Scottsdale, AZ. 1995.
35. Perlin MW. Identifying human remains using TrueAllele® technology. In: Okoye MI, Wecht CH, editors. *Forensic Investigation and Management of Mass Disasters*. Tucson, AZ: Lawyers & Judges Publishing Co; 2007;31-8.
36. Aitken CG, Taroni F. *Statistics and the Evaluation of Evidence for Forensic Scientists*. Second ed. Chicester, UK: John Wiley & Sons, 2004.
37. Essen-Möller E. Die Biesweiskraft der Ähnlichkeit im Vater Schafsnachweis; Theoretische Grundlagen. *Mitteilungen der anthropologischen Gesellschaft in Wien.* 1938;68(9-53).

Figure Legends

Figure 1. Quantitative data at an STR locus showing a four peak, two contributor mixture pattern.

Figure 2. Applying a threshold to quantitative data, reducing it to all-or-none qualitative "allele" events.

Figure 3. Using a quantitative model (gray) to fully explain quantitative data (green).

Figure 4. Three amplifications of the same DNA template at the vWA locus, shown following post-PCR enhancement.

Figure 5. Genotype probability distributions at the vWA locus from before (brown) and after (blue) computer examination of the data. The suspect's genotype is [14 18].

Figure 6. A log-log scatter plot of culprit DNA quantity vs. likelihood ratio for a set of DNA mixture samples. Quantitative (blue) and qualitative (red) mixture interpretation results are shown.

Figure 7. Comparison of quantitative (blue) and qualitative (orange) interpretation of two person mixtures, without a victim reference, conducted on the same data. The y-axis $\log(\text{LR})$ shows identification information.

Figure 8. Comparison of quantitative (blue) and qualitative (gray, green, orange) of 86 case mixture items, conducted on the same data. The y-axis $\log(\text{LR})$ shows identification information.

Figures

Figure 1.

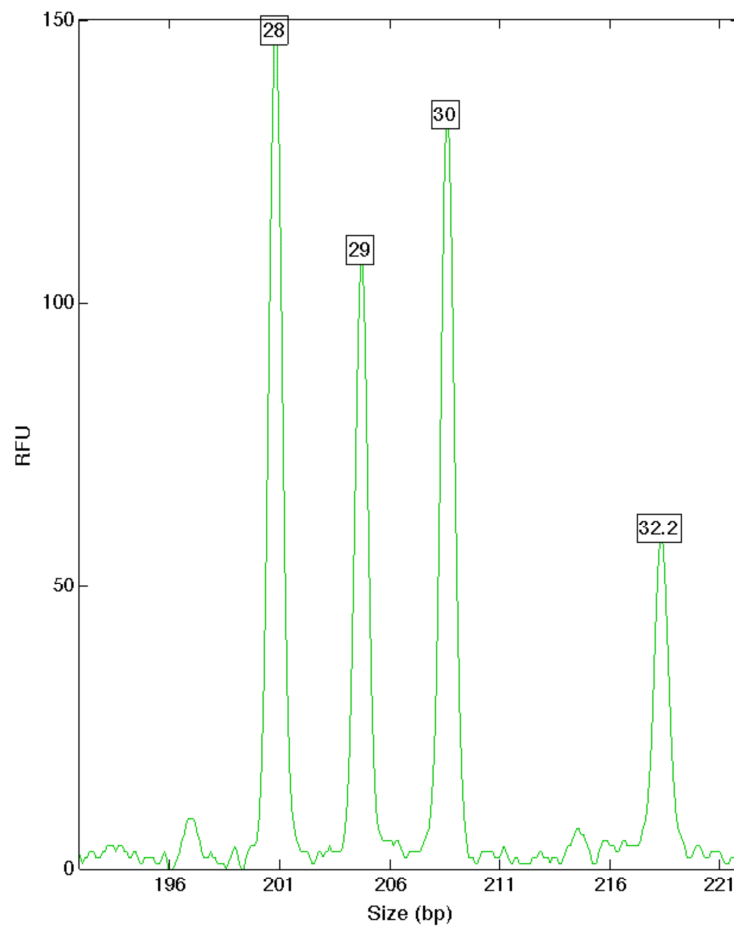


Figure 2.

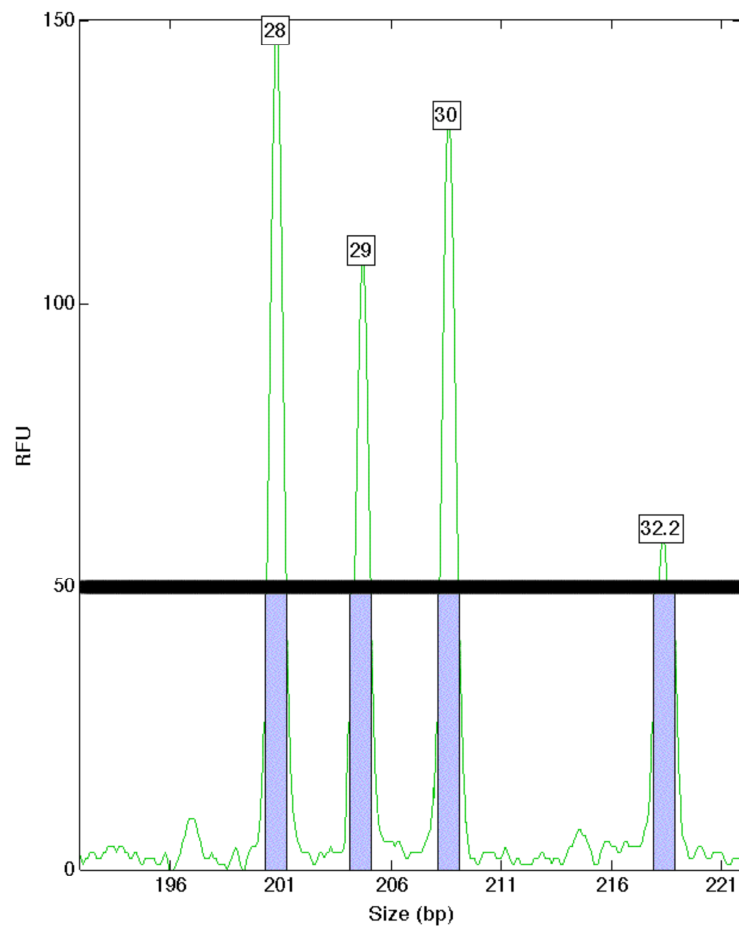
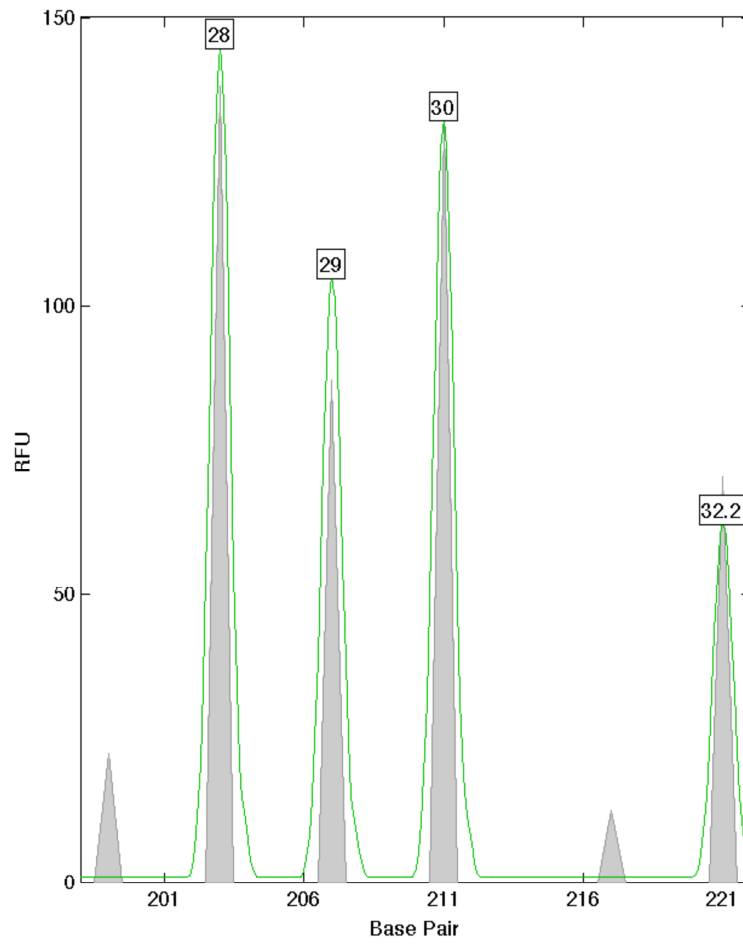


Figure 3.



Explaining the Likelihood Ratio in DNA Mixture Interpretation

Figure 4.

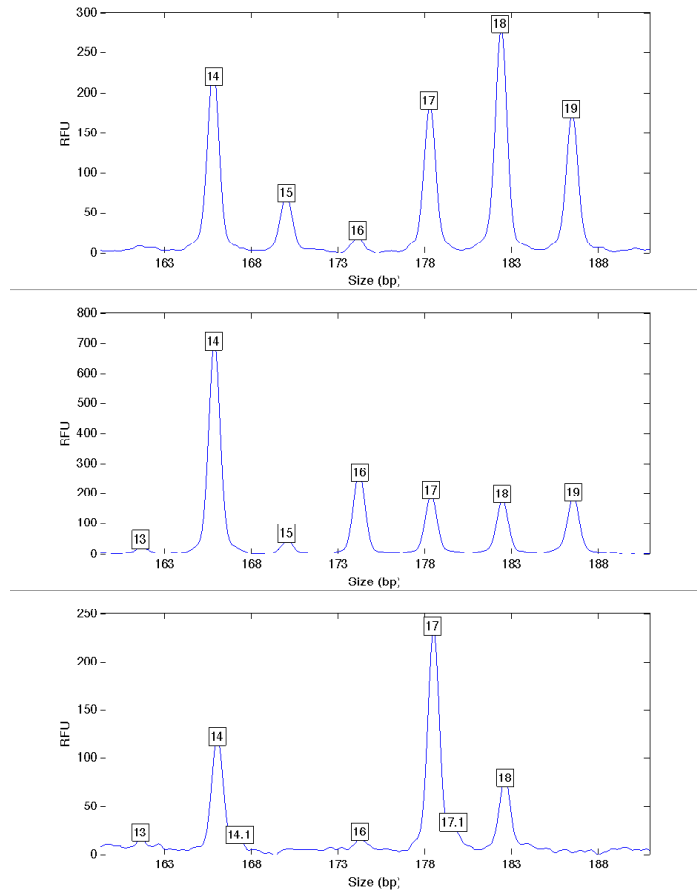


Figure 5.

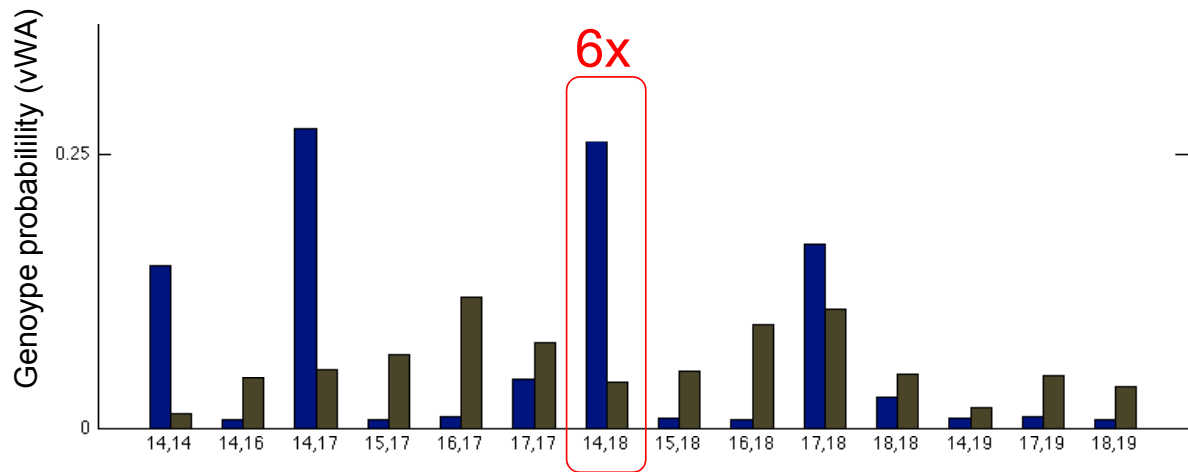


Figure 6.

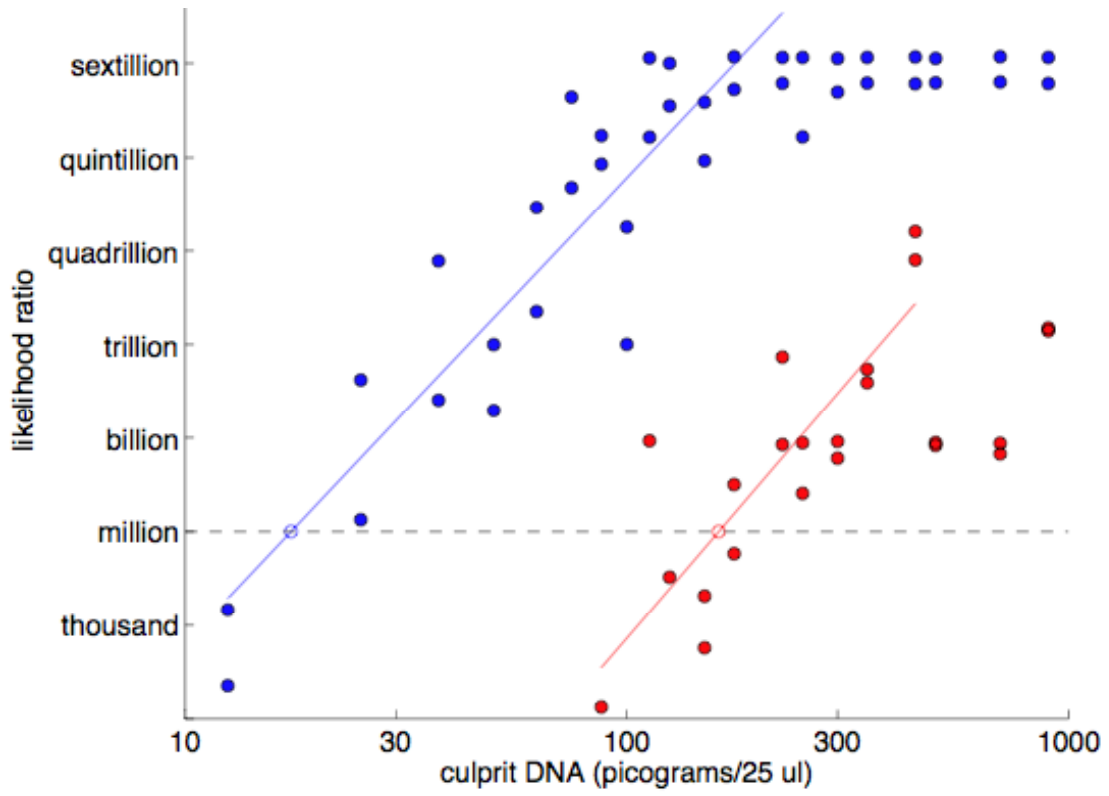


Figure 7.

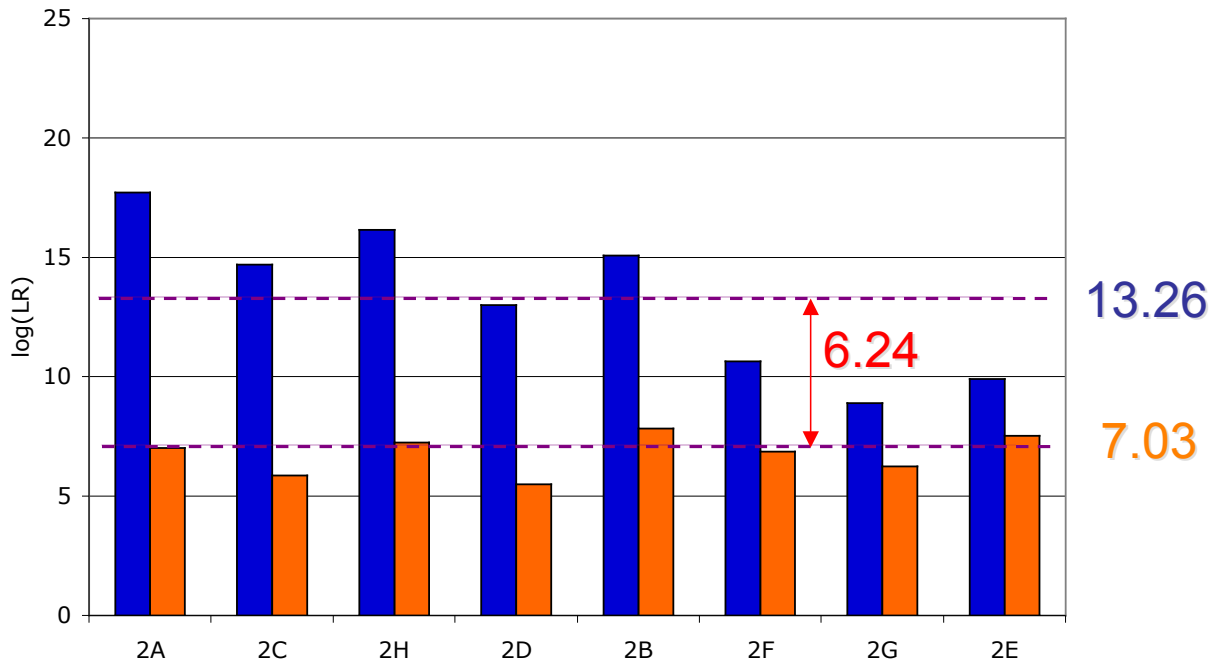


Figure 8.

