

Evaluation of mixed-source, low-template DNA profiles in forensic science

David J. Balding¹

University College London Genetics Institute, University College London, London WC1E 6BT, United Kingdom

Edited by Terence P. Speed, University of California, Berkeley, CA, and accepted by the Editorial Board May 31, 2013 (received for review November 13, 2012)

Enhancements in sensitivity now allow DNA profiles to be obtained from only tens of picograms of DNA, corresponding to a few cells, even for samples subject to degradation from environmental exposure. However, low-template DNA (LTDNA) profiles are subject to stochastic effects, such as “dropout” and “dropin” of alleles, and highly variable stutter peak heights. Although the sensitivity of the newly developed methods is highly appealing to crime investigators, courts are concerned about the reliability of the underlying science. High-profile cases relying on LTDNA evidence have collapsed amid controversy, including the case of Hoey in the United Kingdom and the case of Knox and Sollecito in Italy. I argue that rather than the reliability of the science, courts and commentators should focus on the validity of the statistical methods of evaluation of the evidence. Even noisy DNA evidence can be more powerful than many traditional types of evidence, and it can be helpful to a court as long as its strength is not overstated. There have been serious shortcomings in statistical methods for the evaluation of LTDNA profile evidence, however. Here, I propose a method that allows for multiple replicates with different rates of dropout, sporadic dropins, different amounts of DNA from different contributors, relatedness of suspected and alternate contributors, “uncertain” allele designations, and degradation. R code implementing the method is open source, facilitating wide scrutiny. I illustrate its good performance using real cases and simulated crime scene profiles.

forensic genetics | forensic identification | statistical genetics | weight of evidence

Reliability of Low-Template DNA Profiling

Problems with the courtroom use of low-template DNA (LTDNA) profiles were brought into sharp focus in the United Kingdom in 2007, with the collapse of a trial arising from the Omagh bombing in Northern Ireland in 1998. This crime killed 29 people and injured many more; consequently, early termination of the trial and acquittal of the defendant attracted widespread adverse publicity. The judge gave several reasons, but it was his critical appraisal of the LTDNA evidence that captured headlines. In response to the controversy, a report reviewing LTDNA evidence (1) was commissioned by the UK Forensic Science Regulator. The report found the underlying science to be “sound” and LTDNA profiling to be “fit for purpose,” although admitting that there was lack of agreement “on how LTDNA profiles are to be interpreted.”

I suggest that these comments are somewhat contradictory: Without valid methods of assessing evidential strength, a technique cannot be fit for purpose in the criminal justice system. Fig. 1 shows part of the electropherogram (epg) giving the results from replicate LTDNA profiling runs in a crime investigation. The two epgs show substantial similarity yet also important differences: For example, the 17 allele at locus D19 is detected in Fig. 1 (*Left*) but not in Fig. 1 (*Right*), yet the reverse is true for the 11 allele. Is a technology that produces such variable results reliable? This is often asked by legal commentators, but the term “reliable” is too vague for the question to be useful. What is evident is that there is substantial, but imperfect, information in

these results about the genotypes of individuals contributing DNA to the sample. The important question is whether or not we can extract that information with enough statistical efficiency for it to be useful while avoiding overstatement of evidential strength. Fortunately, progress has been made on this front since publication of the report by Caddy et al. (1), and I propose here a methodology for robust and efficient analyses of LTDNA evidence that is incorporated in a freely available suite of R functions.

Case of Knox and Sollecito

Table 1 shows three interpretations of the DNA evidence at five loci from exhibit 165B of the trial in Perugia, Italy, in 2009. The exhibit includes the clasp of a bra, attached to some apparently blood-stained fabric, that was found near the murdered woman, Meredith Kercher. The report (2), written by two academic experts from the Sapienza Università di Roma, was highly critical of the prosecution’s DNA evidence at trial and led to the convictions of Amanda Knox and Raffaele Sollecito being overturned on appeal. Here, I will use “interpretation” for the process of deciding which epg peaks are allelic and “evaluation” for the calculation of numerical measures of evidential weight for an interpretation.

The interpretation by the Italian Scientific Police presented at trial identified exactly the alleles of the victim and one of the coaccused, Sollecito, in the DNA profiling results. Using methods described below, I computed a weight of evidence (WoE) in favor of the contributors of DNA being Kercher and Sollecito, rather than Kercher and an unknown man, of >15 bans. The ban is the unit of WoE introduced by Alan Turing (3): x bans means $\log_{10}(\text{LR}) = x$, where LR is the likelihood ratio, such that 6 bans means an LR of 1 million. In reviewing the evidence, Vecchiotti and Conti (2) agreed with the alleles originally identified but also reported many additional epg peaks. They cited recommendation 6 of Gill et al. (4) in concluding that all peaks in stutter positions should be regarded as allelic. Of the 24 additional peaks identified by Vecchiotti and Conti (2), of which 6 had heights below the threshold of 50 relative fluorescence units, 9 are included in the profile of the other codefendant, Knox, providing apparent support for the presence of DNA from her. However, four of her alleles were not observed, including two homozygotes, which are less prone to dropout.

These interpretations pose problems for standard methods of evidence evaluation because of the alleles not attributable to any of the profiled individuals, uncertainty over whether or not Knox is a contributor, and the need to allow for the possibility that subthreshold peaks may be allelic. The number of above-threshold alleles recorded at any locus is six or less, which implies three

Author contributions: D.J.B. designed research, performed research, analyzed data, and wrote the paper.

The author declares no conflict of interest.

Freely available online through the PNAS open access option.

This article is a PNAS Direct Submission. T.P.S. is a guest editor invited by the Editorial Board.

¹E-mail: d.balding@ucl.ac.uk.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1219739110/-DCSupplemental.

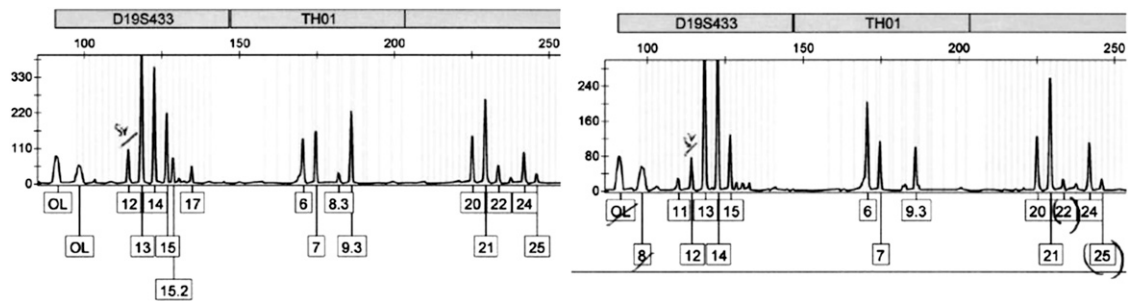


Fig. 1. Illustrative epgs from a swab of a handgun magazine. Two replicates are shown at three loci: D19, TH01, and FGA. Note the different y-axis scales, chosen automatically, in units of relative fluorescence units; the x axis shows fragment length in base pairs. Allele labels in boxes are assigned automatically but can be overridden by a forensic expert taking into account factors like peak morphology and potential stutter. Some manual annotations are shown, indicating subthreshold peaks in () as well as possible artifacts, such as stutter.

or more contributors of DNA. However, if Knox is assumed to be a contributor, the alleles not attributable to her still imply three or more other contributors. I first compare these prosecution (H_p) and defence (H_d) hypotheses for the contributors of DNA:

H_p : Kercher, Knox, Sollecito, and one unknown individual

H_d : Kercher, Knox, and two unknown individuals

I introduce an innovation to likelihood-based analyses to allow for an “uncertain” allele designation. In previous formulations (5–10), the likelihood at a locus in a profiling run is the product over all allelic positions in the epg of one of four possible terms, according to whether or not the corresponding allele is represented in the crime scene profile (CSP) and whether or not it is included in the profiles of any of the hypothesized contributors (*Materials and Methods*). I introduce here a fifth possibility corresponding to an absence of information about whether the allele is present, irrespective of whether or not it is included in the profile of a hypothesized contributor. An assumption of no information is appropriate if there is substantial uncertainty, for example, due to borderline peak height or the possibility that a peak is due to stutter or other artifact.

Using this uncertain designation for the six subthreshold alleles, the estimated dropout rate for Knox is close to 100%. A separate analysis with her as the queried contributor returned an $LR < 1$, also favoring a conclusion of no DNA from her. I reran the analysis excluding Knox from both H_p and H_d , and obtained an LR in favor of H_p of 42 million ($WoE = 7.6$ bans). Thus, although the additional alleles have, by providing evidence for an additional contributor, weakened the evidence implicating Sollecito by a massive 8 bans, this evidence nevertheless remains strong. Moreover, Gill et al. (4) did not consider uncertain

designations for peaks that are potentially due to stutter. After reclassifying as uncertain all peaks below 15% of the height at one extra repeat unit, a common stutter guideline (4), there remain four alleles not attributable to either Sollecito or Kercher and the WoE is increased to 10.7 bans.

Note that I cannot address here issues of how the DNA came to be in the exhibit: Possible contamination was an issue in the trial and appeal. I only consider whether there is DNA from Sollecito for which the evidence remains very strong after allowing for the additional alleles identified by Vecchiotti and Conti (2) and the possibility that apparent stutters are allelic.

LikeLTD Software

The probability model used to calculate these LR_s is implemented in the likeLTD (likelihoods for LTDNA profiles) software, which computes likelihoods for hypotheses, such as H_p and H_d , that specify the contributors to a sample of DNA, some or all of whom may have contributed low levels of possibly degraded DNA. For mixed-source profiles, epg peak heights are potentially informative beyond simply indicating whether or not an allele is present because they can reflect the amount of DNA, which may differ among contributors. However, this information can be difficult to exploit because peak heights for LTDNA are highly variable and this variability can be sensitive to the details of the profiling protocol used. The data input into likeLTD are the reference profiles, together with the CSP, coded as present/uncertain/absent at each allelic position in each replicate. Peak height information is used by the forensic scientist when deciding these classifications, for example, when assessing whether a peak in a stutter position should be regarded as allelic or uncertain. The full set of present/uncertain/absent indicators, combined over alleles, loci, and replicates, is highly informative about the amount of DNA from different contributors, and hence about dropout probabilities, permitting powerful and robust evaluations of evidential weight without the need to use sensitive peak height information.

In this article, I describe the probability model underlying likeLTD and assess its performance on real and artificial CSPs. I show that likeLTD provides a good solution to the problem of evaluating LTDNA profiles with up to two unprofiled contributors in addition to the queried contributor. The DNA evidence for Knox and Sollecito was criticized by Vecchiotti and Conti (2) because only a single DNA profiling run was performed. For any “noisy” scientific process, replicate analyses are desirable. This is broadly true for LTDNA evidence, but replication does have the potential disadvantage of dividing an already minuscule sample, which may adversely affect the results (11). As long as the noise is adequately modeled, which is possible by combining information over alleles and loci, replication is not a prerequisite for valid evaluation of the evidence. In Table 1, the extra uncertainty due to lack of replication has led to a much lower WoE than might have been realized had replicate PCR assays been successfully analyzed. In other words, a penalty for lack of replication

Table 1. Allele calls at 5 of 15 loci in the DNA profile obtained from exhibit 165B (case of Knox and Sollecito)

Locus	Trial*	Appeal [†]	New [‡]
D8	13, 15, 16	<u>11</u> , 12, 13, 14, 15, 16	<u>11</u> , <u>12</u> , 13, 14, 15, 16
D21	30, 32.2, 33.2	29, 30, 32.2, 33.2	29, 30, 32.2, 33.2
D7	8, 11	8, 10, 11	8, <u>10</u> , 11
CSF	10, 12	10, 11, 12	10, <u>11</u> , 12
D3	14, 16, 17, 18	14, <u>15</u> , 16, 17, 18	14, <u>15</u> , 16, 17, 18
LR [§]	7×10^{15} (15.8 bans)	4×10^7 (7.6 bans)	5×10^{10} (10.7 bans)

*Alleles reported at the original trial.

[†]Alleles identified by Vecchiotti and Conti (2); underlined alleles have peak heights <50 relative fluorescent units.

[‡]Apparent stutters are also underlined (peaks with a height <15% of the peak height at one extra repeat unit).

[§]LR for Sollecito to be a contributor of DNA, given that Kercher is a contributor, based on all 15 loci [x bans means $\log_{10}(LR) = x$].

arises automatically in likelihoods that model stochastic phenomena, such as dropin or dropout.

Results

Hammer Case. The profile data in Table 2, consisting of two CSP replicates and reference profiles from a queried contributor, Q, and two victims, K1 and K2, are taken from Table 2 of a study by Gill et al. (7), which did not consider the possibility of uncertain allele calls. There is some variability across the two replicates, a symptom of low-template and/or degraded DNA, such that 12 alleles are observed in only one of the two replicates. There are a total of 6 alleles, ≤ 2 per locus and all of them unreplicated, that are not from Q, K1, or K2. This suggests a comparison of the following two hypotheses for the contributors of DNA to the sample:

$$H'_p: Q + K1 + K2 + U1,$$

$$H'_d: X + K1 + K2 + U1,$$

where X and U1 are both unprofiled individuals. The distinction between them is that X is the alternative to Q; thus, the ethnic backgrounds of X and Q, and the degree of relatedness between them, can have important impacts on the WoE, whereas U1 plays the same role under both H'_p and H'_d .

Every CSP allele attributable to K2 could also come from K1 or Q (Table 2); thus, under H'_p , there is no evidence for DNA from K2. However, under H'_d , the DNA of Q is not present, leaving three CSP alleles that can be attributed to K2 but not to K1. Nevertheless, likeLTD estimates 100% dropout of the alleles of K2 in both replicates under both hypotheses. This is because the three alleles attributable to K2 under H'_d are all replicated, whereas seven other alleles of K2 do not appear at all, indicating very high dropout; thus, likeLTD finds that attribution of the three alleles to K2 is unlikely. Although K2 cannot be excluded from contributing any DNA to the sample, these results indicate that including K2 in the analysis brings no explanatory power and so has a negligible impact on the WoE implicating Q as a contributor.

After removing K2 from H'_p and H'_d , the WoE is 10.6 (SD = 0.10). H'_p is favored over H'_d at every locus except D18 (-0.6 bans); the

most incriminating locus (2.5 bans) was D19, where Q has two rare alleles that appeared in both CSP replicates.

The WoE of 10.6 bans computed here is stronger than that obtained by Gill et al. (the maximum of the blue solid curve in figure 1 of ref. 7 is just over 9 bans). The extra discrimination power of likeLTD results from its extra flexibility, for example, allowing different dropout rates per replicate and per contributor.

I recoded eight CSP alleles as uncertain and observed differing effects at individual loci, depending, for example, on whether the uncertain allele is in the reference profile of Q. The resulting changes in the computed WoE match intuition; for example:

Locus D16, CSPa, allele 11: This is an allele of Q not shared with K1; thus, changing its status from present to uncertain reduces the WoE, but only slightly, because the allele is called in CSPb. The single-locus WoE decreases from 1.17 to 1.14 bans (Table 2, column 4).

Locus D21, CSPb, alleles 29 and 30: These are both alleles of K1 not shared with Q, and they are not called as alleles in CSPa. Thus, changing the allele call to uncertain has a bigger impact, although still modest. The WoE is increased by just over a deciban because of the reduction in possible genotypes for X.

There are also indirect effects on all loci, because the changes in allele calls have an impact on the support for dropout parameter values. Overall, the evidence is weakened but remains very strong at 9.3 bans (Table 2, bottom row).

Simulated Profiles. Detailed results for a range of tests of likeLTD on DNA profiles subject to a number of modifications, such as artificial dropin and dropout, as well as modifications to the modeling assumptions underlying likeLTD, are provided in *SI Text*. I summarize here the main conclusions.

For a one-contributor, two-replicate CSP with no dropin or dropout (*SI Text, section S2*), likeLTD returns almost exactly the same WoE as the standard match probability formula whether or not dropin is explicitly modeled (because the dropin rate is estimated at zero) and whether X is unrelated to Q or is a brother of Q. If I wrongly hypothesize two contributors rather than one, the dropout rate for the additional contributor is estimated at 100% and the WoE is unaffected. If the CSP is modified, the WoE at individual loci changes in line with expectations and the overall WoE is reduced. For a CSP affected by four dropouts and two dropins over the two replicates, repeat likeLTD runs with a search length (n) of 1,000 simulated annealing iterations (*Materials and Methods, Parameter Estimation*) give WoE values with a SD less than half of a deciban (about 11% on the LR scale). For $n = 5,000$ and $n = 10,000$, the WoE is precise to <1% on the LR scale. When the alleged contributor Q was chosen at random, such that the prosecution hypothesis was false, the dropout rate for the noncontributor Q was estimated to be very high and, consequently, the WoE was usually negative and always low.

Proceeding to a two-contributor CSP (*SI Text, section S3*), with neither contributor known, a series of experiments introducing 50% dropout to one or both contributors, as well as uncertain allele calls due to stutter, gave satisfactory results in that parameter estimates and the WoE varied in accord with intuition. With three unknown contributors (*SI Text, section S4*), the larger number of parameters implies less precision in evaluating the WoE. With $n = 1,000$ simulated annealing iterations, the overall WoE has an SD of 0.4 bans (a factor of 2.5 on the LR scale), which is reduced to 0.24 bans (a factor of 1.7 on the LR scale) and 0.03 bans (about 6%) for $n = 5,000$ and $n = 10,000$, respectively. Taking a different three-contributor CSP, this time with one contributor a known and profiled individual, I investigated (*SI Text, section S5*) high and low extremes for the degradation parameters, the variance of the locus-specific parameters, and the dropout model power parameter. The WoE was relatively stable under these extreme perturbations, varying

Table 2. Hammer case DNA profiles and results from two analyses

Locus	D3	D16	D2	D8	D21
CSPa*	14 16	11 ^u , 13 14	20, 23 24, 25	11 ^u , 12 13, 15	28 31
CSPb	14 16	11, 13 14	20, 24 25	11, 12 13, 15 ^u	28, 29 ^u , 30 ^u 31, 31.2
Q [†]	14, 16	11, 14	24, 25	12, 13	28, 31
K1	16, 16	13, 13	20, 20	11, 15	29, 30
K2	<u>15, 17</u>	<u>12, 13</u>	<u>18, 25</u>	11, 13	29, 30
Other [‡]			23		31.2
WoE [§] , bans					
Mean (SD)	1.23 (0.057)	1.17 (0.033)	0.91 (0.084)	0.88 (0.029)	1.48 (0.14)
unc [¶]	1.10	1.14	0.95	0.94	1.59

*The crime scene DNA sample was profiled in duplicate (CSPa and CSPb). Results from 5 of 10 loci are shown.

[†]The profiles of the two uncontested possible contributors of DNA, K1 and K2, and the queried contributor, Q, are shown using the notations: **replicated alleles**, *unreplicated alleles*, and *unobserved alleles*.

[‡]CSP alleles not attributable to any of Q, K1, or K2.

[§]WoE for Q to be a contributor, given that K1 and one unprofiled individual are also contributors. The mean 10-locus WoE from 25 likeLTD runs is 10.6 bans (SD = 0.10).

[¶]Mean WoE based on 10 likeLTD runs when eight alleles were reclassified as "uncertain," of which five were at the displayed loci and are indicated with ^u. The mean 10-locus WoE is 9.3 bans (SD = 0.15).

in the range of 10.3–10.7 bans, compared with a standard analysis WoE of 10.6 bans.

Discussion

There is no “gold standard” test of an LR calculation for LTDNA profiles. Likelihoods reflect uncertainty, and even when the profiles of the true contributors are known in an artificial simulation, this does not tell us what is the appropriate level of uncertainty justified by a given observation affected by stochastic phenomena. Likelihoods depend on modeling assumptions, and there can be no “true” statistical model for a phenomenon as complex as an LTDNA profile.

I have shown here good performance of likeLTD in analyzing a wide range of crime scene DNA profiles involving complex mixtures, uncertain allele designations, dropout and dropout, degradation, stutter, and relatedness of alternative contributors. It behaves consistently over replicate analyses and agrees with well-established formulas in simple settings. The parameter estimates and WoE change in a coherent and interpretable manner under artificial modifications of the CSPs, and they are robust to major modifications of the modeling assumptions. For $n = 5,000$ iterations of the simulated annealing algorithm, the reported WoE values are reasonably precise when the hypotheses involve both U1 and U2 (SD of about 0.25 bans) and very precise when U2 is not required.

The analysis of LTDNA profiles embodied in likeLTD has elements in common with existing methods (6–10, 12, 13). It goes beyond these methods by eliminating nuisance parameters automatically via maximization of penalized likelihoods, avoiding the use of external calibration data specific to the profiling protocol used, as required by other methods (9). Even with extensive calibration data, estimation of dropout and dropout rates for the specific conditions of a crime sample cannot be precise, but precise estimates are not required: “Best fit” (in the sense of maximum penalized likelihood) values under each of the competing hypotheses provide a fair evaluation of the WoE. To achieve this, likeLTD adopts a multidose dropout model (14–16) that uses information across all replicates, loci, and contributors. The model underlying likeLTD is highly flexible, allowing both amounts of DNA and level of degradation to vary over contributors, as well as locus- and replicate-specific dropout rates. In particular, the contribution of DNA from different individuals is estimated and can be zero, such that additional profiled or unprofiled contributors can be proposed with little error arising if, in fact, there is no DNA from those individuals.

As well as providing strong WoE in favor of true contributors in simulation experiments, I showed in examples that likeLTD identified no support for the presence of DNA even when there superficially appeared to be some support and that the WoE declined appropriately as dropins and dropouts were introduced or allele calls were altered to uncertain.

The problem of how to make a numerical expression of the WoE meaningful to judges or jurors is common to all evaluations of complex DNA evidence. The problem is not insurmountable, and illustrative examples can be helpful (17).

An early “consensus” method approach to the analysis of LTDNA profiles took account only of alleles that appear in both of two DNA profiling replicates (12). This method is often claimed to be conservative, but this is not necessarily the case because it allows alleles that are inconvenient for the prosecution case to be “swept under the carpet.” The analysis proposed here makes use of all the results in every DNA profiling run. The consensus method served a useful purpose when few alternative approaches for the analysis of LTDNA profiles were available, but it is no longer best practice.

Methods of analysis that directly use epg peak height information have been developed (18, 19), but software is not currently freely available. These have potential advantages over the method proposed here, in which peak heights are used to classify every allele as present/uncertain/absent in each replicate. However, peak heights can be highly variable, and their statistical

properties can depend sensitively on details of the experimental protocol. Thus, our freely available R code likeLTD may remain useful as a robust and efficient approach to the analysis of LTDNA profiles, even if peak height-based methods can be more statistically efficient, given relevant calibration data. Previous versions of likeLTD have already been used in many criminal investigations, with results presented as evidence in UK and US courts (20).

Materials and Methods

Consider a single crime stain that may have been profiled multiple times from replicate PCR assays of the sample. Forensic DNA profiling predominantly assays autosomal short tandem repeat (STR) loci, using technology in which an allele in the profiled sample is represented by a peak in an epg (5), such as those shown in Fig. 1.

I assume that a reference profile is available for a queried contributor (Q) and the goal is to evaluate the LR for two competing hypotheses, one including Q as a contributor (the “prosecution hypothesis,” H_p), whereas the “defence hypothesis,” H_d , has an unprofiled individual X in place of Q. Both hypotheses may include additional unprofiled contributors [in practice, I can handle up to two (U1 and U2)], as well as profiled possible contributors, for example, the victim or a bystander (K1, K2, ...). The contribution of DNA from each proposed contributor is estimated, and this estimate can be zero, such that including an individual in H_p or H_d does not imply that the individual contributed DNA to the sample.

The LR can depend on, for example, the assumed ethnicity of X and his/her relatedness to Q (the more genetically similar X is to Q, the smaller is the LR). The likeLTD software program allows close relatedness of X to Q, specified with two relatedness coefficients, whereas all other hypothesized contributors must be mutually unrelated and unrelated to X and Q. In addition, remote shared ancestry (“coancestry”) of X with Q is modeled using the population genetics parameter F_{ST} (17). Typically, in US forensic practice, F_{ST} (also called θ) is only used to model intraindividual genetic correlations (i.e., excess homozygosity) (9). However, intraindividual correlations are of little relevance to evidential weight. Only between-individual correlations matter in practice, and failing to model them gives WoE values that are biased against defendants. This deficiency affects some alternative methods for analyzing DNA profiles. Because the relatedness coefficients and F_{ST} account for the positive correlations across loci due to shared ancestry of X and Q, it is reasonable to compute full-profile LR by multiplication of single-locus LR, which is standard practice in the assessment of DNA profile evidence (5). I thus focus below on the single-locus case.

Unless otherwise stated, all analyses reported here use $n = 5,000$ iterations of the simulated annealing algorithm within likeLTD. The allele frequencies from a standard database of ~200 UK Caucasians have undergone sampling and F_{ST} adjustments as described in *SI Text, section S2*.

Single-Locus LR with Dropout. Consider first a single profiling run, with a single contributor who is Q under H_p and X under H_d . If $Q \equiv AB$, where “ \equiv ” denotes “has genotype,” but the CSP shows only A and low epg peak heights suggest that dropout is possible, then the possibility that B has dropped out must be considered. Under a standard model (8, 10), the LR can be written as

$$\frac{D(1-D)}{p_A^2(1-D_2) + 2p_A(1-p_A)D(1-D)}, \quad [1]$$

where D and D_2 denote the probabilities of dropout for heterozygote and homozygote alleles, respectively. The numerator is the probability that the B allele of Q has dropped out (D), whereas the A allele has not ($1 - D$). In the denominator, either X is AA and there has been no dropout (first term) or (second term) X is heterozygous but the non-A allele has dropped out. Logically, D in the numerator of the LR is different from D in the denominator; however, typically similar values are supported under both hypotheses, and they are often taken to be equal for illustrative calculations (7).

Effect of an Uncertain Allele Designation. If I now assume CSP = A[B], where [] denotes an uncertain allele designation, and, again, $Q \equiv AB$, the LR becomes

$$\frac{1-D}{p_A^2(1-D_2) + 2p_A p_B(1-D) + 2p_A(1-p_A-p_B)D(1-D)}. \quad [2]$$

In the numerator, I know that Q’s A allele has not dropped out ($1 - D$) but not whether the B allele has dropped out. In the denominator, the three

terms correspond to $X \equiv AA$, AB , and AZ , respectively, where Z is any allele other than A or B .

Fig. 2 (solid curves) shows LR_s as functions of D for a locus with $p_A = p_B = 0.1$ (after adjustment). As expected, the LR for $CSP = A[B]$ (Fig. 2, red curve) is always intermediate between those for $CSP = AB$ (Fig. 2, black) and $CSP = A$ (Fig. 2, green). When D is high, the red and green curves in Fig. 2 are similar because in the presence of high dropout, both an uncertain designation and an absent designation for B convey little information about whether or not X has a B allele. However, when D is small, the two LR_s differ substantially because $CSP = A$ is inconsistent with $X \equiv AB$, whereas $CSP = A[B]$ is consistent with both $X \equiv AA$ and $X \equiv AB$.

Next, consider the LR_s when there is a second replicate that gives $A[B]$ in each case (Fig. 2, dashed curves). I assume the same D for both replicates. When $CSP = AB + A[B]$, I must have $X \equiv AB$ (I ignore dropin here, as discussed below). When $CSP = A + A[B]$, the LR is

$$\frac{D(1-D)^2}{p_A^2(1-D_2)^2 + 2p_A p_B D(1-D)^2 + 2p_A(1-p_A-p_B)D^2(1-D)^2},$$

whereas for $CSP = A[B] + A[B]$, it is

$$\frac{(1-D)^2}{p_A^2(1-D_2)^2 + 2p_A p_B(1-D)^2 + 2p_A(1-p_A-p_B)D^2(1-D)^2}.$$

Note that I assume the different replicates are independent, conditional on the genotypes of all contributors (6).

I see from Fig. 2 that observing $A[B]$ in the second replicate increases both LR_s when D is small but decreases them when D is large. In fact, when D is very high, observing either A or $A[B]$ in just one replicate yields $LR > 1$, favoring H_p , whereas observing two such observations in independent replicates gives $LR < 1$, against H_p . This is because $X \equiv AA$ under H_d then provides a better explanation of the replicate observations than H_p , because homozygotes are much less likely to drop out than a heterozygote allele.

Additional Contributors. LR_s, such as those in Eqs. 1 and 2, can be rewritten more generally as

$$LR = \frac{P(CSP|Q \equiv AB)}{\sum_{g \in \Gamma} p_g P(CSP|X \equiv g)}, \quad [3]$$

where Γ denotes the set of possible genotypes and p_g denotes the population fraction of genotype g . Eq. 3 makes explicit the requirement to sum over all possible genotypes for the unprofiled contributor X . When there is an

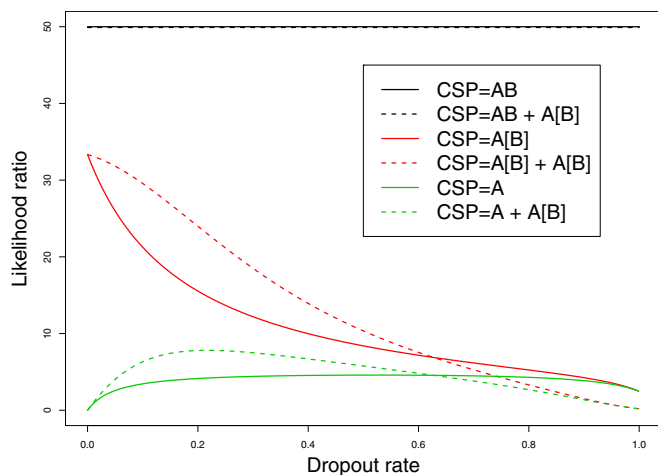


Fig. 2. Single-locus, single-contributor LR_s for three CSPs with one replicate (solid curves) and three CSPs with two replicates (dashed curves). The LR_s are expressed as functions of the dropout rate D , assumed to be the same for all alleles in both the numerator and denominator. In the legend box, “+” separates the two replicates and [] denotes an uncertain allele call. Allele A is observed in every replicate; the designation of allele B is uncertain in the second replicate, whereas it varies over present, uncertain, and absent in the first replicate.

additional unprofiled contributor $U1$, it is necessary to sum over all possibilities for the unknown genotypes, multiplying each term by the genotype probability:

$$LR = \frac{\sum_{g \in \Gamma} p_g P(CSP|Q \equiv AB, U1 \equiv g)}{\sum_{g_1, g_2 \in \Gamma} p_{g_1} p_{g_2} P(CSP|X \equiv g_1, U1 \equiv g_2)}. \quad [4]$$

Each term in these sums follows the same well-established rules used for Q and X above, now applied additionally to the current genotype for $U1$.

Multidose Dropout Model. Individuals contribute different amounts of DNA to a mixed-source sample, and multiple individuals can have one or two copies of a given allele. Thus, given $D(1)$, the dropout probability for a unit “dose” of DNA, it is necessary to evaluate $D(k)$, the dropout probability for dose k of DNA. I adopt the model of Tvedebrink et al. (14), which can be written as

$$\frac{D(k)}{1-D(k)} = (\alpha_s k)^\beta, \quad [5]$$

where s indicates the locus. I choose the scale by fixing the mean over loci of α_s at 1. I take $k=1$ to correspond to a single heterozygote allele of a reference individual, usually X or Q .

The estimates obtained by Tvedebrink et al. (14) from experimental nondegraded LTDNA profiled at the 10 loci of the SGM+ system imply an SD for α_s of 0.141. Because they may depend sensitively on the experimental protocol used, I do not use the estimates of Tvedebrink et al. (14) directly; instead, I estimate the α_s under each hypothesis from the observed CSP. To keep the estimates realistic, I impose a γ -distribution prior on the α_s , with mean = 1 and SD = 0.141 (a different SD may be appropriate, for example, in highly degraded samples). For the global parameter β , I adopt here the estimate $\beta = -4.35$ (14). Fig. 3 illustrates $D(k)$ as a function of $D(1)$ for several values of k , evaluated by substituting $\alpha_s^\beta = D(1)/(1-D(1))$ in Eq. 5. Note, for example, that if $D(1) = 0.5$, $D(1.2) \approx 0.3$ and $D(0.8) \approx 0.7$; thus, a 20% change in DNA dose can have a large impact on dropout probabilities.

The problem of calculating likelihoods for LTDNA profiles was not addressed by Tvedebrink et al. (14); they validated their model by comparing theoretical and empirical dropout rates. To achieve this, they estimated the amount of DNA from each contributor using the heights of peaks due only to that contributor over the whole profile. This is problematic for calculating LTDNA likelihoods because it ignores information from allele peaks with multiple contributors and requires alleles of individual contributors to be distinguished, which is frequently not possible. Here, I directly specify the likelihood for each replicate in terms of $D(k)$ at every allelic position, with k calculated according to the contributions of DNA and the genotypes of all the hypothesized contributors. I thus use present/uncertain/absent information at every allelic position to provide information about amounts of DNA, with the contributions from different contributors being estimated by maximum likelihood.

Degradation Model. DNA degrades over time at a rate that depends on temperature, humidity, and environmental exposure. In forensic DNA profiling, degradation is manifested as higher dropout rates for alleles with large fragment lengths. Our model for the effect of degradation is based on that of Tvedebrink et al. (15), who posited a geometrical distribution for the effective amount of DNA as a function of allele fragment length. Thus, the average allele dose k from the i th contributor subject to dropout is modified at an allele with fragment length l base pairs (centered to have mean zero) according to $k' = k(1+\gamma_i)^{-l}$, where $\gamma_i > 0$. Shorter fragments ($l < 0$) correspond, in effect, to an enhanced allele dose, whereas longer fragments generate a smaller effective allele dose. An STR allele consists of flanking regions in addition to the tandem repeats; thus, the repeat number that characterizes the allele is not a good proxy for fragment length, which can be obtained for many DNA profiling systems at www.cstl.nist.gov/div831/strbase/.

In the spirit of shrinkage regression methods, likeLTD incorporates a weak penalty (exponential, mean = 0.02) on each γ_i . The effect of this penalty is a slight tendency to shrink the parameter estimates toward zero, which is usually negligible but avoids inflated values when there is very little information.

Dropin. Dropin refers to an allele in the CSP that is not included in the genotype of any hypothesized contributor, profiled or unprofiled. Dropin alleles can arise from individuals contributing a very low level of DNA to the sample, for example, via environmental contamination either in the laboratory

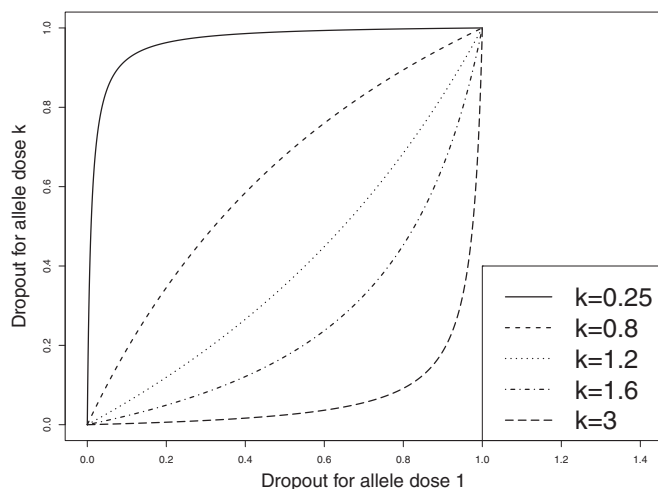


Fig. 3. Dropout probabilities for dose k of DNA (y axis) against those for a unit dose (x axis). The values of k are shown in the legend box.

or at the scene of recovery of the item. They can be generated from tiny fragments of the DNA molecule that persist for some time after the death and decay of a cell. Forensic scientists often restrict “dropin” to laboratory-based alleles, the rate of which can be measured by control runs and is usually found to be low. However, I cannot usually distinguish laboratory-based dropin from alleles arising at the crime scene (12).

Each dropin allele does come from an individual, but it may arise from very few and possibly degraded cells. It is computationally inefficient to sum over all possible genotypes, as in Eq. 4, for such low-level contributors; thus, I allow the possibility of modeling dropin more simply as independent Bernoulli trials (6). Dropin is nondropout of an allele of a low-level contributor; thus, I model the dropin probability as a constant (c) times the nondropout rate for each replicate. As for the γ_i , I impose a weak penalty on c (exponential, mean = 0.5) to discourage solutions with a large c , reflecting background information that dropin is usually rare.

Parameter Estimation. To compute the LR, it is necessary to deal with the “nuisance” parameters under each hypothesis. These are the $D(1)$ (one per replicate), the α_s (one per locus), possibly a dropin parameter c (see above),

the contributions of DNA relative to the reference individual, and the γ_i (one for each contributor subject to dropout). The likeLTD program seeks to maximize a penalized likelihood over these parameters, with penalties on α_s , γ_i , and c as described above. The penalties can be thought of as prior distributions, but I do not use a Bayesian approach because I maximize over unknown parameters rather than integrate. The primary purpose of the penalty is to discourage the maximization algorithm from exploring unrealistic regions of the parameter space.

I use a simulated annealing algorithm (21) to maximize in an approximate manner the penalized likelihood L . Starting with L computed at any set of parameter values, the algorithm repeatedly takes a random step in parameter space; compute the penalized likelihood L' ; and accept the new state with probability $\exp(-(L' - L)/t)$, where t is the temperature, computed here as $t = (1 - i/n)^3$, with i and n being the current and total numbers of iterations. Our goal is to obtain the maximized L under each hypothesis: \widehat{L}_p and \widehat{L}_q . Estimates of the nuisance parameters are available as a byproduct. They may not be precise, because there are some regions of the space of nuisance parameters over which L varies little, particularly when both U1 and U2 are included in the hypotheses being compared, because their genotypes cannot easily be distinguished. However, the assessment of evidential weight uses only $\widehat{L}_p/\widehat{L}_q$ and does not require precise estimates of the nuisance parameters.

Simulated annealing is a well-established algorithm with good properties, but there is no guarantee that it will find the exact maximum value of L . A larger n generally produces better approximations to the maximum; however, beyond a certain value, the improvement may be negligible. In *SI Text, section S2 and S3*, I show that $n = 5,000$ provides good precision for hypotheses involving both U1 and U2, as well as excellent precision when U2 is not required.

Computing Times. Using $n = 5,000$, likeLTD requires a few minutes if neither U1 nor U2 is included in the hypotheses being compared, a few hours for U1 only, and several days for both U1 and U2. Other parameters affecting computing times include the number of replicates and whether dropin is modeled. The runs of likeLTD for the Hammer case and the Knox and Sollecito case reported here required just over 10 min per locus on a standard desktop machine. A much faster implementation of the algorithm is under development.

ACKNOWLEDGMENTS. I thank Ben Lanham of Cellmark Forensic Services for providing Fig. 1 and Professor Carla Vecchiotti of the Sapienza Università di Roma for providing Fig. S1, which is a higher quality version of a figure in ref. 2. I gratefully acknowledge help with computations from Chris Steele and Adrian Timpson, both of University College London, and I also thank Dr. Torben Tvedebrink of Aalborg University and Dr. Norah Rudin, a forensic DNA consultant from California, for helpful comments on a draft of the manuscript.

- Caddy B, Taylor G, Linacre A (2008) *A Review of the Science of Low Template DNA Analysis* (UK Home Office, London).
- Vecchiotti C, Conti S (2011) [DNA evidence in the case against Amanda Knox and Raffaele Sollecito. Corte di Assise di Appello di Perugia], trans komponisto (English translation available at knoxndnareport.wordpress.com, accessed November 12, 2012).
- Good I (1979) Studies in the history of probability and statistics. XXXVII AM Turing's statistical work in World War II. *Biometrika* 66(2):393–396.
- Gill P, et al.; DNA commission of the International Society of Forensic Genetics (2006) DNA commission of the International Society of Forensic Genetics: Recommendations on the interpretation of mixtures. *Forensic Sci Int* 160(2-3):90–101.
- Buckleton J, Triggs C, Walsh S (2004) *DNA Evidence* (CRC, Boca Raton, FL).
- Curran JM, Gill P, Bill MR (2005) Interpretation of repeat measurement DNA evidence allowing for multiple contributors and population substructure. *Forensic Sci Int* 148(1):47–53.
- Gill P, Kirkham A, Curran J (2007) LoComatoN: A software tool for the analysis of low copy number DNA profiles. *Forensic Sci Int* 166(2-3):128–138.
- Balding DJ, Buckleton J (2009) Interpreting low template DNA profiles. *Forensic Sci Int Genet* 4(1):1–10.
- Mitchell AA, et al. (2012) Validation of a DNA mixture statistics tool incorporating allelic drop-out and drop-in. *Forensic Sci Int Genet* 6(6):749–761.
- Gill P, et al. (2012) DNA commission of the International Society of Forensic Genetics: Recommendations on the evaluation of STR typing results that may include drop-out and/or drop-in using probabilistic methods. *Forensic Sci Int Genet* 6(6): 679–688.
- Grisedale KS, van Daal A (2012) Comparison of STR profiling from low template DNA extracts with and without the consensus profiling method. *Investig Genet* 3(1):14.
- Gill P, Whitaker J, Flaxman C, Brown N, Buckleton J (2000) An investigation of the rigor of interpretation rules for STRs derived from less than 100 pg of DNA. *Forensic Sci Int* 112(1):17–40.
- Haned H (2011) Forensim: An open-source initiative for the evaluation of statistical methods in forensic genetics. *Forensic Sci Int Genet* 5(4):265–268.
- Tvedebrink T, Eriksen PS, Mogensen HS, Morling N (2009) Estimating the probability of allelic drop-out of STR alleles in forensic genetics. *Forensic Sci Int Genet* 3(4): 222–226.
- Tvedebrink T, Eriksen PS, Mogensen HS, Morling N (2012) Statistical model for degraded DNA samples and adjusted probabilities for allelic drop-out. *Forensic Sci Int Genet* 6(1):97–101.
- Tvedebrink T, Eriksen PS, Asplund M, Mogensen HS, Morling N (2012) Allelic drop-out probabilities estimated by logistic regression—Further considerations and practical implementation. *Forensic Sci Int Genet* 6(2):263–267.
- Balding DJ (2005) *Weight of Evidence for Forensic DNA Profiles* (Wiley, New York).
- Perlin MW, et al. (2011) Validating TrueAllele® DNA mixture interpretation. *J Forensic Sci* 56(6):1430–1447.
- Cowell RG, Lauritzen SL, Mortera J (2011) Probabilistic expert systems for handling artifacts in complex DNA mixtures. *Forensic Sci Int Genet* 5(3):202–209.
- Lohmueller KE, Rudin N (2013) Calculating the weight of evidence in low-template forensic DNA casework. *J Forensic Sci* 58(Suppl 1):S243–S249.
- Kirkpatrick S, Gelatt CD, Jr., Vecchi MP (1983) Optimization by simulated annealing. *Science* 220(4598):671–680.

Supporting Information

Balding 10.1073/pnas.1219739110

SI Text

S1. Introduction

The likeLTD (likelihoods for low-template DNA profiles) software program is available under a GNU public license from the author.

The file likeLTD-4-4.R (version 4.4) gives R functions for evaluating single-locus likelihoods for an observed crime scene profile (CSP) of DNA, given the profiles of possible contributors of DNA and a specified number of unprofiled contributors. This file is called from likeLTD-4-4-Wrapper.R, which implements a simulated annealing algorithm to maximize the product of single-locus likelihoods, multiplied by penalty terms as discussed in the main text, under both prosecution (H_p) and defense (H_d) hypotheses. The ratio of these maximized, penalized likelihoods is the likelihood ratio (LR). The wrapper code must, in turn, be called from an input file that specifies the CSP and reference profiles and assigns key case-specific parameters, such as numbers of unprofiled contributors and whether dropout is to be modeled. Parameters that tend to be fixed across analyses, such as the penalty hyperparameters and simulated annealing step sizes, are set in the wrapper code.

I report here the performance of likeLTD in computing weight of evidence (WoE) measured in bans [x bans means $\log_{10}(\text{LR}) = x$].

S2. One-Contributor CSP with Dropout and Dropin

First, three CSPs were analyzed with increasing amounts of both dropout and dropin introduced (Table S1). The contributors of DNA to the sample under the prosecution and defense hypotheses are

$$H_p^1: Q \quad \text{and} \quad H_d^1: X,$$

where X is an unprofiled individual unrelated to the profiled, suspected contributor, Q .

CSP1. Exactly the profile of Q (Table S1, row 1) appears in each replicate of CSP1 (rows 2 and 3). There would be no need to use likeLTD for this perfect single-contributor match, but I apply it here to check that likeLTD works correctly in this setting. The WoE (Table S2, row 1) is 11.45 bans. A naive application of the product rule using the allele counts in Table S1 gives 12.9 bans (Table S3); however, using the allele fractions in the bottom two rows of the table, adjusted for coancestry ($F_{ST} = 0.02$) and sampling adjustment ($\text{adj} = 1$), reduces this to 11.43 bans, close to the value generated by likeLTD. There is a small amount of error in optimizing the likelihood under H_d^1 , whereas under H_p^1 , it is evaluated exactly here, such that the WoE is overestimated by about 0.02 bans, which is negligible relative to the reduction in the WoE from sampling and F_{ST} adjustments (1.5 bans). In practical settings there will be some error in estimation of likelihoods under both H_p^1 and H_d^1 .

The adjustment value is added to the database counts for the alleles of Q to reduce the risk of understating the population frequency, particularly for rare alleles. Except in the first column of Table S3, I always use $\text{adj} = 1$. Balding (1) advocated $\text{adj} = 2$ in the absence of an F_{ST} adjustment, but because F_{ST} has a big impact on low-frequency alleles, $\text{adj} = 1$ is adequate when an appropriate F_{ST} adjustment is used. The likeLTD program adjusts for coancestry between X and Q by replacing each population allele fraction estimate p (after sampling adjustment) with:

$$(1 - F_{ST})p / (1 + F_{ST})$$

for alleles not in the profile of Q ,

$$(F_{ST} + (1 - F_{ST})p) / (1 + F_{ST})$$

for a heterozygote allele of Q , and

$$(2F_{ST} + (1 - F_{ST})p) / (1 + F_{ST})$$

for a homozygote allele of Q . I suggest that $F_{ST} = 0.02$ is a conservative value when Q is from a large, well-mixed population, whereas $F_{ST} = 0.03$, or $F_{ST} = 0.05$ in extreme cases, may be more appropriate for small, isolated, or heterogeneous populations, such as many migrant populations. Further discussion and more details of F_{ST} adjustments based on the multinomial-Dirichlet distribution are available elsewhere (2).

Brother alternative. Repeating the one-contributor analysis but now assuming that X is an unprofiled brother of Q gives 4.14 bans (Table S4). The single-locus LRs computed by likeLTD follow very closely the usual formula for a sibling (2), which in the heterozygote (homozygote) case is:

$$\text{LR} = \frac{4}{1 + p_a + p_b + 2p_a p_b} \left(\frac{4}{(1 + p_a)^2} \right).$$

There is no explicit F_{ST} adjustment in this formula because the allele fractions p have been adjusted as described above. This provides a good approximation to the adjustment based on the multinomial-Dirichlet distribution.

Assuming two contributors. If I wrongly guessed that there were two contributors to CSP1, I could compare the hypotheses

$$H_p^2: Q + U1 \quad \text{with} \quad H_d^2: X + U1,$$

where X and $U1$ denote unprofiled individuals who are unrelated to each other and to Q . In this case, likeLTD correctly estimates near 100% dropout for $U1$ under both hypotheses (Table S2, row 2) and the WoE is almost unchanged from assuming one contributor. Similarly, when X is a brother of Q , the single-locus LRs and overall WoE computed by likeLTD are almost identical to the one-contributor case (Table S4).

CSP2. The two replicates of CSP2 differ from those of CSP1 due to one dropin and two dropouts (Table S1, rows 4 and 5). The two dropouts both affect loci with large fragment lengths, consistent with the effects of DNA degradation. Because the dropin allele is at a heterozygous locus where the two alleles of Q are replicated and the two dropout alleles each appear in the other replicate, the evidence implicating Q remains very powerful and the WoE is only modestly reduced to 11.3 bans (Table S2, rows 5 and 6). The dropout parameter estimates are similar over the four replicate/hypothesis combinations. The γ_Q and γ_X estimates are, as expected, moderately large at 0.5%. Because the dropin allele must have come from somebody, it is possible to analyze CSP2 as a two-contributor profile, which implies very heavy dropout for the second contributor, $U1$. I see from (Table S2, rows 7 and 8)

that this analysis gives a slightly stronger WoE (11.4 bans), closer to the WoE for CSP1. The dropout rates for U1 are, as expected, very high, and γ_{U1} is also high (0.8%) because the dropin allele is at a short fragment length locus; thus, when attributed to U1, it gives further support to the pattern of dropout increasing with DNA fragment length.

CSP3. One further dropin and two more dropouts have been introduced into CSP3 relative to CSP2 (Table S1, rows 7 and 8). The extra dropin is at a locus for which Q is homozygous; thus, it must be a dropin under H_p^1 but not under H_d^1 . The two additional dropouts again both affect loci with large fragment lengths, and this time, an allele of Q has dropped out in both replicates (FGA 24), which has a substantial impact on the evidence implicating Q as a contributor. Although the evidence remains powerful, the overall WoE is now reduced by at least two bans compared with CSP2 (Table S2, rows 9–12). CSP3b is affected by three dropouts, whereas CSP3a is affected by only one; consequently, the dropout rate estimates are much higher for CSP3b than for CSP3a. The difference between modeling CSP3 as a one-contributor profile with two dropins and a two-contributor profile with substantial dropout for one of the contributors is now more important (0.7 bans). The two models differ in several respects, and the dropin model can be regarded as a simple approximation that reduces computation time.

Repeatability and simulated annealing search length. To investigate the effectiveness of a different simulated annealing algorithm search length n , I repeated both the one-contributor and two-contributor analyses of CSP3 10 times for $n = 1,000$, $n = 5,000$, and $n = 10,000$ (Tables S5 and S6). Even for $n = 1,000$, the WoE is estimated with $SD < 0.05$ ban, and for $n = 5,000$, the WoE computed by likeLTD is almost exact for both one- and two-contributor analyses. For these examples, the dropout and dropin parameters are well estimated with $n = 5,000$.

Randomly selected suspect. To investigate the performance of likeLTD when Q is not a contributor of DNA (thus, H_p is false), 100 profiles for Q were selected at random, according to population allele fraction estimates. Because random Q usually cannot account for many of the replicated observed alleles in CSP3, the prosecution case must be that there are two contributors of DNA, Q and U1, whereas the defense requires only X with dropin. Comparing these hypotheses, the mean dropout rates estimated for (noncontributor) Q were 90% in run 1 and 95% in run 2. Because of the high dropout rates, the WoE is necessarily low in magnitude, with a mean of -1.1 and a range of -3.6 to 1.7 bans (>0 in 19 of 100 simulations), compared with a WoE around 9 bans for a true contributor Q.

5.3. Two Unprofiled Contributors

CSP4. Table S8 (rows 1 and 2) presents results from a likeLTD analysis of CSP4, which shows exactly the alleles of a two-person mixture in both replicates, with no dropout or dropin. One of the contributors is Q, and the other is treated as unknown (U1); thus, the hypotheses compared are H_p^2 and H_d^2 . The dropout rates, γ_Q and γ_X , are correctly estimated to be close to zero. The WoE is 7.2 bans, reduced by over 4 bans from the case of a single-contributor CSP matching Q (CSP1) because of the additional uncertainty created by the masking effect of the alleles of U1.

CSP5. Table S8 (rows 3 and 4) introduces random 50% dropout for the alleles of U1 not shared with Q. Because of the reduced masking effect, the WoE is much higher than for CSP4, and is now >9 bans. The γ_i is estimated at close to zero, which is correct because the dropout rate was the same for all alleles. The dropout rates for Q are correctly estimated as zero, whereas for X, there can be less certainty because that individual's genotype is unknown and the rates are estimated at 1% and 2%. The actual numbers of dropouts of the alleles of U1 are four in CSP5a and five in

CSP5b, and this is reflected in the dropout rate estimates of around 54% for CSP5a and 72% for CSP5b. These estimates cannot be precise because of the masking of the alleles of U1 by those of Q.

Once again, 100 random profiles for Q were simulated, whereas the individual whose profile is shown in Table S7 (row 1) was treated as a profiled contributor K1. The prosecution hypothesis was that the contributors of DNA were Q, K1, and U1, whereas under the defense hypothesis, they were X and K1. The results were broadly similar as for CSP3: The mean WoE against Q was -1.0 bans, and all values were <0.1 bans this time.

CSP6. Here, the opposite scenario is considered of a random 50% dropout of the alleles of Q not shared with U1. As expected, this reduces the WoE, by 2 bans relative to CSP4. At locus FGA, both alleles of Q have dropped out in both replicates; the WoE is -0.74 at this locus and is >0 at all other loci. There are again four dropouts in CSP6a but five in CSP6b, which is reflected in different dropout rates over the two replicates. Under H_d^2 , both X and U1 are unprofiled, and so are interchangeable. The dropout rates for U1 are correctly estimated to be low under H_p^2 , whereas under H_d^2 , it is X who is assigned a low dropout. In this case, the dropouts were predominantly at the loci with large fragment lengths (D2, D18, and FGA); thus, γ_Q is moderately high, at 0.6%.

CSP7. A new difficulty is introduced in CSP7. In addition to 50% dropout for the alleles of Q, 50% of the alleles of U1 generate stutter peaks that have a nonnegligible probability of masking an allele of Q. Each of these peaks is classified as "uncertain" for the likeLTD analysis, irrespective of whether or not Q has an allele at that position. This additional ambiguity in the CSPs reduces the WoE further, to 4.7. Once again, the dropout and γ_i parameter estimates are broadly reasonable, noting that high precision is not possible here because of the masking effects of the stutters and alleles of U1.

CSP8. Random 50% dropout affects the alleles of both Q and U1. Once again, all the dropout and γ_i estimates are reasonable for the two-contributor analysis, and the WoE is (coincidentally) again 4.7. Even though the CSPs were created assuming two contributors, there are now only three replicate/locus combinations (out of 20) at which more than two alleles are observed. Thus, it would be possible, although not recommended, to analyze this case assuming one contributor plus dropin. The results (Table S8, bottom two rows) again show that this simplified analysis gives a conservative result, with the WoE now being 2.7.

5.4. Three Unprofiled Contributors

Three unprofiled contributors were generated using the profiles from the Hammer case but omitting K2 and treating K1 as unprofiled. The hypotheses compared were then

$$H_p^3: Q + U1 + U2 \quad \text{with} \quad H_d^3: X + U1 + U2.$$

Under both H_p^3 and H_d^3 , a search length of $n = 1,000$ is reasonably precise even when there are three contributors (Table S9). However, Tables S5, S6, and S9 all suggest a slight tendency to overstate the WoE when $n = 1,000$, perhaps because there are more parameters to be estimated under the defense hypothesis than under the prosecution hypothesis.

There is no indication from these results that $n = 5,000$ overstates the WoE, and it gives excellent precision for one and two contributors. For three contributors, an SD of about 0.25 bans should usually be adequate, and if higher precision is required, multiple searches can be run in parallel to reduce total computation time, with the largest values of the penalized likelihoods

under each hypothesis used to form the reported LR. On my desktop computer, with three unknown contributors, likeLTD performs just less than 100 simulated annealing iterations per hour; thus, $n = 5,000$ requires about 2 days.

The estimates of the nuisance parameters from the runs underlying Table S9 are reasonably precise under H_p^3 but not under H_d^3 . This is expected, because there are many different ways to attribute the observed alleles to the three unprofiled contributors and so many plausible combinations of parameter values. As noted above, lack of precision in estimation of the nuisance parameters is unimportant in evaluating the WoE.

55. Effects of Modifying the Dropout Model

The CSP and reference profiles for the Hammer case are given in Table 2, where I report the WoE in favor of H_p over H_d of 10.6 bans for a “standard” analysis that excludes K2 and does not include a dropout term. I also report there the results from an analysis, including dropout, and from an experiment in which some allele calls have been changed to uncertain. Here, I undertake some further analyses of the Hammer case to investigate the effects on the WoE of various modifications to the dropout model, described near Eq. 5 in *Materials and Methods*. Some results are given in Table S10 and are briefly discussed below.

Varying the Degradation Model. Degradation is manifested in higher dropout rates for alleles with a large fragment length. In the SGM+ system, this has the greatest effect for the alleles at the D2, D18, and FGA loci. In Table 2, note a reduced number of observed alleles relative to other loci at D18 and FGA, but not at D2. To investigate the effect of modeling degradation in likeLTD, I repeated the analysis with the degradation parameter γ_i set to zero for every contributor i , which implies no change in dropout probability with fragment length. This increased the

WoE by nearly 3 decibans at locus D2 relative to the standard analysis but decreased it by 0.5 ban at D18.

In contrast to fixing the γ_i , I removed the penalty term (“free” in Table S10). I also included K2 as a profiled possible contributor. The estimates of γ_Q , γ_X , γ_{K1} , and γ_{U1} were almost identical to those given in Table 3 (rows 3 and 6), confirming that the penalty term has little effect in the presence of good information. For K2, I have essentially no information, and γ_{K2} was estimated at 0.000 and 0.000 with the penalty but at 0.001 and 0.003 without it, illustrating that the penalty can prevent inflated γ_i estimates in the presence of little information.

The overall effect on the WoE of both the failure to model degradation and the removal of the penalty term was close to zero (Table S10).

Varying the Locus Adjustments. The locus adjustment parameter α_s allows for differences in dropout rates over loci, beyond the dependence on fragment length that is captured with the degradation model. In the standard analysis, likeLTD imposes a γ -distribution penalty term with both parameters equal to 50, implying a prior variance in α_s of 0.14. I first fixed α_s at one (variance = 0), so that the dropout probability of an allele of a given fragment length is the same for all loci. The overall WoE was 10.5 bans. The converse change was to weaken the γ -penalty by changing its prior SD from 0.14 to 0.32 (the two parameters of the γ -distribution penalty changed from both = 50 to both = 10). This reduced the overall WoE to 10.3 bans.

Varying the Power Parameter. I performed new power parameter (β) analyses assigning $\beta = -3.97$ and $\beta = -4.73$. These values represent 1 SD above and below the central estimate of -4.35 reported by Tvedebrink et al. (3). The overall WoE changes from 10.6 to 10.4 and 10.7, respectively.

1. Balding DJ (1995) Estimating products in forensic identification using DNA profiles. *J Am Stat Assoc* 90(431):839–844.
2. Balding DJ (2005) *Weight of Evidence for Forensic DNA Profiles* (Wiley, Chichester, UK).

3. Tvedebrink T, Eriksen PS, Mogensen HS, Morling N (2009) Estimating the probability of allelic drop-out of STR alleles in forensic genetics. *Forensic Sci Int Genet* 3(4):222–226.

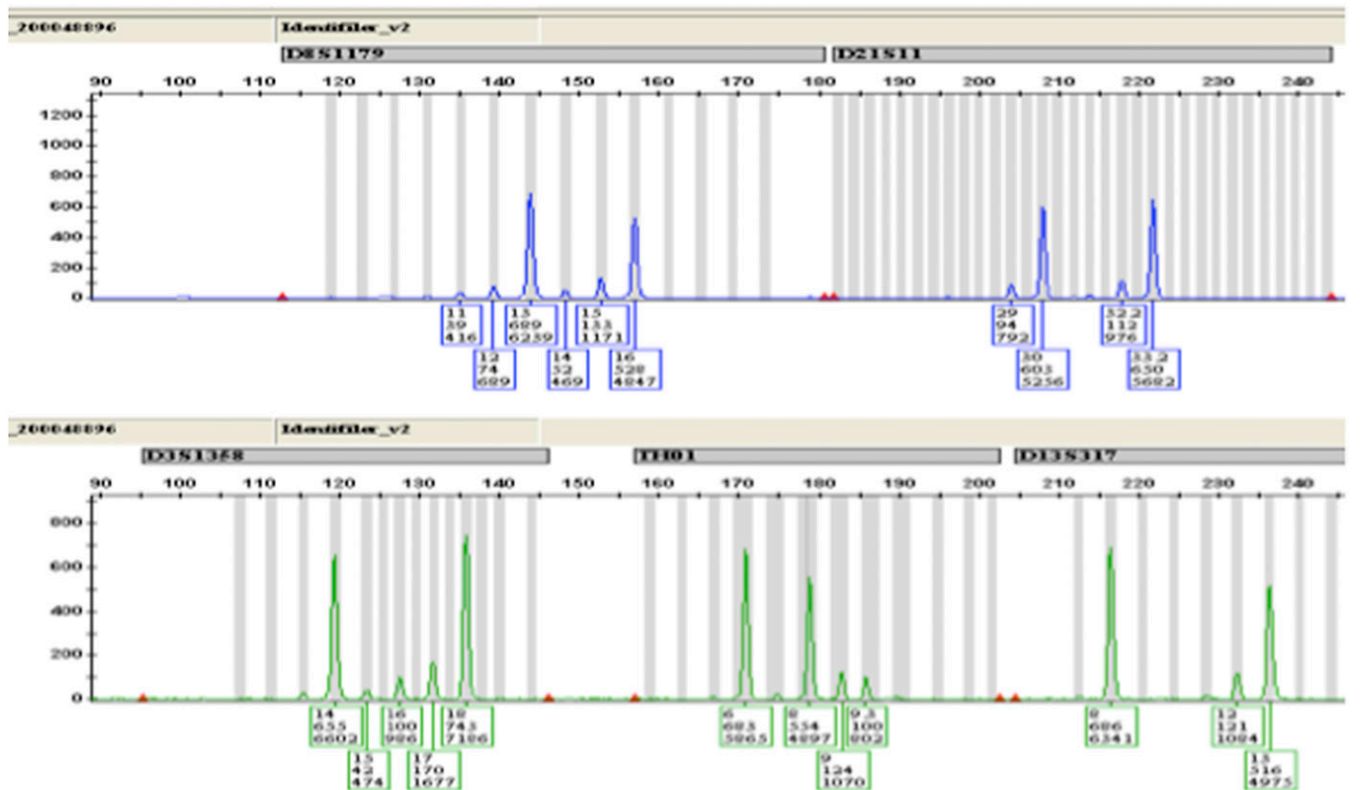


Fig. S1. Electropherogram corresponding to the five loci of Table 1. The y axis shows peak height in relative fluorescence units, and the x axis shows fragment length in base pairs.

Table S1. One-contributor profiles

Locus	D3	vWA	D16	D2	D8	D21	D18	D19	TH01	FGA
Q	15, 16	15, 18	11, 11	17, 24	13, 13	29, 31	10, 16	13, 14	6, 9.3	22, 24
CSP1a	15, 16	15, 18	11	17, 24	13	29, 31	10, 16	13, 14	6, 9.3	22, 24
CSP1b	15, 16	15, 18	11	17, 24	13	29, 31	10, 16	13, 14	6, 9.3	22, 24
CSP2a	15, 16	15, 18	11	17, 24	13	29, 31	10, 16	13, 14	6, 9.3	22
CSP2b	15, 16, 17	15, 18	11	17	13	29, 31	10, 16	13, 14	6, 9.3	22, 24
CSP3a	15, 16	15, 18	11	17, 24	13, 15	29, 31	10, 16	13, 14	6, 9.3	22
CSP3b	15, 16, 17	15, 18	11	17	13	29, 31	16	13, 14	6, 9.3	22
Count*	116, 106	38, 80	130	90, 35	118	74, 28	2, 56	102, 133	77, 151	75, 47
Total†	396	384	418	388	384	384	384	396	384	384
Adjusted	0.302	0.117	0.341	0.244	0.338	0.206	0.027	0.268	0.214	0.209
Fraction‡	0.278	0.221		0.108		0.0918	0.161	0.343	0.398	0.139

The crime scene DNA sample was profiled in duplicate (a and b), and the observed alleles in each replicate are reported in separate rows for each of three hypothetical CSPs: CSP1 is a full single-contributor CSP, whereas CSP2 and CSP3 are the same profile with dropout and dropout introduced.

*Database counts for the alleles of Q.

†Total database allele count.

‡Allele fractions after sampling and F_{ST} adjustments (adj = 1, F_{ST} = 0.02).

Table S2. WoE and some parameter estimates from likeLTD analyses of the Table S1 profiles

CSP, no. of contributors	WoE, bans	Hypothesis	Dropout a	Dropout b	Degradation, γ_i	Dropin
CSP1, 1 contributor, mean (SD)	11.45 (0.01)	H_d^1	0.00 (0.00)	0.00 (0.000)	0.000 (0.0001)	0.00 (0.001)
CSP1, 2 contributors, mean (SD)	11.45 (0.01)	H_d^2	0.00, 1.00 (0.000), (0.000)	0.00, 1.00 (0.000), (0.000)	0.000, 0.000 (0.0001), (0.0001)	—
CSP2, 1 contributor	11.3	H_p^1	0.03	0.03	0.005	0.07
		H_d^1	0.03	0.03	0.005	0.09
CSP2, 2 contributors	11.4	H_p^2	0.05, 1.00	0.02, 0.99	0.005, 0.008	—
		H_d^2	0.06, 0.99	0.03, 0.99	0.004, 0.008	—
CSP3, 1 contributor	8.7	H_p^1	0.03	0.12	0.005	0.15
		H_d^1	0.00	0.16	0.003	0.10
CSP3, 2 contributors	9.4	H_p^2	0.06, 0.95	0.14, 0.98	0.004, 0.007	—
		H_d^2	0.05, 0.96	0.15 0.99	0.003, 0.009	—

Each CSP was analyzed twice, assuming one and then two contributors. The letters a and b denote the two replicate profiling runs. For two contributors, both dropout and γ_i estimates are given (for Q and U1 under H_p^2 and for X and U1 under H_d^2). For CSP1, the mean (SD) based on 10 likeLTD runs are shown for the WoE and parameter estimates under H_d^i ; there are no unknown parameters under H_p^1 . Results for CSP2 and CSP3 are from one likeLTD run.

Table S3. Overall WoE for CSP1

F_{ST} , %	0	0	1	2	3
Sampling adjustment	0	1	1	1	1
WoE, bans	12.9	12.6	11.9	11.43	11.0

The underlying allele counts are given in Table S1.

Table S4. Single-locus LRs and overall WoE for CSP1 when X, the alternative source of the DNA, is a brother of Q

No. of contributors	Single-locus LR										WoE, bans
	D3	vWA	D16	D2	D8	D21	D18	D19	TH01	FGA	
1	2.29	2.88	2.22	2.85	2.23	3.00	3.34	2.23	2.24	2.84	4.14
2	2.29	2.88	2.22	2.85	2.24	3.00	3.35	2.23	2.24	2.85	4.14

Results are shown for both one- and two-contributor analyses. Note that the LR can never exceed 4.

Table S5. Mean (SD) of WoE, the maximized and penalized log-likelihoods, and dropout/dropin parameters over 10 likeLTD analyses of rows 9 and 10 in Table S2 (CSP3, one-contributor)

n	WoE, bans	$\log_{10}(\hat{L})$, mean	Dropout CSP3a	Dropout CSP3b	γ_Q or γ_X	Dropin
1,000	8.78	-3.73 (0.082)	0.03 (0.002)	0.12 (0.005)	0.005 (0.0003)	0.15 (0.007)
	(0.046)	-12.49 (0.044)	0.00 (0.000)	0.15 (0.013)	0.004 (0.0005)	0.10 (0.022)
5,000	8.73	-3.70 (0.001)	0.03 (0.001)	0.12 (0.002)	0.005 (0.001)	0.15 (0.004)
	(0.003)	-12.43 (0.003)	0.00 (0.000)	0.15 (0.006)	0.003 (0.0001)	0.09 (0.004)
10,000	8.72	-3.70 (0.001)	0.03 (0.001)	0.12 (0.003)	0.005 (0.0001)	0.15 (0.002)
	(0.002)	-12.43 (0.003)	0.00 (0.000)	0.15 (0.004)	0.003 (0.0001)	0.09 (0.005)

For each value of the simulated annealing algorithm search length n , the parameter estimates (SD) are given under H_p^1 (top row) and H_d^1 (bottom row).

Table S6. Mean (SD) of WoE, the maximized and penalized log-likelihoods, and dropout/dropin parameters over 10 likeLTD analyses of rows 11 and 12 in Table S2 (CSP3, two-contributor)

n	WoE, bans	$\log_{10}(\hat{L})$, mean	Dropout CSP3a	Dropout CSP3b	γ_Q or γ_X	Contribution from U1
1,000	9.46	-2.96 (0.007)	0.12 (0.183)	0.14 (0.012)	0.004 (0.000)	0.25 (0.048)
	(0.035)	-12.42 (0.033)	0.06 (0.026)	0.15 (0.053)	0.003 (0.001)	0.25 (0.048)
5,000	9.43	-2.95 (0.002)	0.06 (0.002)	0.14 (0.003)	0.004 (0.000)	0.27 (0.003)
	(0.001)	-12.38 (0.001)	0.05 (0.003)	0.14 (0.004)	0.003 (0.000)	0.23 (0.005)
10,000	9.43	-2.95 (0.002)	0.06 (0.002)	0.14 (0.004)	0.004 (0.000)	0.26 (0.003)
	(0.001)	-12.38 (0.001)	0.05 (0.002)	0.14 (0.005)	0.003 (0.000)	0.23 (0.011)

The final column gives the estimated amount of DNA from U1 relative to Q or X, which is the parameter used by likeLTD to estimate dropout rates for U1.

Table S7. Two-contributor profiles

Locus	D3	vWA	D16	D2	D8	D21	D18	D19	TH01	FGA
Q	15, 16	15, 18	11, 11	17, 24	13, 13	29, 31	10, 16	13, 14	6, 9.3	22, 24
CSP4a	15, 16, 17	15, 18	11, 13	17, 20, 24	10, 13	29, 31, 32.2	10, 11, 13, 16	13, 14	6, 8, 9.3	19, 22, 24, 25
CSP4b	15, 16, 17	15, 18	11, 13	17, 20, 24	10, 13	29, 31, 32.2	10, 11, 13, 16	13, 14	6, 8, 9.3	19, 22, 24, 25
CSP5a	15, 16	15, 18	11, 13	17, 24	10, 13	29, 31, 32.2	10, 11, 16	13, 14	6, 8, 9.3	19, 22, 24
CSP5b	15, 16, 17	15, 18	11, 13	17, 24	13	29, 31, 32.2	10, 16	13, 14	6, 9.3	22, 24, 25
CSP6a	15, 16, 17	15, 18	11, 13	17, 20	10, 13	29, 31, 32.2	10, 11, 13	13, 14	6, 8, 9.3	19, 25
CSP6b	15, 16, 17	15, 18	11, 13	17, 20, 24	10, 13	29, 32.2	11, 13	13, 14	6, 8, 9.3	19, 25
CSP7a	15, 17	15, 18	11, 13	17, 20	10, 13	29, 31, 32.2	10, 11, 13	13, 14	6, 8, 9.3	19, 25
	[16]	[14], [17]	[10]	[16], [19]	[12]	—	—	[12]	[5]	[18], [24]
CSP7b	15, 17	15, 18	11, 13	17, 20, 24	10, 13	29, 32.2	11, 13	13, 14	6, 8, 9.3	19, 25
	—	[17]	[12]	[19]	—	[31.2]	[10]	[12]	[5]	[24]
CSP8a	15	15, 18	11	17, 20, 24	—	29, 31	10, 11, 16	13, 14	6, 8	19, 22
CSP8b	16, 17	15, 18	13	17, 20	10, 13	31	10, 16	13	6, 8	22, 24, 25

The CSPs correspond to those of Q and one unprofiled contributor, modified by dropout and (for CSP7) stutter generating uncertain allele calls (denoted by []). The layout is similar to that of Table S1, except there are extra rows for CSP7 to show the uncertain allele calls.

Table S8. WoE and parameter estimates for the CSPs of Table S7

CSP	WoE, bans	Hypothesis	Dropout, a	Dropout, b	Degradation, γ_i	Dropin
CSP4	7.2	H_p^2	0.00, 0.00	0.00, 0.00	0.000, 0.000	—
		H_d^2	0.00, 0.05	0.00, 0.06	0.003, 0.002	—
CSP5	9.7	H_p^2	0.00, 0.54	0.00, 0.71	0.000, 0.000	—
		H_d^2	0.01, 0.55	0.02, 0.73	0.000, 0.001	—
CSP6	5.1	H_p^2	0.22, 0.03	0.30, 0.00	0.006, 0.000	—
		H_d^2	0.00, 0.72	0.00, 0.80	0.000, 0.000	—
CSP7	4.7	H_p^2	0.26, 0.02	0.42, 0.03	0.003, 0.000	—
		H_d^2	0.06, 0.09	0.06, 0.10	0.000, 0.005	—
CSP8	4.7	H_p^2	0.43, 0.70	0.41, 0.68	0.000, 0.000	—
		H_d^2	0.35, 0.74	0.34, 0.73	0.000, 0.000	—
2 contributors						
CSP8	2.7	H_p^1	0.25	0.32	0.000	0.87
1 contributor		H_d^1	0.26	0.21	0.000	0.56

All analyses assume (correctly) two contributors, and their dropout and γ_i values are both shown, except for the final rows, which incorrectly assume one contributor + dropin. The letters a and b refer to profiling replicates.

Table S9. Effect of simulated annealing search length

<i>n</i> *	WoE, bans (SD)										
	Overall	D3	vWA	D16	D2	D8	D21	D18	D19	TH01	FGA
1,000	8.0 (0.40)	0.88 (0.27)	1.02 (0.05)	0.75 (0.12)	0.89 (0.19)	0.15 (0.18)	1.31 (0.26)	0.37 (0.20)	1.79 (0.06)	0.08 (0.19)	0.70 (0.18)
5,000	7.8 (0.24)	0.88 (0.18)	1.03 (0.02)	0.75 (0.10)	0.73 (0.12)	0.15 (0.10)	1.31 (0.14)	0.36 (0.13)	1.85 (0.03)	0.03 (0.10)	0.61 (0.09)
10,000	7.9 (0.03)	0.91 (0.01)	1.05 (0.01)	0.78 (0.02)	0.71 (0.03)	0.18 (0.02)	1.38 (0.01)	0.38 (0.04)	1.86 (0.02)	0.06 (0.02)	0.64 (0.02)

Summary of results from 10 likeLTD analyses for the three-contributor CSP (*SI Text*, section S4).

*The search length of the simulated annealing algorithm.

Table S10. Single-locus and overall WoE values for analyses of the Hammer case imposing various modifications to the dropout model

	WoE, bans										
	D3	vWA	D16	D2	D8	D21	D18	D19	TH01	FGA	Overall
Standard analysis (no dropout, K2 omitted)*											
Mean	1.23	1.07	1.17	0.91	0.88	1.48	-0.59	2.49	0.70	1.27	10.6
SD	(0.057)	(0.017)	(0.033)	(0.084)	(0.029)	(0.14)	(0.078)	(0.034)	(0.037)	(0.089)	(0.10)
Degradation parameter, γ_i											
0	1.23	1.08	1.24	1.18	0.82	1.52	-1.08	2.44	0.67	1.20	10.6
Free [†]	1.21	1.04	1.18	0.84	0.87	1.55	-0.44	2.47	0.69	1.21	10.6
Locus adjustment, α_s , prior variance											
0	1.04	0.98	1.17	1.21	0.88	1.35	-0.37	2.40	0.78	1.03	10.5
0.32	1.13	1.15	1.00	0.68	0.87	1.54	-0.39	2.47	0.70	1.38	10.3
Power parameter, β											
-4.73	1.04	1.00	1.20	1.17	0.87	1.37	-0.67	2.42	0.72	1.25	10.4
-3.97	1.07	1.00	1.22	1.18	0.89	1.40	-0.48	2.45	0.78	1.13	10.7

*These results, based on 25 likeLTD runs, are reproduced from Table 2 for ease of comparison. The other results each correspond to a single likeLTD run.

[†]The exponential penalty has been removed.