

Hidden Multiple Comparisons Increase Forensic Error Rates

Susan Vanderplas^{a,1,2}, Alicia Carriquiry^{b, c}, and Heike Hofmann^{b, c, 1}

This manuscript was compiled on February 5, 2024

When wires are cut, the tool produces striations on the cut surface; as in other forms of forensic analysis, these striation marks are used to connect the evidence to the source that created them. Here, we argue that the practice of comparing two wire cut surfaces introduces complexities not present in better-investigated forensic examination of toolmarks such as those observed on bullets, as wire comparisons inherently require multiple distinct comparisons, increasing the expected false positive rate. We call attention to the multiple comparison problem in wire examination and relate it to other situations in forensics that involve multiple comparisons, such as database searches.

Forensic Evidence | Statistics | Wire cuts | Toolmark analysis

In forensic evaluations, a single conclusion often relies on many comparisons, either implicitly or explicitly. Multiple comparisons arise persistently when developing statistical methods to address scientific problems (1), and greatly increase the probability of false discoveries. Now that vast databases and efficient algorithms are routinely used in forensic evaluations to propose matches to crime scene items, the problem of close non-matches (2) due to multiple comparisons becomes critically important. This often ignored issue increases the false positive rate, and can contribute to the erosion of public trust in the justice system. The multiple comparison problem is not new: it has been raised in the past with regard to DNA (3) and latent print evaluations (4). One of the root causes (5) leading to the wrongful accusation of Brandon Mayfield in the 2004 Madrid train bombing case was that the large size of the IAFIS database used to search for similar prints made it possible to locate ‘unusually’ close non-matches. The probability of finding close non-matches for an item of evidence is proportional to the size of the database available for searching; close non-matches occur more frequently as databases increase in size.

Compounding this issue, the use of algorithms also results in a large number of comparisons that are not obvious to the user. For example, the cross-correlation function (6) was one of the first measures proposed for quantifying the similarity between two patterns in response to the 2009 NRC report (7), and continues to be used in many pattern searching algorithms to find the best alignment between two images and to quantify their overall similarity. Finding the best alignment often consists in sliding one surface across the whole length (for one-dimensional patterns, such as striations) or area (for two dimensional sources, such as impression marks) of the other item while keeping track of the value of a similarity measure. This mirrors the forensic examination process: the examiner visually rotates and shifts items under a comparison microscope to align two surfaces. In order to avoid false accusations and further erosions of public trust in science, we must address the problem of multiple comparisons and control their effect on false discovery rates.

Here, we consider the multiple comparisons problem that arises from a relatively simple toolmark examination: matching a cut wire to a wire-cutting tool. We describe the comparison approach, estimate the (minimal) number of comparisons that are needed to carry out the examination, and discuss how the false discovery rate changes with the number of comparisons involved, using error rates derived from published black-box studies.

Examination Process

A forensics examiner tasked with determining whether a wire in evidence was cut by a specific tool will create one or more blade cuts using the putative tool, which are then compared to the cut surface of the wire. These cuts are made in a sheet of material similar to the wire’s composition, and may be

Author affiliations: ^aStatistics Department, University of Nebraska Lincoln, 350 Hardin Hall, 3310 Holdrege North Wing, Lincoln, NE 68503; ^bDepartment of Statistics, Iowa State University, 1121 Snedecor Hall, 2438 Osborn Dr, Ames, IA 50011; ^cCenter for Statistics and Applications in Forensic Evidence, 195 Durham Center, 613 Morrill Road, Ames, Iowa 50011

Please provide details of author contributions here.

The authors have no competing interests to declare.

¹SVP (Author One) contributed equally to this work with HH (Author Two).

²To whom correspondence should be addressed. E-mail: susan.vanderplas@unl.edu

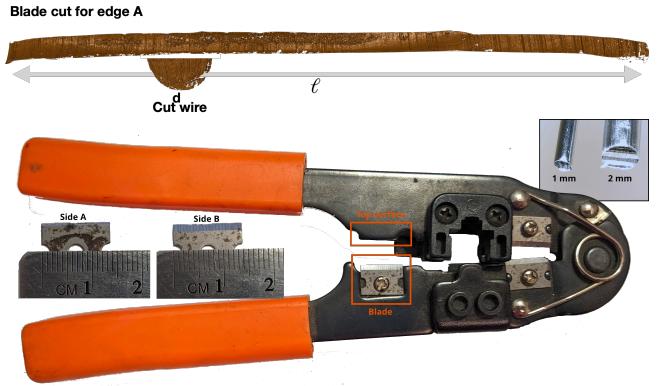
125 performed at multiple angles, as the angle of the tool to
 126 the substrate can affect which striations are recorded on the
 127 substrate surface. The blade cuts will then be compared to
 128 the wire under a comparison microscope, though eventually,
 129 automatic comparison algorithms may also be validated for
 130 lab use. Each side of each blade cut will be compared to each
 131 side of the wire; depending on the tool design, there can be
 132 between two and four cutting surfaces in contact with the
 133 substrate.

134 Calculating the Number of Comparisons

135 In order to calculate the number of comparisons carried out
 136 in the course of one examination, we define ℓ to be the length
 137 of the blade cut, and d to be the diameter of the wire. We
 138 assume that the wire is covered with striations suitable for
 139 comparison across its full diameter d . If this is not the case,
 140 we reduce the value d . Both the blade and the wire are either
 141 digitally scanned at resolution r mm per pixel, or visually
 142 examined using a microscope with a digital resolution that
 143 can be expressed as r equivalent to the digital scan. An
 144 illustration of the sliding comparison process is shown in
 145 Figure 1. Imagine that we move the cut wire along the
 146 blade cut in order to assess whether striations on the blade
 147 cut match the striations on the wire. We can move the wire
 148 unit-by-unit, or we can move the wire by its full length, with
 149 no overlap to the previous comparison.

150 The first option gives us the maximum number of
 151 comparisons $(\ell/r - d/r + 1)$, while the second option gives us
 152 the minimum number of comparisons ℓ/d . In the first case,
 153 sequential comparisons share much of the same physical data
 154 and are highly related; in the second case, no data are shared
 155 between physical comparisons and we can expect that they
 156 are statistically independent, though empirically there will
 157 be nonzero correlations due to physical similarities between
 158 striations. For simplicity, let us consider the number of
 159 comparisons to lie somewhere between these two estimates.
 160 Note that when $\ell/d \approx 1$, as in some toolmark comparisons,
 161 the upper number of comparisons goes to 1. Finally, we must
 162 consider the number of surfaces which must be compared:
 163 the wire may have one or two sets of striae and there may
 164 be two to four blade cut surfaces to examine, depending on
 165 the tool. This results in a multiplier of as much as 8.

166 **A concrete example.** Let us consider a wire-cutting tool with
 167 a 1.5 cm razor blade that meets a cast surface (one such
 168 tool is shown in Figure 1); the wire is held against this
 169 rectangular cast surface as the blade is pushed into the wire,
 170 splitting it in two. This is a minimal scenario - the wire
 171 will acquire striations from one side of the blade, while the
 172 blade itself has two cutting edges, which we will call side A
 173 and side B. A blade cut of a sheet of aluminum will thus
 174 produce two striated edges corresponding to side A and side
 175 B which are compared to cut wires to assess similarity. We
 176 also have a 12 gauge aluminum wire (2 mm diameter) which
 177 may have been cut by the wire-cutting tool described above.
 178 Class characteristics, which are shared by all tools of similar
 179 manufacture, appear to match: there is a flat impression on
 180 one side of the wire corresponding to the cast metal backstop
 181 of the tool, and the wire is cut such that the blade and the
 182 backstop appear to be perpendicular (that is, the wire
 183 appears to have been cut with a tool of similar configuration).



187 **Fig. 1.** (Top) A comparison between a wire and a blade cut requires sliding the wire
 188 along the entire blade cut length to determine the best match (or whether there is a
 189 match). Surfaces shown are rendered 2D topographical scans of a wire and blade
 190 cut taken with a confocal light microscope. (Bottom) RJ45 Crimp tool with a 1.5 cm
 191 razor blade used for cutting. 1 mm and 2 mm diameter aluminum wires cut with the
 192 pliers are shown in a box in the top right corner.

200 In this example, $\ell = 15$ mm, $d = 2$ mm, and there are at least
 201 $\ell/d = 7.5$ comparisons between a wire cut and a blade cut.
 202 As there are two blade cuts (side A and side B), the minimal
 203 number of comparisons overall is 15, as these comparisons
 204 are non-overlapping and independent (on average).

205 Assuming a resolution of $0.645\mu\text{m}$ per pixel, the
 206 maximum number of comparisons per blade cut is around
 207 20,000; thus, we need 40,000 comparisons overall in
 208 order to find the optimal alignment between the wire and
 209 the blade cut. These comparisons are implicit in the
 210 calculation of cross-correlation, which is the first and often
 211 the only step used to quantitatively assess the similarity
 212 between striated evidence such as bullets, aperture shear,
 213 and firing pin impressions. Implicit comparisons are not
 214 unique to algorithms; an examiner would need to physically
 215 align the wire and the blade cut by searching along the
 216 length of the cut to visually match striations, performing
 217 the same process physically that the algorithm performs
 218 computationally. While these sequential comparisons are
 219 highly auto correlated, and we cannot assume sequential
 220 independence when calculating the probability of an error,
 221 they serve as an upper bound on the number of comparisons
 222 which could be performed. As the number of comparisons
 223 increases, the probability of encountering a coincidental
 224 match increases. Statisticians call this the *family-wise*
 225 *error rate E*; it is an important quantity to control when
 226 conducting a series ("family") of tests.

227 Probability of False Identifications

228 There are at least two components of the false positive rate
 229 (FPR): identifying two pieces of evidence that have similar
 230 characteristics but are from different sources (a coincidental
 231 match) and procedural failures (e.g. lab process errors) (2, p
 232 50). In objective disciplines with standardized evaluation
 233 rules (e.g. DNA), these sources can be distinguished.
 234 However, in toolmark examination, no objective evaluation
 235 rules are used; examiners testify based on subjective,
 236 individual rules for how much similarity is sufficient for an
 237 identification. Assuming that lab procedure errors are not
 238 a factor in studies, we use reported error rates from open-

Study	FPR e	False Positives (%) in N comparisons			
		E_{10}	E_{100}	$E_{1,000}$	$E_N < 0.1$
Mattijssen (2021)	7.24%	52.8	99.9	100.0	1
Pooled Error	2.00%	18.3	86.7	100.0	5
Bajic (2020)	0.70%	6.8	50.7	99.9	14
Best (2022)	0.45%	4.5	36.6	98.9	23
	1 in 1,000	1.0	9.5	63.2	105
	1 in 10,000	0.1	1.0	9.5	1,053
	1 in 100,000	10^{-4}	0.1	1.0	10,535

Table 1. Table showing the relationship between false positive rates and the chance of false identifications in N comparisons for a set of different FPRs and different number of comparisons. The last column gives the number of comparisons allowed while ensuring a false identification percentage of at most 10%.

set studies of striated evidence comparisons to estimate the coincidental match rate of a single wire-cut comparison. The three studies we consider (8–10) have FPRs between 0.0045 (10) and 0.072 (9); pooling data from these studies yields an FPR of 0.02. For a single-comparison FPR of e , the family-wise FPR for n comparisons, E_n is $1 - [1 - e]^n$. ?? shows the impact the number of comparisons has on these published error rates. With an error rate of 0.07, as suggested by Bajic (2020), examiners can make up to 14 comparisons, i.e. only the comparisons of one wire to a single-bladed tool are safely covered without exceeding an upper bound of 10% for the family wise false identification error. For a modestly sized database search of size 1000, the initial false positive rate should not exceed 1 in 10,000 to guarantee a family-wise total false identification error of at most 10%.

Under these constraints, the accuracy of an examination involving multiple comparisons between a wire and a tool will be low, as the number of candidate alignments that must be examined is high. Even the most innocuous example (small blade, only 2 cutting surfaces, and a relatively large wire) involves a minimum of 15 comparisons. When we then consider that examiners would make cuts under multiple angles (11), increasing the number of comparisons, the probability of a false identification becomes even more likely. As a result, it is questionable whether wire comparisons made under current protocols are reliable enough to be presented at trial.

Obviously, we need studies for wire evidence. It is also clear, based on the error rates, that larger studies are needed. Going forward, we need to move away from the binary assessment of 'match' or 'non-match' and quantify (a) the similarity of striations expressed between two pieces of evidence, and (b) the frequency of the observed pattern. Unusual striation patterns should be assigned more weight in the process.

Discussion & Conclusions

Forensic practitioners often report the findings from their examinations in the form of a categorical conclusion reflecting a single decision. This is misleading when the decision relies on multiple comparisons which are not individually presented in reports or testimony. In this short contribution, we have shown that the implicit comparisons

performed during forensic analysis of wire cuts increase the family-wise error rate.

We describe a simple scenario where a wire is cut using a two-sided blade, but findings apply to any situation where a forensic evaluation involves multiple comparisons, including, e.g., database searches. Forensic practitioners should understand how the number of comparisons can affect the accuracy of their final conclusion. We propose three strategies to enhance transparency and enable more reliable estimates of examination-specific error rates.

First, examiners should report (or defense attorneys should request) the overall length or area of surfaces generated during the examination process, along with the total consecutive length or area of the recovered evidence. These pieces of information will take the place of ℓ and d and facilitate calculation of examination-wide error rates.

Second, researchers should conduct studies relating the length/area of comparison surface to the error rate. For instance, we have pooled studies looking at bullet striations and firing pin shear marks because we could not find black-box error rate studies of wire cuts. The striated surfaces are of orders of magnitude different lengths, but represent the best estimate of the error rate for striated materials. New studies should be Type II studies as defined by Koehler (12), designed to assess the actual error rate when examiners are making difficult comparisons.

Finally, when databases are used at any stage of the forensic evidence evaluation process (from suitability assessment and triage to reports which will be used at trial), the number of database items searched (or comparisons made) and the number of results returned must be reported. Additionally, the number of results used for further manual comparison should also be reported. For example, if a firearms examiner searches a local NIBIN database with 1000 entries, requests the 20 closest matches to her evidence, and then carries out a physical examination of five exemplars from the list of 20, all of those values should be clearly reported to enable estimation of the examination-wise error rate. This will help make the multiple comparison issue accessible to everyone involved in evaluating the value of forensic evidence, from examiners to lawyers to jurors and judges.

ACKNOWLEDGMENTS. This work was funded (or partially funded) by the Center for Statistics and Applications in Forensic Evidence (CSAFE) through Cooperative Agreements 70NANB15H176 and 70NANB20H019 between NIST and Iowa State University, which includes activities carried out at Carnegie Mellon University, Duke University, University of California Irvine, University of Virginia, West Virginia University, University of Pennsylvania, Swarthmore College and University of Nebraska, Lincoln.

- Y Benjamini, Y Hochberg, Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *J. Royal Stat. Soc. Ser. B (Methodological)* **57**, 289–300 (1995).
- President's council of advisors on science and technology, *Forensic science in criminal courts: Ensuring scientific validity of feature-comparison methods*. (Executive Office of the President of the United States, President's Council), (2016).
- WC Thompson, F Taroni, CGG Aitken, How the probability of a false positive affects the value of dna evidence. *J. Forensic Sci.* **48**, 2001171 (2003).
- JJ Koehler, S Liu, Fingerprint error rate on close non-matches. *J. Forensic Sci.* **66**, 129–134 (2021).
- GA Fine, A review of the fbi's handling of the brandon mayfield case unclassified executive summary, (Washington DC), Technical report (2006).
- T Vorburger, et al., Applications of cross-correlation functions. *Wear* **271**, 529–533 (2011).
- NRC, *National Research Council: Strengthening Forensic Science in the United States: A Path Forward*. (National Academies Press), (2009).

373	8. S Bajic, LS Chumbley, M Morris, D Zamzow, Validation study of the accuracy, repeatability,	435
374	and reproducibility of firearm comparisons. (Ames, IA), Technical report (2020).	436
375	9. EJAT Mattijsen, et al., Firearm examination: Examiner judgments and computer-based	437
376	comparisons. <i>J. Forensic Sci.</i> 66 , 96–111 (2021) _eprint:	438
377	https://onlinelibrary.wiley.com/doi/pdf/10.1111/1556-4029.14557 .	439
378	10. BA Best, EA Gardner, An assessment of the foundational validity of firearms identification	439
379	using ten consecutively button-rifled barrels. <i>AFTE J.</i> 54 , 28–37 (2022).	440
380	11. M Baker, Toolmark variability and quality depending on the fundamental parameters:	440
381	Angle of attack, toolmark depth and substrate material. <i>Forensic Sci. Int.</i> p. 10 (2015).	441
382	12. JJ Koehler, Intuitive Error Rate Estimates for the Forensic Sciences. <i>Jurimetrics</i> 57 ,	442
383	153–168 (2017).	443
384		444
385		445
386		446
387		447
388		448
389		449
390		450
391		451
392		452
393		453
394		454
395		455
396		456
397		457
398		458
399		459
400		460
401		461
402		462
403		463
404		464
405		465
406		466
407		467
408		468
409		469
410		470
411		471
412		472
413		473
414		474
415		475
416		476
417		477
418		478
419		479
420		480
421		481
422		482
423		483
424		484
425		485
426		486
427		487
428		488
429		489
430		490
431		491
432		492
433		493
434		494
		495
		496