

Three-dimensional data of wirecut surface scans under the confocal microscope (110 character maximum, inc. spaces)

Yuhang Lin^{1,2} and Heike Hofmann^{1,2}

¹Iowa State University, Department of Statistics, Ames,

²Center for Statistics and Applications in Forensic Evidence (CSAFE), Iowa State University, Ames,

*corresponding author(s): Yuhang Lin (ymlin@iastate.edu) ???????it.corresponding returns nothing

*corresponding author(s): Heike Hofmann (hofmann@iastate.edu) ???????it.corresponding returns nothing

ABSTRACT

Wire cut data is important in forensic investigations but lacks a systematic way of analyzing the data. We created a data set of 120 scans of aluminum wire cut in $\times 3p$ format, using 5 wire cutters and 3 locations along the 4 blades, with 2 replicates for each combination. A systematic pipeline with multiple analysis plots was developed to analyze the data and draw conclusions based on numerical measures.

(maximum 170 words) This is a manuscript template for Data Descriptor submissions to *Scientific Data* (<http://www.nature.com/scientificdata>). The abstract must be no longer than 170 words, and should succinctly describe the study, the assay(s) performed, the resulting data, and the reuse potential, but should not make any claims regarding new scientific findings. No references are allowed in this section.

Please note: Abbreviations should be introduced at the first mention in the main text no abbreviations lists or tables should be included. Structure of the main text is provided below.

Background & Summary

Wire cut data is a type of forensic tool mark data used to identify the source of a wire cutter based on the striations left on the surface. There have been cases ([don't mention the name?](#)) where the wire cut data played a crucial role in the criminal investigation. However, there is a lack of a standardized method to analyze it except for visual comparison. Here, we provide a data set of 120 scans of aluminum wire cut in $\times 3p$ format, and we conduct a technical validation to introduce a systematic pipeline for analyzing these types of data based on numerical measures. We hope this pipeline developed using this data set can be further generalized and applied to real crime scenes to help investigators draw conclusions based on real wire cut data.

(unlimited length) An overview of the study design, the assay(s) performed, and the created data, including any background information needed to put this study in the context of previous work and the literature. The section should also briefly outline the broader goals that motivated the creation of this dataset and the potential reuse value. We also encourage authors to include a figure that provides a schematic overview of the study and assay(s) design. The Background & Summary should not include subheadings. This section and the other main body sections of the manuscript should include citations to the literature as needed.

Methods

In this study, aluminum wire was used to create an optimal scenario where the most amount of information could be transferred from the tool to the substrate, despite the wire in some real cases being made of lead. The aluminum wire used was 16 Gauge/1.5 mm, anodized.

In order to cut the wire, 4-inch pieces were unspooled and cut using Kaiweets wire cutters, model KWS-105, as shown in Figure 1(a), for 1 blade location, either inner, middle, or outer, which gives us 1 replicate. Each piece was then cut into half to create 2-inch pieces for each side, AB and CD, with a sharpie line marking the cut ends, giving us 4 samples. Here, we are showing AB sides only in Figure 1(b) ([need a different tent figure](#)), and the CD sides are similar from the other side of the cut, with the back of A being C and the back of B being D. Both AB and CD sides form tent structures on the tips of the wire, and we can separate each side of the tent into 2 pieces along the bending position, resulting in 8 scans. We repeated this process for all 3 locations along the blade and 5 wire cutters, with 2 replicates for each tool-edge-location combination, resulting

37 in 120 scans. Each piece was labeled with the naming conventions, T(ool) 1/2/3/4/5 (Edge) A/B/C/D W(ire) - L(ocation)
38 I(nner)/M(iddle)/O(uter) - R(epetition) 1/2, with T1AW-LI-R1 being the piece cut by tool 1 on the A edge at the inner location
39 for the first repetition.

40 A more detailed procedure, including standard scanning protocols for the confocal microscope in Figure 1(c) (need an
41 extra pic of the very tip), can be found in the [README of the GitHub repository heike/Wirecuts](#) (High-res pics needed
42 in the README). The scanned surfaces are saved in a resolution of $0.645\mu\text{m} \times 0.645\mu\text{m}$ per square pixel in an x3p file
43 format.

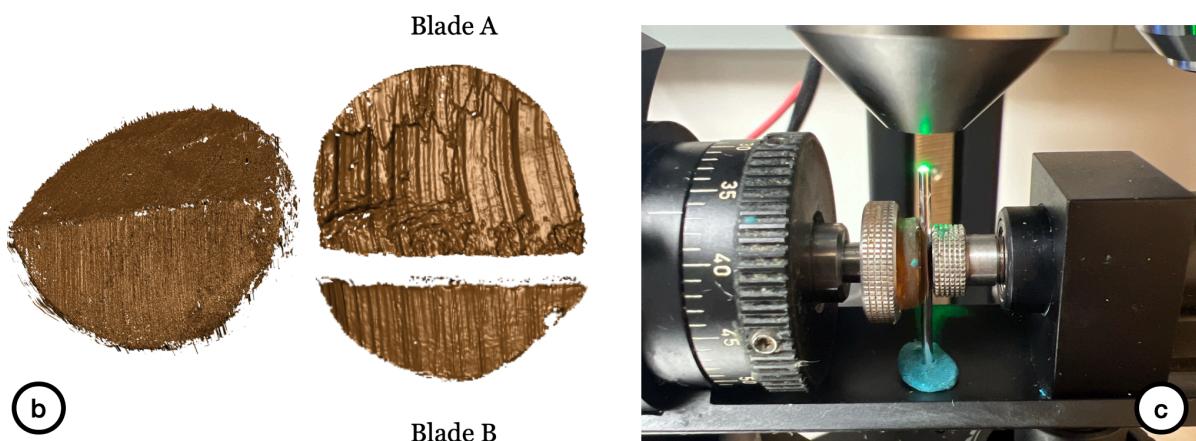
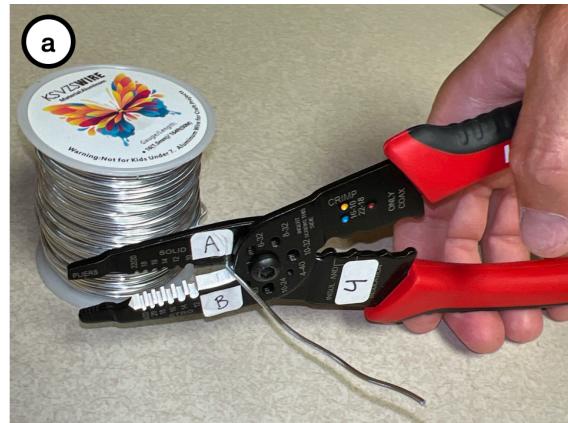


Figure 1. (a) A Kaiweets wire cutter of model KWS-105 was used to cut the wire. (b) A tent structure created by blade AB. After separating 2 tent structures by the connecting position, we obtained 2 samples - 1 sample from blade A and B. (c) A confocal microscope was used to scan the wire surfaces.

44 (unlimited length) The Methods should include detailed text describing any steps or procedures used in producing the data,
45 including full descriptions of the experimental design, data acquisition assays, and any computational processing (e.g. normalization,
46 image feature extraction). See the detailed section in our submission guidelines for advice on writing a transparent
47 and reproducible methods section. Related methods should be grouped under corresponding subheadings where possible, and
48 methods should be described in enough detail to allow other researchers to interpret and repeat, if required, the full study.
49 Specific data outputs should be explicitly referenced via data citation (see Data Records and Citing Data, below).

50 Authors should cite previous descriptions of the methods under use, but ideally the method descriptions should be complete
51 enough for others to understand and reproduce the methods and processing steps without referring to associated publications.
52 There is no limit to the length of the Methods section. Subheadings should not be numbered.

53 Authors should review the transparent methods checklist below, and ensure that their manuscript complies with any relevant
54 points. Authors are also encouraged to search FAIRsharing.org for community reporting standards that may be relevant
55 to their specific data-type.

56 **Transparent Methods Checklist**

- 57 • Materials & reagents: Identify commercial suppliers of reagents, instrumentation or kits, when the source is critical to

58 the outcome of the experiments. Declare any restrictions on the availability of unique materials (more information here).
59 Provide catalogue or clone numbers for all antibodies (if available). For primary antibodies, provide proof of validation
60 for the relevant species and applications.

- 61 • Exclusion criteria: If any data or samples were excluded, explain the exclusion criteria and state in the methods whether
62 the criteria were established before the study was conducted.
- 63 • Randomization & blinding: For any studies that involve assigning samples, animals or participants into different groups:
64 State clearly whether randomization methods were used. If randomization was not employed, this should be clearly
65 stated. State clearly whether blinding was employed during data collection. If blinding was not employed, this should
66 be clearly stated.
- 67 • Animal & human studies (full journal policies here): Experiments involving human participants must identify the
68 committee approving the experiments, and include a statement confirming that informed consent was obtained from all
69 participants. Studies employing nonhuman animals should ensure that methods descriptions comply with the ARRIVE
70 checklist.
- 71 • Cell lines: For each eukaryotic cell line used, state the source and whether the cell line has been authenticated or
72 otherwise tested for integrity. If any commonly misidentified cell lines were used (see ICLAC or NCBI Biosample),
73 justify their use. Report whether the cell lines were tested for mycoplasma contamination.
- 74 • Chemistry & materials science: Manuscripts describing chemical syntheses, or characterizing new chemicals or materi-
75 als should refer to the guidance at Nature Chemistry.

76 Data Records

77 The complete data set is available on the ISU DataShare repository at <https://iastate.figshare.com/>, which is public and open
78 access for every interested researcher. The data set consists of 120 scans in the x3p file format with the naming convention
79 as described before. (Explain the x3p header info?)

80 (unlimited length) The Data Records section should be used to explain each data record associated with this work, in-
81 cluding the repository where this information is stored, and to provide an overview of the data files and their formats. Each
82 external data record should be cited numerically in the text of this section, for example ?, and included in the main reference
83 list as described below. A data citation should also be placed in the subsection of the Methods containing the data-collection
84 or analytical procedure(s) used to derive the corresponding record. Providing a direct link to the dataset may also be helpful
85 to readers (<https://doi.org/10.6084/m9.figshare.853801>).

86 Tables should be used to support the data records, and should clearly indicate the samples and subjects (study inputs), their
87 provenance, and the experimental manipulations performed on each (please see 'Tables' below). They should also specify the
88 data output resulting from each data-collection or analytical step, should these form part of the archived record.

89 Technical Validation

90 For the data collection process, two team members did the cutting and labeling together, then one person did the scanning and
91 named according to the naming convention above. The scanning was done in a specific order to ensure consistency across all
92 scans. The data was saved in a consistent format to ensure they could be easily accessed and analyzed. A third person then
93 checked the data to ensure that the data was consistent in naming and accurate.

94 We also conduct numerical comparisons between 2 replicates made by the same side of the same tool at the same location.
95 For example, we compare T1AW-LI-R1 to T1AW-LI-R2, T1CW-LI-R1 to T1CW-LI-R2, and so on. That is comparing each
96 row in Figure 2. The comparison is done by

- 97 1. manually extracting profiles from scans,
- 98 2. obtaining the filtered signals from scans,
- 99 3. aligning signals and computing the cross-correlation function (CCF) values.

100 Extract profiles

101 A representative function with the most information is considered as a signal for one scan, which can be used for comparison
102 later. To obtain this function, we first need a profile of the scan, which is a sequence of values along a user-drawn line on
103 the surface. The profile should capture most features of the scan, and be orthogonal to the striation marks of the scan, which
104 are formed by ups and downs of grooves. So, we draw the line across the wide region of the scan to maximize the feature
105 captured, as shown in dark blue in Figure 3(a). We can then investigate the values under this profile line. The profile function
106 is along the line is plotted in Figure 3(b).

107 Filtered signals

108 With the profile extracted, we can then obtain the signal. Two Gaussian filters are applied to these resulting profiles. In
109 particular, we first used a large low-pass filter with bandwidths of 400 microns to remove large trend, as it can overwhelm the
110 signals, and then used a small high-pass filter of 40 microns to average across noise and remove spikes, as shown in Figure
111 3(c). (add reference: W. S. Cleveland, E. Grosse and W. M. Shyu (1992) Local regression models. Chapter 8 of Statistical
112 Models in S eds J.M. Chambers and T.J. Hastie, Wadsworth & Brooks/Cole.). Finally, the extreme tail values are removed.

113 Align signals and compute the cross-correlation function (CCF) values

114 Signals extracted from different scans can be put together for comparison, and we maximize the overlapping region of the
115 signals by aligning them. We know that signals from two replicates with the same tool-edge-location combination should
116 yield similar signals as in the first and second column of Figure 4, which will give alignments of massive overlapping and high
117 CCF values close to 1. The alignments and values we got in the rightmost column of Figure 4 fulfill our expectations.

118 For validation of all other tools and locations of scan replicates, see the detailed [report](#).

119 (unlimited length) This section presents any experiments or analyses that are needed to support the technical quality of
120 the dataset. This section may be supported by figures and tables, as needed. This is a required section; authors must present
121 information justifying the reliability of their data.

- 122 • Measurement of data quality?

- 123 – Numeric measurements / tests: ?
- 124 – Visualizations: ?
- 125 – Check with existing data: ?
- 126 – Questionable / slur procedures:

127 * [AidDatas Geospatial Global Chinese Development Finance Dataset](#): Second, all data collected is reviewed
128 by at least two individuals. Although this is not a double-blind review procedure, the use of satellite imagery
129 to verify project locations results in far less uncertainty when compared to previous approaches to geocoding
130 where locations were selected entirely based on text descriptions.

131 * [A large open access dataset of brain metastasis 3D segmentations on MRI with clinical and imaging information](#): A medical student (D.R.) double checked and adjusted the revised NIfTI segmentation masks and
132 manually counted the number of lesions with contrast-enhancement, necrosis, and peritumoral edema for
133 each patient.

135 * [Time series of freshwater macroinvertebrate abundances and site characteristics of European streams and](#)
136 [rivers](#): Technical validation of the TREAM dataset was achieved through exclusion of time series data that
137 did not match our inclusion criteria and data standardisation steps (outlined in Methods above). Any noted
138 issues that did not adhere to the outlined standardisation within the datasets from the 41 independent projects
139 included in this dataset were checked with data providers and corrected or removed when standardisation
140 was not achievable (e.g., when collection methods changed over the course of the time series).

141 * [3D surgical instrument collection for computer vision and extended reality](#): The main issue...Since we
142 store our models in a standard format (STL), they are compatible with a large variety of visualisation and
143 processing software.

144 * [Three-dimensional reconstruction of high latitude bamboo coral via X-ray microfocus Computed Tomography](#) : Regular quality assurance inspections are carried out on the μ -CT scanner to verify its metrological and
145 geometrical (alignments) accuracy for conducting the scans. The geometry of source to object and source
146 to detector distances are verified whenever there is any significant physical interaction with the source such

as re-alignment, change of filament, or source anode change. This calibration process involves scanning a specially designed phantom known as an hourglass36, which consists of three pairs of high-sphericity spheres. The sphere sizes are as follows: two spheres with a diameter of 3.000 mm, two spheres with a diameter of 6.000 mm, and two spheres with a diameter of 9.525 mm, and each sphere is kept in contact with its size-counterpart. By using this phantom, it becomes possible to accurately determine a known distance, specifically the centre-to-centre distance of the spheres, in a threshold-independent manner. If the measured distance deviates beyond the acceptable limits of metrological accuracy, the systems calibration parameters are adjusted to ensure agreement between the measured distance and the actual distance.

Usage Notes

The data set can be easily accessed with the CRAN R package `x3ptools`. Further analysis can be conducted with the GitHub R package `wire` and the GitHub R shiny app `wireShiny` ([citation](#)). We already conduct between-replicate comparisons in the technical validation section, and we can also conduct across-replicate comparisons to establish error rates threshold and produce other analysis plots.

Suppose we put the CCF values in a tilemap with different tool, location and edge combinations. In that case, we expect only the diagonal to have high CCF values, close to 1 and marked as orange in the tilemap, as the diagonal represents the same source, and the rest of the matrix to have low CCF values, close to 0 and marked as gray. In Figure 5, the behavior is consistent with our expectation overall, except for some rare cases with tool 5 edge D, which is caused by [?????](#). We also put the resulting CCFs in the boxplot, as in Figure 6. We can see that the CCF values for the same sources are close to 1, while the CCF values for different sources are much lower than expected. This difference can be used to establish a threshold for CCF and help us draw conclusions about the similarity between wire cut scans numerically, which can be used in real crime scenes. The density plot in Figure 7 shows the distribution of the CCF values with the same sources and different sources. The overlapping points between the tails of these two distributions can be a rough threshold. Furthermore, the receiver operating characteristic (ROC) curve in Figure 8 shows the sensitivity / true positive rate against the false positive rate (FPR) (1 - specificity). The curve is very close to the upper left corner, which is excellent for classification and drawing conclusion. It gives us a true threshold of 0.589 to control the FPR to be less than 0.05 with false negative rate (FNR) to be 0, ([false positive rate \(FPR\) / false discovery rate \(FDR\)](#) -> define the H0 or call it [false identification rate \(FIR\)](#)??), and 0.658 to control the FPR to be less than 0.01, with FNR to be 0.02.

(unlimited length) The Usage Notes should contain brief instructions to assist other researchers with reuse of the data. This may include discussion of software packages that are suitable for analysing the assay data files, suggested downstream processing steps (e.g. normalization, etc.), or tips for integrating or comparing the data records with other datasets. Authors are encouraged to provide code, programs or data-processing workflows if they may help others understand or use the data. Please see our code availability policy for advice on supplying custom code alongside Data Descriptor manuscripts.

For studies involving privacy or safety controls on public access to the data, this section should describe in detail these controls, including how authors can apply to access the data, what criteria will be used to determine who may access the data, and any limitations on data use.

Code availability

The data set can be easily accessed with the CRAN R package `x3ptools`. Further analysis can be conducted with the GitHub R package `wire` and the GitHub R shiny app `wireShiny` ([citation](#)) ([again](#)??).

For all studies using custom code in the generation or processing of datasets, a statement must be included under the heading "Code availability", indicating whether and how the code can be accessed, including any restrictions to access. This section should also include information on the versions of any software used, if relevant, and any specific variables or parameters used to generate, test, or process the current dataset.

1 End of Body

(Note that the bibliography style and the name of the bib-file are hard coded in the template file right now.)

LaTeX formats citations and references automatically using the bibliography records in your .bib file, which you can edit via the project menu. Use the cite command for an inline citation, e.g. [????](#). For data citations of datasets uploaded to e.g. `figshare`, please use the `howpublished` option in the bib entry to specify the platform and the link, as in the `Hao : gidmaps : 2014` example in the sample bibliography file. For journal articles, DOIs should be included for works in press that do not yet have volume or page numbers. For other journal articles, DOIs should be included uniformly for all articles or not at all. We recommend that you encode all DOIs in your bibtex database as full URLs, e.g. <https://doi.org/10.1007/s12110-009-9068-2>.

198 **Acknowledgements**

199 Acknowledgements should be brief, and should not include thanks to anonymous referees and editors, or effusive comments.
200 Grant or contribution numbers may be acknowledged.

201 **Author contributions statement**

202 Y.L. did all of the work, H.H. made him do the work. But seriously, this paper is the one where we need to cite everybody:
203 Eden Amin, Curtis Mosher, Jeff Salyards. Alicia? Must include all authors, identified by initials, for example: A.A. conceived
204 the experiment(s), A.A. and B.A. conducted the experiment(s), C.A. and D.A. analysed the results. All authors reviewed the
205 manuscript.

206 **Competing interests**

207 (mandatory statement)

208 The corresponding author is responsible for providing a [competing interests statement](#) on behalf of all authors of the paper.
209 This statement must be included in the submitted article file. H.H. is a technical advisor to AFTE (Association of Firearms and
210 Toolmarks Examiners), fellow of the ASA (American Statistical Association), and committee member of the ASA Forensic
211 Science Committee. H.H. has testified as court witness on behalf of judge April Neubauer, NY State Supreme Court Criminal
212 Term in New York City.

213 **Figures & Tables**

214 Figures, tables, and their legends, should be included at the end of the document. Figures and tables can be referenced in
215 L^AT_EX using the ref command, e.g. Figure 9 and Table 1.

216 Authors are encouraged to provide one or more tables that provide basic information on the main inputs to the study
217 (e.g. samples, participants, or information sources) and the main data outputs of the study. Tables in the manuscript should
218 generally not be used to present primary data (i.e. measurements). Tables containing primary data should be submitted to an
219 appropriate data repository.

220 Tables may be provided within the L^AT_EX document or as separate files (tab-delimited text or Excel files). Legends, where
221 needed, should be included here. Generally, a Data Descriptor should have fewer than ten Tables, but more may be allowed
222 when needed. Tables may be of any size, but only Tables which fit onto a single printed page will be included in the PDF
223 version of the article (up to a maximum of three).

224 Due to typesetting constraints, tables that do not fit onto a single A4 page cannot be included in the PDF version of the
225 article and will be made available in the online version only. Any such tables must be labelled in the text as Online-only tables
226 and numbered separately from the main table list e.g. Table 1, Table 2, Online-only Table 1 etc.

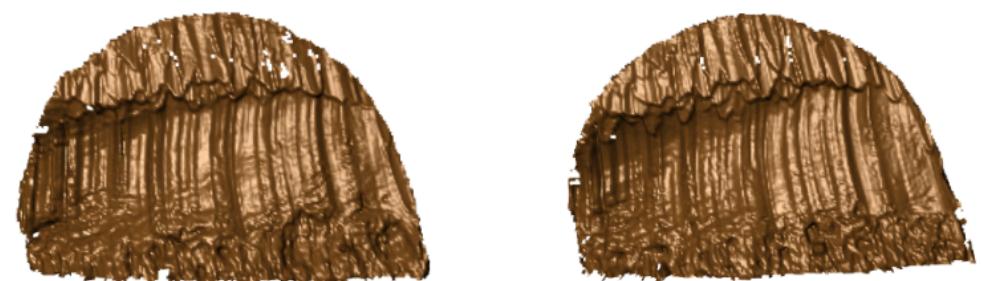
| Condition | n | p |
|-----------|----|------|
| A | 5 | 0.1 |
| B | 10 | 0.01 |

227 **Table 1.** Legend (350 words max). Example legend text.

Edge A



Edge C



Edge B



Edge D



Figure 2. Scans from different sides of tool 1 at the inner location.

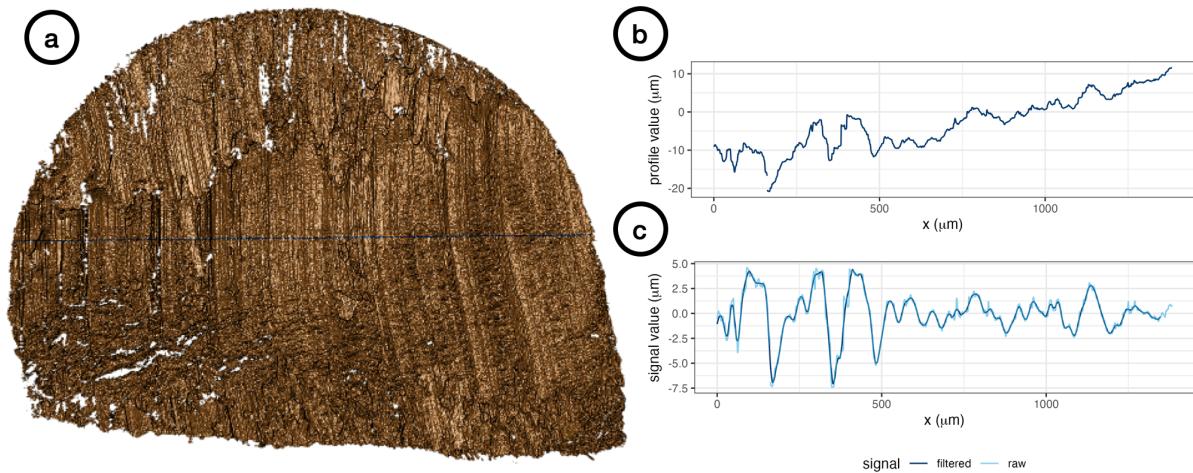


Figure 3. (a) A profile line in dark blue was drawn across the striations of the scan. (b) The profile function extracted along the profile line in (a). (c) The raw signal in light blue is obtained by using the low-pass filter on the profile function in (b) and the filtered signal is obtained by using the high-pass filter on the raw signal.

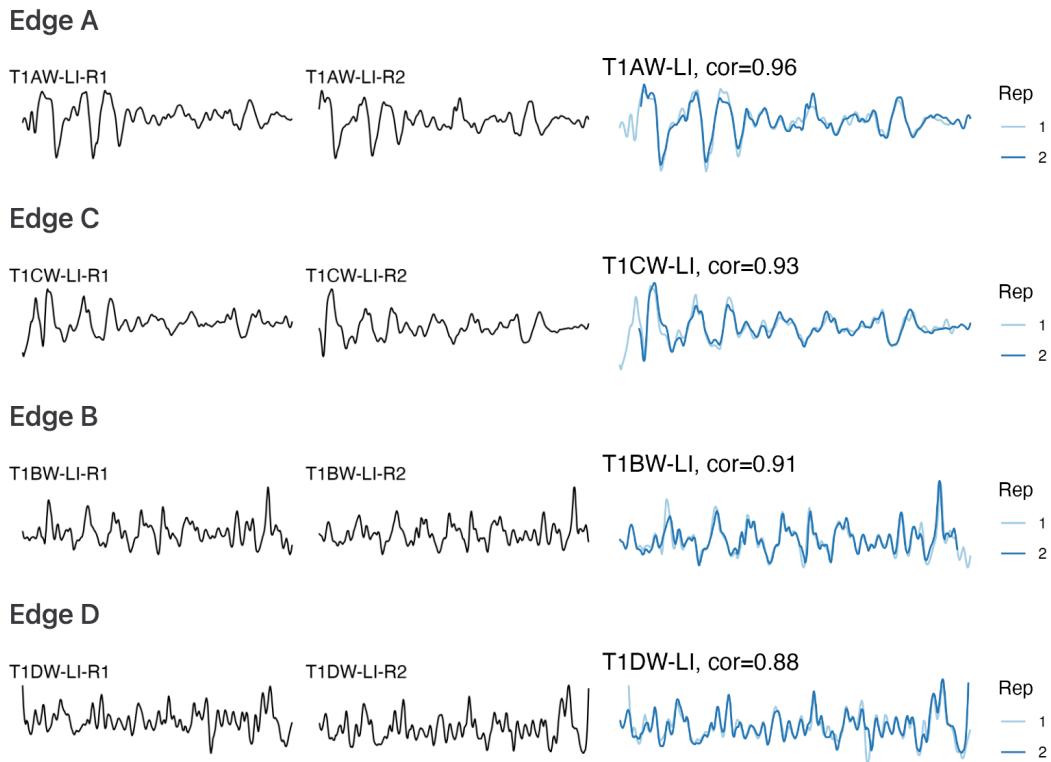


Figure 4. The first and second columns show the signals extracted from Figure 2, and the third column shows the alignments and CCF values between pairs of signals.

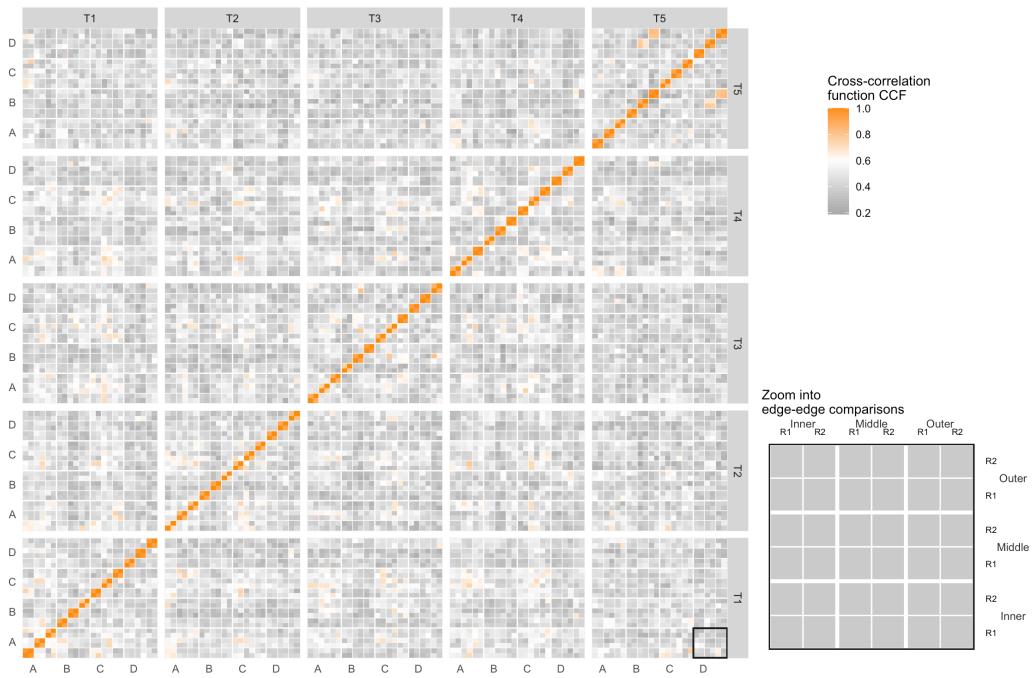


Figure 5. The tilemap shows signals from the same source have CCFs close to 1.

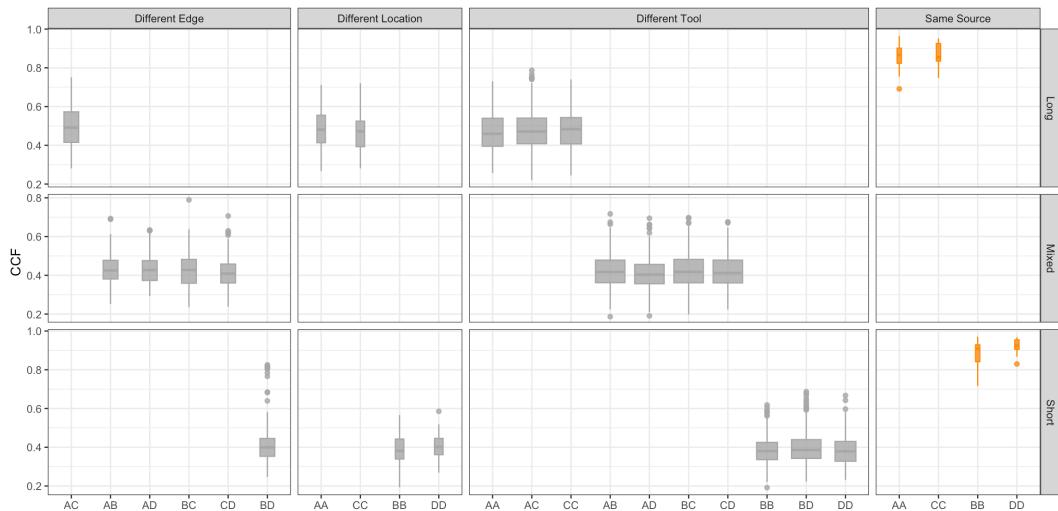


Figure 6. The boxplot shows that signals from the same sources have higher CCFs than those from different sources.

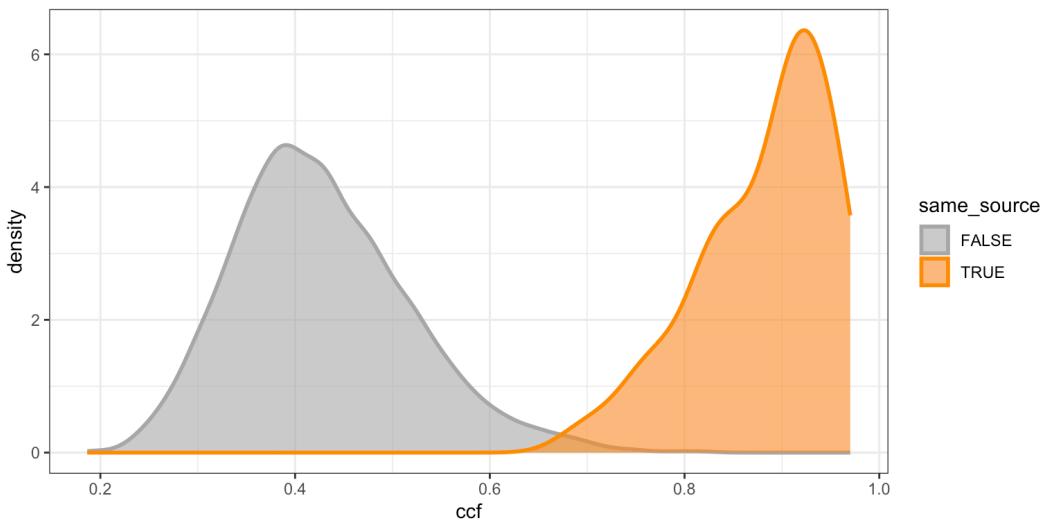


Figure 7. The density plot shows tails of distributions overlap, which can be used as a rough threshold for drawing conclusions.

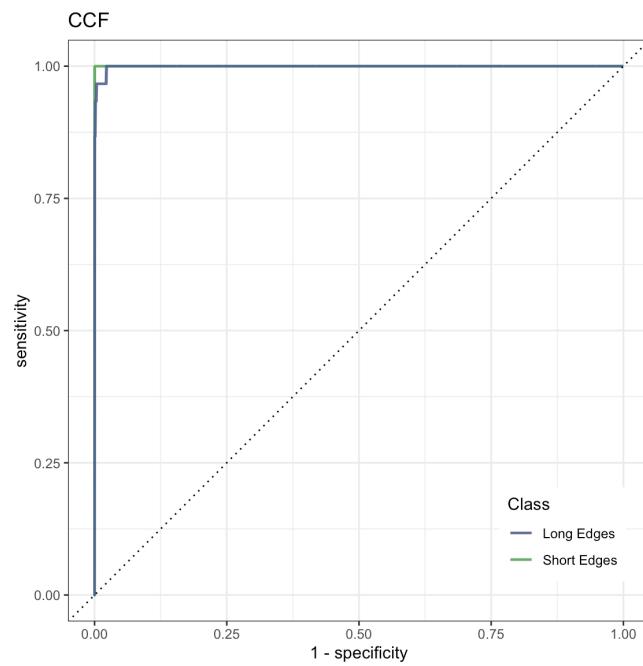


Figure 8. The ROC curve is bending very close to the upper left corner, which means excellent in classification and drawing conclusions.

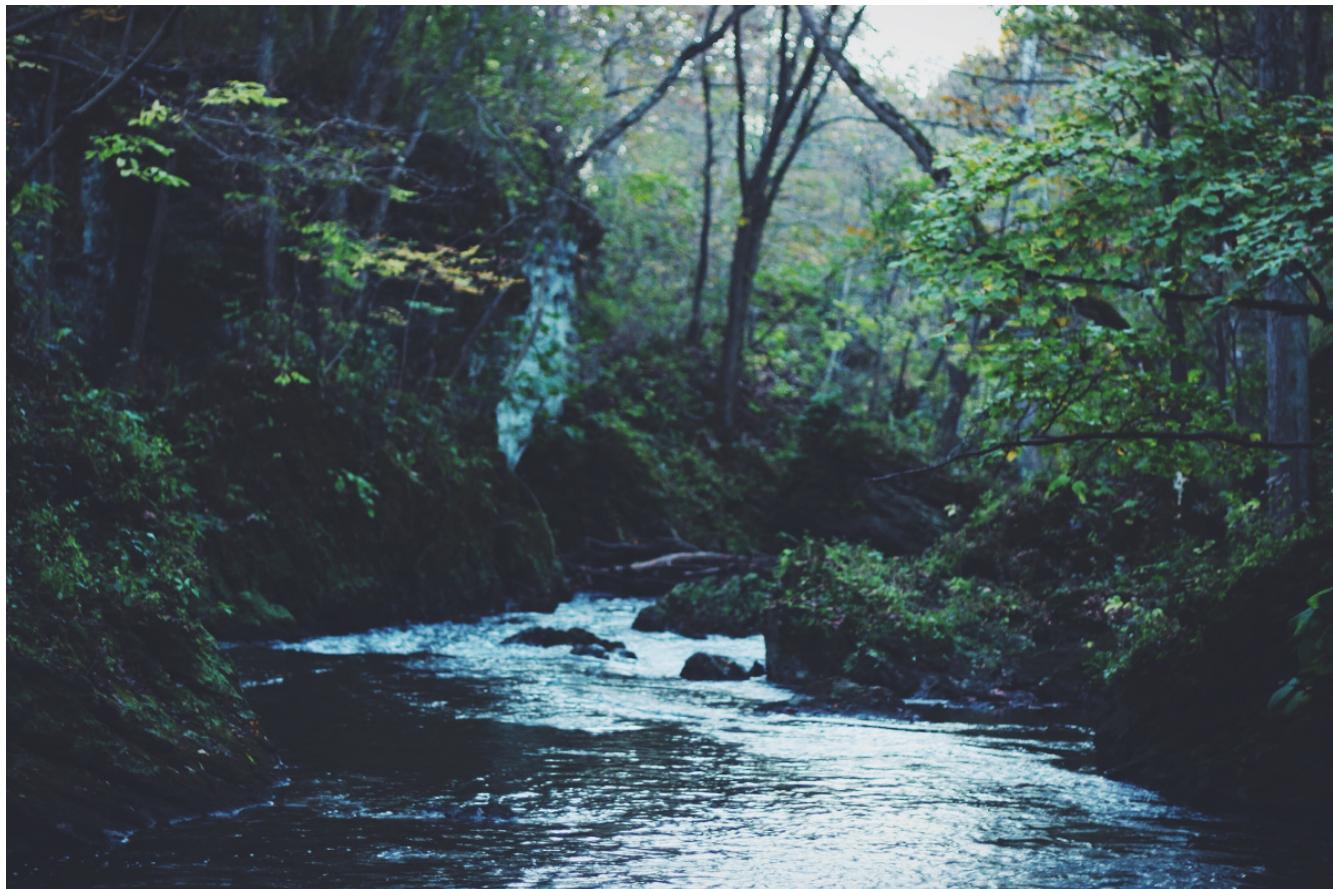


Figure 9. Legend (350 words max). Example legend text.