

# **Data Descriptor Title (110 character maximum, inc. spaces)**

**Yuhang Lin<sup>1,2</sup> and Heike Hofmann<sup>1,2</sup>**

<sup>1</sup>Iowa State University, Department of Statistics, Ames,

<sup>2</sup>Center 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**

(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**

(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 James Genrich's case being made of lead. The aluminum wire used was 16 Gauge/1.5 mm, anodized.

To cut the wire, 4-inch pieces were unspooled and cut using Kaiweets wire cutters, model KWS-105, as shown in Figure 1, 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. Both AB and CD sides form tent structures, as shown in Figure 2, and we can separate each side of the tent into 2 pieces, as shown in Figure 3, 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 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) 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 for the first repetition.

A more detailed procedure, including standard scanning protocols for the confocal microscope in Figure 4, can be found in the [README of the GitHub repository heike/Wirecuts](#) (High-res pics needed 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 format.

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

Authors should cite previous descriptions of the methods under use, but ideally the method descriptions should be complete enough for others to understand and reproduce the methods and processing steps without referring to associated publications.



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

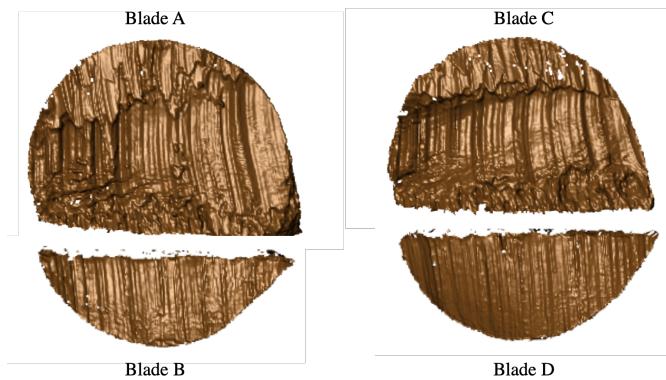


**Figure 2.** Legend (350 words max). Example legend text.

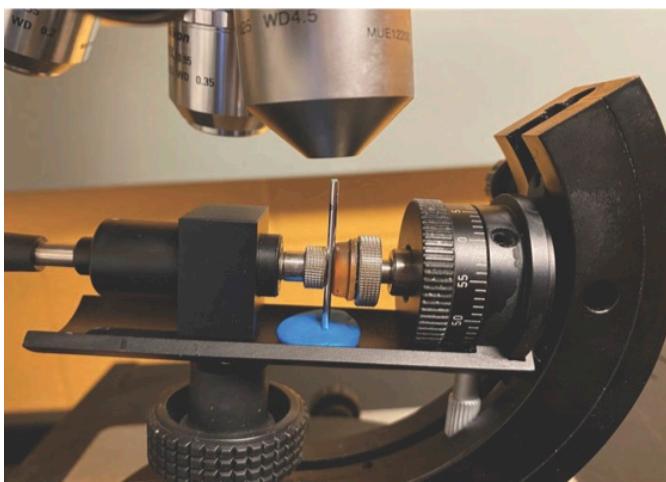
42 There is no limit to the length of the Methods section. Subheadings should not be numbered.  
43 Authors should review the transparent methods checklist below, and ensure that their manuscript complies with any relevant  
44 points. Authors are also encouraged to search FAIRsharing.org for community reporting standards that may be relevant  
45 to their specific data-type.

46 **Transparent Methods Checklist**

- 47 • Materials & reagents: Identify commercial suppliers of reagents, instrumentation or kits, when the source is critical to  
48 the outcome of the experiments. Declare any restrictions on the availability of unique materials (more information here).  
49 Provide catalogue or clone numbers for all antibodies (if available). For primary antibodies, provide proof of validation  
50 for the relevant species and applications.
- 51 • Exclusion criteria: If any data or samples were excluded, explain the exclusion criteria and state in the methods whether  
52 the criteria were established before the study was conducted.
- 53 • Randomization & blinding: For any studies that involve assigning samples, animals or participants into different groups:  
54 State clearly whether randomization methods were used. If randomization was not employed, this should be clearly  
55 stated. State clearly whether blinding was employed during data collection. If blinding was not employed, this should  
56 be clearly stated.
- 57 • Animal & human studies (full journal policies here): Experiments involving human participants must identify the  
58 committee approving the experiments, and include a statement confirming that informed consent was obtained from all  
59 participants. Studies employing nonhuman animals should ensure that methods descriptions comply with the ARRIVE  
60 checklist.
- 61 • Cell lines: For each eukaryotic cell line used, state the source and whether the cell line has been authenticated or  
62 otherwise tested for integrity. If any commonly misidentified cell lines were used (see ICLAC or NCBI Biosample),



**Figure 3.** Legend (350 words max). Example legend text.



**Figure 4.** Legend (350 words max). Example legend text.

justify their use. Report whether the cell lines were tested for mycoplasma contamination.

- Chemistry & materials science: Manuscripts describing chemical syntheses, or characterizing new chemicals or materials should refer to the guidance at Nature Chemistry.

## 66 Data Records

The complete data set is available on the ISU DataShare repository at <https://iastate.figshare.com/>, which is public and open access for every interested researcher. The data set consists of 120 scans in the x3p file format with the naming convention as described before. **Explain the x3p header info to some extent? x3ptools?**

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

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

## 79 Technical Validation

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

84 We also conduct numerical comparisons between 2 replicates made by the same side of the same tool at the same location.  
85 For example, we compare T1AW-LI-R1 to T1AW-LI-R2, T1BW-LI-R1 to T1BW-LI-R2, and so on, as shown in Figure 5.  
86 The comparison is done by aligning the two scans and the cross-correlation function (CCF) between the two scans. Explain  
87 the manual extraction in details? explain striations, profiles, signals, etc? We know that signals from two replicates with the  
88 same tool-edge-location combination should yield similar signals, which results in high CCF values close to 1, and the values  
89 we got in Figure 6 achieve our expectation. For validation of all other pairs of scan replicates, see the detailed report.

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

- 93 • Measurement of data quality?

- 94     – Numeric measurements / tests: ?

- 95     – Visualizations: ?

- 96     – Check with existing data: ?

- 97     – Questionable / slur procedures:

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

- 102         \* [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  
103        manually counted the number of lesions with contrast-enhancement, necrosis, and peritumoral edema for  
104        each patient.

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

- 112         \* [3D surgical instrument collection for computer vision and extended reality](#): The main issue...Since we  
113        store our models in a standard format (STL), they are compatible with a large variety of visualisation and  
114        processing software.

- 115         \* [Three-dimensional reconstruction of high latitude bamboo coral via X-ray microfocus Computed Tomography](#) : Regular quality assurance inspections are carried out on the  $\mu$ -CT scanner to verify its metrological and  
116        geometrical (alignments) accuracy for conducting the scans. The geometry of source to object and source  
117        to detector distances are verified whenever there is any significant physical interaction with the source such  
118        as re-alignment, change of filament, or source anode change. This calibration process involves scanning  
119        a specially designed phantom known as an hourglass36, which consists of three pairs of high-sphericity  
120        spheres. The sphere sizes are as follows: two spheres with a diameter of 3.000 mm, two spheres with a di-  
121        ameter of 6.000 mm, and two spheres with a diameter of 9.525 mm, and each sphere is kept in contact with  
122        its size-counterpart. By using this phantom, it becomes possible to accurately determine a known distance,  
123        specifically the centre-to-centre distance of the spheres, in a threshold-independent manner. If the measured  
124        distance deviates beyond the acceptable limits of metrological accuracy, the systems calibration parameters  
125        are adjusted to ensure agreement between the measured distance and the actual distance.

## 127 Usage Notes

128 The data set can be easily accessed with the CRAN R package `x3ptools`. Further analysis can be conducted with the  
129 GitHub R package `wire` and the GitHub R shiny app `wireShiny`. We already conduct between-replicate comparisons  
130 in the technical validation section, and we can also conduct across-replicate comparisons to establish error rates threshold and  
131 produce other analysis plots. The resulting CCFs are shown in the boxplot in Figure 7. We can see that the CCF values for  
132 the same sources are close to 1, while the CCF values for different sources are much lower. This difference can be used to  
133 establish a threshold for CCF and help us draw conclusion about the similarity between wire cut scans, which can be used  
134 in real crime scenes. The density plot in Figure 8 shows the distribution of the CCF values. The overlapping points between  
135 the two distributions can be a rough threshold. The receiver operating characteristic (ROC) curve in Figure 9 shows the true  
136 positive rate against the false positive rate, gives us a true threshold of 0.589 to control the false positive rate (FPR) to be less  
137 than 0.05, and 0.658 to control the FPR to be less than 0.01.

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

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

## 146 Code availability

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

## 151 1 End of Body

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

## 159 Acknowledgements

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

## 162 Author contributions statement

163 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:  
164 Eden Amin, Curtis Mosher, Jeff Salyards. Must include all authors, identified by initials, for example: A.A. conceived the  
165 experiment(s), A.A. and B.A. conducted the experiment(s), C.A. and D.A. analysed the results. All authors reviewed the  
166 manuscript.

## 167 Competing interests

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

174 **Figures & Tables**

175 Figures, tables, and their legends, should be included at the end of the document. Figures and tables can be referenced in  
176 L<sup>A</sup>T<sub>E</sub>X using the ref command, e.g. Figure 10 and Table 1.

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

181 Tables may be provided within the L<sup>A</sup>T<sub>E</sub>X document or as separate files (tab-delimited text or Excel files). Legends, where  
182 needed, should be included here. Generally, a Data Descriptor should have fewer than ten Tables, but more may be allowed  
183 when needed. Tables may be of any size, but only Tables which fit onto a single printed page will be included in the PDF  
184 version of the article (up to a maximum of three).

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

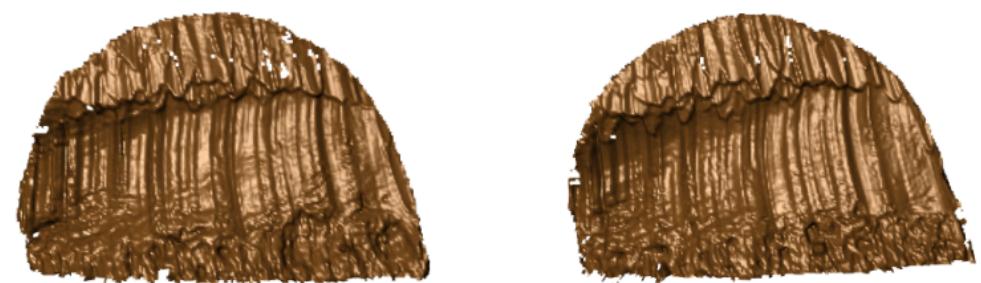
Condition	n	p
A	5	0.1
B	10	0.01

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

Edge A



Edge C



Edge B



Edge D



**Figure 5.** Legend (350 words max). Example legend text.

### Edge A



### Edge C



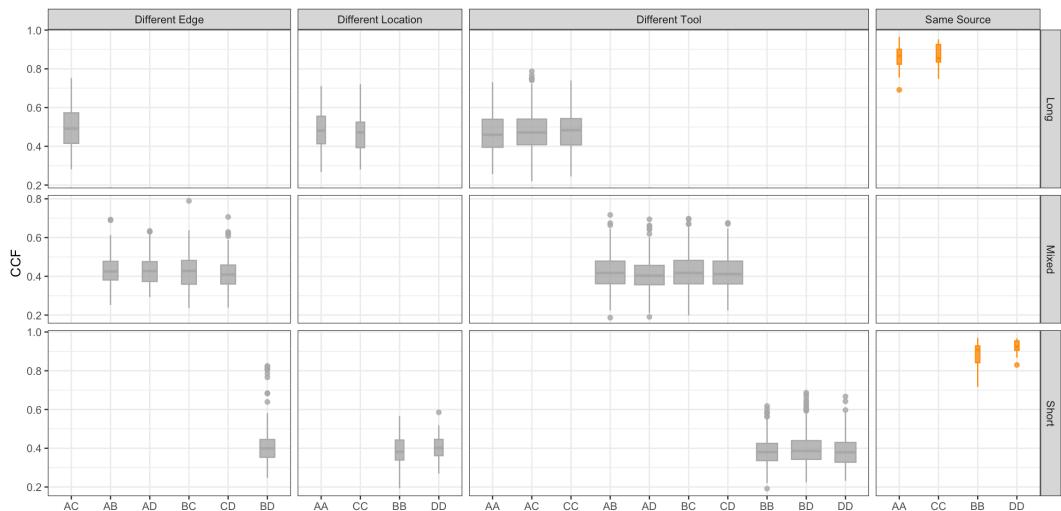
### Edge B



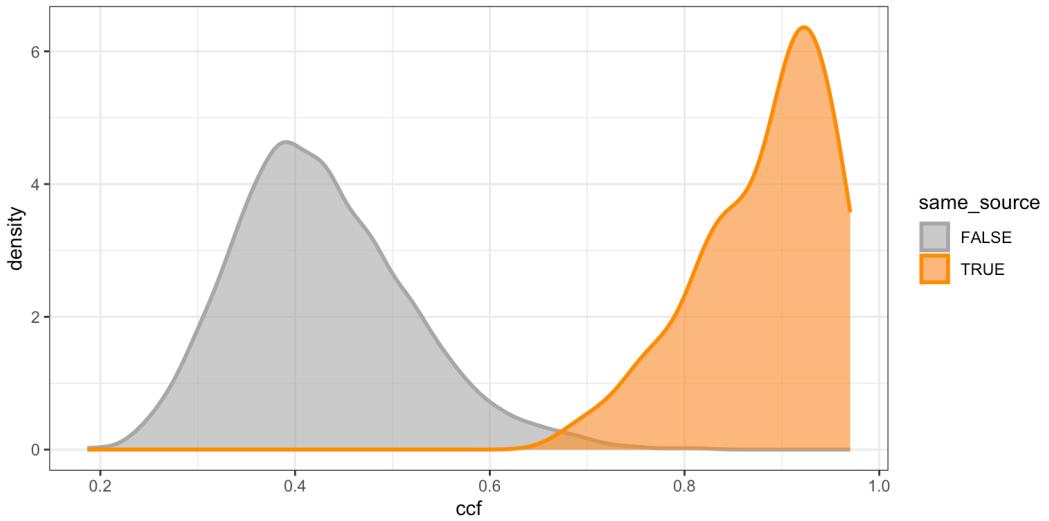
### Edge D



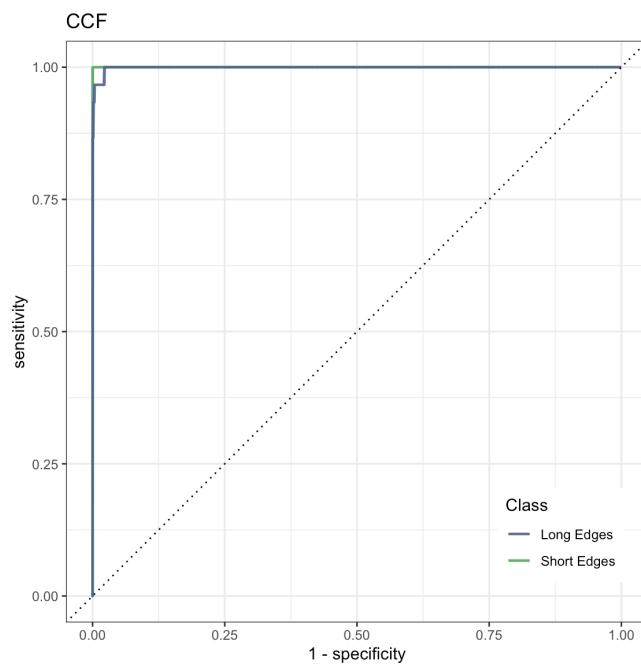
**Figure 6.** Legend (350 words max). Example legend text.



**Figure 7.** Legend (350 words max). Example legend text.



**Figure 8.** Legend (350 words max). Example legend text.



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



**Figure 10.** Legend (350 words max). Example legend text.