

1 **Disclaimer:**

2
3 As a condition to the use of this document and the information contained herein, the
4 Facial Identification Scientific Working Group (FISWG) requests notification by e-mail
5 before or contemporaneously to the introduction of this document, or any portion
6 thereof, as a marked exhibit offered for or moved into evidence in any judicial,
7 administrative, legislative, or adjudicatory hearing or other proceeding (including
8 discovery proceedings) in the United States or any foreign country. Such notification
9 shall include: 1) the formal name of the proceeding, including docket number or similar
10 identifier; 2) the name and location of the body conducting the hearing or proceeding;
11 and 3) the name, mailing address (if available) and contact information of the party
12 offering or moving the document into evidence. Subsequent to the use of this document
13 in a formal proceeding, it is requested that FISWG be notified as to its use and the
14 outcome of the proceeding. Notifications should be sent to: chair@fiswg.org
15

16 **Redistribution Policy:**

17
18 FISWG grants permission for redistribution and use of all publicly posted documents
19 created by FISWG, provided that the following conditions are met:
20
21 Redistributions of documents, or parts of documents, must retain the FISWG cover
22 page containing the disclaimer.
23
24 Neither the name of FISWG, nor the names of its contributors, may be used to endorse
25 or promote products derived from its documents.
26
27 Any reference or quote from a FISWG document must include the version number (or
28 creation date) of the document and mention if the document is in a draft status.
29



Understanding and Testing for Face Recognition Systems Operation Assurance: Identity Ground Truth

30 Purpose

31 This document provides guidelines and techniques to help administrators of
32 automated face recognition systems (FRS) produce advanced and accurate recognition
33 statistics on the face recognition systems.

34 The intended audience of this document is system owners, system users, and system
35 administrators of existing automated face recognition systems. Outside the scope of this
36 document include, but not necessarily limited to, system setup, system tuning, workflow
37 management and improvement, and proof of concept pilots.

38 This document is a follow on from the FISWG document: “Understanding and Testing
39 for Face Recognition Systems Operation Assurance” (version 1.0, 2020.12.11)

40 The issues presented in this document form a base for other considerations and
41 advanced topics when testing (e.g., system setup and tuning) which will be covered in
42 future FISWG documents.

43 **Scope**

44 The scope of this document is to provide a detailed process and examples of testing
45 and repairing identity ground truth in facial data sets which is a critical initial step before
46 recognition statistics are created and reviewed. This document does not address facial
47 accuracy but is focused solely on testing and correcting identity ground truth. Any facial
48 algorithm can be used with these processes. It is assumed that all facial images create
49 proper templates. The facial data set used in this document is the “Labeled Faces in
50 the Wild” (LFW) but conceptually any other facial data set with identity ground truth can
51 be used.

52 **Background**

53 When doing accuracy profiling, there is always one key aspect which must be
54 addressed first: what is the identity ground truth within the images? All data sets will
55 potentially have some corruption in the identity ground truth with the data. Detecting
56 and correcting this so pristine results can be reviewed is always a critical portion of
57 profiling.

58 This type of identity ground truth verification will potentially exist in all data sets no
59 matter where the data sets come from. This is an iterative process as the agency
60 learns the algorithms, the data, and how the two interact with each other. If proper care
61 is not given in these early stages, then incorrect assumptions on the outcomes will be
62 made. It’s critical to understand this process with an investigative mindset before the

63 agency gets to the operational data sets which may have identity corruption and image
64 quality issues that may be large but not uncommon in operational deployments. If the
65 agency gets to the operational data set without a firm awareness and knowledge base
66 on the how the core algorithms work with verified data, then the agency could be
67 incorrectly assessing and measuring performance of the FRS. Agencies need to lay the
68 groundwork to know and trust the algorithms before they get to the possibly unmanaged
69 and unknown operational data.

70 Most of the work in these processes is on creating the testing frameworks and
71 understanding how to repeatedly run tests, make corrections, and do retesting with what
72 has been learned. Once the frameworks and the processing are understood, then the
73 agency can make diligent progress, but it takes time and focus. The outcomes are
74 worth the time spent as you begin to understand how the data interacts with the
75 algorithms which give the agency the ability to trust the solution and not just assume the
76 data is invalid.

77 Setting up frameworks to do enrollment and searching while recording results is fairly
78 mechanical as you learn the facial algorithms and the data sets to develop proper
79 profiling. Understanding the data and building frameworks to analytically qualify the
80 results is not trivial but must be done so effective operational metrics can be derived
81 and applied.

82 Before doing this analysis on operational data, it is recommended that the agency
83 develop and test the framework on experimental datasets. After some experience is
84 gained in this process and confidence that the process is correct, one could then assess
85 the operational dataset.

86 This document describes procedures to assess an experimental dataset, which can
87 be replicated before assessing operational datasets.

88 **Data Set**

89 Care should be taken in selecting data sets to profile experimentally. It is
90 recommended to select data sets which:

- 91 • Have operational relevancy
- 92 • Have consistent image quality aspects: type of capture, size of images, subject
93 poses, etc.
- 94 • Have sufficient identities and images to test with. This decision will be agency
95 specific.
- 96 • Includes associated identity ground truth information which links each image to a
97 unique identity

98 The data set used for this document is the LFW data set available at: [http://vis-
www.cs.umass.edu/lfw/](http://vis-
99 www.cs.umass.edu/lfw/) See Appendix 2: “**LFW Data Set Information**“ for more
100 details. Conceptually any other facial data set with identity ground truth can be used.

101 LFW is a widely used open source data set which will work well for this specific
102 document serving as an introductory data set. Information on this data set includes:

- 103 • Has smaller but consistent image sizes and file formats
- 104 • Has over 5,700 identities and over 13000 images
- 105 • Has a wide range of subjects: sex, pose, lighting, etc
- 106 • Has stated identity ground truth errors

107 A key point in the LFW errata is that there are known errors in the LFW data set.
108 While the LFW URL addresses these errors, this document will show how to locate
109 them and give examples on how determining identity ground truths in relating images to
110 identities is critical and needs to be addressed in any operational testing scenario.

111 **Ground Truth Process**

112 **Step 1:** Enroll the facial images into a facial gallery for searching.

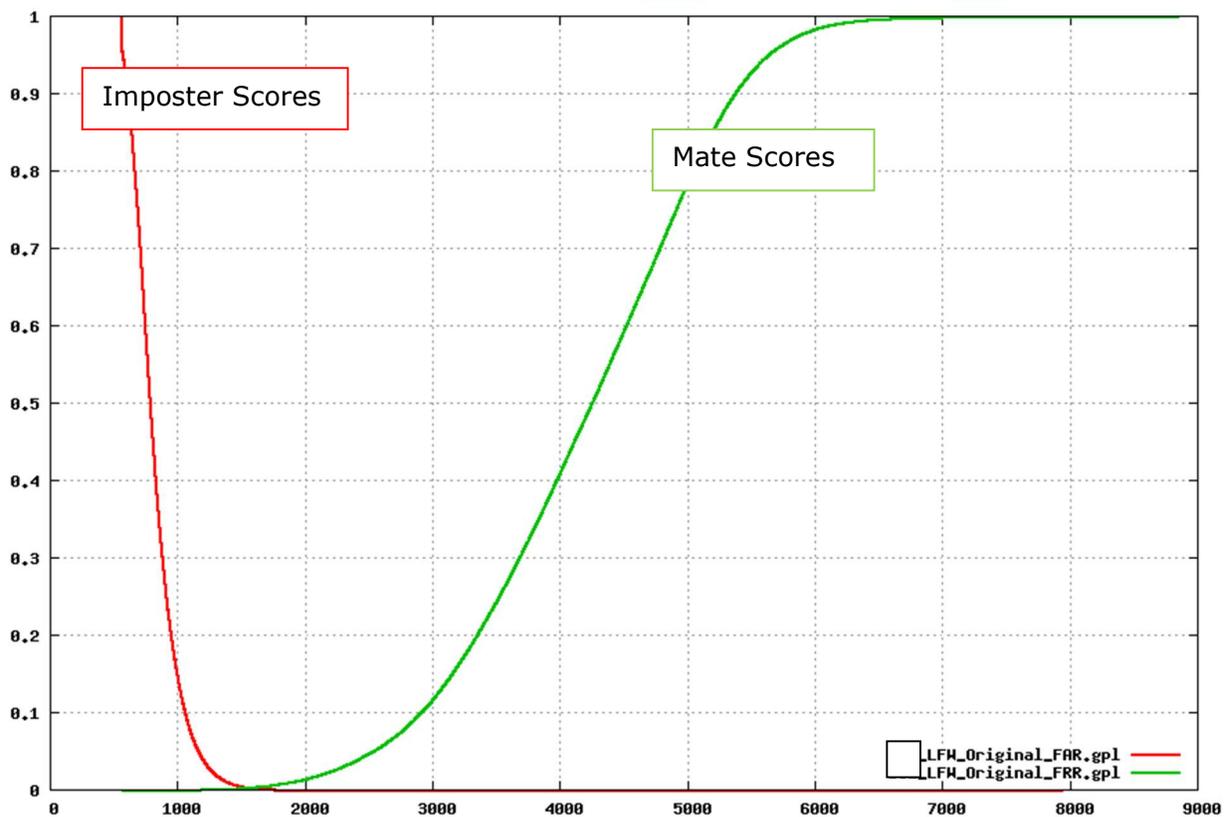
113 **Step 2:** Search the facial images against the facial gallery. The number of
114 candidates returned for this document was 100. This number may vary with agency
115 specifics and the biometric algorithm deployed. It is recommended to test with a larger
116 number of candidates than what may be operationally used so that deeper accuracy
117 investigations can be analyzed. Do not use scoring thresholds.

118 **Step 3:** Analyze the scoring to delineate every candidate in all 1:N searches:

- 119 • Image file name

- 120 • Image identity
- 121 • Score
- 122 • Rank
- 123 • Mate scoring
- 124 • Imposter scoring

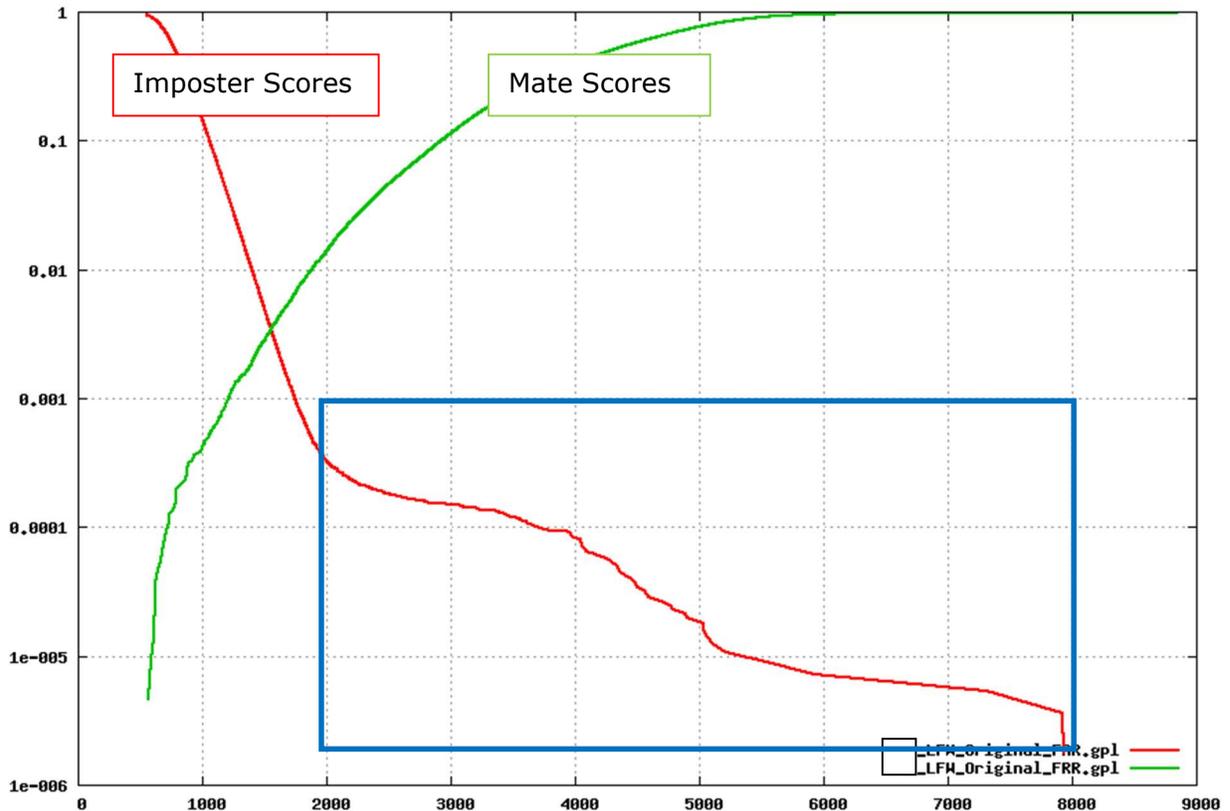
125 When this was done for this document the following accuracy curves were obtained:



126
127

FIG. 1 FAR and FRR from uncorrected imagery

128 In Figure 1 the FAR and FRR scores are presented on a linear Y-axis. This can tend
129 to hide the identity errors.



130
131 **FIG. 2 FAR and FRR from uncorrected imagery**

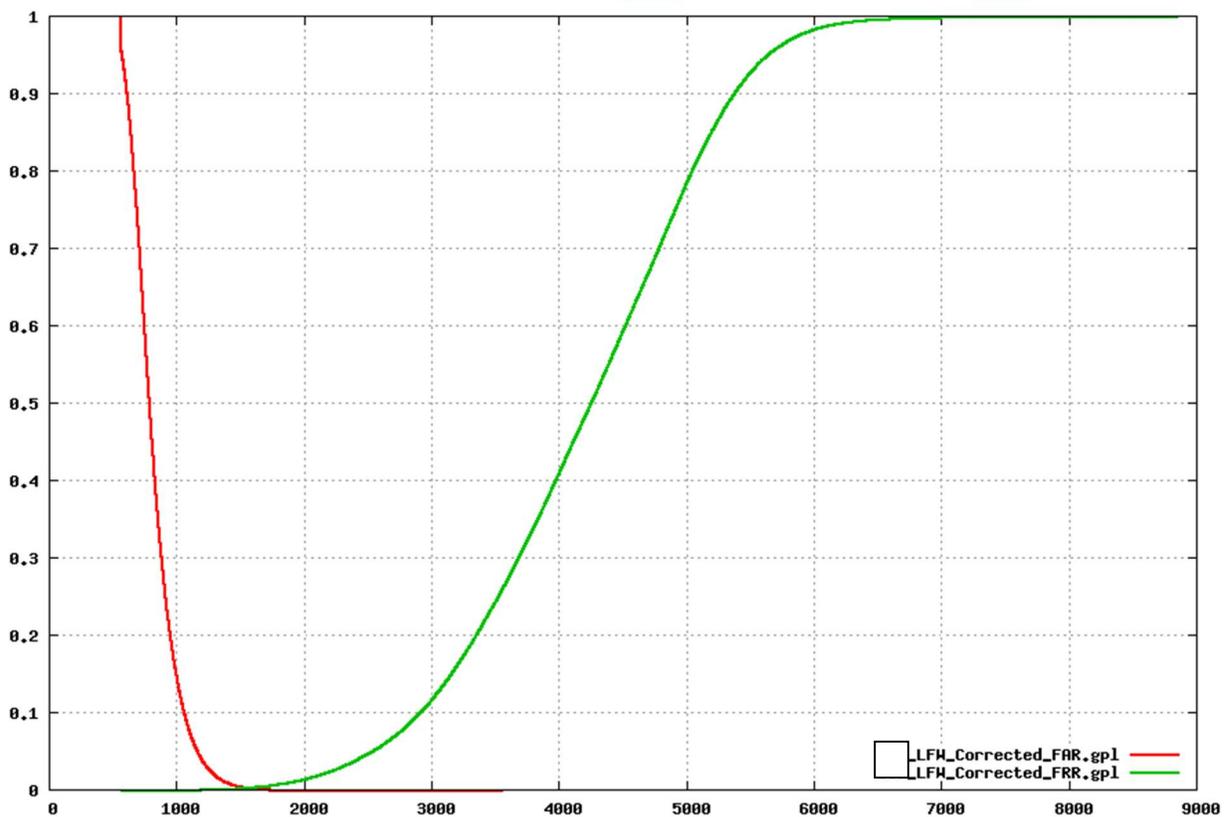
132 In Figure 2 the FAR and FRR scores are presented on a logarithmic Y-axis. This
133 shows that there may be identity errors in the high scoring FAR scores (blue area).
134 The sudden increase in FAR scores indicates potential mates that are incorrectly
135 labeled imposters.

136 Other methods to select incorrectly labeled imposters are publicly available. See
137 Appendix 1: “**Alternative Methods**”.

138 **Step 4:** Analyze the scoring to resolve high scoring imposters to see if identity errors
139 do exist in the data.

140 **Step 5:** Iterate between Steps 1-4 as long as identity errors are suspected. The
 141 expectation is that several passes will be needed to achieve an objective measure of
 142 correct ground truth.

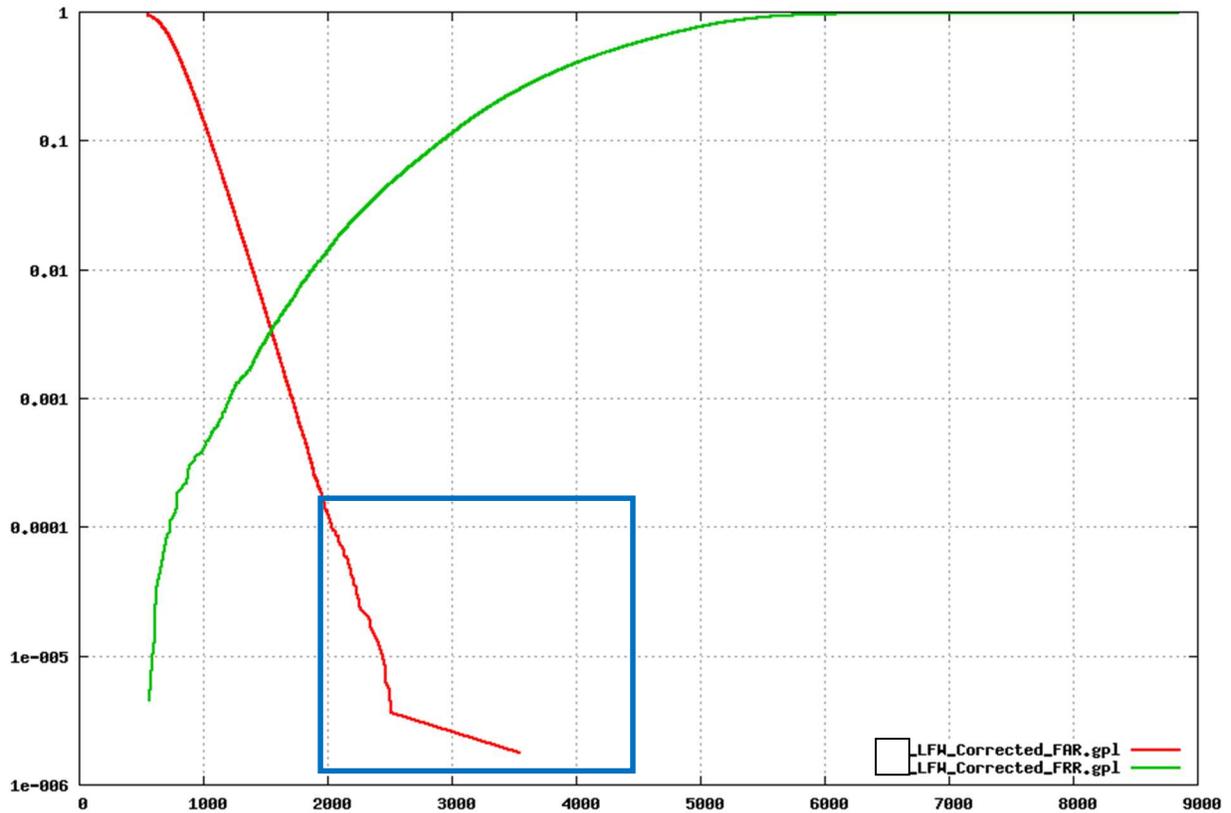
143 NOTE: In the LFW data set several iterations were done and approximately 90
 144 identities were modified. After these corrections were made the following accuracy
 145 plots were derived.



146
 147

FIG. 3 FAR and FRR from corrected imagery

148 In Figure 3 the FAR and FRR scores are presented on a linear Y-axis. This can tend
 149 to hide the identity errors.



150
151

FIG. 4 FAR and FRR from uncorrected imagery

152 In Figure 4 the FAR and FRR scores are presented on a logarithmic Y-axis. This
153 shows there may be several more identity errors in the FAR scoring (blue area).

154

155 One can find the following images by investigating the high score imposters.



156
157



The images above are siblings



The images above are close appearances



The images above are close appearances

159 As these high scoring imposters were investigated, it became apparent that the high
 160 scores in the corrected FAR were heavily influenced by:

- 161 • Twins
- 162 • Siblings
- 163 • What are referred to as “doppelgangers”

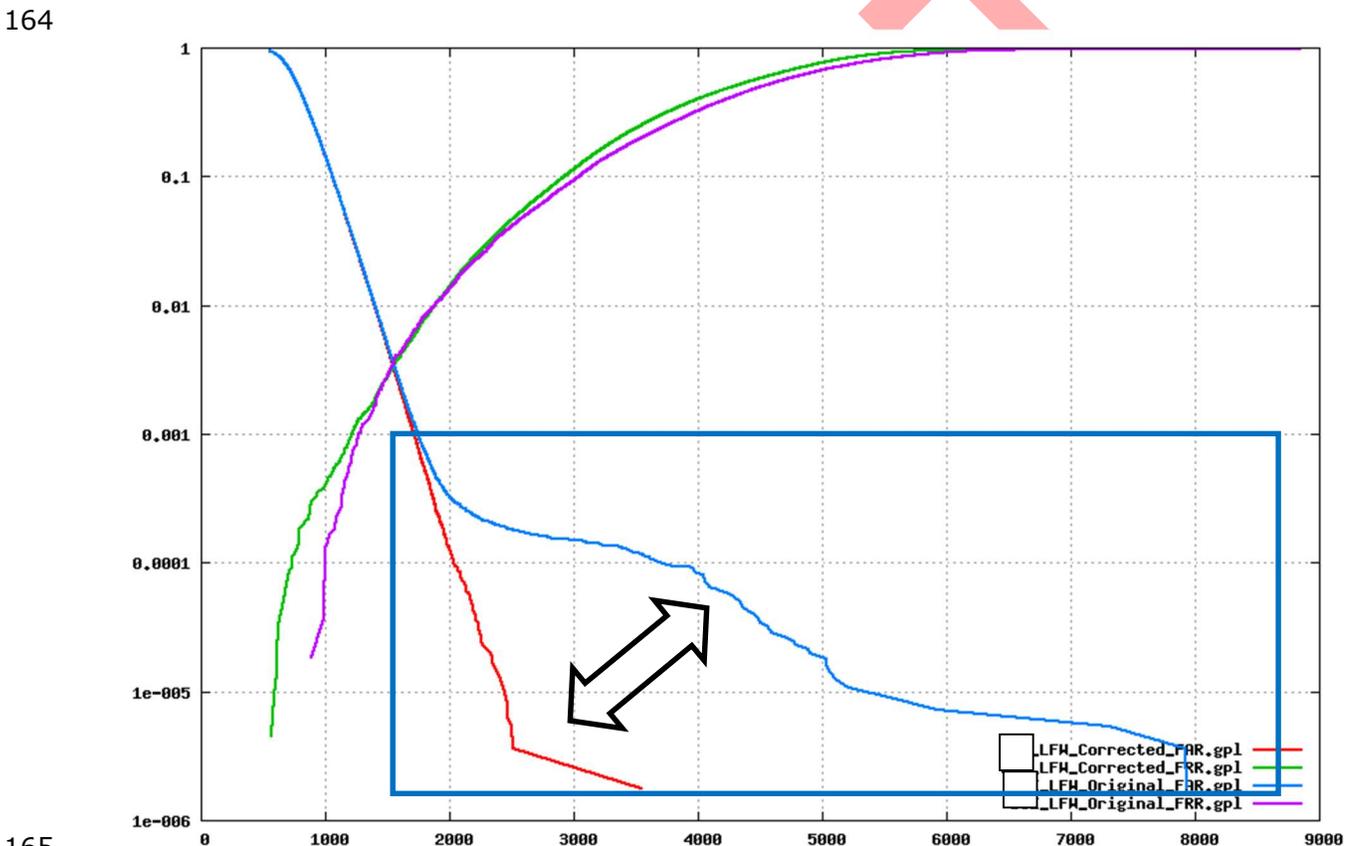


FIG. 5 FAR and FRR from original and corrected LFW imagery

167 In Figure 5 the FAR and FRR from original vs. corrected imagery are presented on a
 168 logarithmic Y-axis. This shows the scoring variances after the identity errors were
 169 corrected (blue area). This shows how a few improper identities can dramatically affect
 170 the FAR and FRR scoring profiles.

171

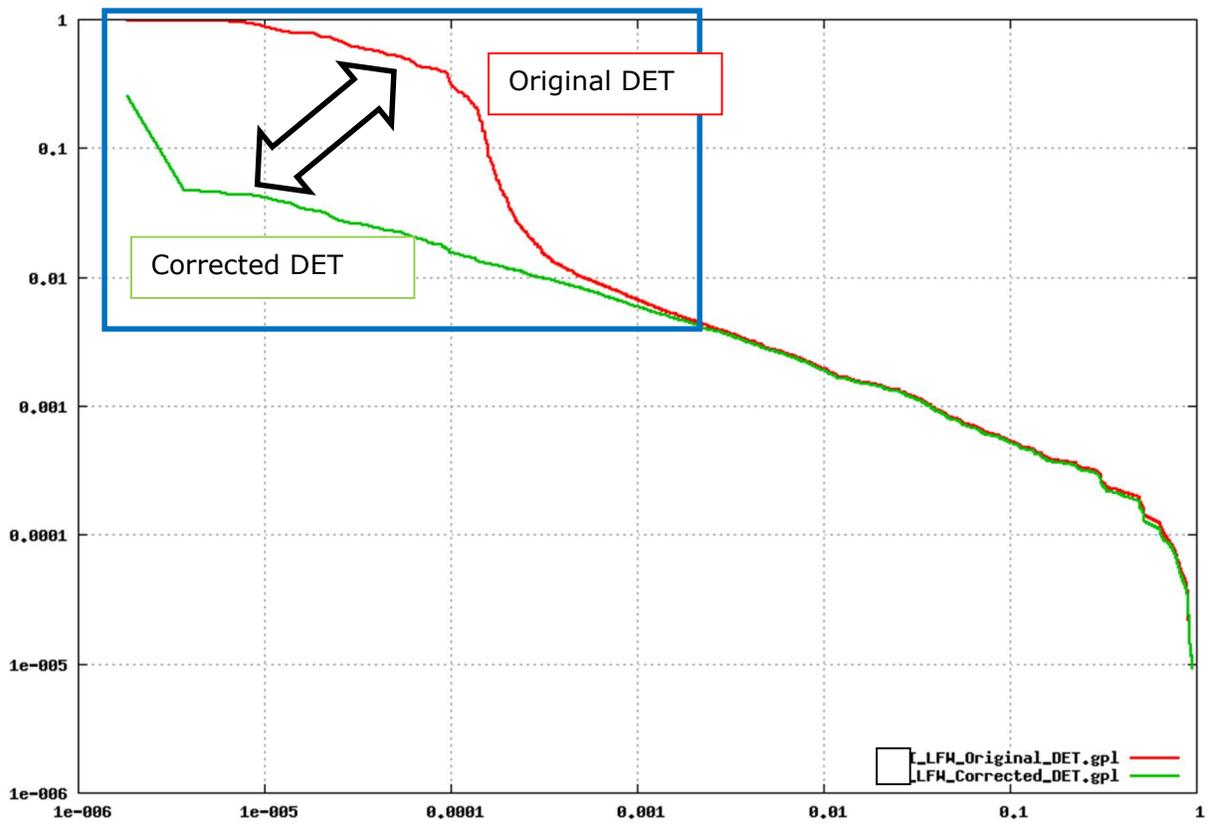
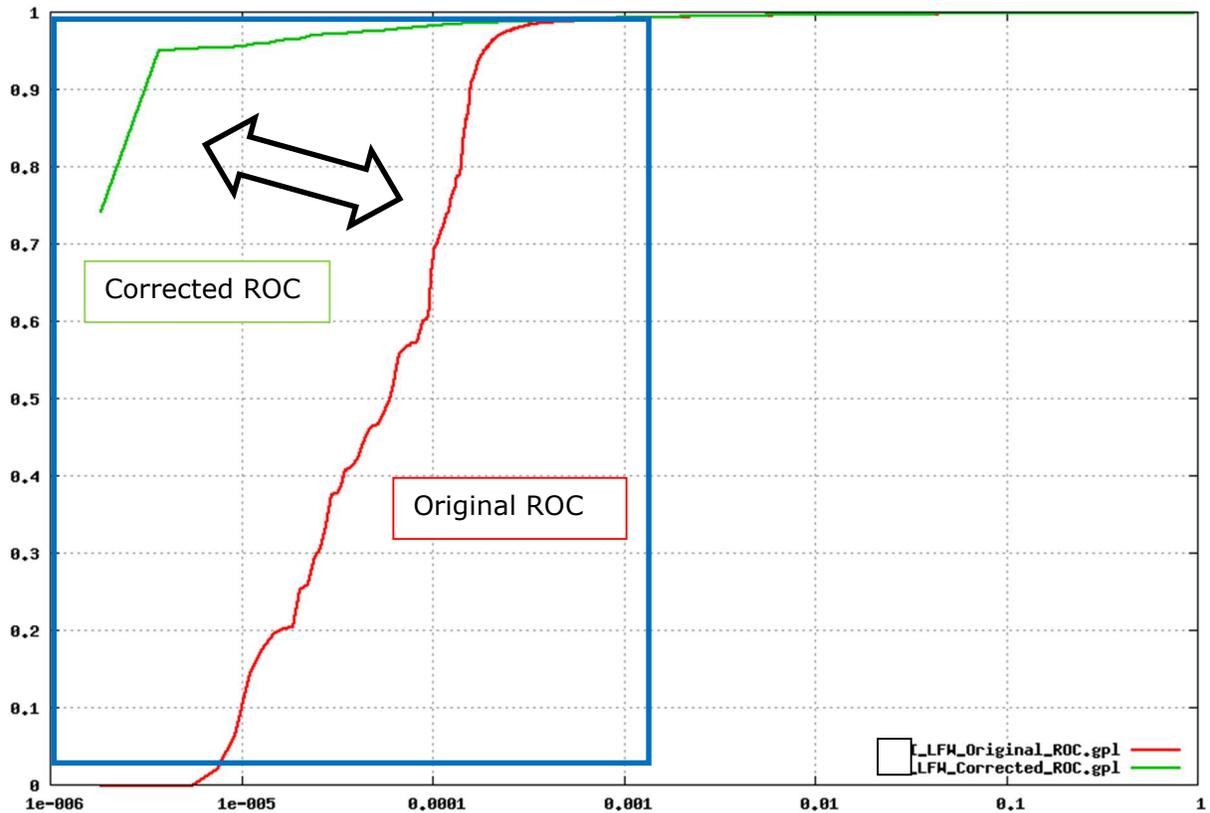
172
173

FIG 6 DET comparison between original and corrected LFW imagery

174 In Figure 6 the DET curve from original vs. corrected imagery are presented. This
 175 shows the scoring variances after the identity errors were corrected (blue area). This
 176 again shows how a few improper identities can dramatically affect the FAR and FRR
 177 scoring profiles.



178
179

FIG. 7 ROC comparison between original and corrected imagery

180 In Figure 7 the ROC curve from original vs. corrected imagery are presented with a
181 logarithmic X-axis. This shows the scoring variances after the identity errors were
182 corrected (blue area).

183

184 Investigating low scoring mates - the following images were confirmed identities but
185 scored very low.





186 Outcomes

187 Based on this data set and the testing process documented here:

- 188 • The LFW data set tested had identity errors which need to be adjusted to ensure
- 189 proper scoring analysis could be done. From the default data set downloaded ~90
- 190 identities should have been corrected.
- 191 • Iterative processes to properly analyze and modify this data set were required
- 192 which focused on using the FAR/FRR curves to locate the identity errors.

- 193 • Correcting the identity errors improved the FAR/FRR scoring analysis
- 194 • Twins, siblings, and doppelgangers did affect the scoring analysis. Two twins
- 195 should have been located with several siblings and doppelgangers causing the highest
- 196 FAR scores.
- 197 • FAR, FRR, DET and CMC curves were utilized in these processes.
- 198 • Critical scoring analysis required the presentation of the scoring in both linear and
- 199 logarithmic presentations to see the imposter scores which had identity errors.
- 200 • A variety of methods can be used to resolve identity errors in facial datasets. FRS
- 201 administrators should be aware of the advantages and disadvantages of each method
- 202 before selecting and applying a method, especially on operational datasets.

203 **Glossary**

- 204 • FR: face recognition
- 205 • FRS: face recognition systems
- 206 • ROC: Receiver Operating Characteristics
- 207 • DET: Detection error tradeoff
- 208 • FAR (False Accept Rate): the measure of the probability that the biometric system
- 209 will incorrectly accept an access attempt by an unauthorized user
- 210 • FRR (False Reject Rate): the measure of the probability that the biometric system
- 211 will incorrectly reject an access attempt by an authorized user
- 212 • Doppelganger: an apparition or double of a living person.

213 **Appendix 1: Alternative Methods**

214 **General/Initial assessment**

215 The initial assessment objectives are twofold: first, to acquire a general sense of how
216 the FRS interacts with the data, and second, to select the identities or images for
217 detailed review.

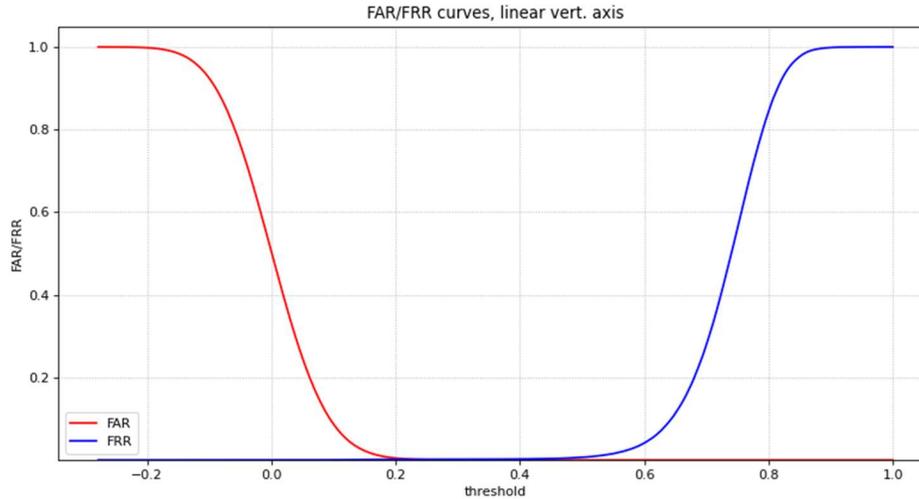
218 The initial assessment assumes that similarity scores, both genuine and impostor,
219 are distributed according to predictable and smooth monomodal curves and that errors
220 in the dataset are expected to introduce "unnatural" multimodal variations in the curves
221 obtained from the sets of scores.

222 Because the specific shape of this curve (e.g., Gaussian, exponential, log-normal)
223 varies from system to system and also depends on the dataset itself, it is important to
224 develop a general sense of the interaction between the FRS and the data.

225 It is important to note that the effectiveness of these procedures, which are ultimately
226 based on the distribution of scores, is fundamentally dependent on the accuracy of
227 algorithm used to generate these scores. Algorithms whose scores are highly
228 discriminative are more effective in identifying outliers.

229 An initial assessment could then be based on curves derived from genuine and
230 impostor scores distributions as, e.g., the FAR/FRR curve, or directly inspecting both
231 distributions of scores.

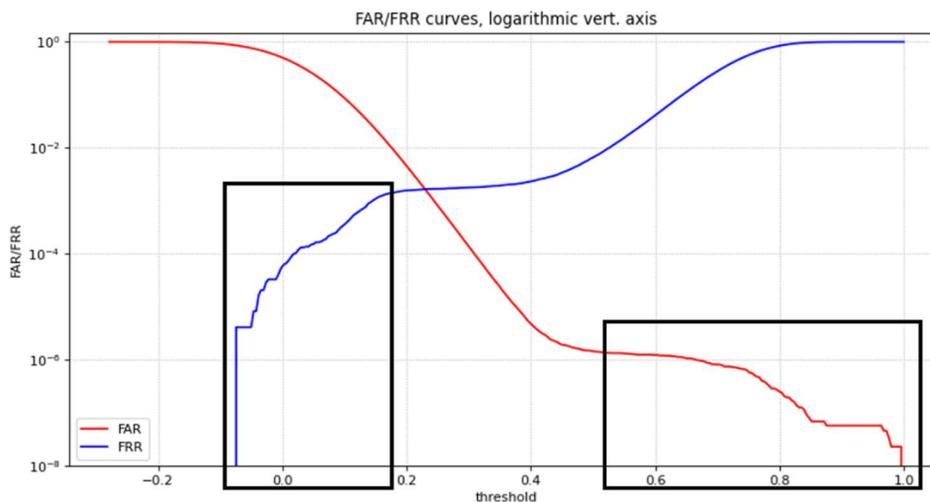
232 FAR/FRR curve plots both of these rates as a function of threshold. An example of
 233 such a curve for the LFW dataset is shown below.



234
 235

FIG. A1.1 FAR and FRR from uncorrected imagery

236 In this plot, the vertical axis is linear, but plotting the same data with the vertical axis
 237 in a logarithmic scale facilitates the inspection of the curves in the lowest values.

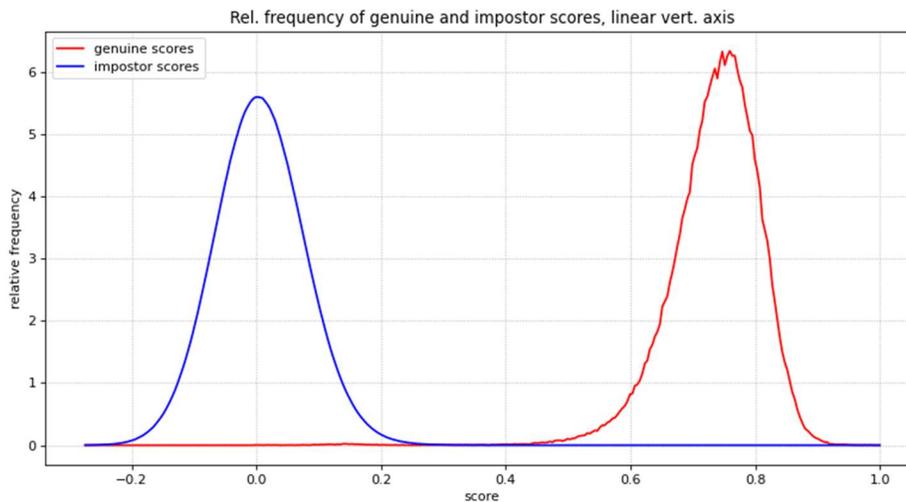


238
 239

FIG. A1.2 FAR and FRR from uncorrected imagery (Logr axis)

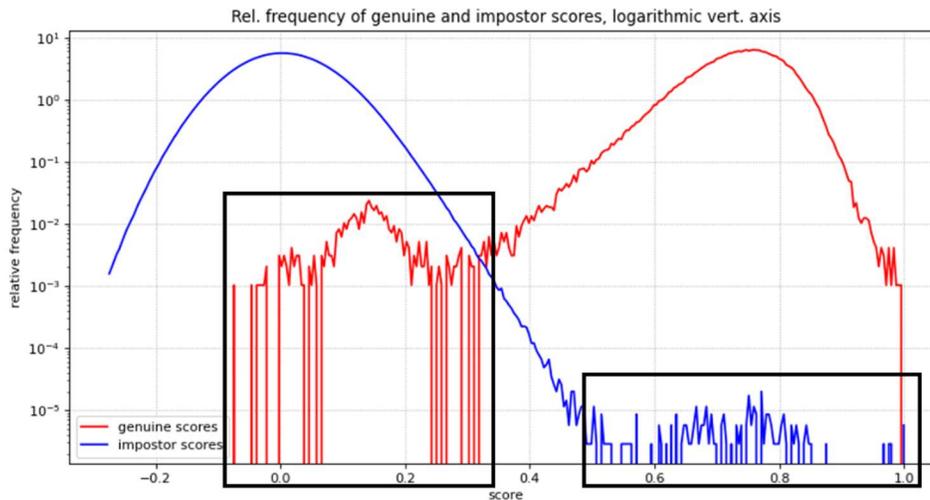
240
 241 In this plot, it is now possible to observe variations on the shape of the FAR and FRR
 242 curves that violate the assumption of smoothness of a monomodal distribution curve
 243 (framed portions). This raises suspicions that errors exist in the dataset, resulting in
 244 more impostor scores with high values and more genuine scores with low values.

245 This assessment can be further explored using the distributions of scores directly.
 246 The plot below is based on the same data.



247
 248 **FIG. A1.3 Genuine and Imposter Scoring Histograms uncorrected imagery**

249 The vertical axis's linear scale again makes it challenging to inspect the lower values,
 250 but using a logarithmic vertical axis allows one to observe a violation of the assumption
 251 of smoothness of a monomodal distribution curve (framed portions).



252
253

FIG. A1.4 Genuine and Impostor Scoring Histograms uncorrected imagery (Logr axis)

254 Both the higher impostor scores and the lower genuine scores could be selected from
255 this plot to a detailed review of the images involved. Impostor scores should be selected
256 for further inspection in highest to lowest order, and the opposite for genuine scores.

257 Apart from the visual assessment of the curves derived from the distribution of
258 scores, statistical methods for identifying outliers in each set of scores can be used.

259 One such method involves selecting outliers based on multiples of the median
260 absolute deviation from median (MAD) of each set of scores. In univariate statistics, the
261 MAD is a robust dispersion measure in the presence of outliers, contrary to the more
262 commonly used standard deviation around the mean. [1]

¹ Leys, C., et al., Detecting outliers: Do not use standard deviation around the mean, use absolute deviation around the median, *Journal of Experimental Social Psychology* (2013), <http://dx.doi.org/10.1016/j.jesp.2013.03.013>,

263 In practical terms, for genuine scores, those scores that lie below a multiple of MAD
264 from the median should be considered as outliers and then manually verified. For
265 impostor scores, the outliers are those higher than a multiple of MAD from the median.
266 Larger multiples of MAD will result in smaller sets of outliers to be manually verified,
267 while smaller multiples of MAD will result in a larger number of scores being considered
268 as outliers.

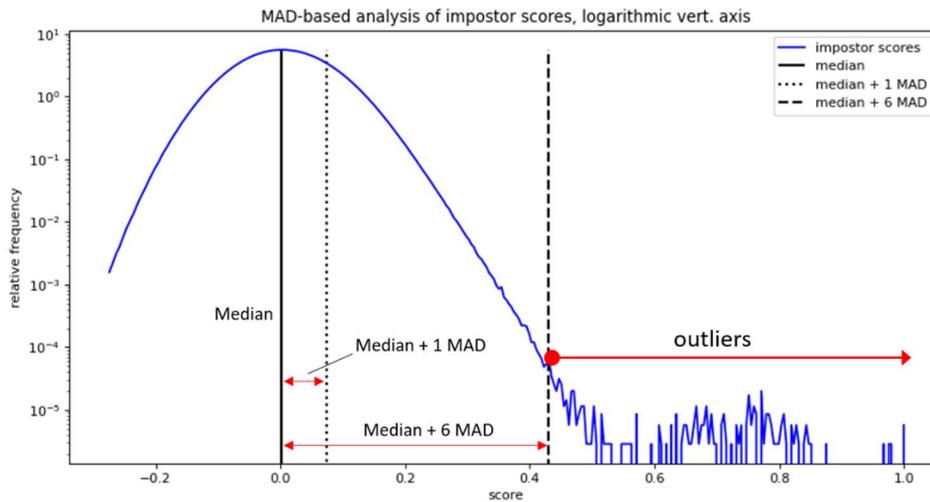
269 It is essential to notice that both the MAD and median statistics should be calculated
270 independently for each set of scores, genuine and impostors.

271 The specific multiple of MAD for each set should be determined experimentally, on a
272 case by case basis. For the worked example, it was determined through a controlled
273 experiment with the algorithm and the LFW dataset.

274 For this specific purpose, controlled errors were introduced in the dataset, which
275 added 764 errors in the genuine scores and the same quantity in the impostor scores.
276 These experimental sets were verified using the criteria of multiples of MAD, which
277 resulted in the number of six MADs as being adequate to identify the controlled errors
278 as outliers in both genuine and impostor scores.

279 The remainder of the analysis was conducted in the original LFW dataset, without the
280 artificially introduced errors.

281 The figure below illustrates the procedure for selecting outliers in the impostor scores,
282 in the original LFW dataset, after the appropriate multiple of MAD was determined. The
283 procedure for genuine scores is analogous.



284
285

FIG. A1.5 MAD Based Imposter Scores uncorrected imagery (Logr axis)

286 In the case of the original LFW dataset, without prior corrections of the known errors,
287 the number of six MADs for the genuine and impostor scores sets resulted in the
288 selection of 474 and 220 image pairs as outliers in each set, respectively.

289 The amount of detailed inspection is subject to the available workforce. Ideally, it
290 should be done until all image pairs identified as outliers are reviewed, and the FRS
291 administrator is assured of the dataset integrity and correctness.

292 **Specific Assessments**

293 The general/initial assessment will result in two sets of image pairs (allegedly genuine
294 and allegedly impostors) that must be manually reviewed to check for identity or other
295 kinds of errors.

296 Depending on the number of image pairs, a direct assessment of all pairs can be a
297 viable approach. However, if the number of pairs is relatively large, some strategies can
298 be employed to optimize this reviewing process.

299 One such strategy is to verify if some images appear more frequently in each of the
300 sets selected for review. These images should be inspected first.

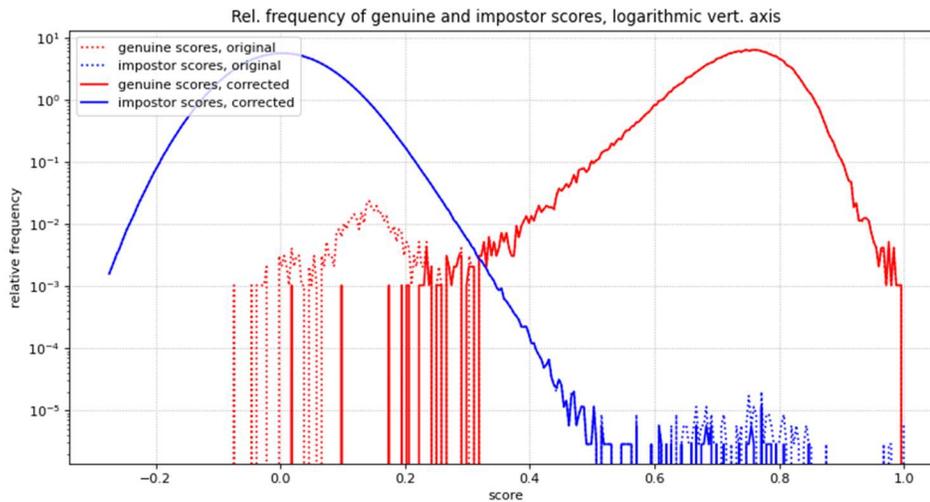
301 **Results**

302 The remaining pairs of images were inspected, and all the known errors published on
303 the LFW website were identified.

304 Some of the high impostor scores were confirmed to be from different persons. In
305 most cases, this was caused by the presence of twins, siblings, or close appearances.

306 For the low genuine scores, some were verified as being the same person, but
307 variations in pose, make-up, facial expression, and age were factors that negatively
308 affected the scores.

309 After correcting the errors, the FAR/FRR and the scores distributions plots show a
310 smoother monomodal behavior.



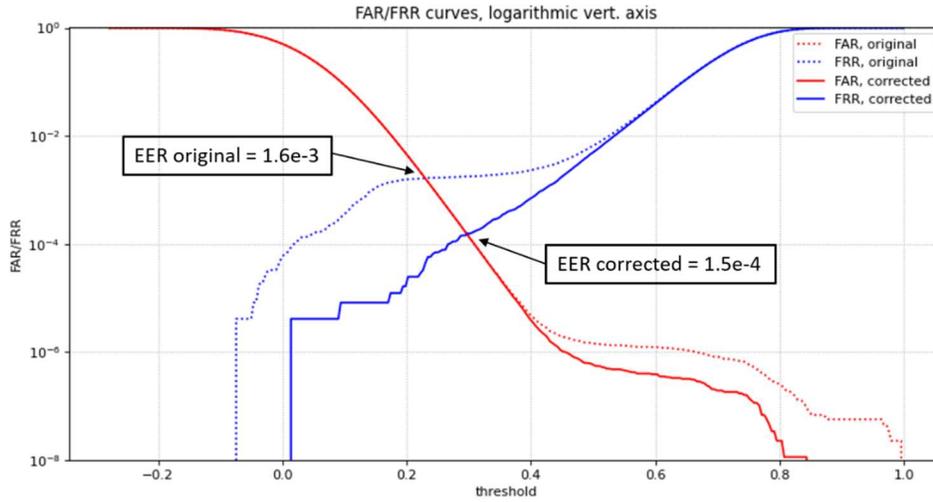
311 **FIG. A1.6 Genuine and Impostor Scoring Histograms corrected imagery (Logr axis)**
312

313 The distributions of scores improved noticeably, with fewer outliers both in the
314 genuine and impostor sets after the corrections were applied.

315 This improvement can also be observed in the FAR/FRR plot, which is directly related
316 to the scores' distributions. After the corrections, EER improved by an order of
317 magnitude.

318

319

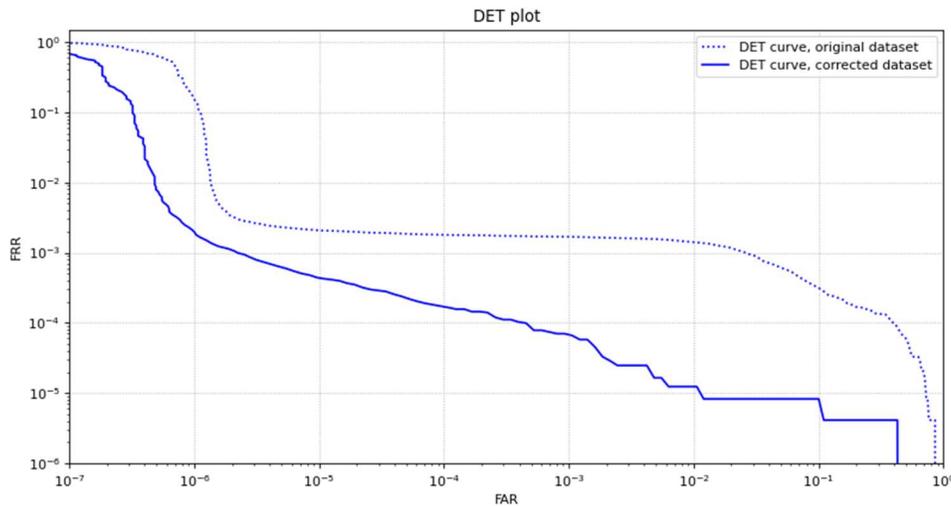


320
321

FIG. A1.7 FAR and FRR from uncorrected/corrected imagery (Logr axis)

322
323

Finally, the DET plot also reveals improvement in all operational points.



324
325

FIG. A1.8 DET from uncorrected/corrected imagery

326
327

The whole process can be repeated until the FRS administrator is assured of the

328

dataset integrity and correctness or while there are resources available.

329 Other Methods

330 The verification of integrity and correctness of biometric datasets is an evolving
331 area in the scientific literature. In this section, some methods and approaches that
332 can be used to tackle this problem are referenced.

333 The chapter on Exploratory Data Analysis (EDA) of the NIST/SEMATECH e-
334 Handbook of Statistical Methods [2] is recommended. The book describes EDA as
335 "(...) an approach/philosophy for data analysis that employs a variety of techniques
336 (mostly graphical) to (i) maximize insight into a data set; (...) (iv) detect outliers and
337 anomalies; (...)". Although not specific to inspecting facial datasets, many of the
338 concepts presented in this chapter can be used for this purpose.

339 In [3], the authors describe a method specifically designed to sort out identity
340 errors in large facial datasets, similar to the MAD analysis described in the worked
341 example. Their method is based on a two-layered thresholding process to select
342 outliers. First, identity outliers are selected, and, for each selected identity, images
343 considered as outliers are manually reviewed.

² NIST/SEMATECH e-Handbook of Statistical Methods, <https://doi.org/10.18434/M32189>

³ V. Varkarakis and P. Corcoran, "Dataset Cleaning — A Cross Validation Methodology for Large Facial Datasets using Face Recognition," 2020 Twelfth International Conference on Quality of Multimedia Experience (QoMEX), Athlone, Ireland, 2020, pp. 1-6, doi: 10.1109/QoMEX48832.2020.9123123.

344 In [4], the authors present the concept of Biometric Menagerie and propose a
345 framework for evaluating biometric systems. This framework is based on the
346 relationship between a person's genuine and impostor scores and could be explored
347 to select identities for further inspection.

348 Apart from methods that aid in selecting identities or images for manual review,
349 some methods are proposed in the literature to automatically clean large facial
350 datasets, with a large number of images assigned to each identity.

351 In one such approach, presented in [5], a clustering algorithm, Density-Based
352 Spatial Clustering of Applications with Noise – DBSCAN, clusters similar images in
353 each identity set, retaining the cluster with most images.

⁴ Yager, Neil and Ted Dunstone. "The Biometric Menagerie." IEEE Transactions on Pattern Analysis and Machine Intelligence 32 (2010): 220-230.

⁵ Gallo, Ignazio & Nawaz, Shah & Calefati, Alessandro & Piccoli, Gabriele & Zamberletti, Alessandro. (2018). A Pipeline to Improve Face Recognition Datasets and Applications. 10.1109/IVCNZ.2018.8634724.

354 Another automatic method [6] employs the community detection algorithm to
355 identify and delete mislabeled images while preserving diversity in each identity's
356 images.

357 The reader should be aware that these fully automated methods will remove some
358 images from the dataset without human reviewing.

359 Although there is a limited number of commercially available software specialized
360 for this task, most scientific literature methods can be implemented in software
361 devoted to statistics and mathematics.

362

DRAFT

⁶ Chi Jin, Ruochun Jin, Kai Chen, Yong Dou, "A Community Detection Approach to Cleaning Extremely Large Face Database", Computational Intelligence and Neuroscience, vol. 2018, Article ID 4512473, 10 pages, 2018. <https://doi.org/10.1155/2018/4512473>

363 **Appendix 2: LFW Data Set Information**

364 This is widely used open source data set which will work well for this specific
365 document. Information on this data set includes:

366

DISCLAIMER:

Labeled Faces in the Wild is a public benchmark for face verification, also known as pair matching. No matter what the performance of an algorithm on LFW, it should not be used to conclude that an algorithm is suitable for any commercial purpose. There are many reasons for this. Here is a non-exhaustive list:

- Face verification and other forms of face recognition are very different problems. For example, it is very difficult to extrapolate from performance on verification to performance on 1:N recognition.
- Many groups are not well represented in LFW. For example, there are very few children, no babies, very few people over the age of 80, and a relatively small proportion of women. In addition, many ethnicities have very minor representation or none at all.
- While theoretically LFW could be used to assess performance for certain subgroups, the database was not designed to have enough data for strong statistical conclusions about subgroups. Simply put, LFW is not large enough

to provide evidence that a particular piece of software has been thoroughly tested.

- Additional conditions, such as poor lighting, extreme pose, strong occlusions, low resolution, and other important factors do not constitute a major part of LFW. These are important areas of evaluation, especially for algorithms designed to recognize images “in the wild”.

For all of these reasons, we would like to emphasize that LFW was published to help the research community make advances in face verification, not to provide a thorough vetting of commercial algorithms before deployment.

While there are many resources available for assessing face recognition algorithms, such as the Face Recognition Vendor Tests run by the USA National Institute of Standards and Technology (NIST), the understanding of how to best test face recognition algorithms for commercial use is a rapidly evolving area. Some of us are actively involved in developing these new standards and will continue to make them publicly available when they are ready.

Welcome to Labeled Faces in the Wild, a database of face photographs designed for studying the problem of unconstrained face recognition. The data set contains more than 13,000 images of faces collected from the web. Each face has been labeled with the name of the person pictured. 1680 of the people pictured have two or more

distinct photos in the data set. The only constraint on these faces is that they were detected by the Viola-Jones face detector. More details can be found in the technical report below.

Information:

- 13233 images
- 5749 people
- 1680 people with two or more images

Reference:

Please cite as:

[Gary B. Huang](#), Manu Ramesh, [Tamara Berg](#), and [Erik Learned-Miller](#).

Labeled Faces in the Wild: A Database for Studying Face Recognition in Unconstrained Environments.

University of Massachusetts, Amherst, Technical Report 07-49, October, 2007.

[\[pdf\]](#)

BibTeX entry:

@TechReport{LFWTech,

author = {Gary B. Huang and Manu Ramesh and Tamara Berg and
Erik Learned-Miller},

title = {Labeled Faces in the Wild: A Database for Studying
Face Recognition in Unconstrained Environments},

institution = {University of Massachusetts, Amherst},

year = 2007,

number = {07-49},

month = {October}}

[Gary B. Huang](#) and [Erik Learned-Miller](#).

Labeled Faces in the Wild: Updates and New Reporting Procedures.

University of Massachusetts, Amherst, Technical Report UM-CS-2014-003,

May, 2014. [\[pdf\]](#)

@TechReport{LFWTechUpdate,

author = {Gary B. Huang Erik Learned-Miller},

title = {Labeled Faces in the Wild: Updates and New Reporting Procedures},

institution = {University of Massachusetts, Amherst},

year = 2014,

number = {UM-CS-2014-003},

month = {May}}

368 Critical to this document's purpose is the errata found on the LFW URL:

Errata:

The following is a list of known errors in LFW. Due to the small number of such errors, the database will be left as is (without corrections) to avoid confusion.

It is important that users of the database provide their algorithms with the database as is, i.e. without correcting the errors below, since previous results published for the database did not have the advantage of correcting for these errors.

Currently, there are six incorrectly labeled matched pairs in View 2. While we do not believe this should have a significant effect on accuracy, we do encourage researchers to be aware of these errors when producing any visualizations (e.g. matched pairs most confidently predicted as mismatched, as the matched pair may actually be mismatched).

369 A list of errors can be viewed at: <http://vis-www.cs.umass.edu/lfw/>

370 **Reference List**

371 [1] ANSI/NIST-ITL Standard Homepage:

372 http://www.nist.gov/itl/iad/ig/ansi_standard.cfm

373 [2] P. Grother, M. Ngan, K. Hanaoka “NISTIR 8271 DRAFT SUPPLEMENT Face
374 Recognition Vendor Test (FRVT) Part 2: Identification”

375 https://pages.nist.gov/frvt/reports/1N/frvt_1N_report.pdf

376 FISWG documents can be found at: www.FISWG.org