

2 **Disclaimer:**

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

17 **Redistribution Policy:**

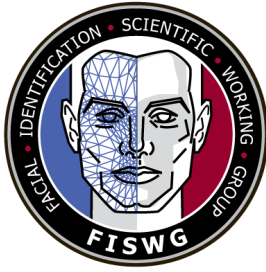
18  
19 FISWG grants permission for redistribution and use of all publicly posted documents  
20 created by FISWG, provided that the following conditions are met:

21  
22 Redistributions of documents, or parts of documents, must retain the FISWG cover  
23 page containing the disclaimer.

24  
25 Neither the name of FISWG, nor the names of its contributors, may be used to endorse  
26 or promote products derived from its documents.

27  
28 Any reference or quote from a FISWG document must include the version number (or  
29 creation date) of the document and mention if the document is in a draft status.

30



# Facial Recognition Systems Operation Assurance: Scoring Thresholds

## 31 1. Scope

32 1.1 The scope of this document is to provide a detailed process and examples of  
33 how to evaluate scoring thresholds when adjusting operational workflows. Properly  
34 executing biometric profiling producing appropriate charts to inspect and review can be  
35 done for many reasons. This document will focus on how this analysis can support a  
36 systematic process to determine appropriate facial scoring thresholds to support end  
37 user requirements. This document is relevant to systems that operate with automated  
38 workflows as well as investigative systems requiring a human practitioner to review a  
39 candidate list.

40 1.2 Understanding how to evaluate facial biometric scoring is critical for both system  
41 accuracy and workflows of the human practitioners.

42 1.3 Topics outside of this document include, but are not necessarily limited to setup,  
43 system tuning, workflow management and improvement, and proof-of-concept pilots.

## 44 2. Referenced Documents

45 2.1 *ASTM Standards:*<sup>1</sup>

---

<sup>1</sup> For referenced ASTM standards, visit [www.nist.gov/osac/astm-launch-code](http://www.nist.gov/osac/astm-launch-code), or the ASTM website, [www.astm.org](http://www.astm.org), or contact ASTM Customer Service at [service@astm.org](mailto:service@astm.org). For Annual Book of ASTM Standards volume information, refer to the standard's Document Summary page on the ASTM website.

46 E2916 Terminology for *Digital* and Multimedia Evidence Examination

47 E2825 Standard Guide for Forensic Digital Image Processing

48 2.2 *Other Standards:*

49 ANSI/NIST- ITL-1-2011 Data Format for the Interchange of Fingerprint, Facial &  
50 Other Biometric Information<sup>2</sup>

51 NISTIR 8271 Face Recognition Vendor Test (FRVT) Part 2: Identification

### 52 3. Terminology

53 3.1 *Definitions:*

54 3.1.1 For terms relating to digital and multimedia evidence, refer to Terminology  
55 E2916.

56 3.2 *Definitions of Terms Specific to This Standard:*

57 3.2.1 *Doppelganger*—an apparition or double of a living person.

58 3.3 *Acronyms*

59 3.3.1 *FR*—Face Recognition

60 3.3.2 *FRS*—Facial Recognition Systems

61 3.3.3 *CMC*—Cumulative Match Characteristic

62 3.3.4 *ROC*—Receiver Operating Characteristics

63 3.3.5 *DET*—Detection Error Tradeoff

---

<sup>2</sup> National Institute of Standards and Technology (NIST) standards available from website <https://www.nist.gov>.

64 3.3.6 *FMR*–False Match Rate proportion of the completed biometric non-mated  
65 comparison trials that result in a false match. This will be referred to as FAR (false  
66 acceptance rate) and does not include errors from images which do not create valid  
67 templates.

68 3.3.7 *FNMR*–False Non-Match Rate proportion of the completed biometric mated  
69 comparison trials that result in a false non-match. This will be referred to as FRR (false  
70 reject rate) and does not include errors from images which do not create valid  
71 templates.

72 3.3.8 *IPD*: Interpupillary Distance

#### 73 **4. Summary of Guide**

74 4.1 This document provides guidelines and techniques to help administrators of  
75 automated face recognition systems (FRS) produce recognition statistics on the face  
76 recognition systems which can be used to improve overall biometric performance.

77 4.2 The intended audience of this document is system owners, system users, and  
78 system administrators of existing automated face recognition systems.

79 4.3 This document is a continuation of the FISWG documents:

80 4.3.1 “Understanding and Testing for Face Recognition Systems Operation  
81 Assurance”

82 4.3.2 “Facial Recognition Systems Operation Assurance: Part 2, Identity Ground  
83 Truth”

84 4.3.3 “Facial Recognition Systems Operation Assurance: Part 3, Image Quality  
85 Assessment”

86 4.3.4 “Facial Recognition Systems Operation Assurance: Part 4, Manual Facial  
87 Localization”

88 4.4 The issues presented in this document form a foundation for other  
89 considerations and applications when testing such as system setup and tuning.

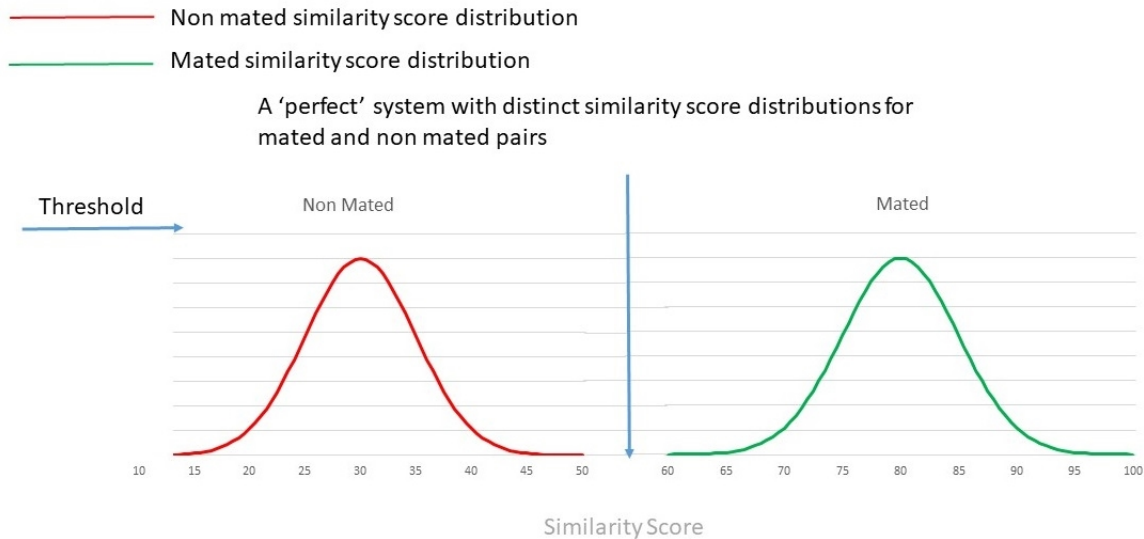
## 90 **5. Significance and Use**

### 91 5.1 Introduction

92 5.1.1 For 1:N searches, a probe image is searched against a collection of images  
93 stored in a gallery. For threshold based workflows, the search returns candidate(s) that  
94 ‘match’ the probe image above a pre-defined threshold. The length of the candidate list  
95 is usually set by the user with longer candidate lists requiring more human review effort  
96 or resource.

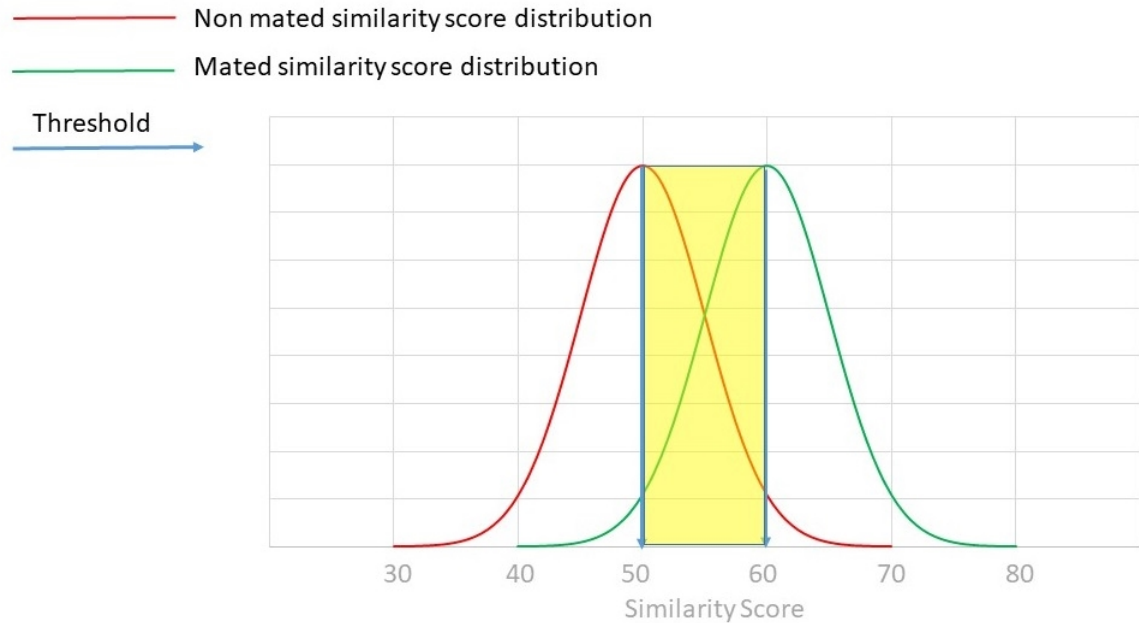
97 5.1.2 The candidates are normally sorted to have the highest score as the first  
98 (Rank 1) candidate and the lowest score as the last candidate (Rank N).

99 5.1.3 Ideally, a probe with a mate in the gallery (mated pair) will return a high score  
100 and a probe without a mate in the gallery (non-mated pair) will return a low score,  
101 allowing a clear discrimination in the meaning and inference of the results in the  
102 candidate list. If you had a perfect biometric algorithm and searched a large number of  
103 probes against a gallery with known mates and non-mates and plotted the candidate  
104 scores you would come up with a distribution shown below in Figure 1.



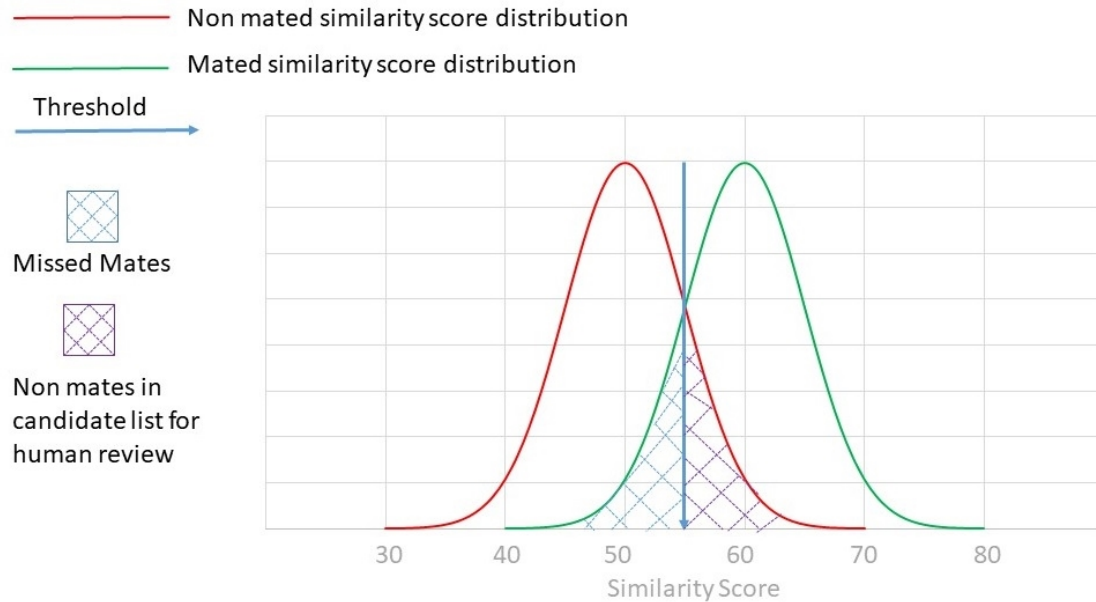
**Figure 1: Distribution of similarity scores for mated & non-mated image pairs for a 'perfect' biometric system**

- This plot shows that all the mates have a high score and all the non-mates have a low score. For this, a simple threshold score can be used to determine a mate and this would support an automated workflow scenario with no manual review.
- Unfortunately, there is no such thing as a perfect biometric algorithm. Some mates score low and some non-mates score high, which result from various conditions. The result of this is that there is a cross-over in the distributions of scores for mated and non-mated pairs and a resulting range of scores, sometimes referred to as "Yellow Resolve" where the score alone can't accurately and reliably be used to determine a mate from a non-mate see Figure 2). A facial practitioner must therefore review the candidate list and make a potential mate/non- mate decision.



122 **Figure 2: Distribution of similarity scores for mated & non-mated image pairs**  
 123 **showing 'yellow resolve' zone**

- 124
- 125 • The numeric ranges and meaning of the score distribution within this
- 126 'yellow resolve' range need to be understood by the operators of the
- 127 solution.
- 128 • The solution will have specific requirements for accuracy in terms of human
- 129 practitioner resource availability balanced against the acceptability of
- 130 incorrect matches (the False Accept Rate, requiring more time and effort to
- 131 review) and the risk of missing a mated image (the False Reject Rate) -see
- 132 Figure 3



**Figure 3: Distribution of similarity scores for mated & non-mated image pairs showing the proportion of False Rejects and False Accepts for a particular threshold setting**

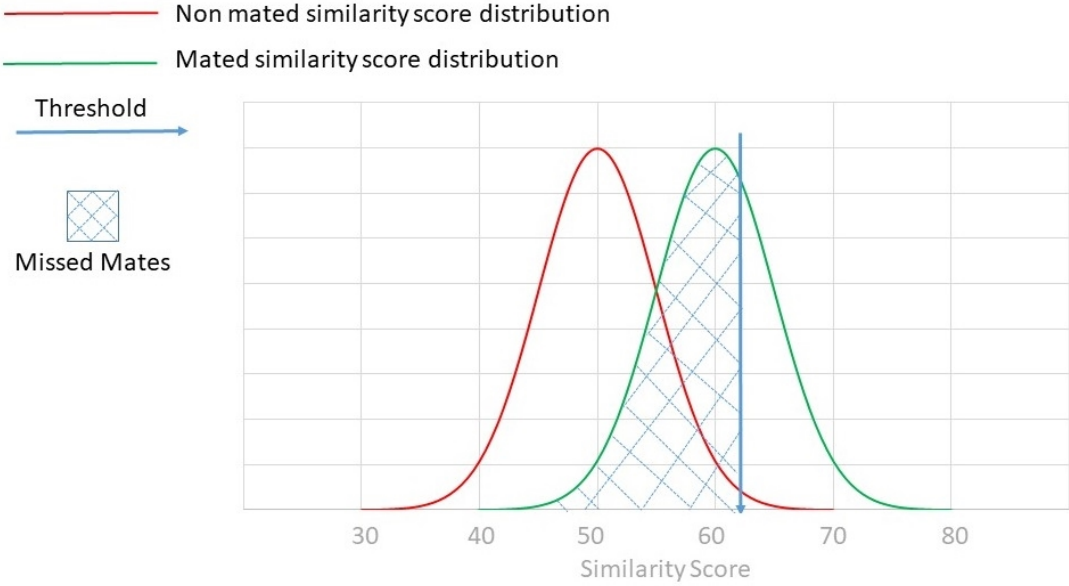
134  
135  
136

137

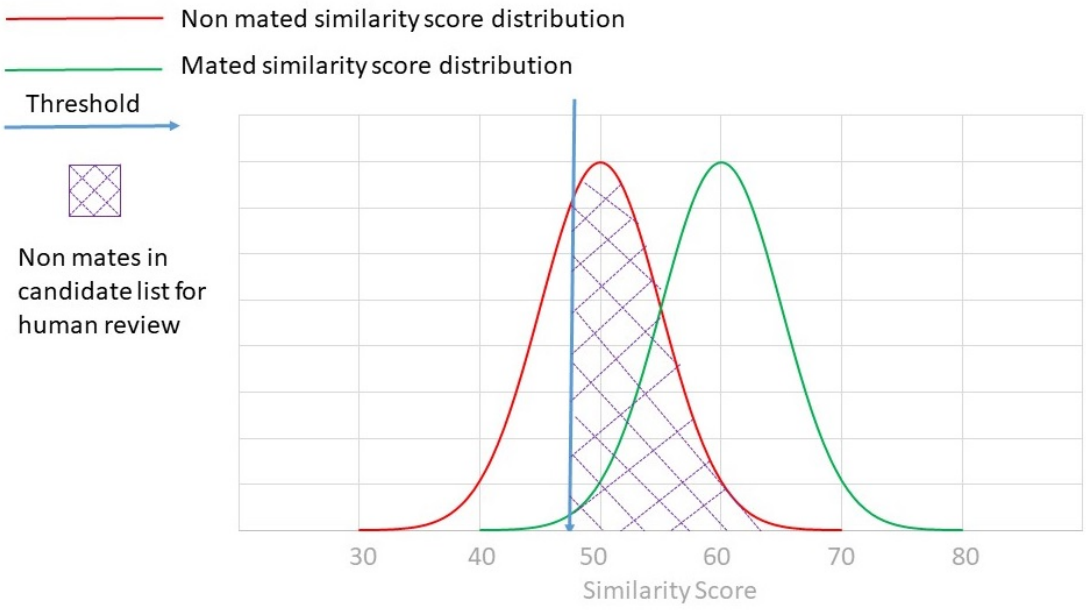
- Setting the system-operating threshold should be aligned to the Concept of Operations (ConOps) and the level of human resource available. For a high throughput ConOps where human review must be minimized, a high threshold setting will eliminate False Accepts, but there is a trade off as the proportion of False Rejects is increased (see Figure 4). For a high security ConOps, a low system threshold will ensure that mated pairs are returned in the candidate list but the trade-off is that the proportion of False Accepts is increased requiring a high level of human review resource (see Figure 5).

145





146  
 147 **Figure 4: High Throughput ConOps where a high threshold eliminates False Accepts**  
 148 **but results in a high proportion of missed mates**



150 **Figure 5: High Security ConOps where a low threshold maximizes the return of True**  
 151 **Mates but results in a high proportion of False Accepts for human review**

## 152 5.2 Considerations

### 153 5.2.1 The tuning and operational implementation of the mate, non-mate, and yellow

154 resolve areas are dependent on the areas mentioned previously:

- 155 • The facial data being evaluated
- 156 • The performance of the algorithm
- 157 • The requirements of the solution

### 158 5.2.2 Factors involving the data include:

- 159 • Were the images captured in controlled or uncontrolled environments?
- 160 • Were different capture systems used to populate the gallery?

### 161 5.2.3 Factors involving the facial algorithm include:

- 162 • How do the varying types of image quality interact with each other when
- 163 doing 1:N based searching?
- 164 • How selective is the algorithm in properly discriminating the mates and
- 165 non-mates?

### 166 5.2.4 Factors involving the requirements of the solution include:

- 167 • What is least favorable? A false accept (e.g., a high scoring non-mate) or a
- 168 false reject (e.g., a low scoring mate)?
- 169 • How will the yellow resolve area be managed? This usually involves
- 170 manual review of the search to determine if a mate exists. This manual
- 171 review takes a varying amount of time to determine and depending on
- 172 agency protocols; two separate examinations may be required by distinct
- 173 practitioners. This adds time to the search, which can then be translated
- 174 into a labor cost and a cost in terms of operational search throughput.

175 5.2.5 A key reference document to review is:

176 NISTIR 8271 DRAFT SUPPLEMENT Face Recognition Vendor Test (FRVT) Part 2: Identification

177 **Low similarity scores:** In thousands of mugshot cases the correct gallery image is returned at rank 1 but its similarity score is nevertheless low, below some operationally required score threshold. This is not so important when face recognition is used for “lead generation” in investigational applications because human reviewers are specifically required to review potentially long candidate lists and the threshold is effectively 0. In applications where search volumes are higher and labor is not available to review the results from searches, a higher threshold must be applied. This reduces the length of candidate lists and false positive identification rates at the expense of increased false negative miss rates. The tradeoff between the two error rates is reported extensively later.

178 From: [https://pages.nist.gov/frvt/reports/1N/frvt\\_1N\\_report.pdf](https://pages.nist.gov/frvt/reports/1N/frvt_1N_report.pdf) or

179 <https://doi.org/10.6028/NIST.IR.8271>

### 180 5.3 Important Notes

181 5.3.1 Care should be taken in selecting data sets to assess. It is recommended to select data sets which:

- 182 • Have operational relevancy
- 183 • Have consistent image quality aspects: type of capture, size of images, subject poses, etc
- 184 • Have sufficient identities and images to test with. This decision will be agency specific. This includes associated identity ground truth information.

185 5.3.2 The data set used for this document is the LFW (Labeled Faces in the Wild) data set available at: <http://vis-www.cs.umass.edu/lfw/> See section “LFW Data Set

189 Information“ for more details in referenced document [2]. Conceptually any other facial  
190 data set with identity ground truth can be used.

191 5.3.3 LFW is a widely used open source data set which will work well for this  
192 specific document. Information on this data set includes:

- 193 • Smaller but consistent image sizes and file formats
- 194 • Over 5,700 identities and over 13,000 images Has a wide range of  
195 subjects: sex, pose, lighting, etc.
- 196 • Stated identity ground truth errors

## 197 **6. Procedure**

### 198 6.1 Yellow Resolve Determination Process

199 6.1.1 Ensure the data set to use has verified ground truth and that any manual  
200 localization to the facial images was performed.

201 6.1.2 Enroll the facial images into a facial gallery for searching.

202 6.1.3 Search all the facial images against the facial gallery (NxM). For this test 50  
203 candidates were returned for each search but this number may vary with agency  
204 specifics and the biometric algorithm deployed. It is recommended to test with a larger  
205 number of candidates than what may be operationally used so that deeper accuracy  
206 investigations can be analyzed. Do not apply a score threshold for this test phase.

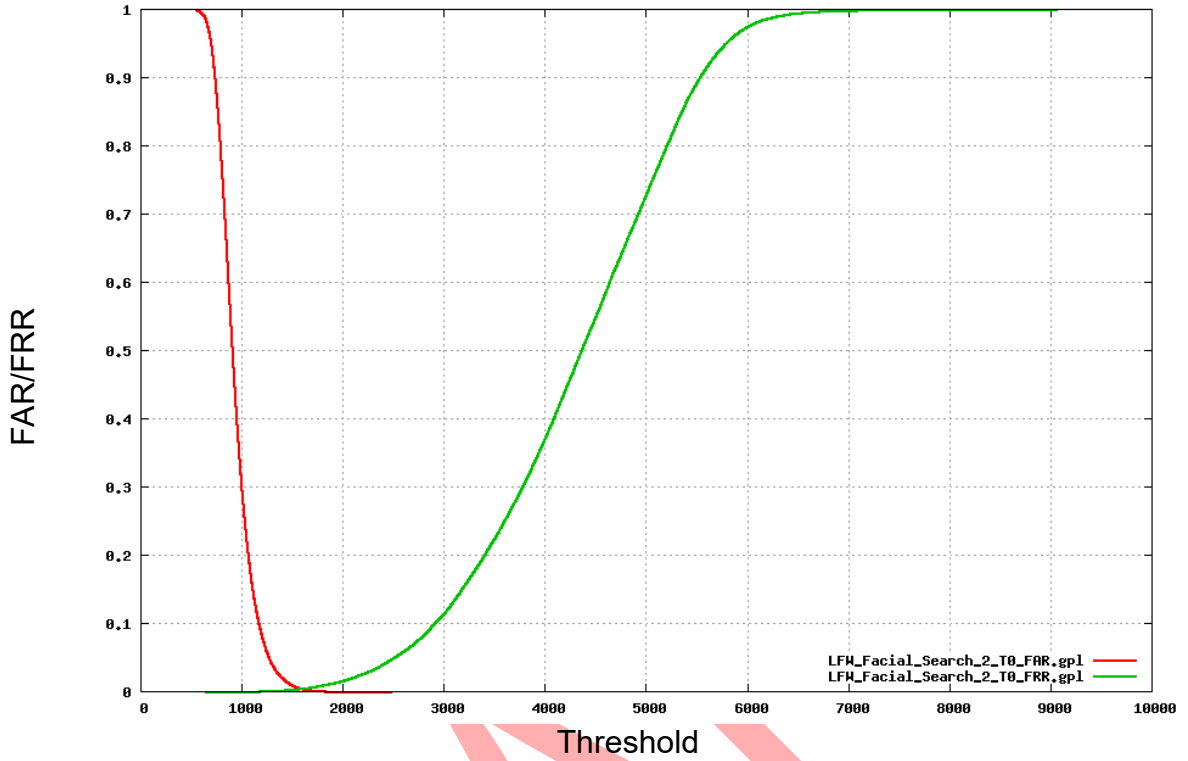
207 6.1.4 Plot the mated (FRR) and non-mated scores (FAR)

- 208 • Determine an approximate high score imposter from the FAR. This should  
209 be the highest true imposter score plus some small percentage added for  
210 safety.

- 211 • Determine a low mate score by analyzing the FRR. There are numerous  
212 ways to do this. For this document the following approach is being  
213 demonstrated.
- 214 ○ Determine the approximate equal error rate where the FRR score  
215 crosses the FAR. This score should be lowered by some small  
216 percentage added for safety.
  - 217 • The yellow resolve scoring range is therefore slightly higher than the  
218 highest scoring imposter and slightly lower than the equal error rate score.
  - 219 • Determine how many searches have a candidate score which is in the  
220 yellow resolve range. This represents the number of “searches” which  
221 would need manual review.
  - 222 • Determine how many images have a true mate lower score than the yellow  
223 resolve range. This represents the number of “mates” which would have  
224 been missed.

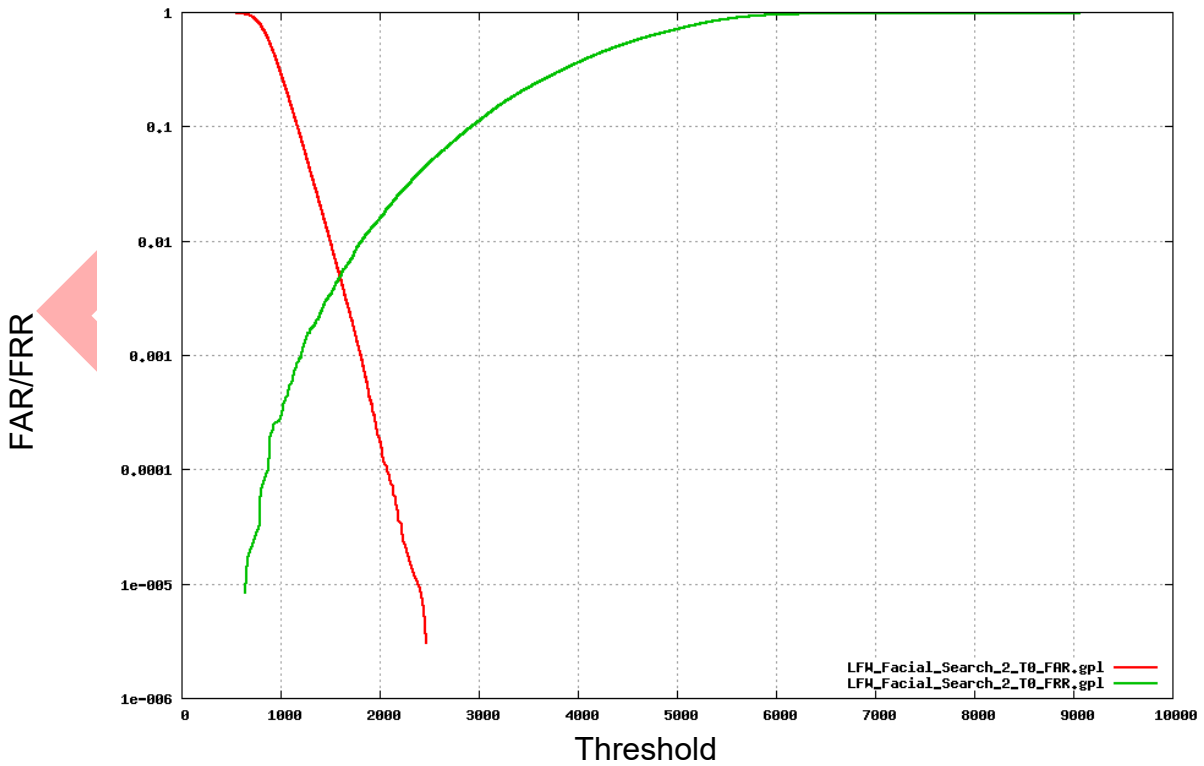
## 225 6.2 Process Outcomes

### 226 6.2.1 Step 4 Outputs:



227

Figure 3: FAR/FRR Scoring Linear Axis



228

Figure 4: FAR/FRR Scoring Logarithmic Axis

229

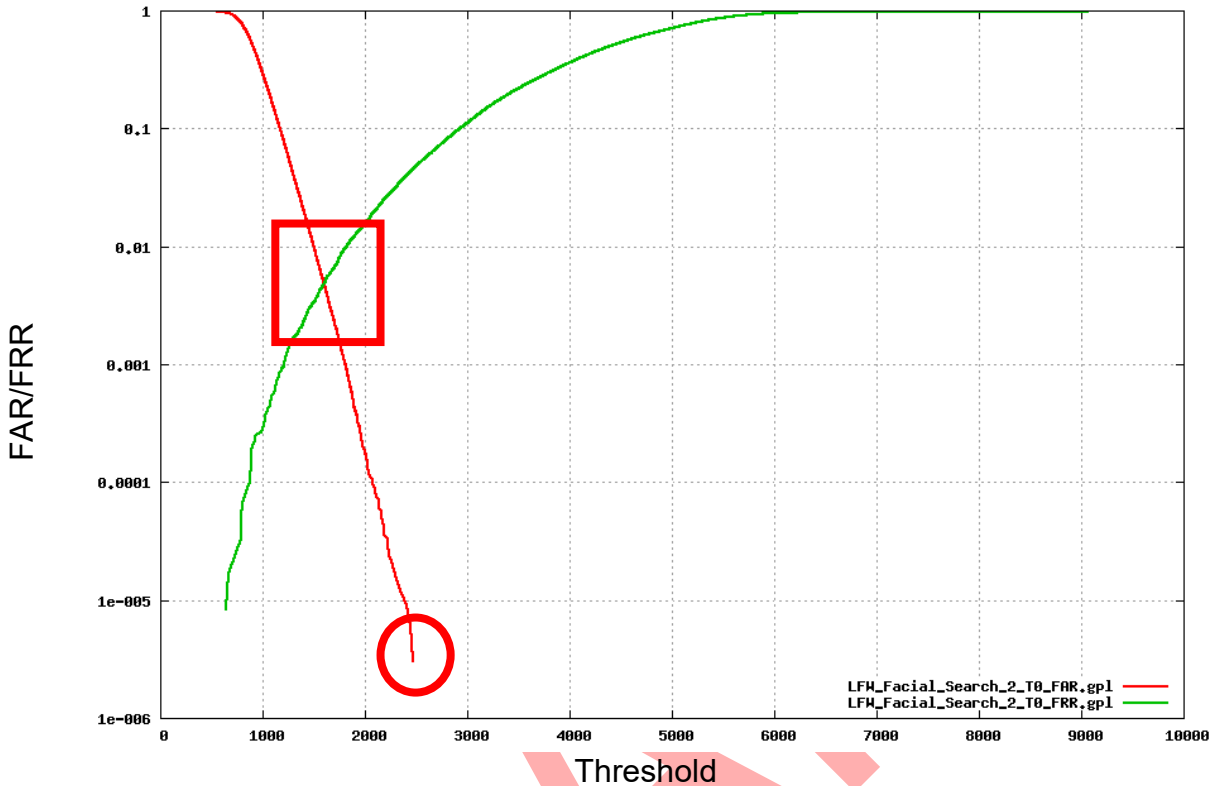


Figure 5: FAR/FRR Scoring Logarithmic Axis

230

231

232

## Notes:

233

- The highest FAR score (red circle) is ~2460.

234

- The FAR/FRR crossover score (red square) is ~1550.

235

- So given these values a yellow resolve range could be assumed to be:

236

- Low yellow: 1395 (1550 – 10%)

237

- High yellow: 2700 (2460 + 10%)

238

- This score range can then be applied to all the search results given the

239

searches with 50 candidates.

240

- Based on the sample set if we set the score range of 1395 to 2700 the

241

results were as follows:

- 242           ▪ 250 searches were in the yellow resolve range of 1395 and 2700
- 243           ▪ 289 mates were missed with a score less than 1395
- 244           ○ However with a tighter score range of 1550 to 2460 the results showed
- 245           many more mates were missed:
- 246           ▪ 136 searches were in the yellow resolve range of 1550 and 2460
- 247           ▪ 512 mates were missed with a score less than 1550
- 248           • The yellow resolve score range can be modified to gauge how many
- 249           searches would need manual adjudication by human practitioners versus
- 250           how many false positives and false negatives occurred.

## 251 6.3 Outcomes

### 252 6.3.1 Based on this data set and the testing process documented here:

- 253           • FAR and FRR curves were utilized in these processes.
- 254           • Analyzing the behavior and interactions of the FAR and FRR is an effective
- 255           evidence based method to determine yellow resolve scoring thresholds
- 256           • Once yellow resolve thresholds are determined the number of searches
- 257           which would need manual review and the number of missed mates can be
- 258           compared.
- 259           • The yellow resolve thresholds can then be modified to refine how the
- 260           agency can balance manual reviews versus missed mates.

261



262

263

FISWG documents can be found at: [www.FISWG.org](http://www.FISWG.org)

DRAFT