



Disclaimer:

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

Redistribution Policy:

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

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

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

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



Understanding and Testing for Facial Recognition Systems Operation Assurance

1

1. Scope

2

3 1.1. This document provides guidelines and techniques to help administrators of
4 automated facial recognition systems (FRS) test them in an operational setting
5 to provide assurance of facial recognition accuracy, system integrity,
6 configuration, and data storage.

7 1.2. The intended audience of this document is system owners, system users, and
8 system administrators of existing automated facial recognition systems.

9 Outside the scope of this document are types of testing including, but not
10 necessarily limited to, system setup, system tuning, workflow management and
11 improvement, and proof of concept pilots.

12 2. Referenced Documents

13 2.1. ISO/IEC (International Organization for Standardization/International
14 Electrotechnical Commission). 2012. ISO/IEC 19795-6:20121 — Information

¹ For the referenced ANSI standard, visit the ANSI webstore at
<http://webstore.ansi.org/RecordDetail.aspx?sku=ISO/IEC+19795-6:2012>.

15 technology — Biometric performance testing and reporting — Part 6: Testing
 16 methodologies for operational evaluation.

17 2.2. NIST/ITL (National Institute of Standards and Technology/ Information
 18 Technology Laboratory). 2011. NIST Special Database 32–Multiple Encounter
 19 Dataset (MEDS)².

20 2.3. NIST publication "Relating ROC and CMC Curves³," 2012, by Arun Ross.

21 2.4. "NIST Interagency Report 8271 DRAFT SUPPLEMENT Face Recognition
 22 Vendor Test (FRVT) Part 2: Identification,⁴" by P. Grother, M. Ngan, K.
 23 Hanaoka.

24 3. Terminology

25 3.1. Definitions:

26 3.1.1. *FR*, acronym for facial recognition.

27 3.1.2. *FRS*, acronym for facial recognition systems.

28 3.1.3. *CMC*, acronym for cumulative match characteristic.

29 3.1.4. *ROC*, acronym for Receiver Operating Characteristics.

30 3.1.5. *Exact binary copy*, *n* - two images that are byte-identical, i.e., they
 31 yield the same MD5 checksum. Note: An image encoded as a PNG is not an
 32 exact binary copy of the same image encoded as a JPG.

33 3.1.6. *Mate*, *n* - a separately captured image of the same subject
 34

² For the referenced NIST document visit <http://www.nist.gov/itl/iad/ig/sd32.cfm>.

³ For the referenced NIST document visit https://www.nist.gov/system/files/documents/2016/12/06/12_ross_cmc-roc_ibpc2016.pdf.

⁴ For the referenced NIST document visit https://pages.nist.gov/frvt/reports/1N/frvt_1N_report.pdf,

35 3.1.7. *Retrieval rank, n* - the order with which a particular image in the gallery
36 was retrieved respective to a probe and all other images in the gallery

37 4. Related Work

38 4.1. The overarching goal of ISO/IEC 19795-6:2012 — Information technology —
39 Biometric performance testing and reporting — Part 6: Testing methodologies
40 for operational evaluation⁵ “to measure or monitor operational biometric system
41 performance” parses into sub-goals of:

42 4.1.1. Determining if performance meets expectations or may be improved
43 through system tuning or reconfiguring;

44 4.1.2. Predicting expected performance for increases in number of enrollments
45 and/or systems;

46 4.1.3. Obtaining information that affects system performance (e.g., changes in
47 target population and environmental parameters); and/or

48 4.1.4. Obtaining performance data from a pilot implementation or to benchmark
49 future systems.

50 4.2. This document serves as a qualitative introduction to a subset of the topics and
51 concepts that are covered in greater technical detail in ISO/IEC 19795-6:2012.

⁵ ISO/IEC (International Organization for Standardization/International Electrotechnical Commission). 2012. ISO/IEC 19795-6:2012 — Information technology — Biometric performance testing and reporting — Part 6: Testing methodologies for operational evaluation. <http://webstore.ansi.org/RecordDetail.aspx?sku=ISO/IEC+19795-6:2012>.

52 **5. Operational Testing Techniques**

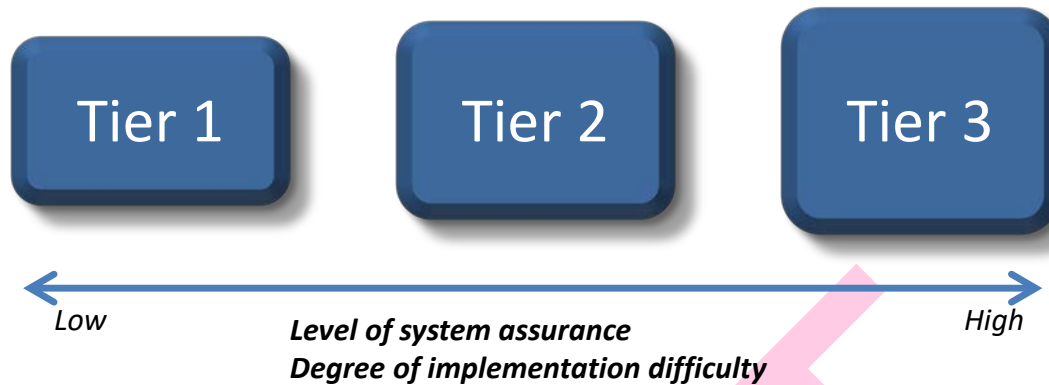
53 5.1. The operational testing techniques provided in this document can be performed
54 manually, automatically, or semi-automatically, depending on operational
55 factors and the administrator's technical knowledge.

56 5.2. Note: System administrators familiar with performance tests run by the National
57 Institute of Standards and Technology (NIST) Information Technology
58 Laboratory (ITL) like the Face Recognition Vendor Tests (FRVT) and the
59 Multiple Biometrics Grand Challenge (MBGC) should be aware that biometric
60 test results may not be applicable to other datasets and operational systems
61 with different processing constraints. These test results should be reviewed
62 within the context in which they are obtained and might not generalize to
63 operational data.

64 **6. Tier-based System Evaluation Guide**

65 6.1. A tier-based approach to operational system evaluation is proposed which
66 provides a tradeoff between the difficulty of implementing the testing strategy
67 and the level of system assurance that can be achieved through such testing.

68 Figure 1 illustrates the tier-based testing paradigm.



69

70

71

72

Figure 1. Recommendations for performing an operational evaluation of a facial recognition system are tiered, which allows administrators a tradeoff between the degree of testing difficulty and the level of system assurance

73

6.2. Tier 1 – Basic:

74

6.2.1. The simplest level of testing is Tier 1 testing, or self-identification, which

75

is a basic check to ensure there are not egregious errors in the system.

76

Tier 1 testing can be performed by the system administrator, but may

77

also be performed by system users in certain circumstances.

78

6.2.2. Tier 1 testing is performed as follows:

79

6.2.2.1. Query the facial recognition system with images whose exact

80

binary copies have already been enrolled in the database.

81

6.2.2.2. Ensure the Rank-1 match candidate is the same image as the

82

probe.

83

6.2.2.3. Continue to perform the first two steps with additional images as

84

often as availability of computing resources and human effort

85

allow.

86

6.2.3. Tier 1 testing is performed to ensure:

87

6.2.3.1. Gallery images are properly enrolled and their corresponding

88 templates are both valid and accessible.

89 6.2.3.2. The Facial recognition (FR) system's network access is not
90 interrupted for any distributed resources.

91 6.2.3.3. Software running on local and network resources is not
92 exhibiting any failures.

93 6.2.4. Tier 1 testing provides minimal assurance regarding the recognition
94 accuracy of the facial recognition system. It provides a basic level of
95 assurance of the system integrity, configuration, and data storage.

96 **6.3. Tier 2 – Intermediate:**

97 6.3.1. Tier 2 testing provides intermediate retrieval accuracy statistics on the
98 facial recognition system using test subjects whose mates are known to
99 be in the gallery. Depending on operational factors and the
100 administrator's technical knowledge, Tier 2 testing may be performed by
101 the system administrator and/or the FR algorithm vendor/integrator.

102 6.3.2. Tier 2 testing must be undertaken using the operating or vendor
103 recommended default threshold in order that the results are truly
104 representative of expected operational performance.

105 6.3.3. Testing should be undertaken with operationally representative images.
106 If the query images are from a number of different sources (e.g.
107 passport, CCTV, social media) it may be necessary to categorize the
108 images into different query sets based on these sources. It may be
109 acceptable to adjust the threshold in order to determine the retrieval rate

110 for lower quality images. However, these tests should not be run with a 0
111 threshold.

112
113

6.3.4. Tier 2 testing is performed as follows:

114 6.3.4.1. Query system with images whose mates are known to be in the
115 system. The query image should be from a separate encounter
116 than at least one gallery mate. Query images should be
117 representative of operational distributions. For example, for a
118 database of 10,000 images comprised of 80% white males and
119 20% black females, the system administrator might query
120 images of 80 white males and images of 20 black females.

121 6.3.4.2. Record the top rank in which the probe's mate was retrieved. If
122 an image from same encounter as probe image is contained in
123 the gallery (e.g., the exact binary copy), then do not consider
124 same encounter candidate list in the retrieval results.

125 6.3.4.3. Continue to perform the first two steps with additional images
126 as often as availability of computing resources and human
127 effort allow.

128 6.3.4.4. Use recorded rank retrievals from all retrieval tests to generate
129 Cumulative Match Characteristic (CMC) accuracies.

130 6.3.5. When testing on galleries that continually increase in size, the recorded
131 accuracies for this test generally will decrease over time. However, for a
132 properly functioning recognition system no major negative changes in

133 retrieval accuracy should occur for major system updates or high
134 frequency fixed-interval tests.

135 6.3.6. How to generate CMC scores:

136 6.3.6.1. The CMC scores list what percentage of image queries had
137 their mate returned at a particular retrieval rank or better. The
138 CMC scores typically would store accuracies up to the first N
139 ranks. The value of N would be determined based on how
140 many retrievals are typically examined in the system's
141 operational use. Thus, for an application where analysts
142 examine the top 20 matches, N=20.

143 6.3.6.2. The CMC scores contain N values, which correspond to the
144 Rank-1 accuracy, the Rank-2 accuracy, and all the way up to
145 the Rank-N accuracy. These accuracies are interpreted as
146 follows. The Rank-1 accuracy is what percentage of queries
147 had their mate at the first retrieval rank (i.e, the closest match).
148 The Rank-2 accuracy is what percentage of queries had the
149 mate returned at the second retrieval rank, or better. Similarly,
150 the Rank-3 accuracy is the percentage of queries that had their
151 mate retrieved at the third rank, or better.

152 6.3.6.3. It is important to note the "or better" portion of the above
153 description. If an image is matched at Rank-1, then it is also
154 included in the Rank-2 (and beyond) CMC scores. Thus, when
155 plotting the CMC scores respective the retrieval ranks, the

156 graph will not decrease.

157 6.3.6.4. When computing CMC scores in operational testing, the user
 158 would submit an image query and record the highest rank
 159 pertaining to a match. This process would be continued k
 160 times, where (as discussed above) k is determined based on
 161 available resources.

162 6.3.6.5. As an example, suppose k = 50 (i.e., 50 different images
 163 submitted to the system). For N = 7 ranks, the test may yield the
 164 results in Table 1:

Rank	1	2	3	4	5	6	7
Number at Rank	32	7	1	3	1	0	1

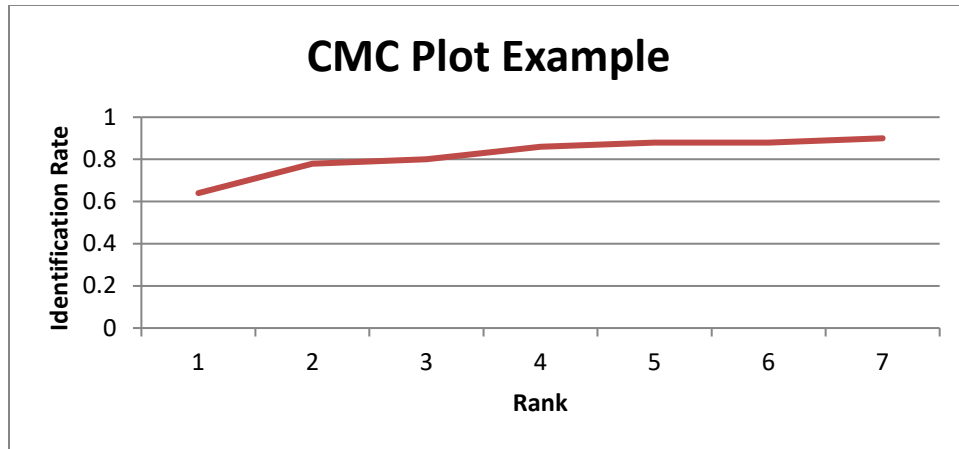
165 **Table 1. Example of Candidates for Rank Order 1 - 7**

166 The remaining 5 of the 50 submissions do not return within the
 167 designated threshold of 7 candidates, thus are not reported in
 168 the table. Given these results, the CMC scores are as shown in
 169 Table 2:

Rank	1	2	3	4	5	6	7
CMC score	32/50 =0.64	39/50 =0.78	40/50 =0.8	43/50 =0.86	44/50 =0.88	44/50 =0.88	45/50 =0.9

170 **Table 2. Example of CMC Scores for Candidates for Rank Order 1 - 7**

171 Thus, when performing Tier 2 testing, these CMC scores would
 172 be logged each time the test was performed. With a fixed set of
 173 query images, one should not expect major deviations from
 174 these CMC scores.



175

176

Figure 2. Example CMC Plot

177

Figure 2 is a plot of the CMC curve of the same example

178

previously discussed. One can visualize certain aspects of

179

CMC curves, and how they can be interpreted. One key point is

180

that the Identification Rate is never decreasing with increasing

181

Rank. When reading this plot, it can be interpreted as follows:

182

at Rank=3, the rate is 0.8. Thus, in the test example, 80% of

183

the time the subjects are matched within the top three ranks.

184

Similarly, at Rank=7, the identification rate is 0.9. Thus, in the

185

test example, 90% of the subjects are matched within the top

186

seven ranks. Finally, because only results of the top seven

187

ranks are recorded, the CMC curve⁶ never reaches 100%

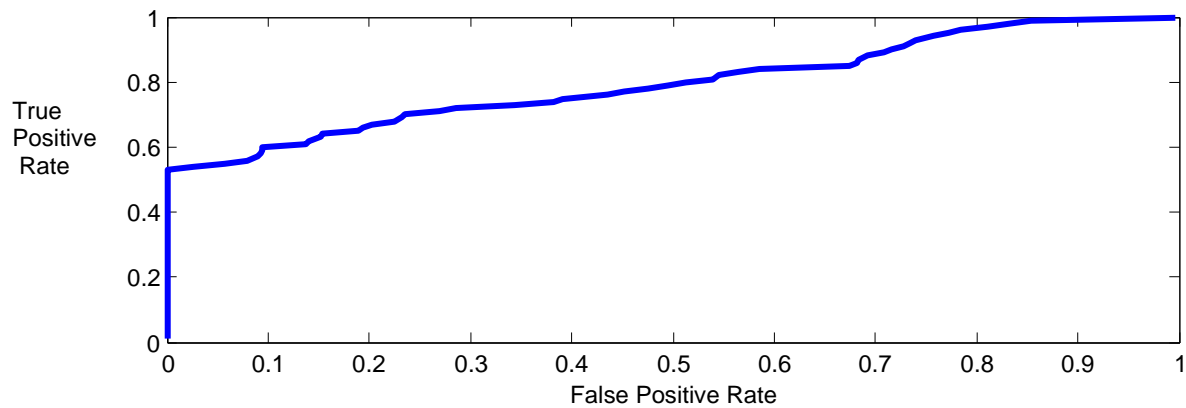
⁶ For more detail about creating a CMC curve, readers should refer to the NIST publication "Relating ROC and CMC Curves," 2012, by Arun Ross at https://www.nist.gov/system/files/documents/2016/12/06/12_ross_cmc-roc_ibpc2016.pdf.

188 accuracy. Measuring results up to a higher rank, such as 50,
189 may or may not allow one to record an identification rate of 1.0.

190 6.4. **Tier 3 – Advanced:**

191 6.4.1. Tier 3 testing provides advanced recognition accuracy statistics on the
192 facial recognition systems. This section introduces and briefly discusses
193 the Receiver Operating Characteristics (ROC) measurement and how to
194 interpret this result. For more in depth discussions on how to generate
195 and interpret advanced results such as ROC, readers are referred to
196 other documents that provide these details [3][4].

197 6.4.2. Generally, these accuracy measurements cannot be computed by the
198 FR system administrator. Instead, these measurements will often be
199 made available through software provided by the FR algorithm
200 vendor/integrator. Additionally, these results can be computed in non-
201 operational scenarios and offline environments using sufficient and
202 representative ground truth operational data. For example, the National
203 Institute of Standards and Technology routinely measures the accuracy
204 of FR vendor algorithms on various image matching scenarios [4]. By
205 understanding how to read advanced accuracy measurements, FR
206 system administrators will have a better idea of which algorithms will
207 best suit their respective organization's needs.



208

209

Figure 3 Example ROC curve

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

6.4.3. Figure 3 is an example of an ROC curve. A ROC curve plots on the x-axis the false positive rate (i.e., the percentage of impostor (or different subject) comparisons that exceed the match threshold) versus the true positive rate (i.e., the percentage of genuine (or same subject) comparisons that exceed the match threshold) on the y-axis. Any particular point on the ROC plot corresponds to the measured true positive and false positive rate at a particular match score threshold. The value of visualizing FR algorithm accuracy in the form of an ROC plot is that an organization (with the help of its integrator) can better determine which match threshold should be used as a function of how many false positive identifications will require human resource and effort to process, or the minimum true positive rate acceptable by its' mission.

6.4.4. Note that while these performance metrics can be computed for most FR systems, they may not be appropriate for certain FRS applications. For example, Rank based systems, used primarily in human adjudicated applications may not benefit from such testing. Further, when tracking

226 ROC accuracy over time in systems with highly variable gallery
227 characteristics, it may be the case that results are not stable over time.
228 By contrast, FRS applications that threshold match scores will generally
229 have more use for such ROC analysis, as they can use these
230 measurements to tune such decision thresholds over time.

231 **7. Additional Considerations and Best Practices:**

232 7.1. When to perform testing:

233 7.1.1. Operational testing should be performed before and after all major
234 system updates (e.g., software, hardware, network) to ensure success of
235 any such updates and measure changes in biometric matching
236 accuracy. Such testing is imperative due to the volatility of system
237 updates.

238 7.1.2. Fixed-interval testing should also be performed. If the system is
239 susceptible to attacks (e.g., cyber-security related), or has experienced
240 recent errors, then fixed-interval testing should occur at a higher
241 frequency.

242 7.1.3. The system administrator should set operational and fixed-interval
243 testing schedules based on the availability of computing resources and
244 human effort. There is an expectation that Tier 1 and Tier 2 testing will
245 be performed much more frequently than Tier 3 testing, which might be
246 performed only before and after major system updates.

247 7.2. How to select test images:

248 7.2.1. Operational test images should target enrollments that span different

249 periods of time from the near term to the long term. This may help
250 uncover any defects in system performance that result from template
251 corruption, errors in updates, algorithm changes, or patches to extraction
252 and/or matching algorithms.

253 7.2.2. For Tier 2 and Tier 3 testing, when selecting testing images, subjects
254 contained in those images who are recidivate (i.e., enrolled frequently in
255 the database), should be avoided, as they could reduce the consistency
256 of the accuracy reports.

257 7.2.3. Subject to agency policies around retention of images, deceased
258 subjects are ideal for enrolling into the database as test images, as
259 deceased state mostly ensures no new encounters with test subjects.
260 The MEDS database [2] contains deceased subjects, and is publicly
261 available.

262 7.2.4. Data may be collected from an uncontrolled set of test subjects that are
263 reflective of the system's target population or a test crew that is
264 representative of the system's target population, but that has not been
265 enrolled in the system.

266 7.3. Other:

267 7.3.1. If available, a non-operational system may be configured to use in an
268 "evaluation mode" to collect information not available during normal
269 system operation.

270 7.3.2. It is important that test data is representative of the operational database
271 in terms of both subjects and image quality.

- 272 7.3.3. Due to the introductory nature in this document, there are advanced
273 configurations and operational concerns not addressed which include:
274 7.3.3.1. Use of score thresholds when testing
275 7.3.3.2. Use of other accuracy statistics: Detection Error Tradeoff
276 (DET), False Accept Rate (FAR), False Reject Rate (FRR)
277 7.3.3.3. Demographic differentials in FRS algorithms
278 7.3.3.4. Selection and continuity of imagery used for probe and
279 database
280 7.3.3.5. How to address identity ground truth in mate/imposter imagery
281 7.3.4. These advanced issues and other considerations when testing (e.g.,
282 system setup and tuning) will be covered in future FISWG documents.
283
284
285
286
287
288

289 FISWG documents can be found at: www.fiswg.org
290