



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

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: FISWG@yahoogroups.com

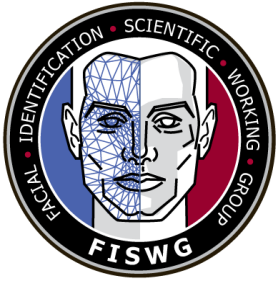
Redistribution Policy:

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 Face Recognition Systems Operation Assurance

35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52

53
54
55
56
57
58

59
60
61
62
63
64
65
66
67

68
69
70

The document provides guidelines and techniques to help administrators of automated face recognition systems test them in an operational setting to provide assurance of face recognition accuracy, system integrity, configuration, and data storage.

The intended audience of this document is system users and administrators of existing automated face recognition systems. Outside the scope of this document are types of non-operational testing including, but not necessarily limited to, system setup, system tuning, and proof of concept pilots.

Related Work

The overarching goal of *ISO/IEC 19795-6:2012 — Information technology — Biometric performance testing and reporting — Part 6: Testing methodologies for operational evaluation* [1] “to measure or monitor operational biometric system performance” parses into sub-goals of:

- ▶ Determining if performance meets expectations or may be improved through system tuning or reconfiguring;
- ▶ Predicting expected performance for increases in number of enrollments and/or systems;
- ▶ Obtaining information that affects system performance (e.g., changes in target population and environmental parameters); and/or
- ▶ Obtaining performance data from a pilot implementation or to benchmark future systems.

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

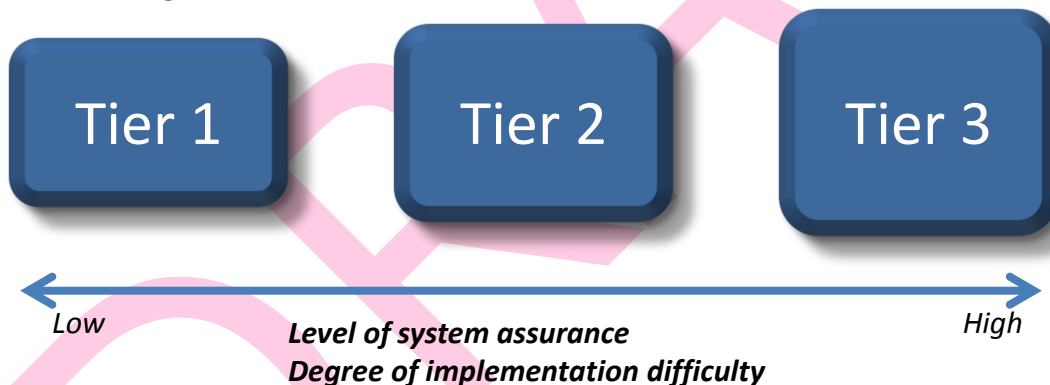
71 **Operational Testing Techniques**

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

75
76 Note: System administrators familiar with performance tests run by the
77 National Institute of Standards and Technology (NIST) Information
78 Technology Laboratory (ITL) like the Face Recognition Vendor Tests (FRVT)
79 and the Multiple Biometrics Grand Challenge (MBGC) should be aware that
80 biometric test results may not be applicable to other datasets and
81 operational systems with different processing constraints. These test results
82 have limited relevance outside of the context in which they are obtained and
83 caveats apply to the quantitative results and conclusions of the biometric
84 tests.

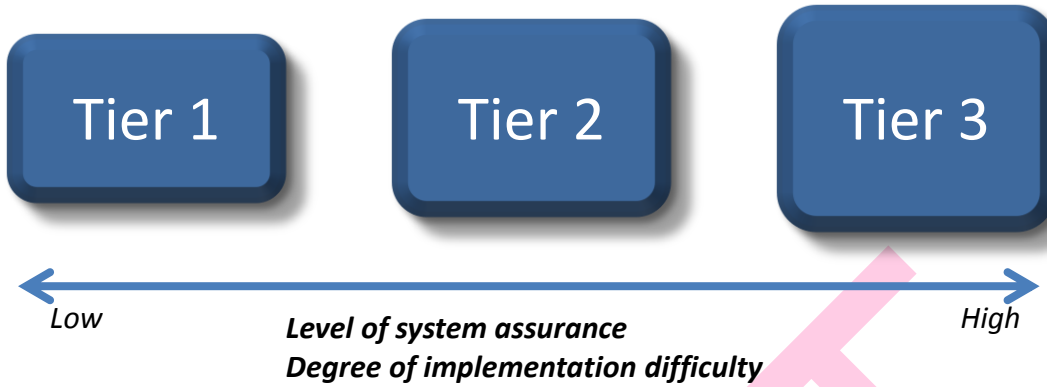
85 **Tier-based System Evaluation Guide**

86
87 A tier-based approach to operational system evaluation is proposed which
88 provides a tradeoff between the difficulty of implementing the testing
89 strategy and the level of system assurance that can be achieved through
90 such testing.



91 Figure 1 illustrates the tier-based testing paradigm.
92

93
94
95
96



97

98

99

100

101

102

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

103 **Tier 1 – Basic:**

104

105 The simplest level of testing is Tier 1 testing, or self-identification, which is a
106 basic check to ensure there are not egregious errors in the system. Tier 1
107 testing can be performed by the system administrator, but may also be
108 performed by system users in certain circumstances.

109 Tier 1 testing is performed as follows:

- 110 ▶ Query the face recognition system with images whose exact binary
- 111 copies have already been enrolled in the database.
- 112 ▶ Ensure the Rank-1 match candidate is the same image as the
- 113 probe.
- 114 ▶ Continue to perform the first two steps with additional images as
- 115 often as availability of computing resources and human effort allow.

116

117 Tier 1 testing is performed to ensure:

- 118 ▶ Gallery images are properly enrolled and their corresponding
- 119 templates are both valid and accessible.
- 120 ▶ The Face Recognition (FR) system's network access is not
- 121 interrupted for any distributed resources.
- 122 ▶ Software running on local and network resources is not exhibiting
- 123 any failures.

124 Tier 1 testing provides minimal assurance regarding the recognition accuracy
125 of the face recognition system. It provides a basic level of assurance of the
126 system integrity, configuration, and data storage.

127 **Tier 2 – Intermediate:**

128

129 Tier 2 testing provides intermediate retrieval accuracy statistics on the face
130 recognition system using test subjects whose mates are known to be in the
131 gallery. Depending on operational factors and the administrator's technical
132 knowledge, Tier 2 testing may be performed by the system administrator
133 and/or the FR algorithm vendor/integrator. Tier 2 testing is performed as
134 follows:

- 135 ▶ Query system with images whose mates are known to be in the
- 136 system. The query image should be from a separate encounter than
- 137 at least one gallery mate. Query images should be representative of
- 138 operational distributions. For example, for a database of 10,000
- 139 images comprised of 80% white males and 20% black females, the
- 140 system administrator might query images of 80 white males and
- 141 images of 20 black females.

- 142 ▶ Record the top rank in which the probe's mate was retrieved. If an
- 143 image from same encounter as probe image is contained in the
- 144 gallery (e.g., the exact binary copy), then do not consider same
- 145 encounter candidate list in the retrieval results.
- 146 ▶ Continue to perform the first two steps with additional images as
- 147 often as availability of computing resources and human effort allow.
- 148 ▶ Use recorded rank retrievals from all retrieval tests to generate
- 149 Cumulative Match Characteristic (CMC) accuracies.

150 When testing on galleries that continually increase in size, the recorded
 151 accuracies for this test generally will decrease over time. However, for a
 152 properly functioning recognition system no major negative changes in
 153 retrieval accuracy should occur for major system updates or high frequency
 154 fixed-interval tests.

155 **How to generate CMC scores:**

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

162 The CMC scores contain N values, which correspond to the Rank-1 accuracy,
 163 the Rank-2 accuracy, and all the way up to the Rank-N accuracy. These
 164 accuracies are interpreted as follows. The Rank-1 accuracy is what
 165 percentage of queries had their mate at the first retrieval rank (i.e, the
 166 closest match). The Rank-2 accuracy is what percentage of queries had the
 167 mate returned at the second retrieval rank, or better. Similarly, the Rank-3
 168 accuracy is the percentage of queries that had their mate retrieved at the
 169 third rank, or better.

170 It is important to note the "or better" portion of the above description. If an
 171 image is matched at Rank-1, then it is also included in the Rank-2 (and
 172 beyond) CMC scores. Thus, when plotting the CMC scores respective the
 173 retrieval ranks, the graph will not decrease.

174 When computing CMC scores in operational testing, the user would submit
 175 an image query and record the highest rank pertaining to a match. This
 176 process would be continued k times, where (as discussed above) k is
 177 determined based on available resources.
 178

179

Table 1 Example of Candidates for Rank Order 1 through 7

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

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

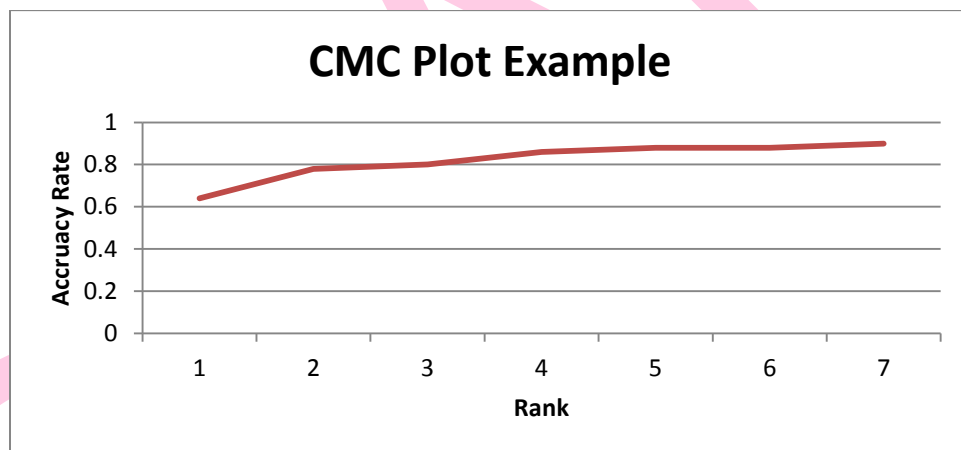
183

Table 2 Example of CMC Scores for Candidates for Rank Order 1 through 7

184

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

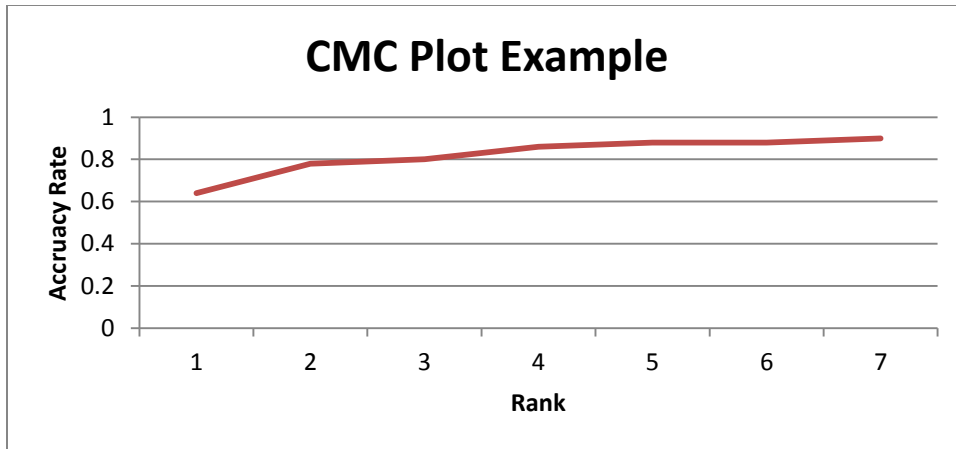
185 Thus, when performing Tier 2 testing, these CMC scores would be logged
 186 each time the test was performed. With a fixed set of query images, one
 187 should not expect major deviations from these CMC scores.



188

189

Figure 2 Example CMC Plot



190

191 Figure 2 is a plot of the CMC curve of the same example previously
 192 discussed. One can visualize certain aspects of CMC curves, and how they
 193 can be interpreted. One key point is that the Accuracy Rate is never
 194 decreasing with increasing Rank. When reading this plot, it can be
 195 interpreted as follows: at Rank=3, the rate is 0.8. Thus, in the test example,
 196 80% of the time the subjects are matched within the top three ranks.
 197 Similarly, at Rank=7, the retrieval rate is 0.9. Thus, in the test example,
 198 90% of the subjects are matched within the top seven ranks. Finally,
 199 because only results of the top seven ranks are recorded, the CMC curve
 200 never reaches 100% accuracy. Measuring results up to a higher rank, such
 201 as 50, may or may not allow one to record an accuracy rate of 1.0.

202 For more detail about creating a CMC curve, readers should refer to the Face
 203 Recognition Vendor Test 2002 [3].

204 **Tier 3 – Advanced:**

205

206 Tier 3 testing provides advanced recognition accuracy statistics on the face
 207 recognition system, such as the Receiver Operating Characteristics (ROC)
 208 and the False Positive Identification Rate (FPIR). This section briefly
 209 discusses some such measurements, and how to interpret these results. For
 210 more in depth discussions on how to generate and interpret such results,
 211 readers are referred to other documents that provide these details [3] [4].

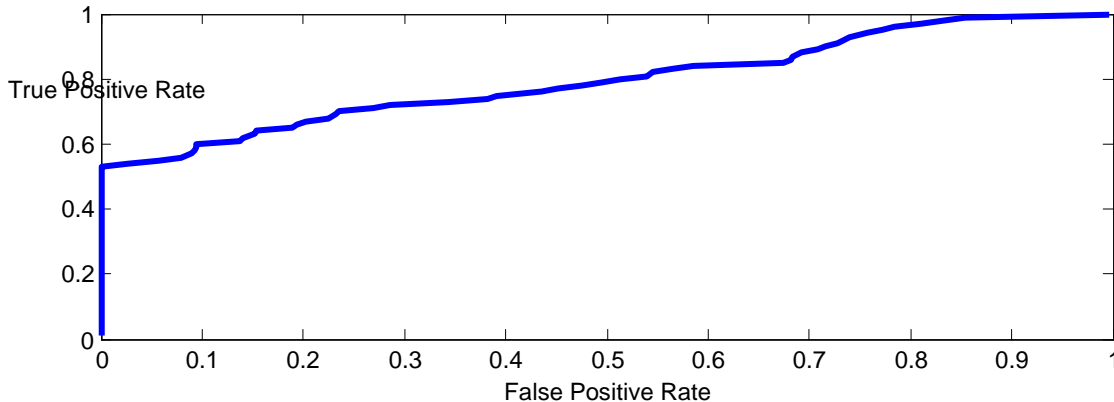
212

213 Generally these accuracy measurements cannot be computed by the FR
 214 system administrator. Instead, these measurements will often be made
 215 available through software provided by the FR algorithm vendor/integrator.
 216 Additionally, these results can be computed in non-operational scenarios and
 217 offline environments using sufficient and representative ground truth
 218 operational data. For example, the National Institute of Standards and

219 Technology routinely measures the accuracy of FR vendor algorithms on
220 various image matching scenarios [5]. By understanding how to read
221 advanced accuracy measurements, FR system administrators will have a
222 better idea of which algorithms will best suit their respective organization's
223 needs.
224

DRAFT

225



226

227

228

Error! Reference source not found. is an example of an ROC curve. An 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 it should use as a function of how many false positive identifications it has the human effort to process, or the minimum true positive rate acceptable by its mission.

241

Note that while these performance metrics can be computed for most any FR system, they may not be appropriate for certain FR applications. For example, FR systems used primarily in human adjudicated applications may not benefit from such testing. Further, when tracking ROC accuracy over time in systems with highly variable gallery characteristics, it may be the case that results are not stable over time. By contrast, FR applications that threshold match scores will generally have more use for such ROC analysis, as they can use these measurements to tune such decision thresholds over time.

251 **Additional Considerations and Best Practices**

252

253 When to perform testing:

254

255

- ▶ Operational testing should be performed before and after all major system updates (software, hardware, network) to ensure success of

256 any such updates and measure changes in biometric matching
 257 accuracy. Such testing is imperative due to the volatility of system
 258 updates.

- 259 ▶ Fixed-interval testing should also be performed. If the system is
 260 susceptible to attacks (e.g., cyber-security related), or has
 261 experienced recent errors, then fixed-interval testing should occur
 262 at a higher frequency.
- 263 ▶ The system administrator should set operational and fixed-interval
 264 testing schedules based on the availability of computing resources
 265 and human effort. There is an expectation that Tier 1 and Tier 2
 266 testing will be performed much more frequently than Tier 3 testing,
 267 which might be performed only before and after major system
 268 updates.

269 How to select test images:

- 270 ▶ Operational test images should target enrollments that span
 271 different periods of time from the near term to the long term. This
 272 may help uncover any defects in system performance that result
 273 from template corruption or errors in updates or patches to
 274 extraction and/or matching algorithms.
- 275 ▶ For Tier 2 and Tier 3 testing, when selecting testing images,
 276 subjects contained in those images who are recidivate (i.e., enrolled
 277 frequently in the database), should be avoided, as they could
 278 reduce the consistency of the accuracy reports.
- 279 ▶ Deceased subjects are ideal for enrolling into the database as test
 280 images, as deceased state mostly ensures no new encounters with
 281 test subjects. The MEDS database [1] contains deceased subjects,
 282 and is publicly available.
- 283 ▶ Data may be collected from an uncontrolled set of test subjects that
 284 are reflective of the system's target population or a test crew that is
 285 representative of the system's target population, but that has not
 286 been enrolled in the system.

287 Other:

- 288 ▶ If available, a system may be configured to use an "evaluation
 289 mode" to collect information not available during normal system
 290 operation.
- 291 ▶ It is important that test data is representative of, and ages the
 292 same as, the operational database.
- 293 ▶ There are other types of testing like system setup and tuning and
 294 this will be covered in a future FISWG document.

295 **Glossary**

296 FR = face recognition

297

298 CMC = cumulative match characteristic

299

300 exact binary copy = two images that are byte-identical, i.e., they yield the
 301 same MD5 sum. Note: An image encoded as a PNG is **not** an exact binary
 302 copy of the same image encoded as a JPG.

303

304 mate = a separately captured image of the same subject

305

306 retrieval rank = the order with which a particular image in the gallery was
 307 retrieved respective to a probe and all other images in the gallery

308 **References**

309 [1] ISO/IEC (International Organization for Standardization/International
 310 Electrotechnical Commission). 2012. *ISO/IEC 19795-6:2012 — Information
 311 technology — Biometric performance testing and reporting — Part 6: Testing
 312 methodologies for operational evaluation.*

313 <http://webstore.ansi.org/RecordDetail.aspx?sku=ISO/IEC+19795-6:2012>.

314

315 [2] NIST/ITL (National Institute of Standards and Technology/Information
 316 Technology Laboratory). 2011. *NIST Special Database 32—Multiple Encounter
 317 Dataset (MEDS)*. <http://www.nist.gov/itl/iad/ig/sd32.cfm>.

318

319 [3] P. J. Phillips, P. Grother, R. Micheals, D. Blackburn, E. Tabassi, J. M.
 320 Bone, "Face Recognition Vendor Test 2002: Evaluation Report," NIST, 2003.

321

322 [4] P. Grother, et al. "IREX III: Performance of Iris Identification
 323 Algorithms." NIST Interagency Report 7836, National Institute of Standards
 324 and Technology (2012).

325

326 [5] P. J. Grother, G. W. Quinn, and P. J. Phillips, "MBE 2010: Report on the
 327 evaluation of 2D still-image face recognition algorithms," NIST Interagency
 328 Report, National Institute of Standards and Technology, vol. 7709, 2010

329

330

331 **Reference List**

332

333

334

335

FISWG documents can be found at:

www.FISWG.org

DRAFT