



Guideline for Facial Recognition System End Users

European Network of Forensic Science Institutes
Digital Imaging Working Group

V1.0 June 2022

Project Team

Margreet Aerts-Bruintjes (Center for Biometrics, Netherlands Police)

Zsuzsanna Bartha (Hungarian Institute for Forensic Sciences)

Elisabet Leitet (National Forensic Centre, Swedish Police Authority)

Eszter Orsolya Lévai (Hungarian Institute for Forensic Sciences)

Sergio Castro Martinez (Comisaria General de Policía Científica. Spanish National Police)

Reuben Moreton (The Open University)

Johanna Morley (Interpol)

Elisabeth Pickersgill (Germany BKA)

Arnout Ruifrok (Netherlands Forensic Institute)

GUIDELINE FOR FACIAL RECOGNITION SYSTEM END USERS			
DOCUMENT TYPE: DRAFT	REF. CODE:	ISSUE NO:	ISSUE DATE:



28

29 **Table of Contents**

30 **1 INTRODUCTION 1**

31 **2 AIMS 1**

32 **3 SCOPE..... 1**

33 **4 DEFINITIONS AND TERMS 1**

34 **5 General introduction to FR systems 5**

35 **5.1 Criminal investigative use case6**

36 **5.2 FR search output.....6**

37 **5.3 Know your system7**

38 5.3.1 The FR algorithm.....7

39 5.3.2 Image enrollment quality7

40 5.3.3 Algorithm performance7

41 5.3.4 Rank based systems.....8

42 **5.4 Database of reference facial images10**

43 **6 METHODOLOGY..... 12**

44 **6.1 Flowchart12**

45 **6.2 Source material12**

46 **6.3 Assess suitability for FR search13**

47 6.3.1 Basic criteria13

48 6.3.2 Imaging acquisition factors affecting FR.....13

49 6.3.3 Subject (Human) Factors affecting FR14

50 **6.4 Request for additional materials15**

51 **6.5 Enrollment of probe image.....15**

52 **6.6 Initial image processing.....16**

53 **6.7 Run database search16**

54 6.7.1 Rank based approach17

55 6.7.2 Threshold based approach17

56 **6.8 Analysis of candidate list.....18**

57 6.8.1 Exclusion of candidates20

58 6.8.2 Comparison of candidates20

59 6.8.3 Comparing images of children20

60 6.8.4 Other relevant information21

61 **6.9 Image processing21**

GUIDELINE FOR FACIAL RECOGNITION SYSTEM END USERS			
DOCUMENT TYPE: DRAFT	REF. CODE:	ISSUE NO:	ISSUE DATE:



62	6.10	Reducing the search space using metadata filtering	22
63	6.11	Potential candidate	22
64	6.12	No potential candidate	23
65	6.12.1	Save probe to “unresolved cases database”	23
66	6.13	Multiple potential candidates	23
67	6.14	Verification process	24
68	6.15	Search aborted	24
69	6.16	Reporting and follow-up	25
70	6.16.1	Reporting a potential candidate	25
71	6.16.2	Reporting no potential candidate	25
72	6.16.3	Reporting multiple potential candidates	25
73	6.16.4	Auditing trail	26
74	6.16.5	Follow-up	26
75	7	TRAINING AND COMPETENCY	27
76	7.1	Types of FR users	27
77	7.1.1	Facial reviewer	27
78	7.1.2	Facial examiner	28
79	7.2	Selection of FR users	29
80	7.3	Training	30
81	7.4	Competency Testing	30
82	8	VALIDATION	32
83	9	REFERENCES	33
84	10	AMENDMENTS AGAINST PREVIOUS VERSION	36
85			
86			

87 **1 INTRODUCTION**

88 Facial Recognition (FR) systems can be used to search and match faces (extracted from
89 images or videos) against a database of facial images. The accuracy of FR systems is
90 nowadays high for a diverse range of image quality, mainly due to the introduction of
91 Artificial Intelligence (AI) or convolutional neural networks. In both the public and private
92 sector, the technology has many uses, such as one to many (1:N) searches for
93 Identification of an unknown subject, one to one (1:1) comparison for the Verification of
94 a claimed identity, or N:M for Clustering individuals together under a single identity.

95 The vast potential of FR, mainly in relation to the analysis of large volumes of data,
96 combined with the speed of analysis and the recent improvement in accuracy rates, has
97 led to the increased use of FR for law enforcement and investigative purposes, in Europe
98 and around the globe.

99

100 **2 AIMS**

101 This document aims to provide a framework for end-users of FR systems and deliver:

102 Basic knowledge and information on FR system set up

103 Best practices with respect to search optimization, analysis and decision making

104 Guidance with respect to the training and competency of users

105

106 **3 SCOPE**

107 This guideline addresses the processing and examination of facial images for FR
108 searching, optimizing the search parameters, the understanding and interpretation of FR
109 search results and reporting outcomes.

110 The scope is specific to retrospective 1:N identification searches for law enforcement,
111 investigatory and forensic applications.

112 Other FR technology applications, such as live (real time) FR, border control and access
113 to secure sites are out of scope.

114 Human 1:1 facial image comparison and other uses of facial images are also out of scope
115 of these guidelines. For 1:1 comparison please refer to the document "Best practice
116 Manual for Facial Image Comparison" published by the ENFSI [1].

117

118 **4 DEFINITIONS AND TERMS**

119 The terms used throughout this guide are described below:

120 **Accuracy:** A measure of how well the facial recognition process performs in terms of
121 false positive and false negative errors. It should be noted that the FR process combines
122 the automated FR system and the human review.

123 **False positive:** when an FR algorithm and/or human reviewer/examiner
124 incorrectly associates (matches) two faces together but the ground truth is that
125 the two faces belong to different individuals (referred to as non-mated search).

126 **False negative:** when an FR algorithm and/or human reviewer/examiner fails to
127 associate (match) two faces of the same person, but the ground truth is that the
128 two faces are the same individual (referred to as mated searches).

129 **False Positive Identification Rate (FPIR):** The proportion of the total number of
130 non-mated searches where one or more potential candidates are returned as a
131 mate.

132 **False Negative Identification Rate (FNIR):** The proportion of the total number
133 of mated searches where the true mate is not returned as a potential candidate.
134 This may be because it is found below a configured threshold, outside the top
135 rank candidates or not selected as a potential candidate by the human reviewer.

136 **True Negative Identification Rate (TNIR):** the complement of the FPIR, giving
137 a statement of how often non mated searches are not returned as a potential
138 candidate.

139 **True Positive Identification Rate (TPIR):** the complement of FNIR, giving a
140 positive statement of how often mated searches are successful.

141 **Audit trail:** A chronological record or set of records that provide documentary evidence
142 of the sequence of activities that have affected a specific procedure, material, system,
143 analysis or decision.

144 **Biometric data:** EU directive 2016/680, definitions: Biometric data means personal data
145 resulting from specific technical processing relating to the physical, physiological or
146 behavioural characteristics of a natural person, which allow or confirm the unique
147 identification of that natural person, such as facial images or dactyloscopic data.

148 **Candidate list:** List of images returned by the FR system, ranked according to similarity
149 score. The number of candidates shown will depend on system configuration settings
150 such as threshold and/or a defined candidate length.

151 **Cognitive bias:** A broad term that includes a variety of processes that may lead to
152 inaccurate judgments or interpretations; cognitive biases can affect memory, reasoning,
153 and decision-making.

154 **Confirmation bias:** is a type of cognitive bias, whereby people test hypotheses by
155 looking for confirming evidence rather than for potentially conflicting evidence.

156 **Cumulative Match Characteristic (CMC):** Summarizes the accuracy of mated-
157 searches and plots the proportion of mated searches returning the mate at rank R or
158 better.

159 **Deep convolutional neural network:** A class of machine learning based on artificial
160 neural networks, most commonly applied to analyze imagery.

161 **Enrollment:** The process of localizing and aligning the face from an image or video and
162 encoding the facial features to generate a template.

163 **Facial Examiner:** A trained facial comparison practitioner that conducts the task of facial
164 examination (see **Facial image comparison; Examination**).

165 **Facial image comparison:** Is a manual process undertaken by a human to identify
166 similarities and differences between facial images. Facial image comparison is used in
167 different applications, involves different levels of evaluation according to the purpose of
168 the comparison:

169 **Examination:** Detailed and methodological process to compare one image of a
170 face to another (1:1) in accordance with scientific recommendations for the
171 purpose of effecting a conclusion. It is often used in forensic applications.

172 **Review:** Is a comparison of image-to-image often used in either investigative or
173 operational lead generating applications. Review encompasses a broad range of
174 purposes and levels of rigor involved in the analysis. An independent technical
175 review or verification by at least one additional reviewer should be conducted.

176 **Facial Recognition system (FR system):** Software, which is able to detect, enroll and
177 compare faces from digital images or a video frame against a database of enrolled
178 reference facial images.

179 **Identification (1:N):** The automated comparison of an unknown biometric
180 sample against a database of (N) reference images in order to return a
181 corresponding identity.

182 **Verification (1:1):** The automated comparison of a biometric sample against the
183 reference biometric sample corresponding to the claimed identity, resulting in a
184 computer-evaluated similarity score.

185 **Clustering (N:M):** The automated grouping of biometric samples, for example,
186 represented within a collection of facial images, based on computer evaluated
187 similarity.

188 **Facial Reviewer:** A trained facial comparison practitioner that conducts the task of facial
189 review (see **Facial image comparison; Review**).

190 **Holistic comparison:** The innate human ability of comparing faces by looking at the
191 face as a whole without specifically considering the component parts in isolation.

192 **Image enrollment quality:** An algorithm or supplier proprietary metric of 'photo
193 acceptance', that provides a measure of subject characteristics (for example pose,
194 expression) and environmental/capture characteristics (for example illumination,
195 resolution, focus).

196 **Intra-variability:** Differences in (biometric) samples taken from the same person.

197 **Inter-variability:** Differences in (biometric) samples taken from different people.

198 **Metadata:** Additional (non-biometric) alpha-numeric information associated with the
199 facial images. Common examples of metadata include sex, age/date of birth, offence
200 type, date of image capture.

201 **Morphological comparison:** The direct comparison of class and individual facial
202 characteristics without explicit measurement. It is the method of facial image comparison
203 in which the features and components of the face are compared.

204 **Potential candidate:** A person in the candidate list that is judged by the reviewer to
205 show significant similarities to the person depicted in the probe image.

206 **Probe image:** Imagery in which the identity of the depicted subject is unknown. The
207 probe image may be captured under either controlled or uncontrolled conditions. Other
208 terms used as synonyms for probe image are questioned image and query image.

209 **Rank:** The position of a reference image in the candidate list as based on the similarity
210 score, where rank 1 is deemed by the FR system to have the highest level of
211 correspondence to the probe image.

212 **Reference database:** The combination of reference imagery and associated biographic
213 and other relevant information (such as name, date of birth, crime reference etc.)

214 **Reference image:** Imagery in which the identity of the depicted subject is usually known
215 and has been verified. Reference imagery is often captured under controlled conditions,
216 for example mugshot images.

217 **Similarity score:** The degree of correspondence between two facial templates as
218 judged by the FR system. The range of possible similarity scores are proprietary to the
219 algorithm. This is sometimes just referred to as 'score'.

220 **Template:** A digital representation, created by the FR system, of features extracted from
221 a biometric (face) sample. Unlike fingerprints, there are no face template standards.
222 Templates are specific to the algorithm.

223 **Threshold:** Configurable setting of the minimum image enrollment quality metric or
224 similarity score which must be reached in order for the FR system to return a candidate.

225 **True mate:** Two images that are taken from the same person. Another term used as a
226 synonym for true mate is true source.

227 **Verification:** The review and/or independent analysis of the search by another FR
228 system user.

229 **Blind verification:** A type of verification in which the subsequent examiner(s)
230 has no knowledge of the original examiner's decisions, conclusions or observed
231 data used to support the conclusion.

232 **Non-blind verification:** A type of verification in which the subsequent
233 examiner has access to the original examiner's decisions, conclusions or
234 observed data used to support the conclusion.

235

236

237 5 General introduction to FR systems

238 Facial recognition systems have been deployed since the 1990s. Early systems were
239 limited to use with high quality passport style images with an evenly lit frontal pose face
240 and neutral expression. However, the adoption of deep convolutional neural networks
241 means that FR systems are increasingly tolerant of poorly illuminated or ill posed
242 subjects and low quality images.

243 The generalized accuracy of FR systems on 'good quality' images has improved
244 dramatically and even over just the last three years there has been a circa 10% reduction
245 in false match rates, which are now well below 1% at a false non match rate of 0.0001%
246 [2].

247

248 Automated facial recognition involves four steps as illustrated in Figure 1.

249

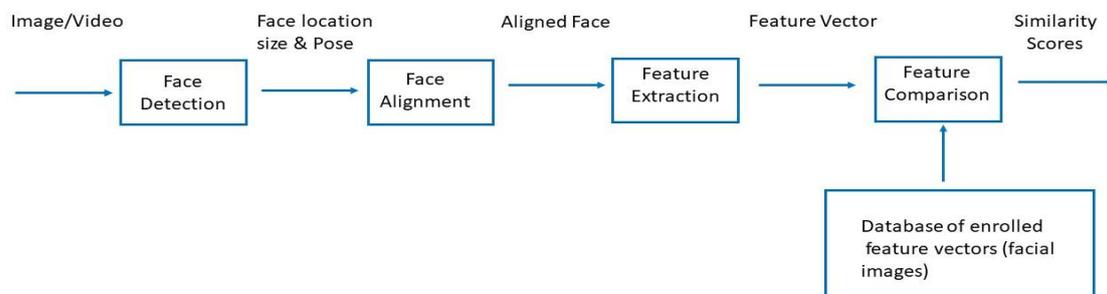


Figure 1: The steps involved in enrollment and search of a facial image using a FR system.

250

- 251 • **Face detection.** Locate one or more faces in the image and mark with a
252 bounding box
- 253 • **Face alignment.** Normalize the face to a standard frontal pose
- 254 • **Feature extraction.** Extract features¹ from the face to create a feature vector,
255 which can be thought of as a 'template' that is used in the recognition task
- 256 • **Face recognition/search.** Compare the face template against the collection of
257 (known) face templates in the database to generate a similarity score² against
258 each face

259 A given system may have a separate module for each step, which was traditionally the
260 case, or may combine some or all of the steps into a single process.

¹ Feature extraction as part of a convolutional neural network (CNN) should not be confused with facial features as observed by humans. Early FR systems generated a template using geometric measurements of facial features. However CNNs use all the information available in a facial image and use machine learning to learn the face representation that correspond to different levels of abstraction, building the set of definition data, referred to as the feature vector.

² The majority of commercially available FR systems generate a similarity score, where a higher number indicates a greater level of correspondence between the two facial images. There are some FR systems that generate a difference score (a lower number indicates a greater level of correspondence). Throughout this document, the term 'similarity score' will be used.

261 It should be noted that in the European Union, the face template is considered as
262 'Biometric data', which is a special category of personal data. Searches using biometric
263 data for identification purposes by law enforcement is only allowed when strictly
264 necessary and requires certain safeguards [3].

265

266 **5.1 Criminal investigative use case**

267 The major use case for investigative applications is 1:N searching of an unknown subject
268 against a reference database for the purposes of progressing an identification. Although
269 this is referred to as 'Identification', it should be noted that, unlike fingerprint searches,
270 the outcome of an FR search is not a positive identification (that can be, for example,
271 produced as evidence) but a (list of) potential candidate(s) that can be proposed for
272 further investigation.

273 An unknown facial image, referred to as the probe or query image, is searched against
274 a database, generally containing reference images. The outcome by the FR system is in
275 most cases a candidate list, ordered from highest to lowest similarity score according to
276 the criteria of the algorithm. Similarity scores are proprietary to the FR system, generally
277 with a higher similarity score indicating a greater degree of correspondence between the
278 facial images. They are sometimes reported as a percentage. This **must not** be
279 considered as a probability score that two images depict the same individual. A human
280 review of the candidate list is required in order to determine if any potential candidate is
281 present.

282

283 **5.2 FR search output**

284 A 1:N search using an FR system normally outputs a list of the highest ranked matches,
285 where the highest similarity score is at rank 1. FR systems can be configured as
286 threshold based, rank based, or a combination of both.

- 287 • For threshold based systems, all facial images with a similarity score above a
288 fixed similarity score threshold (set by the system users) are returned in the
289 candidate list
- 290 • For rank based systems, no threshold is set but candidates are returned in order
291 of high to low similarity score with the number of candidates set by default or by
292 the system user
- 293 • For combination systems, the similarity score threshold and the maximum
294 number of candidates returned is fixed

295 For investigative purposes, the most common deployment configuration is rank based,
296 which returns a candidate list to be reviewed by a (trained) human.

297 The output configuration may be adjusted to the use or objective of the system. For
298 example, for high throughput systems such as real time FR, a combined threshold & rank
299 configuration can be deployed and only the highest scoring candidate above the
300 threshold returned. For these systems, the majority of transactions with the system will
301 not result in any potential candidate.

302

303 **5.3 Know your system**

304 It is important to know the properties of the data within your FR system, as well as the
305 performance of your algorithm.

306

307 **5.3.1 The FR algorithm**

308 The face representation of the algorithm is trained on large amounts of labelled data.
309 The composition of the training data set in terms of the distribution of demographic
310 variations can directly impact on the “fairness” of deep models, i.e. the models should be
311 similar in the accuracy rates for different sexes or ethnicities. There are a number of
312 different methods to mitigate demographic bias in FR systems. However, both the
313 hierarchical architecture of the algorithm and composition of the training data are
314 considered to be commercially sensitive information and are in general not shared with
315 system users.

316 Therefore, it is incumbent upon system owners, administrators and users to ensure that
317 they ‘know their algorithm’. Reference should be made to large scale, independent and
318 transparent testing of FR algorithms such as those undertaken by the National Institute
319 of Standards & Technology (NIST) e.g. [2] as well as undertaking due diligence testing
320 with operationally representative data.

321

322 **5.3.2 Image enrollment quality**

323 Standardized facial image quality assessment is currently a topic of intensive research
324 [4] and may become integrated into FR systems in the near future. Currently, quality
325 scores, which incorporate factors such as resolution, sharpness and face localization,
326 are proprietary to the algorithm. Depending on the system, quality thresholds can be
327 applied, resulting in enrollment of the face only when the quality of the facial image meets
328 the threshold.

329 It is recommended that initially, image enrolment quality thresholds are set according to
330 specifications by the system providers. Agencies should undertake a quality assessment
331 to profile their face data in order to ensure that any quality thresholds set are realistic
332 based on the types of images that are received. Guidance on how to undertake such as
333 assessment are provided in the Facial Identification Scientific Working Group (FISWG)
334 document “Facial Recognition Systems Operation Assurance: Image Quality
335 Assessment” [5].

336

337 **5.3.3 Algorithm performance**

338 The performance or accuracy of an FR system can be described using two key metrics;
339 the False Positive Identification Rate (FPIR) and the False Negative Identification Rate
340 (FNIR). The methodology for generating these metrics requires a large quantity of ground
341 truth data and is outside the scope of this document. However end users should be

342 familiar with these terms and ask their system supplier to provide evidence of the
343 supplied algorithm performance.

344 Further consideration should be given to the 'equitability' of the system such that similar
345 algorithm performance should be observed across different demographics (such as
346 ethnicity, sex and age). A detailed discussion of demographic effects can be found in the
347 NIST publication 'Face Recognition Vendor Test (FRVT) Part 3: Demographic Effects'
348 [6].

349

350 **5.3.4 Rank based systems**

351 It is recommended that prior to operational implementation, testing with ground truth data
352 of known mated pairs is undertaken. Ground truth data should contain images across a
353 range of quality scores such as mugshot - mugshot, CCTV - mugshot etc. It is essential
354 that test data should contain imagery that is representative of that expected in casework.
355 The output from this testing can be used to plot a Cumulative Match Curve (CMC), which
356 summarizes the accuracy of mated-searches and plots the proportion of mated searches
357 returning the mate at rank R or better. This is not dependent on similarity scores (only
358 the rank at which the mated pair is returned), so does not distinguish between strong
359 (high similarity score) and weak mates.

360 An example CMC plot is provided in Figure 2. In this example, for mugshot probe images,
361 100% of mated pairs are returned within the top rank 1-3 positions. The number of
362 candidate images that would need to be reviewed to ensure that every mated pair is
363 returned for social media and high definition CCTV probe images is 12 and 14
364 respectively. However, it can be seen that the number of candidate images that need to
365 be reviewed for low quality CCTV images is 38. Information like this can help provide
366 guidance regarding resource implications and policy setting for use of the FR system.

367 It is important that agencies run similar tests with their specific algorithm, database and
368 case relevant probe images. This testing will assist with similarity score threshold setting
369 (where required) and determining the optimal candidate list size for review in order to
370 provide reliable search result without putting an unnecessary amount of workload on
371 reviewers³.

372 Details on how to run system tests can be found in the FISWG document 'Understanding
373 and Testing for Face Recognition Systems Operation Assurance' [7].

374

³ The table under the "Investigation" tab at <https://pages.nist.gov/frvt/html/frvt1N.html> shows the error rates on whether the mated pair appears in the first 50 rank. Many algorithms with different image types were tested. This information may also provide guidance on optimum candidate list length.

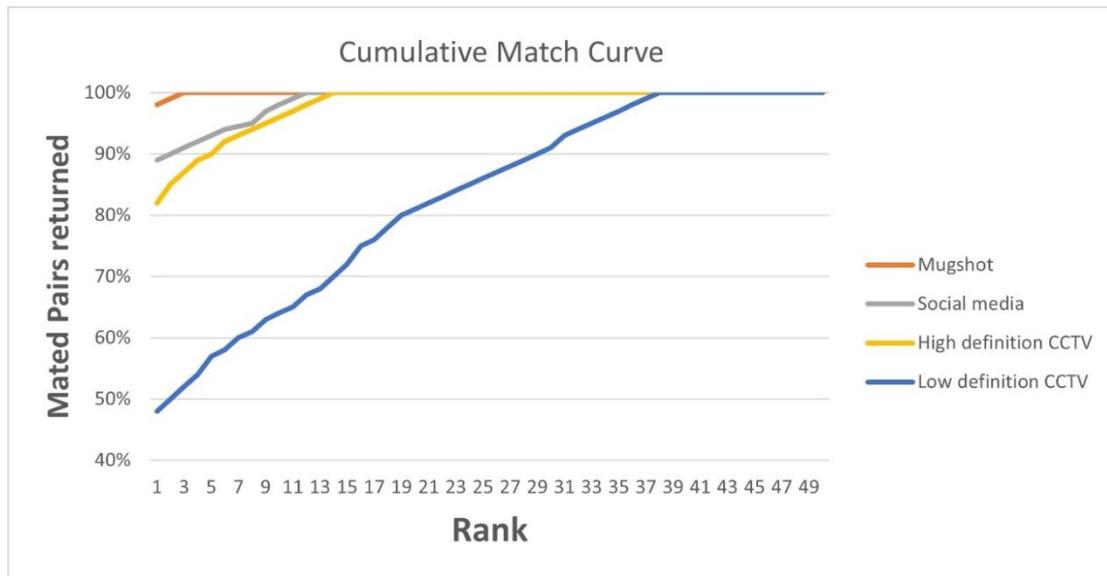


Figure 2: Example CMC for different quality probe images.

375

376 It should be noted that there is a relationship between the image quality (of the probe
 377 and gallery images) and the reviewer training requirements and time taken to review a
 378 candidate list. The following text is a summary from the FISWG document 'Facial
 379 Recognition System Methods and Techniques' [8], describing possible scenarios
 380 encountered when using FR systems.

381

382 **Comparison of:**

383

384 **1) A high-quality probe against the high-quality portions of the facial gallery**

385 Optimal images for facial comparison are high resolution and have sufficient focus to
 386 resolve features of interest, such as facial marks and facial lines, with minimal
 387 compression artifacts or distortion. The obvious advantage of comparing a high-quality
 388 probe against a high-quality gallery image is that the practitioner will be able to clearly
 389 view features, on each image, to support the morphological analysis of the face. The
 390 higher the quality of the probe image, the better the chance the system will return a
 391 potential candidate at a high ranking position in the list of reference images returned.

392

393 **2) A low-quality probe against the high-quality portions of the facial gallery and**
 394 **vice-versa**

395 Each agency and practitioner will have his/her own definition of what constitutes a low-
 396 quality probe image. These include, but are not limited to, distorted photos, low resolution
 397 face, and limited dynamic range, each of which may impede the practitioner's ability to
 398 clearly discern the subject's facial features. An FR system may accept a less-than-
 399 optimal probe image, but the lack of discernible facial features means that it will be more
 400 time consuming to review the list of returned images or that the examiner will be unable
 401 to validate a potential candidate.

402

403 **3) A low-quality probe against the low-quality portions of the facial gallery**

404 The most-challenging scenario, the submission of a low-quality probe image against a
405 collection of low quality gallery images for search by an FR system may be
406 disproportionately impacted by 'pose, expression, illumination' factors. Images returned
407 in the candidate list may be influenced by similar imaging conditions or candidates may
408 be returned on the basis of similar pose, rather than the face similarity. The time taken
409 to review such images will be significant.

410

411 **5.4 Database of reference facial images**

412 Generally, FR systems for investigative purposes consist of a reference dataset of known
413 individuals against which unknown probe images are searched. The reference dataset
414 contains image(s) and metadata associated to the individual, according to each agency's
415 policies, such as name, Date of Birth, sex, offence for which they were arrested etc.

416 It is recommended that reference images are captured under controlled conditions and
417 that, **as a minimum**, they meet appropriate quality standards for automated facial
418 recognition as depicted in Figure 3 [9]. Guidance on how to take images that meet this
419 criteria can be found in 'Standard Guide for Capturing Facial Images for Use with Facial
420 Recognition Systems' [10]. Further information can also be found in Police Standard for
421 Still Digital Image Capture and Data Interchange of Facial/Mugshot and Scar, Mark &
422 Tattoo Images [11].



Figure 3: An example of ANSI/NIST-ITL or ISO/IEC 19794-5 compliant image with distances expressed in pixels.

423

424 Additional high resolution frontal and non-frontal images (for example 45° or 90° profile)
425 will support post FR search human review of the candidate list.

426 Due to recidivism, many subjects in the reference database will have more than one
427 image associated to them, where an image is taken every time they are arrested. Studies
428 undertaken by the National Institute of Standards and Technology have demonstrated
429 that FR accuracy can be improved *when* all reference (controlled) images associated to
430 an individual are enrolled and available for searching [12]. Improvements may be a result
431 of early 'template' level fusion or score level fusion. System suppliers should be
432 consulted on the most appropriate strategy for consolidated multiple encounter
433 enrollment. Multiple enrollments of a subject can also assist the human review process.

434

435 Where it is necessary to search a probe image against other collections of images, for
436 example those taken under semi controlled conditions or collections of 'unresolved' crime
437 images, it is recommended that these images are contained in a separate partitioned
438 database. Based on agency policy, probe images can be searched against each
439 database partition separately.

440

DRAFT DOCUMENT

441 **6 METHODOLOGY**

442 The following sections describe the recommended methodology for 1:N FR searches as
443 led out in the flowchart in Figure 4.

444

445 **6.1 Flowchart**

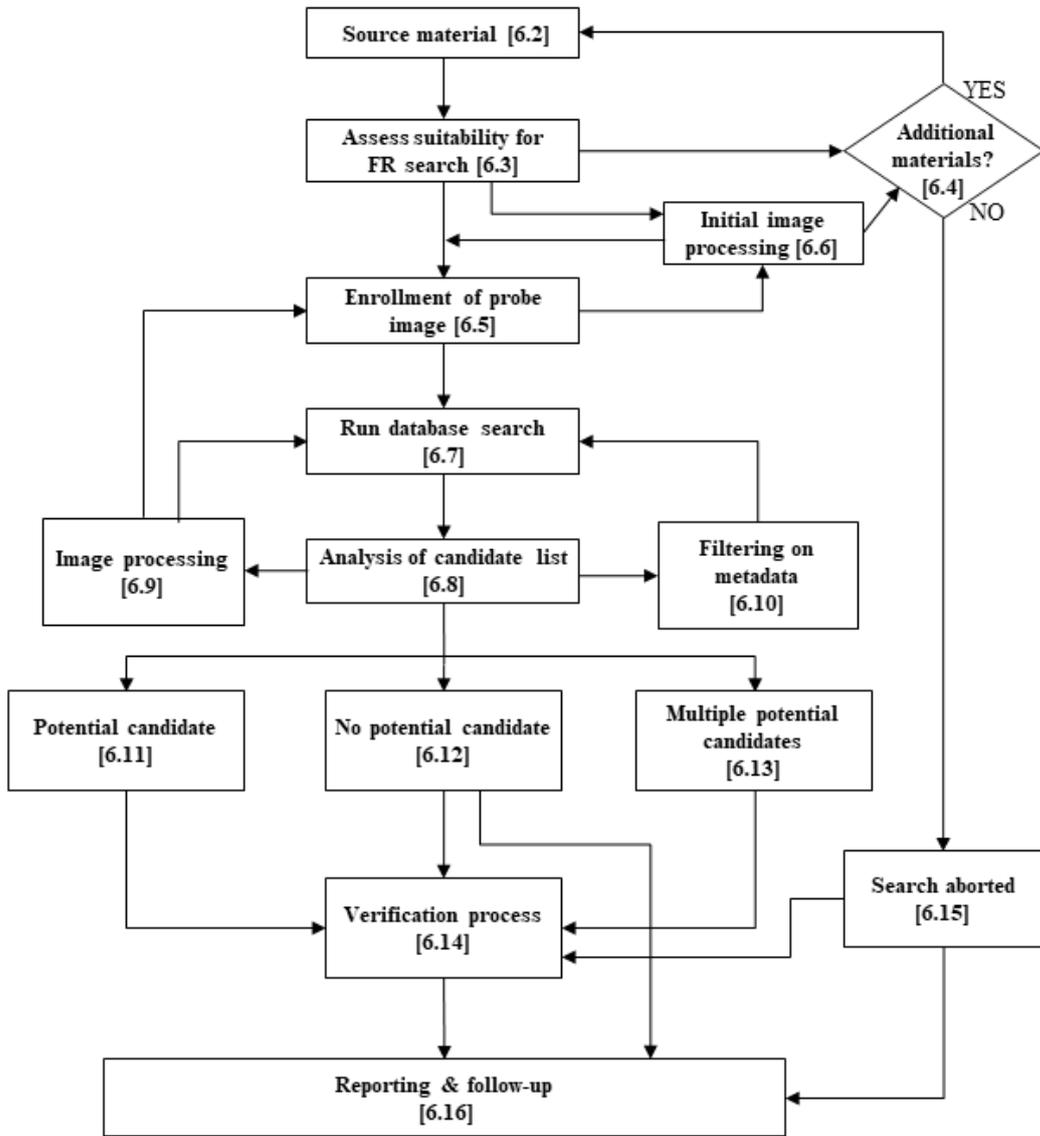


Figure 4: Flowchart showing the workflow.

446

447 **6.2 Source material**

448 Every FR search starts with the judgment of the material received. Image quality should
449 be judged by the multiple criteria that influence the FR result. FR systems are designed
450 to work best with ANSI/NIST-ITL or ISO/IEC 19794-5 compliant images (see Figure 3
451 above), and deviation from these requirements will degrade the performance of the FR

452 system. However, it is possible to achieve good results also from lower quality probe
453 images.

454 It is recommended that, where available, still images from the original data are used.
455 When the provided material is video, extraction of several frames from the original
456 footage is advised. Frames should be extracted according to best practice digital imaging
457 procedures [13]. Non-original material, for example screenshots, may add distortions,
458 lower quality and quantity of details, resulting in an overall decrease of quality. Some FR
459 systems can directly use video or multiple image search and have various mechanisms
460 for optimizing the search.

461 The criteria and imaging acquisition factors described below can be used as a general
462 guideline for assessment of the quality of facial images for use with FR systems.

463 For the remainder of this text it is assumed that one image is used as input to the FR
464 system.

465

466 **6.3 Assess suitability for FR search**

467 **6.3.1 Basic criteria**

468 According to the ISO/IEC 19794-5 standard, some basic criteria should be met in a facial
469 image. These include that the resolution of the image should be at least 60 pixels
470 between the center of the eyes, and the eyes, nose and mouth should all be visible in a
471 frontal image.

472 In reality, probe images are seldom depicted under these controlled conditions. Although
473 FR systems work best with eyes, nose and mouth all visible in a frontal image, current
474 systems may also work with off angle poses and significant parts of the face covered
475 e.g. by a facemask. Many system vendors also claim that results with resolutions down
476 to around 15 pixels between the eyes are feasible.

477 Apart from the basic criteria, the outcome of the FR system is influenced by other image
478 properties. No strict criteria can be set for these properties, but the factors listed below
479 should help in the judgement of the suitability of the image for FR. Some of these
480 properties can be 'enhanced' using image processing, but it is important to have an
481 understanding for how image processing might modify certain characteristics of the face.
482 Image processing of the probe image is described in sections 6.6 and 6.9.

483

484 **6.3.2 Imaging acquisition factors affecting FR**

485 The following factors are known to negatively influence FR performance, including, but
486 not limited to:

- 487
- 488 • Very low resolution images
 - 489 • Compression: (Re-) compression should as much as possible be prevented
490 because it will result in quality loss, unless lossless compression is used
 - 491 • Lighting: Over or under exposure or hard shadows present in the image.
492 Saturated or black areas over the face should preferably be avoided
 - Low/high contrast

- 493 • Sharpness: The face should ideally be in focus. If blurred, details may get lost.
494 Also, movement during acquisition resulting in motion blur may decrease
495 sharpness
- 496 • Artifacts due to, for example, compression, motion, signal error or data corruption
497 or re-acquisition artifacts like scanning of ID documents with security elements or
498 photographing an image displayed on a screen
- 499 • Noise (e.g. images taken at low light). Some noise reduction using image
500 processing may be performed, but will influence sharpness of the image.
- 501 • Distortion due to close distance to the camera, lens properties (e.g. fish-eye) or
502 aspect ratio changes
- 503 • Occlusions of parts of the face due to e.g. dark glasses or mouth-covering
- 504 • Pose: High viewing angle and the degree of pitch and/or yaw. FR will work well
505 on images with up to 45° yaw. Although many FR systems are now developing
506 capability for recognition from profile images (up to 90° yaw), this is still nascent
507 technology and caution should be exercised if these are the only probe images
508 available for searching

509 Visualization of pitch/roll/yaw is depicted in Figure 5.

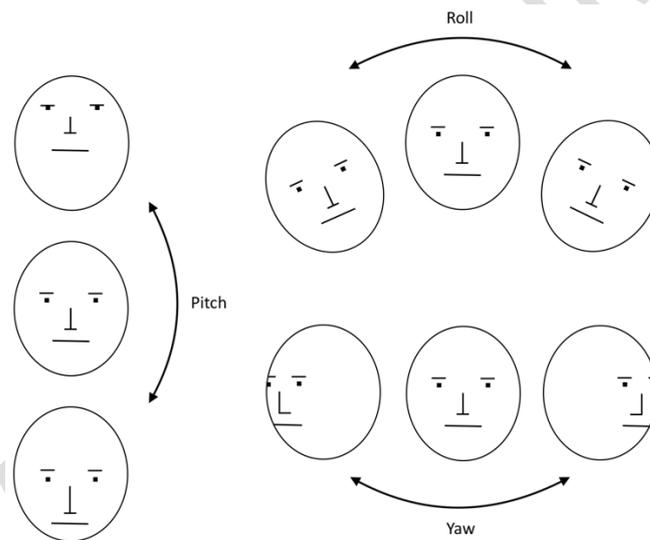


Figure 5: Visualization of Pitch, Roll & Yaw of the face

- 510
- 511 • Mirroring: Many selfie images secured from mobile phones are mirrored.
512 Information in the background might in some cases help the assessment
- 513 • Image manipulation: Intentional manipulation (including morphing or “beauty
514 filters”)
- 515 • Environmental factors: Weather (rain, etc.) windows/glass/reflections/insects
516 between subject and camera

517

518 6.3.3 Subject (Human) Factors affecting FR

519 Changes in a person’s face may negatively affect FR performance, including, but not
520 limited to:

- 521 • Expression
- 522 • Ageing
- 523 • Weight/ weight change
- 524 • Health/medical - Illness/hydration/drug use
- 525 • Post mortem changes
- 526 • Intentional alteration (makeup, surgery)
- 527 • Loss or change of features through self-mutilation or inflicted

528 While some facial features are more stable over time, others can change radically. The
529 stability of facial features in adults is listed in [14].

530

531 FR system performance on images of children (under 16 years) is known to be more
532 challenging because of the rapid changes in craniofacial morphology for infants through
533 to adolescents [15]. There are two factors which need to be considered;

- 534 • The absolute age of the subject in the test images, where generally, young
535 children (aged 0-13) are both harder to correctly recognize (high intra-variability)
536 and harder to distinguish between (low inter-variability) [16]
- 537 • Age variation, the age difference between the probe and reference image(s) with
538 a general trend of decreasing recognition rates with an increasing time gap
539 between images, where the reference image was taken up to 16 years previously
540 when the subject was a child [17]

541

542 **6.4 Request for additional materials**

543 If the submitted material fails to meet requirements in terms of quality, it might be relevant
544 that the reviewer contact the requester and ask for supplementary materials. Such action
545 should typically be balanced according to the importance of the case, the human
546 resources available and agency policies. If no additional materials are available, the
547 examination should be aborted according to section 6.15.

548

549 **6.5 Enrollment of probe image**

550 The probe image should preferably be enrolled in its original format along with relevant
551 metadata. However, in some cases it may be warranted to start with initial image
552 processing according to section 6.6, before the enrollment.

553 On enrollment, the face detection algorithm typically displays the automatically
554 determined eye positions. Only if these positions obviously deviate from the correct
555 positions, should the reviewer change these. The correct eye position should be
556 positioned according to the technical specifications of the FR system.

557

558 Manual modification of the face quality thresholds (see section 5.3.2) on an image by
559 image basis may be allowed according to agency policy in order to enroll the image.

560 If the probe image fails to enroll, the operator could proceed with adjustment of the eye
561 positions or initial image enhancements according to section 6.6. If the image still fails to
562 enroll, the reviewer should either ask additional materials, or the examination should be
563 aborted.

564 If there is more than one image of the unknown person, for example multiple still images
565 or frames from a video etc., it can be beneficial to submit more than one image for
566 searching instead of only picking one, which appears to be the best. The reason is that
567 a single image, which might seem the best for the human eye, is not always the best
568 according to the algorithm criteria. Choosing the best image by the algorithm quality
569 scores based on alignment and reliability of the face after detection are not necessarily
570 the best predictors of success. There is a chance that the mate image in the database is
571 captured in a less than ideal pose, corresponding better with an image of the unknown
572 person that is not considered to be the 'best image'.

573 For video, some systems are able to select the most appropriate facial image(s) of the
574 unknown person (based on internal quality metrics) for searching, or are capable of
575 fusing several frames together. For the latter, this may not be the most effective strategy
576 as lower quality images might be included in the search that adversely impact the
577 outcome.

578

579 **6.6 Initial image processing**

580 If the face in the image can not be found, the enrollment will fail and the facial image
581 search will not take place. To facilitate the enrollment, initial image processing might be
582 performed.

583 It is recommended that intital image processing should not significantly alter the biometric
584 data and should be kept to a minimum. Processing could include for example: cropping
585 (to remove background and /or isolate the relevant face if there are multiple faces),
586 marking the centre of the eyes, horizontal flip/mirroring (this should be utilized if the probe
587 image might have been taken as a reflection, in case of a 'selfie' image, or if the image
588 may have been flipped in transmission), enlarging using interpolation or aspect ratio
589 corrections. No (additional) compression should take place and caution should be
590 exercised when adjusting the aspect ratio as it may result in altering the geometry of the
591 face.

592

593 Image processing is further discussed in section 6.9.

594

595 **6.7 Run database search**

596 After the face (on the probe image) has been correctly enrolled, a 1:N search can be run
597 by the FR system against one or more selected reference database(s). The algorithm
598 compares the template generated from the face on the probe image to the templates
599 generated from the facial images in the database and returns a candidate list of reference
600 images.

601 As was mentioned above under section 5.2, candidate lists can be rank-based,
602 threshold-based or both.

603

604 **6.7.1 Rank based approach**

605 The rank of the true match may be affected by the quality of the probe image and
606 reference images (all the factors described in sections 6.3.2 and 6.3.3) as well as the
607 overall size of the reference database. Despite the possible effects of these factors,
608 modern algorithms have significantly improved in the rate of returning a true mate at a
609 high rank position, although its similarity score may be low.

610 In case of a rank-based approach, where the similarity score threshold is set to 0, the
611 number of the candidates in the list (candidate list size) is configured before the search
612 run. In most systems, there is a default candidate length, but in most cases this can be
613 changed manually by the user. Law Enforcement Agencies in EU countries usually return
614 searches with between 10-200 candidates [18]. For low quality images, the longer the
615 candidate list, the greater the chance that a true mate is present. However, it is worth
616 keeping in mind that the human review of the list needs time and long candidate lists will
617 affect the workload. Additionally, research has demonstrated that long candidate lists of
618 100 or more images can result in increased false alarms, lower detection of true matches,
619 lower decision confidence and increased response times [19]. These factors and the
620 severity of the crime should be taken into consideration when determining candidate list
621 size.

622

623 **6.7.2 Threshold based approach**

624 With threshold-based searches, the number of candidates returned in the candidate list
625 depends on the similarity score threshold. This threshold can be default as
626 recommended by the system supplier or manually set by the system administrator or the
627 user. The threshold should be set according to the operational requirements and
628 resources available. It should be noted that a high threshold setting minimizes the
629 workload for human review (for example eliminating non-mated candidates being
630 returned) but may result in true mates being excluded from the candidate list. This is
631 depicted in Figure 6. The converse is also true; a low threshold increases the chance
632 that a true mate is returned in the list but is resource intensive from a review perspective.

633

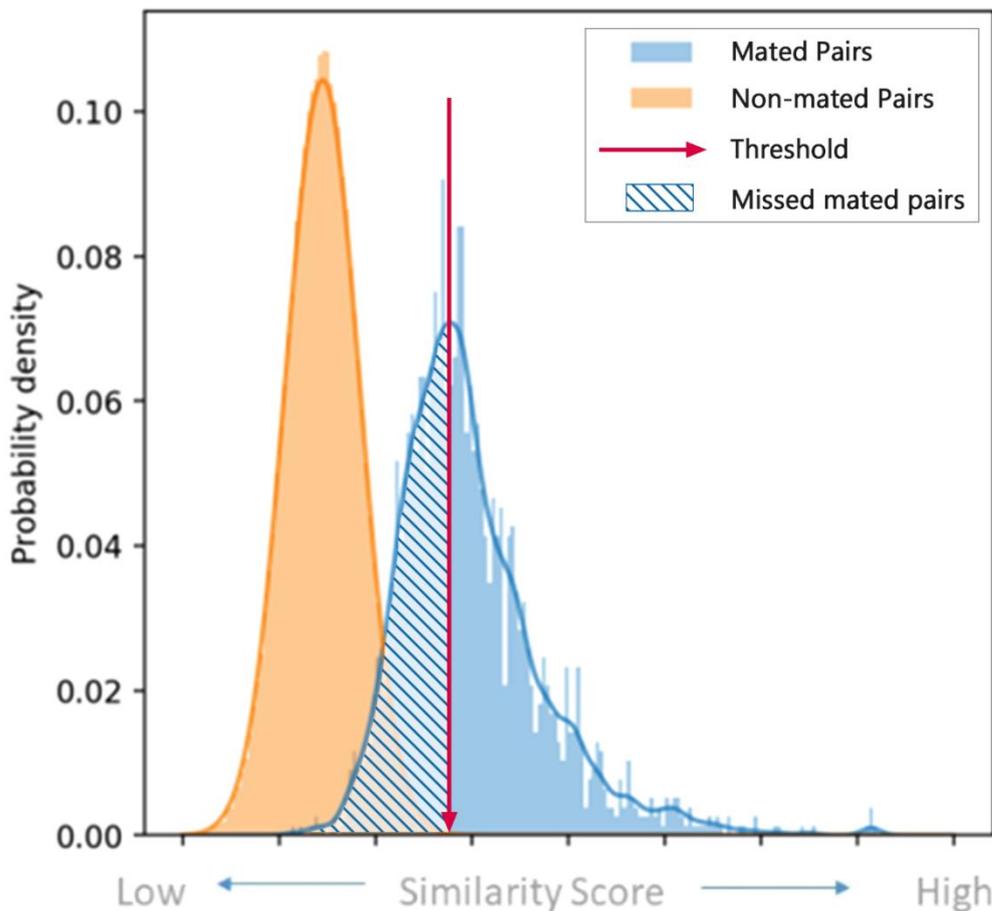


Figure 6: Similarity score distribution showing impact of setting a high threshold – everything (including a significant proportion of mated pairs) to the left of the threshold setting would be excluded from a candidate list.

634

635 The similarity score between two facial image templates is affected by both the quality
 636 of the probe image and the quality of the reference image. For example, when the probe
 637 image and one of the reference images has a similar low resolution or a similar type of
 638 noise, these may gain a high similarity score regardless of facial morphology. These
 639 reference images of candidates can outrank the true mate and may hinder the search
 640 results. Therefore, it is recommended to store different quality reference images in
 641 separate databases and run separate FR searches against each collection. Examples
 642 for quality separation are ICAO standard versus non-ICAO standard image database;
 643 mugshot gallery versus records of uncontrolled facial images; etc.

644

645 **6.8 Analysis of candidate list**

646 An image is merely a representation of reality, and a depiction of an individual is not
 647 perfectly identical to its original. Also, the appearance of a person will change over time.

648 This means that there will be a range of similarity scores returned by any FR system
649 upon a comparison of two images of the same individual (intra-variability).

650 Conversely, there is also a risk that the system perceives different people as very similar
651 to each other, even beyond obvious similarities such as close kinship, and returns a high
652 similarity score. As can be seen in Figure 6, the similarity scores of the mated and the
653 non-mated searches both show a range of values, which may be different depending on
654 the data sets at hand (e.g. mugshots versus low-quality CCTV images). It is expected
655 that there will be an overlap in the similarity score distribution for both same person
656 (mated) and different person (non-mated) comparisons by the system.

657 When using an FR system with good quality images, the difference between images of
658 the same person is small, resulting in high similarity scores, while the difference between
659 images of different persons is generally larger, resulting in low similarity scores. Hence,
660 the scores usually display a low intra-variability and high inter-variability when using good
661 quality images, resulting in good discrimination between mated and non-mated score
662 distributions. However, for low definition CCTV images the situation may be different.
663 The poor quality (uncontrolled) images of the same person may show a large difference,
664 e.g. due to differences in pose, resulting in a large variation in similarity scores (high
665 intra-variability), while images from different persons may show a lot of similarity, e.g.
666 due to similarity in pose, resulting in relatively high similarity scores for images of different
667 persons.

668 The imaging or subject factors described in sections 6.3.2 and 6.3.3, may also hamper
669 and affect the manual analysis of the candidate list and the reviewer should never
670 presume that the correct identity will be at the top of the search list even if he or she
671 exists in the database. Nor can it be presumed that the correct identity exists in the
672 database at all.

673 Nevertheless, the rank and similarity score might give valuable information to the
674 reviewer. Examples for such scenarios are (i) when the difference in the similarity scores
675 between the persons at rank 1 and rank 2 is large, (ii) if different images of the same
676 person appears more than once amongst the reference images in the candidate list or
677 (iii) if the same person appears in different candidate lists from multiple database
678 searches of the same probe image or (iv) if reference image(s) of the same person
679 appears in different candidate lists from multiple searches with different probe images of
680 the same person of interest.

681 Some agencies use the rank and score information, but other agencies prefer not to
682 show the similarity score and/or present the candidate list in random order so as to not
683 bias the reviewer by the FR system outcome. Whether information from the FR system
684 is used or not in the evaluation is therefore highly agency specific.

685 The analysis of the candidate list normally includes a first quick assessment to exclude
686 candidates, followed by a more in-depth comparison of the remaining prospective
687 candidates. There might also be other relevant information that can be used in the
688 evaluation of candidates.

689

690 **6.8.1 Exclusion of candidates**

691 The FR database search results in a list of reference images, generally ranked by
692 similarity score. Depending on the quality of both the probe and the reference images, a
693 first pass of exclusion of unlikely candidates can be performed using a holistic
694 comparison method by a reviewer. In the holistic comparison, an exclusion of unsuitable
695 candidates is made based on the basic features of the persons, such as sex, **obvious**
696 skin tone differences⁴ or other distinctive differences. Holistic comparison can therefore
697 be regarded as a first and quick evaluation to eliminate candidates and identify
698 candidates for further review using a more detailed morphological comparison.

699

700 **6.8.2 Comparison of candidates**

701 For a more in-depth morphological comparison (when the visible facial components and
702 sub-components are compared), the remaining reference images in the candidate list
703 should be viewed side by side with the probe image. Most FR systems will include a set
704 of tools, such as linked zoom to facilitate analysis. However some of these tools such as
705 superimposition (the placement of one image or video over another and adjusting the
706 transparency) [20], wiping (to slide or fade the visible part of a superimposed image to
707 illustrate similarities (slow wiping) and differences (fast wiping)) [21] or photo-
708 anthropometry (the measurement of dimensions and angles of anthropological
709 landmarks and other facial features) [22] are not reliable and should not be used [1].

710 When comparing facial features between the probe and the reference images, any
711 factors that might have an impact on the appearance of a person need to be considered.
712 These include technical and environmental conditions (e.g. resolution, lighting,
713 reflections) and subject factors (e.g. ageing, face expression). Keeping these influencing
714 factors in mind, facial features are analyzed that either support or oppose a possible
715 match. If a person exists with corresponding features to the person in the probe image,
716 this person might be regarded as a potential candidate depending on the level of details
717 observed. More on the decision for potential candidate/no potential candidate/multiple
718 potential candidates can be found in sections 6.11, 6.12 and 6.13. A potential candidate
719 listed by the first reviewer during an FR run should be verified (according to agency
720 SOPs) by a second FR reviewer, ideally in an independent review process according to
721 section 6.14.

722 Depending on the analysis of the candidate list, image processing and/or metadata
723 filtering according to sections 6.9 and 6.10 respectively might be applied and the search
724 rerun.

725

726 **6.8.3 Comparing images of children**

727 Comparing or matching facial images of children is a challenging task for humans, just
728 as for FR algorithms. People from novices to experts generally demonstrate significantly

⁴ Care should be exercised when using skin tone to exclude an image as apparent differences may be due to image factors (capture, color tone etc.)

729 lower performance with images of children in face matching tasks compared to their
730 performance with images of adults [23]. Both absolute (chronological) age and age
731 variation of the children depicted on the images impacts the human performance [24].
732 Images of younger children, as well as greater age variation make these tasks more
733 difficult, which is likely due to the rapid and great amount of facial changes happening
734 throughout childhood and the less discriminating facial features of younger children [25].
735 Facial morphological development during childhood is not evenly paced over time, facial
736 features have growth spurts and maturation occurs at different ages [25]. Studies on
737 human performance have indicated that there is a false positive response bias for 1:1
738 comparisons of child images, and that this increases with wider age variation [24].

739 In summary, this means that reviewing searches with images of children should be done
740 with extra care. According to research, methods recommended for the comparison of
741 adults have their limitations with children's faces and to date there is no international
742 standard methodology accepted for the latter [24]. Understanding the early facial
743 development patterns and child-specific training of reviewers, who are expected to do
744 work with images of children is recommended [24].

745

746 **6.8.4 Other relevant information**

747 When there is more than just the face visible of the person in the probe image, it is
748 recommended to base the comparison on all visible features of the person. Some FR
749 systems allow direct access to the police records where more reference images of the
750 candidates might be stored. These images can include, for example, face profile
751 (left/right), full body or scars/marks/tattoos. These images should be added to the
752 comparison, especially when the reference facial image in the FR system does not
753 provide all aspects of the person the probe image offers.

754 Other available information may also be taken into account, for instance additional probe
755 images, capture date or written descriptions of the person in the database or police
756 records.

757

758 **6.9 Image processing**

759 If the analysis of the candidate list does not yield a potential candidate, the reviewer may
760 choose to process the image further so as to refine the search results. Such processing
761 should be made in accordance with [26] by using either the system's own functions or a
762 separate image editing tool (after which the image can be enrolled into the system again).
763 It should be noted that the purpose of image processing is to optimize the image for
764 searching by the FR system, not to create an aesthetically pleasing image.

765 Image processing techniques that can be applied to influence the FR system performance
766 and may include, but are not limited to:

- 767 1. Histogram equalization
- 768 2. Brightness or contrast adjustment
- 769 3. Color/tint corrections
- 770 4. Grayscale conversion

- 771 5. Noise reduction
- 772 6. Red eye reduction
- 773 7. De-blurring or sharpening

774 If no potential candidate is found after a search using an image processed in this way,
775 the reviewer may choose to proceed with processing the image in a way that has a higher
776 risk of altering its biometric/geometric contents, such as:

- 777 8. Aspect Ratio correction
- 778 9. Lens distortion correction. Some images, such as those from smart phones,
779 automated teller machines (ATM's) and Body Worn Video cameras that use
780 wide-angle lenses typically exhibit significant perspective ('barrel')
781 distortion.
- 782 10. 3D pose correction

783 A log or audit trail should be kept of the image processing techniques and settings
784 applied to the image.

785 **Caution:** Image processing should always be performed on a copy of the original image.
786 The copy image should not be further compressed when image processing is applied.

787 The reviewer should primarily use methods that do not significantly alter the
788 biometric/geometric data in the image.

789 The effect of image processing will vary with different FR systems and may in some
790 cases even degrade performance rather than improve it.

791 The user should also use the original image (as well as the processed image) to aid
792 review of the candidate list and decision making

793

794 **6.10 Reducing the search space using metadata filtering**

795 The reviewer should primarily search the probe image against the entire gallery. If such
796 a search does not yield the desired results, refined searches can be performed by filtering
797 the database on metadata such as sex or approximate age span.

798 By using filtering, the database size is reduced making it possible to generate a more
799 case specific and relevant candidate list. The risks of applying such filters should be
800 noted however, as metadata may have been filled out inaccurately during the booking
801 procedure. Certain metadata may be used differently at different booking stations (for
802 example the sex of transgender persons), while others may contain subjective
803 assessments (for example age estimates).

804 **Caution:** It is important to be aware that in using metadata filtering, there is a risk that
805 the correct identity will never be included in the candidate list despite existing in the
806 database.

807

808 **6.11 Potential candidate**

809 When a candidate is considered, the reviewer must evaluate both the observed
810 similarities and differences of the persons in the images and, given the different

811 imaging or subject factors, make a final conclusion. Also, other relevant information
812 (see section 6.8.4) might be weighted into the decision.

813 *A potential candidate* means that the person in the reference image shows significant
814 similarities to the person in the probe image. A potential candidate does not mean that
815 the two persons must share the same identity.

816 In order for a potential candidate to be reported, no differences other than what is
817 considered to be caused by imaging or subject factors (see section 6.3.2 and 6.3.3),
818 should be allowed.

819 In some cases it might be appropriate to provide more than one candidate, see section
820 6.13.

821

822 **6.12 No potential candidate**

823 If none of the persons in the candidate list show significant similarities to the person in
824 the probe image, or if all persons in the candidate list show dissimilarities that cannot be
825 attributed to imaging or subject factors, the outcome of the search is *no potential*
826 *candidate*.

827 That no potential candidate is found does not necessarily mean that the correct identity
828 does not exist in the database. There are a number of reasons why an FR search might
829 not result in a match against a person even if that person exists in the database, including
830 imaging and subject factors as well as FR algorithm or manual review limitations.

831

832 **6.12.1 Save probe to “unresolved cases database”**

833 Depending on the agency specific procedures, the probe images resulting in *No*
834 *candidate* might be enrolled into a database of unresolved cases. These images might
835 also be searched regularly against new entries to the reference database, and give
836 retrospective matches. Typically, when an unresolved cases probe image is matched to
837 a known identity, the probe image is removed from the unresolved cases database.

838 Also, it is possible to match unresolved probe images against other unresolved probe
839 images, in an attempt to detect links between different crimes. It should be noted that
840 this is a challenging use case of FR technology, and that a low quality image matched
841 against another low quality image may erroneously result in a high similarity score. The
842 challenges of this scenario are discussed further in section 5.3.4.

843 The unresolved cases database use case is not further discussed in this guideline.

844

845 **6.13 Multiple potential candidates**

846 In some cases it might be appropriate to report more than one potential candidate.
847 Examples when this might occur are when persons in the reference database are very
848 similar to each other (such as identical twins), and the probe image quality does not
849 provide the reviewer sufficient details to separate between them. Another example could

850 be when the reviewer suspects that the same person has been enlisted in the reference
851 database under multiple identities (aliases).

852 Another scenario for reporting more than one candidate can be when no potential
853 candidate was concluded, but when it was not possible to exclude all persons on the
854 candidate list due to imaging or subject factors.

855 It is not recommended that the candidate list returned by the FR system is reported
856 without having been subject to a human review process.

857

858 **6.14 Verification process**

859 To mitigate the risk of a false positive result, it is highly recommended that **all** potential
860 candidates go through a verification process before the results are reported.

861 Some agencies also use a verification process for searches which have not resulted in
862 a potential candidate, a strategy which is especially recommended in an operational
863 setting where false negative results can lead to negative consequences.

864 The verification process can start at any point in the flow chart according to agency
865 specific procedures (for example a second reviewer could run the search process
866 independently) and may be conducted by;

- 867 • A second facial reviewer undertaking a blind verification process, or
- 868 • A second facial reviewer verifying the results of the first reviewer (non-blind), or
- 869 • One or more facial examiner(s) performing a 1:1 facial image comparison of the
870 probe image and potential candidate(s). The process of such a comparison is
871 described in detail in [1].

872 In blind verification, the result from the first reviewer is not known to the second reviewer,
873 while in non-blind verification the result of the first reviewer is known. Non-blind
874 verification will have a higher risk of confirmation bias.

875 If there is no consensus about the search result, the agency should have in place a policy
876 for how disagreements will be handled. Each agency should consider the benefits and
877 risks to the investigation when setting their policy, for example, reporting an incorrect
878 potential candidate versus not reporting anything. A potential strategy is that one (or
879 more) additional reviewer(s) provide their opinion(s).

880

881 **6.15 Search aborted**

882 The FR search may need to be aborted for a number of reasons, including:

- 883 • The probe image quality does not meet expected standard during the pre-
884 assessment
- 885 • The probe image fails to be enrolled into the FR system
- 886 • The FR system fails to find the face in the probe image
- 887 • The FR system returns reference images that are matched on non-facial features
888 such as image distortions
- 889 • The request is withdrawn by the requester

890 Whether any of these reasons to abort include a verification process or not will be
891 determined by agency specific procedures.

892

893 **6.16 Reporting and follow-up**

894 Results after the FR search process can be either:

- 895 • A potential candidate
- 896 • No candidate
- 897 • Multiple potential candidates
- 898 • Search aborted
- 899 • Other agency specific answer

900 For all outcomes, the results must be communicated to the requester. Normally, the
901 outcome of a FR search is reported as an investigative lead report.

902 The information to be included in an investigative lead report is normally agency or use
903 case specific. FISWG recommend that reports should include any agency disclaimer,
904 identifying the limitations of the method used and the recommended usage of the report
905 [27].

906

907 **6.16.1 Reporting a potential candidate**

908 When a potential candidate is reported, it is recommended that the following information
909 is provided:

- 910 • That a potential candidate does not indicate a positive identification of the person,
911 and a summary of the reasons why.
- 912 • That important decisions such as fundamental rights limitations (for example
913 arrests, crimes imputation, freedom of movement, etc.) should not be adopted
914 based exclusively on the potential candidate reported.
- 915 • That the investigative lead report is not intended as evidence in court
916 proceedings. The main purpose for the investigative lead report is to provide
917 criminal investigation with intelligence.

918

919 **6.16.2 Reporting no potential candidate**

920 When no potential candidate is reported, it is recommended that the following information
921 is provided:

- 922 • That a “No potential candidate” report is no guarantee that the correct identity is
923 not in the database, and a summary of the reasons why.

924

925 **6.16.3 Reporting multiple potential candidates**

926 When multiple potential candidates are reported, it is recommended that the following
927 information is provided:

- 928 • Clarification that the manual analysis could not conclude a single potential
929 candidate.
930 • Whether the potential candidates provided are listed according to any specific
931 ranking and if so which (for example the FR ranking or random).

932

933 **6.16.4 Auditing trail**

934 During the whole FR search process, from receiving the request, through the FR search,
935 the human review of the candidate list and reporting, it is recommended that an audit
936 trail of the case is recorded. Some topics that could be recorded, include;

- 937 • Administrative details
938 • FR search details
939 • Human review and verification details

940 Which information is recorded should follow agency specific protocols. The level of
941 documentation should be proportionate to the task and balance the workload, resources
942 and intended use of the report.

943

944 **6.16.5 Follow-up**

945 If the investigation against a reported potential candidate is pursued and the image
946 materials are to be used as evidence, it is recommended that a forensic 1:1 facial image
947 comparison is requested, to evaluate the imagery evidential support or not of the
948 proposed candidate. The process of such a comparison is described in detail in [1].

949 Some agencies do not report a potential candidate when such is found, but instead
950 require that a forensic 1:1 facial image comparison between the probe and reference
951 images is performed and reported instead.

952

953

954

955 **7 TRAINING AND COMPETENCY**

956 For the workflow described in this document, a human intervention is required to review
957 the results from the FR system.

958 Studies have consistently demonstrated that human performance in comparing
959 unfamiliar faces is highly varied, and on average much poorer than our ability to
960 recognize familiar faces [28].

961 Research has also shown that human operator performance can substantially impact on
962 the reliability of results from an FR search [29], which is often overlooked when
963 evaluating the performance of FR systems [30].

964 Therefore, the selection, training and testing of FR users requires careful consideration
965 to mitigate the risk of error from human review. This section provides recommendations
966 on the following aspects of FR users:

- 967 • Types of FR users
- 968 • Selection of individuals for FR user roles
- 969 • Formalized training
- 970 • Ongoing competency assessment

971 Untrained individuals, who have not received any specific training in FR review beyond
972 basic vendor training for their FR system and have not been specially selected for the
973 role, are not recommended as FR users. This is due to their lack of formalized training
974 and absence of demonstrable competency.

975 Untrained individuals are unlikely to apply any specific processes or methodology when
976 reviewing candidate lists, therefore the accuracy of the review will be largely dependent
977 upon their innate ability at comparing unfamiliar faces. In the absence of formalized
978 procedures basic users may be particularly susceptible to sources of cognitive bias, such
979 as contextual information, potentially increasing the risk of error.

980

981 **7.1 Types of FR users**

982 FR users are responsible for operating the FR system, including enrolling unknown
983 images, searching against a database, reviewing the candidate list to determine potential
984 matches and verifying results.

985 FR users can be classified as follows, in accordance with the extent of the user's training,
986 knowledge and demonstrable competency:

987

988 **7.1.1 Facial reviewer**

989 Facial reviewers are specialist FR users that have received formalized training and
990 should be able to demonstrate ongoing competency and proficiency in FR review. 'Facial
991 reviewers' typically do not receive as extensive training as facial examiners and work in
992 high throughput environments with greater time constraints than examiners.

993 Due to the large number of comparisons that a reviewer must make when reviewing
994 candidate lists it is expected that the processes used will be less rigorous and with

995 comparably less documentation of observations, than a Facial Image Comparison (FIC)
996 undertaken for forensic or evidential purposes.

997 FISWG [31] defines a facial reviewer as a FR user who:

998 *“Performs a comparison of image(s)-to-image(s) generally resulting from the adjudication*
999 *of a candidate list generated by an FRS. The comparison results are often used in either*
1000 *investigative and operational leads or intelligence gathering applications.”*

1001 Facial reviewers represent a diverse set of users and the results from studies of facial
1002 reviewer performance on facial comparison tests are similarly diverse, with some groups
1003 of reviewers demonstrating superior performance to lay persons and others not [32].
1004 Approaches to training also vary from agency to agency, from one day courses to long-
1005 term training programmes that can last for a year or more. The content of training
1006 materials are also highly varied [33].

1007 At the time of writing there is no published data to support a recommended approach to
1008 training for facial reviewers, however, when selecting and training reviewers the following
1009 general principles should be adhered to:

- 1010 • Facial reviewers typically receive shorter durations of training compared to facial
1011 examiner [33], therefore reviewers should also be selected based on innate face-
1012 matching ability, using ecologically-valid tests that are representative of the
1013 operational work they will undertake [34]
- 1014 • Facial reviewer training should be evidence-based and validated as suitable for
1015 the intended use case [35]
- 1016 • Facial reviewers should undergo continuous professional development, including
1017 ongoing competency testing using operationally-representative images and
1018 following agency specific policies and procedures

1019

1020 **7.1.2 Facial examiner**

1021 Highly trained specialists in forensic FIC, ‘facial examiners’ typically work in small teams
1022 and operate within a controlled quality management system that is sometimes accredited
1023 to an international standard (e.g. ISO 17025). Facial examiners are considered to be
1024 experts in FIC and can provide opinion-based evidence in a court of law.

1025 FISWG [31] defines the role of a facial examiner as:

1026 *“performs a comparison of image(s)-to-image(s) using a rigorous morphological analysis,*
1027 *comparison, and evaluation of images for the purpose of effecting a conclusion, often*
1028 *used in a forensic application.”*

1029 Facial examiners follow detailed procedures to perform FIC, generally based on the
1030 ACE-V framework (Analysis, Comparison, Evaluation and Verification) [36], according to
1031 international recommendations such as the *ENFSI BPM for Facial Image comparison*
1032 [1]. Studies have consistently demonstrated that facial examiners have superior
1033 performance in FIC when using their standard policies and procedures [32]. However,
1034 given that the procedures used for forensic FIC are overly time-consuming they are
1035 unlikely to be directly applicable to the task of FR review.

1036 If working in the role of an FR user, facial examiners may apply a less rigorous approach
1037 to the review of the candidate list, followed by a more detailed 1:1 comparison of viable
1038 candidates. In some situations, the facial examiner may only conduct a 1:1 comparison
1039 of viable candidates, with the candidate list review being conducted by a facial reviewer.

1040 Regardless of how the facial examiner operates as an FR user it is recommended that,
1041 as for facial reviewers, they have received validated training in the task of FR review and
1042 undergo continuous professional development and ongoing proficiency testing.

1043 Given the different task demands between FR review and forensic FIC, it is not
1044 appropriate for a facial examiner to carry out FR review without relevant training and
1045 testing to demonstrate competency in the task, even if they have demonstrated
1046 competency in forensic FIC.

1047

1048 **7.2 Selection of FR users**

1049 Human innate ability in unfamiliar facial comparison is highly varied and falls onto a wide
1050 distribution of performance. At the upper end of the distribution, a small number of
1051 individuals consistently perform exceptionally well at the task, often referred to as *super*
1052 *recognizers* in the academic literature [37]. At the bottom end of the distribution
1053 individuals perform exceptionally badly and may meet the diagnostic definition of
1054 *prosopagnosia* (or face blindness) [38]. The majority of individuals are somewhere in
1055 between.

1056 Given this heterogeneity in facial comparison ability, agencies should consider testing
1057 personnel prior to selection as FR users and enrollment in formalized training, to identify
1058 higher performing individuals.

1059 There are numerous laboratory-based tests for unfamiliar face comparison ability that
1060 are freely available, and many have normalized control data for comparison of
1061 performance [39]–[42].

1062 Whilst such tests may provide an initial indication of face comparison ability, there is
1063 limited evidence that performance on laboratory-based tests directly correlates with
1064 improved operational performance in real-world settings [43]. Additionally, individual
1065 performance can vary across different, related face-processing tasks [44] and even
1066 within the same task due to extraneous factors such as fatigue and motivation.

1067 Therefore, a single laboratory-based test may not provide an accurate picture of an
1068 individual's consistency of performance in facial comparison, and may not be directly
1069 applicable to applied tasks, like reviewing FR candidate lists in operational settings.

1070 Given that many laboratory-based tests are freely available online there is also a risk that
1071 individuals may repeatedly take tests or employ some other means to fake a high score.
1072 So, whilst online laboratory-based tests can provide an initial indication of facial
1073 comparison ability, agencies should also use a variety of screening tests that are not
1074 available to the public [42].

1075 Agency screening tests should be task-specific and representative of the types of images
1076 that an FR user will encounter operationally [34]. Potential FR users should ideally

1077 undergo multiple screening tests, taken on different days to give a measure of
1078 consistency in performance.

1079 FR user screening tests should be validated for their intended use-case and provide an
1080 accurate measure of an individual's performance by having a criterion score or cut-off
1081 for superior performance.

1082

1083 **7.3 Training**

1084 Facial reviewers and facial examiners undergo formalized training to achieve
1085 competency in their role. There has been limited research into the efficacy of formalized
1086 training for FR review, however studies have shown that short, one-off training courses
1087 of three days or less are largely ineffective at improving facial comparison performance
1088 [35], [45].

1089 Approaches to training FR users vary substantially between agencies. In some cases
1090 training lasts one day or less, whereas other agencies provide months of training that
1091 includes one to one mentoring [33]. There is some indication that extensive on the job
1092 training and mentoring is a source of expertise for facial examiners carrying out forensic
1093 FIC [46], however the benefits of such training has not, to date, been evaluated for FR
1094 users in the review of candidate lists.

1095 In the absence of empirical studies of FR training efficacy agencies should validate their
1096 training to ensure that the approach is effective at improving operational performance
1097 and reducing the risk of error.

1098 Training should also be evidence-based and incorporate exercises that have been
1099 demonstrated to improve facial comparison performance. Feed-back training is one such
1100 approach for improving facial comparison performance [47]. Feedback should be
1101 provided to trainees on their performance during training tasks, which allows learning
1102 from mistakes and indicates when and why trainees have performed well at a task. There
1103 is also some evidence that collaborative working on facial comparison exercises can
1104 provide a training benefit, particularly for lower performers [48].

1105 Given the limited effectiveness of short, one-off training courses and the absence of data
1106 supporting the use of longer-term training and mentoring for FR users, agencies should
1107 supplement training with screening tests for selection prior to training (section 7.2 and
1108 task-specific competency testing for completion after training (section 7.4).

1109

1110 **7.4 Competency Testing**

1111 In addition to screening tests and training, FR users should also participate in ongoing
1112 competency testing. Competency testing is used to demonstrate that an individual can
1113 reliably and accurately complete a particular task, in this case FR review and related
1114 sub-processes (e.g. processing of probe image, review of candidate list, and verification
1115 of results).

1116 FR users should be competency tested after selection and training, and prior to
1117 undertaking FR reviews.

1118 It is preferable for competency tests to be conducted using test items of known ground-
1119 truth, however, in some instances competency may be demonstrated through on-the-job
1120 training, such as during workplace mentoring.

1121 Competency tests should encapsulate all of the processes an FR user is expected to
1122 undertake and should be conducted according to local policies and procedures. When
1123 designing competency tests, in addition to ensuring they are task-relevant, agencies
1124 should also ensure that the purpose of the test is clearly specified and is testable (i.e.
1125 can be assessed as pass or fail). Tests may also require multiple measures of accuracy
1126 for evaluation of results. The following measures of accuracy for human review may be
1127 useful:

- 1128 • True and false positive identification rates
- 1129 • True and false negative identification rates

1130 The stimuli used in a test should be representative of the material FR users will encounter
1131 operationally with consideration given to the difficulty of the test, ensuring that the test is
1132 sufficiently challenging to provide a meaningful test of competency but not so difficult that
1133 a pass is unachievable.

1134 For further advice on the design of human performance tests in forensic disciplines see
1135 [49].

1136 Agencies should have a documented procedure for competency testing that defines at
1137 what intervals ongoing competency testing should occur and what action should be taken
1138 if an FR user's competency has lapsed.

1139

1140 **8 VALIDATION**

1141 The final output from an FR process is derived from multiple interactions between
1142 automated computer components and human users, which can all have an impact on
1143 the accuracy and reliability of the result [30].

1144 To ensure that the FR process is fit for purpose agencies should consider carrying out
1145 an end-to-end validation of the entire FR process, including an evaluation of the
1146 performance of the algorithm, the competency of the users and the impact of any
1147 interactions between the automated components and the users. Such a validation study
1148 encompasses the entire FR method, including the testing processes discussed in section
1149 5.3 (Know your system) and section 7.4 (Competency testing) of this guideline.

1150 For agencies intending to gain accreditation for their FR process under ISO 17025:2017
1151 standard, it is a requirement of the standard that all methods are formally validated prior
1152 to implementation [50].

1153 ENFSI Guidelines for the single laboratory Validation of Instrumental and Human
1154 Based Methods in Forensic Science [51] provides guidance on validation for both
1155 quantitative and qualitative methods.

1156

1157 In addition to the requirements of ISO 17025:2017 [50] and validation guidance
1158 provided by ENFSI [51], when designing validation studies for FR processes the
1159 following should also be considered:

- 1160 • Validation of FR processes should only be conducted using ground truth
1161 material that is representative of the types of images that will be encountered in
1162 cases.
- 1163 • A formal procedure for the planning, carrying out, reporting and approval of the
1164 validation exercise should be documented prior to starting the validation.
- 1165 • When documenting the validation plan the following should be clearly
1166 established:
 - 1167 ○ The requirements of the end-users of the method, which may include
1168 the FR user, the investigator and the wider criminal justice system.
 - 1169 ○ The specifications for how the method meets the end-user's
1170 requirements. Specifications should be single testable statements that
1171 can be tested during the validation exercise.
 - 1172 ○ Acceptance criteria that determine whether testing has satisfied the
1173 specifications of the end-user requirements.
- 1174 • Validation tests should only be conducted by competent FR users.

1175

1176 Any major changes to agency procedures or updates to software, in particular the FR
1177 algorithm may require all or part of the validation exercise to be repeated.

1178

1179

1180

1181 9 REFERENCES

- 1182 [1] European Network of Forensic Science Institutes, "ENFSI Best Practice Manual
1183 for Facial Image Comparison," 2018.
- 1184 [2] P. Grother, M. Ngan, and K. Hanaoka, "Face Recognition Vendor Test (FRVT)
1185 Part 2: Identification," 2019.
- 1186 [3] EU Directive, *2016/680 Article 10*. Official Journal, 2016.
- 1187 [4] P. Grother, A. Hom, M. Ngan, and K. Hanaoka, "Ongoing Face Recognition
1188 Vendor Test (FRVT) Part 5 : Face Image Quality Assessment," 2021.
- 1189 [5] Facial Identification Scientific Working Group, "Facial Recognition Systems
1190 Operation Assurance : Image Quality Assessment," 2021.
- 1191 [6] P. Grother, M. Ngan, and K. Hanaoka, "Face Recognition Vendor Test (FRVT)
1192 Part 3 : Demographic Effects," 2019.
- 1193 [7] Facial Identification Scientific Working Group, "Understanding and Testing for
1194 Face Recognition Systems Operation Assurance," 2020.
- 1195 [8] Facial Identification Scientific Working Group, "Facial Recognition Systems
1196 Methods and Techniques," 2013.
- 1197 [9] *Information technology - Biometric data interchange formats - Part 5: Face image
1198 data*. ISO/IEC 19794-5:2011, 2011.
- 1199 [10] Facial Identification Scientific Working Group, "Standard guide for capturing facial
1200 images for use with facial recognition systems," 2019.
- 1201 [11] National Police Improvement Agency, "Police Standard for Still Digital Image
1202 Capture and Data Interchange of Facial / Mugshot and Scar , Mark & Tattoo
1203 Images," no. May, 2007.
- 1204 [12] P. Grother, G. Quinn, and J. Phillips, "Multiple-Biometric Evaluation (MBE) 2010:
1205 Report on the Evaluation of 2D Still-Image Face Recognition Algorithms," *NIST
1206 Interag. Rep. 7709*, 2011.
- 1207 [13] *Standard Guide for Forensic Digital Image Processing*. ASTM E2825-21, 2021.
- 1208 [14] Facial Identification Scientific Working Group, "Physical Stability of Facial
1209 Features of Adults," 2021.
- 1210 [15] K. Ricanek, S. Bhardwaj, and Michael Sodomsky, "A Review of Face
1211 Recognition against Longitudinal Child Faces," *BioSiG 2015*, pp. 15–26, 2015.
- 1212 [16] P. Grother and M. Ngan, "Face Recognition Vendor Test (FRVT)," *Natl. Inst.
1213 Stand. Technol.*, 2014.
- 1214 [17] P. Grother, M. Ngan, and K. Hanaoka, "NISTIR 8271 DRAFT SUPPLEMENT Face
1215 Recognition Vendor Test (FRVT) Part 2: Identification," 2022.
- 1216 [18] TELEFI Project, "Summary Report of the project 'Towards the European Level of
1217 Exchange of Facial Images,'" 2021.
- 1218 [19] R. Heyer, C. Semmler, and A. T. Hendrickson, "Humans and Algorithms for Facial
1219 Recognition: The Effects of Candidate List Length and Experience on
1220 Performance," *J. Appl. Res. Mem. Cogn.*, vol. 7, no. 4, pp. 597–609, 2018.
- 1221 [20] A. Strathie, A. Mcneill, and D. White, "In the Dock: Chimeric Image Composites

- 1222 Reduce Identification Accuracy,” *Appl. Cogn. Psychol.*, vol. 26, no. 1, pp. 140–
1223 148, 2012.
- 1224 [21] A. Strathie and A. McNeill, “Facial Wipes don’t Wash: Facial Image Comparison
1225 by Video Superimposition Reduces the Accuracy of Face Matching Decisions,”
1226 *Appl. Cogn. Psychol.*, vol. 30, no. 4, pp. 504–513, 2016.
- 1227 [22] R. Moreton and J. Morley, “Investigation into the use of photoanthropometry in
1228 facial image comparison,” *Forensic Sci. Int.*, vol. 212, no. 1–3, pp. 231–237, 2011.
- 1229 [23] R. S. S. Kramer, J. Mulgrew, and M. G. Reynolds, “Unfamiliar face matching with
1230 photographs of infants and children,” *PeerJ*, vol. 2018, no. 6, pp. 1–25, 2018.
- 1231 [24] D. Michalski, R. Heyer, and C. Semmler, “The performance of practitioners
1232 conducting facial comparisons on images of children across age,” *PLoS One*, vol.
1233 14, no. 11, pp. 1–17, 2019.
- 1234 [25] C. Wilkinson, “Juvenile facial reconstruction,” in *Craniofacial Identification*, C.
1235 Wilkinson and C. Rynn, Eds. Cambridge University Press, 2012, pp. 254–260.
- 1236 [26] Facial Identification Scientific Working Group, “Standard Practice / Guide for
1237 Image Processing to Improve Automated Facial Recognition Search
1238 Performance,” 2020.
- 1239 [27] Facial Identification Scientific Working Group, “Minimum Guidelines for Facial
1240 Image Comparison Documentation,” 2020.
- 1241 [28] A. W. Young and A. M. Burton, “Are We Face Experts?,” *Trends Cogn. Sci.*, vol.
1242 22, no. 2, pp. 100–110, 2018.
- 1243 [29] D. White, J. D. Dunn, A. C. Schmid, and R. I. Kemp, “Error rates in users of
1244 automatic face recognition software,” *PLoS One*, vol. 10, no. 10, 2015.
- 1245 [30] A. Towler, R. I. Kemp, and D. White, “Unfamiliar Face Matching Systems in
1246 Applied Settings,” in *Face Processing: Systems, Disorders and Cultural
1247 Difficulties*, M. Bindemann and A. M. Megreya, Eds. Nova Science Publishers,
1248 2017, pp. 21–40.
- 1249 [31] Facial Identification Scientific Working Group, “Guide for Role-Based Training in
1250 Facial Comparison,” 2020.
- 1251 [32] D. White, A. Towler, and R. I. Kemp, “Understanding professional expertise in
1252 unfamiliar face matching,” in *Forensic Face Matching: Research and practice*, M.
1253 Bindemann, Ed. Oxford University Press, 2021.
- 1254 [33] R. Moreton, C. Havard, A. Strathie, and G. Pike, “An International Survey of
1255 Applied Face-Matching Training Courses,” *Forensic Sci. Int.*, vol. 327, p. 110947,
1256 2021.
- 1257 [34] R. Moreton, G. Pike, and C. Havard, “A task- and role-based perspective on super-
1258 recognizers: Commentary on ‘Super-recognizers: From the laboratory to the world
1259 and back again,’” *Br. J. Psychol.*, p. bjop.12394, Mar. 2019.
- 1260 [35] A. Towler *et al.*, “Do professional facial image comparison training courses work?,”
1261 *PLoS One*, vol. 14, no. 2, pp. 1–17, 2019.
- 1262 [36] R. Moreton, “Forensic face matching: Procedures and application,” in *Forensic
1263 Face Matching*, M. Bindemann, Ed. Oxford University Press, 2021.
- 1264 [37] R. Russell, B. Duchaine, and K. Nakayama, “Super-recognizers: people with
1265 extraordinary face recognition ability.,” *Psychon. Bull. Rev.*, vol. 16, no. 2, pp. 252–

- 1266 257, 2009.
- 1267 [38] S. Bate and J. J. Tree, "The definition and diagnosis of developmental
1268 prosopagnosia," *Q. J. Exp. Psychol.*, vol. 70, no. 2, pp. 193–200, 2017.
- 1269 [39] A. M. Burton, D. White, and A. McNeill, "The Glasgow Face Matching Test.,"
1270 *Behav. Res. Methods*, vol. 42, no. 1, pp. 286–291, 2010.
- 1271 [40] M. C. Fysh and M. Bindemann, "The Kent Face Matching Test," *Br. J. Psychol.*,
1272 pp. 1–13, 2017.
- 1273 [41] D. White, D. Guilbert, V. P. L. Varela, R. Jenkins, and A. M. Burton, "GFMT2: A
1274 psychometric measure of face matching ability," *Behav. Res. Methods*, 2021.
- 1275 [42] J. D. Dunn, S. Summersby, A. Towler, J. Davis, and D. White, "UNSW Face Test:
1276 A screening tool for super-recognizers," pp. 1–19, 2020.
- 1277 [43] M. Ramon, A. K. Bobak, and D. White, "Super-recognizers: From the lab to the
1278 world and back again," *Br. J. Psychol.*, vol. 110, no. 3, pp. 461–479, 2019.
- 1279 [44] M. C. Fysh, L. Stacchi, and M. Ramon, "Differences between and within
1280 individuals, and subprocesses of face cognition: implications for theory, research
1281 and personnel selection," *R. Soc. Open Sci.*, vol. 7, no. 9, p. 200233, 2020.
- 1282 [45] R. Moreton, "Expertise in Applied Face Matching: Training, Forensic Examiners,
1283 Super Matchers and Algorithms," The Open University, 2021.
- 1284 [46] A. Towler, D. White, and R. Kemp, "Can face identification ability be trained?
1285 Evidence for two routes to expertise," in *Forensic Face Matching*, M. Bindemann,
1286 Ed. Oxford University Press, 2021.
- 1287 [47] D. White, R. I. Kemp, R. Jenkins, and A. M. Burton, "Feedback training for facial
1288 image comparison," *Psychon. Bull. {&} Rev.*, vol. 21, no. 1, pp. 100–106, 2014.
- 1289 [48] A. J. Dowsett and A. M. Burton, "Unfamiliar face matching: Pairs out-perform
1290 individuals and provide a route to training," *Br. J. Psychol.*, vol. 106, no. 3, pp.
1291 433–445, 2015.
- 1292 [49] K. A. Martire and R. I. Kemp, "Considerations when designing human performance
1293 tests in the forensic sciences," *Aust. J. Forensic Sci.*, pp. 1–17, Nov. 2016.
- 1294 [50] *General requirements for the competence of testing and calibration laboratories.*
1295 ISO/IEC 17025:2017, 2017.
- 1296 [51] European Network of Forensic Science Institutes, "Guidelines for the single
1297 laboratory Validation of Instrumental and Human Based Methods in Forensic
1298 Science," no. 001, pp. 1–31, 2014.
- 1299
- 1300
- 1301

1302

1303 **10 AMENDMENTS AGAINST PREVIOUS VERSION**

1304 None.

1305

1306

1307

1308

DRAFT DOCUMENT