

Rank-1 Identity Consistency: Gallery Membership Classification for Operational 1:N Contexts

Anonymous CVPR submission

Paper ID ****

Abstract

001 *In real-world deployments, 1:N face identification systems*
 002 *face a fundamental question: is a probe image’s identity*
 003 *enrolled in the gallery or not? We propose the first ap-*
 004 *proach to using consistency of rank-1 identity across mul-*
 005 *tiple matchers as a method to classify the result of 1:N*
 006 *search as in-gallery (ING) or out-of-gallery (OOG). Our*
 007 *“1-consistency” method classifies a probe as ING if all*
 008 *matchers return the same rank-1 identity, and OOG oth-*
 009 *erwise. We compare its performance to two threshold-based*
 010 *methods: score-thresholding (using raw similarity scores)*
 011 *and gap-thresholding (using the score difference between*
 012 *rank-1 and rank-2 identities).*

013 *We evaluate these methods across 12 experimental con-*
 014 *figurations that systematically vary image quality, gallery*
 015 *enrollment structure, and demographic composition. On*
 016 *average, 1-consistency achieves the highest ING recall*
 017 *(92.0% vs. 72.1% score, 79.6% gap), highest overall accu-*
 018 *racy (92.8% vs. 81.4% score, 89.6% gap), and lowest de-*
 019 *mographic disparities (4.6 pp vs. 10.5 pp score, 8.5 pp gap).*
 020 *Under degraded-probe conditions—the most operationally*
 021 *relevant scenario—1-consistency achieves win margins av-*
 022 *eraging 16.1 pp for ING recall and 6.1 pp for overall accu-*
 023 *racy. With this combination of quality robustness and dem-*
 024 *ographic fairness, 1-consistency is well-suited for real-world*
 025 *contexts like law enforcement investigations—and has the*
 026 *potential to reduce wrongful arrests and enable better allo-*
 027 *cation of investigative efforts.*

028 1. Introduction

029 In operational contexts, 1:N face identification systems con-
 030 tend with two major challenges. The first is image quality.
 031 Probes frequently originate from surveillance footage and
 032 exhibit poor resolution, blur, and other low-quality condi-
 033 tions. Galleries have traditionally been controlled-capture
 034 datasets (e.g., databases of state arrest mugshots or driver’s
 035 licenses), but are increasingly web-scraped. For example,

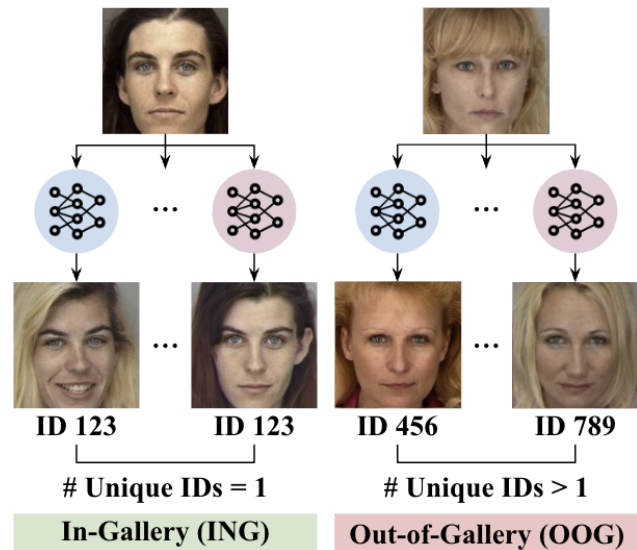


Figure 1. For each probe, multiple matchers perform a 1:N search and return a rank-1 identity. If all rank-1 identities are the same, the probe is labeled in-gallery (ING), else, out-of-gallery (OOG).

036 Clearview AI’s web-scraped database of 70+ billion face
 037 images has been used in over 2 million law enforcement
 038 searches [1, 8]. Since these images are sourced from “news
 039 media, mugshot websites, public social media, and many
 040 other open sources,” their quality is unknown [1]. Quality
 041 has a drastic effect on identification performance: state-of-
 042 the-art systems with near-perfect accuracy on high-quality
 043 images exhibit not only higher FPIRs (upwards of 20%) but
 044 also substantial demographic disparities when probe and/or
 045 gallery quality is poor [4, 20, 21].

046 The second challenge is uncertainty about gallery mem-
 047 bership. An investigator searching surveillance footage
 048 cannot know in advance whether a subject of interest has
 049 ever been photographed and enrolled. This open-set iden-
 050 tification problem—determining if a probe’s true identity
 051 is present in the gallery—is fundamental to investigative
 052 workflows. If a probe is out-of-gallery (OOG), any re-

053	turned candidate is by definition a false positive, which	104
054	can lead to time being wasted investigating the wrong per-	105
055	son and, in the worst case, a wrongful arrest. Conversely,	106
056	if the probe is in-gallery (ING) but misclassified as OOG,	107
057	valid investigative leads may be missed. Under ideal con-	108
058	ditions (high-quality imagery with guaranteed probe en-	109
059	rollment), modern matchers achieve over 99.9% accuracy	110
060	even in galleries containing 12–26 million images of 12	111
061	million identities [12]. However, distinguishing ING from	112
062	OOG probes under realistic operational conditions—where	113
063	quality is degraded and gallery enrollment is unbalanced—	114
064	remains challenging yet essential.	115
065	Existing approaches to gallery membership classification	116
066	often rely on similarity scores to determine whether to ac-	117
067	cept or reject returned candidates, e.g., by applying thresh-	118
068	olds [13, 15] or estimating probabilities [22–24]. However,	119
069	since score distributions vary substantially across matcher	120
070	architectures, quality conditions, and gallery compositions,	121
071	and exhibit significant demographic differentials, a single	122
072	global threshold may produce disparate outcomes. That is,	123
073	a threshold calibrated on one system or dataset may per-	124
074	form poorly if gallery demographics shift, image quality	125
075	degrades, or the matcher is retrained.	126
076	1.1. Contribution of this Work	127
077	In this work, we investigate a multi-matcher fusion strategy	128
078	based on rank stability. While each independently trained	129
079	matcher instance operates in a distinct embedding space—	130
080	and may therefore produce varying similarity scores or can-	131
081	didate orderings—consensus on the top-ranked identity sig-	132
082	nals high certainty. The proposed metric, 1-consistency,	133
083	classifies a probe based on whether all matcher instances	134
084	return the same rank-1 identity. Our hypothesis is that	135
085	ING probes should produce stable rank-1 predictions, while	136
086	OOG probes will exhibit greater variability as different	137
087	matcher instances select different gallery impostors.	138
088	To our knowledge, this is the first work to use cross-	139
089	matcher identity consensus for ING/OOG prediction in op-	140
090	erational 1:N identification. It is also one of few works to	141
091	evaluate ING/OOG prediction under controlled variations	142
092	in image quality and gallery composition, and to report	143
093	demographic-specific performance differentials. Across the	144
094	tested conditions, 1-consistency offers a better balance of	145
095	in-gallery retention and out-of-gallery rejection <i>and</i> more	146
096	consistent cross-demographic performance than threshold-	147
097	-based methods—properties critical for high-stakes inves-	148
098	tigative deployments.	149
099	2. Related Work	150
100	Prior work on distinguishing in-gallery from out-of-gallery	151
101	probes falls into several categories.	152
102	Score-thresholding approaches apply fixed or adaptive	153
103	thresholds on similarity scores to accept or reject candi-	
	dates [13, 15]. These methods are conceptually simple but	
	highly sensitive to matcher architecture, quality distribu-	
	tions, and gallery composition.	
	Statistical modeling and Extreme Value Theory (EVT)	
	methods model tail distributions of match scores to estimate	
	the probability that a score originates from a known versus	
	unknown identity [22–24]. Approaches such as W-SVM,	
	Compact Abating Probability (CAP), and Extreme Value	
	Machine (EVM) produce calibrated decision boundaries but	
	require distributional assumptions that may not hold across	
	operational conditions.	
	Deep open-set recognition methods modify neural net-	
	work architectures or loss functions to enable rejection of	
	unknown identities [2, 3]. These methods calibrate final-	
	layer activations using EVT but remain dependent on score-	
	-based thresholding that is sensitive to quality variation.	
	Feature-fusion and multi-image aggregation methods	
	stabilize matching by aggregating scores or features across	
	multiple enrolled images [7, 16, 25]. Quality-aware score	
	fusion improves robustness under pose or illumination vari-	
	ation but remains fundamentally limited by reliance on sim-	
	ilarity score distributions.	
	Rank-pattern learning approaches train classifiers to dis-	
	tinguish in-gallery from out-of-gallery probes using rank	
	distributions [5]. This approach shows promise under de-	
	graded conditions but is only applicable to identities with	
	more than one image enrolled in the gallery.	
	In contrast to the approaches described above, we pro-	
	pose a rank-level fusion method that operates on identity	
	consensus across multiple independently trained matcher	
	instances. By leveraging rank stability rather than score cal-	
	ibration, our approach demonstrates improved robustness to	
	quality degradation and more consistent performance across	
	demographic groups.	
	3. Methodology	
	In this section, we describe our experimental methodol-	
	ogy, including the datasets and degradations used, the	
	matcher configurations, the per-probe processing pipeline,	
	and the classification and evaluation framework. Together,	
	these components allow us to evaluate the robustness of	
	1-consistency across a range of gallery configurations and	
	quality conditions.	
	3.1. Datasets	
	All experiments use images from the MORPH dataset,	
	which consists of mugshot images captured under con-	
	trolled conditions (i.e., nominally frontal pose, neutral ex-	
	pression, and consistent lighting), along with annotations	
	for race, gender, and age. The images vary in native reso-	
	lution, ranging from 116×154 to 507×631 pixels, with a	
	modal resolution of 400×480 . Each image is close-cropped	

154 around the face region using dlib [17], then aligned and re-
155 sized to 112×112 pixels to meet CNN input requirements.

156 We focus on four demographic cohorts within MORPH:
157 Black Female (BLF), Black Male (BLM), White Female
158 (WHF), and White Male (WHM). From each cohort, we se-
159 lect 3,400 identities, forming a “base identity set” of 13,600
160 total identities used throughout the experiments. For each
161 identity, the most recent image serves as the probe.

162 All remaining (non-probe) images for each identity con-
163 stitute the pool from which gallery sets are constructed.
164 Across all MORPH-based experiments, every gallery is ei-
165 ther an equal set or a strict superset of the probe identities.
166 Thus, any MORPH probe’s identity is always present in the
167 gallery it is matched against; what differs across experi-
168 ments is the demographic composition of the gallery and
169 the degree of control over per-identity enrollment.

170 3.2. Degradations

171 To study quality robustness under controlled and demo-
172 graphically consistent settings, we apply synthetic degrada-
173 tions to MORPH images at both the probe and gallery levels
174 in the experiments where quality is varied. Two degradation
175 types are used, each corresponding to a common problem in
176 probe images taken from surveillance video. First, Gaussian
177 blur is applied to simulate defocus; we choose $\sigma = 4$ as a
178 mid-range severity within the standard NIST 1–7 blur scale
179 [14]. Second, downsampling–upsampling simulates native
180 low-resolution capture: an image is reduced to 18×18 pix-
181 els and then resized back to 112×112 .

182 Fig. 2 illustrates the relative quality of the original ver-
183 sus degraded images. Previous work [21] has shown that
184 $\sigma = 4$ blur and 18×18 resolution are equivalent in terms
185 of impact on similarity scores (when performing 1:1 self-
186 matching of each degraded version to the original using Arc-
187 Face and AdaFace, the matchers used in this work) and on
188 face image quality assessment scores using MagFace [18],
189 a top-performing metric in NIST evaluations [11].

Table 1. Gallery composition for Experiments 3, 4, and 5.

Exp	Dem	# IDs	# Imgs	Mean	Med
3	BLF	3,400	14,664	4.3	3
	BLM	3,400	19,045	5.6	4
	WHF	3,400	10,228	3.0	2
	WHM	3,400	14,173	4.2	3
4	BLM	21,106	155,715	7.4	5
5	BLF	710	2,909	4.1	3
	BLM	9,414	69,660	7.4	5
	WHF	665	1,891	2.8	2
	WHM	8,846	35,306	4.0	3
	HSM	6,563	24,730	3.8	3
	HSF	494	1,467	3.0	3

190 3.3. Matchers

191 We conducted preliminary experiments with both AdaFace
192 and ArcFace. However, ArcFace exhibited substantially
193 worse performance in the quality-degraded scenarios that
194 are central to this work. Given that our focus is on quality-
195 robust gallery membership prediction, we present results
196 only for AdaFace.

197 We use 10 independently trained instances of AdaFace,
198 each denoted as a *run*. While the runs share the same base
199 architecture (a ResNet-100 backbone with adaptive mar-
200 gin loss trained on Glint360K [10]), their learned embed-
201 ding spaces are distinct. Consequently, the templates they
202 produce are not interchangeable. For example, perform-
203 ing matching with run X on templates generated by run Y
204 would yield authentic (mated) similarity scores resembling
205 those typically observed for impostor (non-mated) pairs, as
206 demonstrated in [6].

207 3.4. Per-Probe Processing

208 In all experiments, each MORPH probe is enrolled in the
209 gallery it is matched against, producing both mated and
210 non-mated candidates for that probe.

211 To obtain a probe’s in-gallery (ING) outcome, we use the
212 full candidate list (mated + non-mated). To simulate an out-
213 of-gallery (OOG) scenario, we simply discard the mated re-
214 sults and retain only the non-mated candidates—effectively
215 treating the probe as if it had never been enrolled in the
216 gallery. This provides a controlled setup in which ING and
217 OOG cases can be compared directly while holding all other
218 factors constant.

219 3.4.1. Mean Rank-1 Score (Score-Thresholding)

220 For each probe in a given configuration, the rank-1 re-
221 sult consists of a single image and its similarity score—the
222 “most similar” image to the probe for that run. Across the
223 10 independently trained runs, each probe therefore has 10
224 rank-1 scores. We average these to obtain a single Mean
225 Rank-1 Score per probe per configuration.

226 3.4.2. k-Consistency (1-Consistency)

227 From each rank-1 result we also extract the identity of
228 the corresponding image, which we refer to as the *rank-1*
229 *ID*—the “most similar identity” to the probe for that

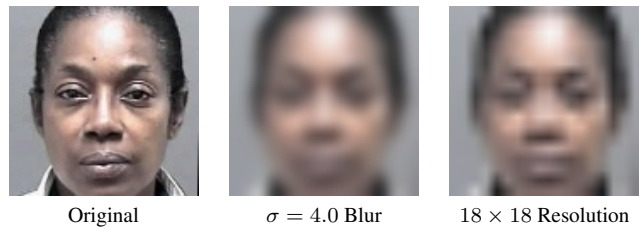


Figure 2. Example MORPH image and degraded variants.

230 run. Each probe thus has 10 rank-1 IDs across the 10
231 runs. We count the number of unique identities appearing
232 among these 10 values and denote this count as k . A probe
233 with k unique rank-1 IDs is said to be k -consistent, where
234 $1 \leq k \leq 10$.

235 3.4.3. Mean ID Gap (Gap-Thresholding)

236 We also measure how strongly the rank-1 identity is pre-
237 ferred over all other identities in the gallery. To do this,
238 we identify the “second most similar identity”—the iden-
239 tity (different from the rank-1 ID) that yields the highest
240 similarity score for that run.

241 This requires special handling when a gallery contains
242 multiple images of the same identity. In such cases, sev-
243 eral top-ranked images may all belong to the rank-1 identity.
244 For example, the rank-2, rank-3, ..., rank- X results may all
245 correspond to the same identity as the rank-1 result. Only
246 at rank $X + 1$ do we encounter an image from a different
247 identity. We refer to that identity as the *rank-2 ID*.

248 Our goal is to quantify the score “distance” between
249 the rank-1 identity and the strongest competing identity.
250 Several definitions are possible (e.g., comparing averages
251 across all images of each identity, or comparing extremal
252 image pairs). For simplicity and consistency, we use the
253 following process: (1) take the best (highest) score among
254 all images belonging to the rank-1 identity (i.e., the rank-1
255 score), (2) take the best score among all images belonging
256 to the rank-2 identity, (3) subtract the latter from the for-
257 mer. We record this value as the ID Gap for that run. Each
258 probe produces 10 such gap values (one per run), which we
259 average to obtain a single Mean ID Gap per probe per con-
260 figuration.

261 3.5. Classification and Evaluation

262 We perform binary classification of ING/OOG. In each ex-
263 perimental configuration, we calculate the three values— k -
264 consistency, mean rank-1 score, and mean ID gap—for each
265 probe. We determine values of k , mean rank-1 score, and
266 mean ID gap that best separate probes in the baseline con-
267 ditions of Experiment 1. These thresholds are then used
268 across all other configurations, and our classification rule
269 is as follows: if a probe has a below-threshold value, it is
270 classified as OOG; if above, it is classified as ING.

271 Because this task involves two classes that are perfectly
272 balanced within each configuration, performance is most
273 naturally expressed in terms of recall for each class:

$$274 \quad R(\text{ING}) = \frac{\# \text{ ING probes correctly labeled ING}}{\# \text{ ING probes}} \quad (1)$$

$$275 \quad R(\text{OOG}) = \frac{\# \text{ OOG probes correctly labeled OOG}}{\# \text{ OOG probes}} \quad (2)$$

276 In a balanced binary setting, these recall values share
277 identical denominators and are therefore directly compara-

278 ble. Because the task is binary and the ING and OOG sets
279 are perfectly balanced, recall alone fully characterizes clas-
280 sification performance: the two recall values share iden-
281 tical denominators, directly reflect the misclassification rates,
282 and determine overall accuracy. Precision, while definable,
283 does not provide additional insight in this setting, since all
284 relevant error patterns are already captured by recall.

285 Misclassification rates follow immediately from the re-
286 calls. The fraction of ING probes misclassified as OOG is
287 $1 - R(\text{ING})$, and the fraction of OOG probes misclassified
288 as ING is $1 - R(\text{OOG})$. Because the classes are balanced,
289 overall accuracy for any method is simply the average of the
290 two recalls.

291 4. Results and Analysis

292 To evaluate ING/OOG classification robustness under in-
293 creasingly realistic operational conditions, we design five
294 experiments that systematically vary quality degradation,
295 cross-demographic matching, or gallery enrollment size.

296 In this section, we describe the setup and results for each
297 experiment individually. For each, we specify the experi-
298 mental conditions, provide the gallery configuration details,
299 and analyze classification performance.

300 4.1. Experiment 1: Baseline

301 We establish classification thresholds in ideal conditions.
302 Each probe has two mated gallery instances, and is only
303 matched against same-demographic images. Probe and
304 gallery images are original mugshot-quality.

305 Fig. 3 shows the k -consistency, mean rank-1 score, and
306 mean ID gap distributions for ING and OOG probes. For
307 each method, we determine the threshold that best separates
308 the pair of distributions. Thresholds are indicated by solid
309 blue lines: $k = 1$ for 1-consistency, mean score = 0.454
310 for score-thresholding, and mean gap = 0.102 for gap-
311 thresholding. For the score and gap plots, the means of the
312 ING/OOG distributions are shown in solid green/red.

313 These thresholds are applied to all subsequent experi-
314 ments without modification. For subsequent experiments,
315 we quantify distribution shifts and classification perfor-
316 mance using tabular results rather than distribution plots to
317 conserve space. Tab. 2 provides per-experiment perfor-
318 mance metrics (ING recall, OOG recall, and overall accu-
319 racy) for each method.

320 **Results.** Under baseline conditions, all three methods
321 achieve near-perfect performance, as expected given that
322 thresholds were selected to optimally separate ING and
323 OOG in this setting. ING recall exceeds 99% for all three
324 methods, with 1-consistency marginally best at 99.8%.
325 OOG recall is 99.1% for both thresholding methods, while
326 1-consistency is lower at 91.8%. Overall accuracy reflects
327 this pattern: threshold-based methods achieve 99.2–99.4%,

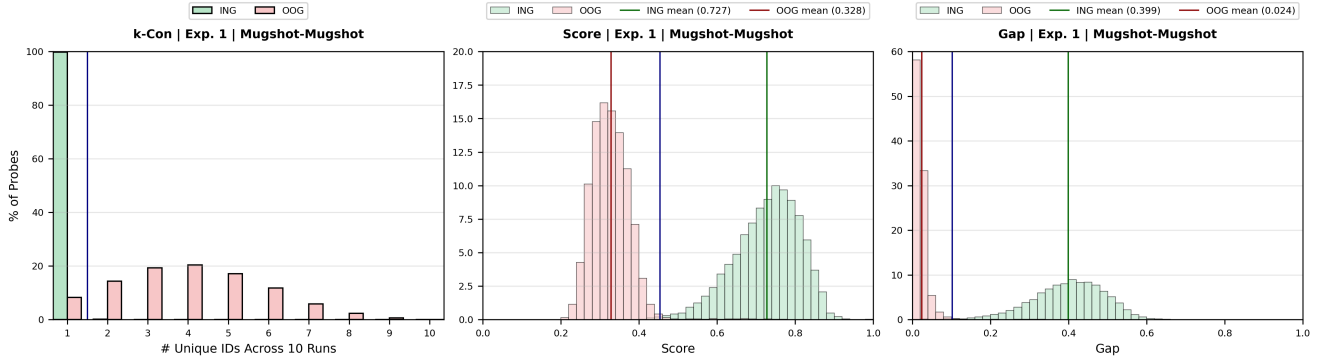


Figure 3. Baseline comparison for good-quality images and single-demographic galleries with balanced enrollment.

Table 2. Average performance metrics for the three classification methods. Best method(s) per configuration and metric highlighted in bold green. Blur (BL) and resolution (RE) degradations indicated where applicable.

Exp	Gal Composition		Degradation		ING Recall			OOG Recall			Accuracy		
	N Imgs	N Dems	Probe	Gal	1-Con	Score	Gap	1-Con	Score	Gap	1-Con	Score	Gap
1	2	1	—	—	99.8	99.7	99.4	91.8	99.1	99.1	95.8	99.4	99.2
2	2	1	BL	—	87.9	44.7	71.4	95.8	99.9	99.6	91.9	72.3	85.5
			RE	—	87.2	47.1	70.6	95.6	99.9	99.6	91.4	73.5	85.1
			—	BL	92.1	45.7	78.6	95.9	99.9	99.7	94.0	72.8	89.2
			—	RE	91.6	49.7	77.2	95.8	99.9	99.7	93.7	74.8	88.5
			BL	BL	83.6	97.0	71.0	90.0	31.0	99.3	86.8	64.0	85.2
			RE	RE	83.8	91.2	70.8	91.1	63.6	99.4	87.5	77.4	85.1
3	2-17	4	—	—	99.9	99.9	99.5	92.1	98.8	99.1	96.0	99.4	99.3
			RE	—	89.3	55.9	73.9	96.1	99.8	99.7	92.7	77.9	86.8
4	1-108	1	—	—	100	100	99.9	94.2	96.1	99.6	97.1	98.1	99.8
			RE	—	87.9	71.9	67.2	96.9	99.7	99.9	92.4	85.8	83.6
5	1-108	6	RE	—	91.3	62.3	75.5	97.2	100	99.9	94.3	81.2	87.7
Averages Across All Configurations					92.0	72.1	79.6	94.4	90.7	99.6	92.8	81.4	89.6

328 while 1-consistency achieves 95.8%. Demographic differ- 345
 329 ences in accuracy are minor across all three methods (1–2 346
 330 percentage points). 347

331 These results demonstrate that under ideal quality and 348
 332 gallery enrollment conditions, distinguishing ING from 349
 333 OOG probes is straightforward using rank-1 scoring infor- 350
 334 mation alone, whether raw similarity scores or score gaps 351
 335 relative to competing identities. However, such ideal con- 352
 336 ditions rarely hold in operational deployments. In the fol- 353
 337 lowing experiments, we evaluate classification performance 354
 338 under increasingly realistic and challenging scenarios: de- 355
 339 graded image quality, varied gallery structures, and cross- 356
 340 demographic matching.

341 4.2. Experiment 2: Variable Image Quality 357

342 **To what extent does quality degradation—in probe im- 358
 343 ages, gallery images, or both—impact classification per- 359
 344 formance?** As in experiment 1, the gallery is single- 360
 361
 362

demographic with two enrolled images per probe. 345

346 We evaluate three quality scenarios: (1) probe-only 347
 348 degradation, representing the most operationally relevant 349
 350 scenario where surveillance probes are matched against 350
 351 controlled-capture mugshot databases, (2) gallery-only 351
 352 degradation, an operationally unlikely scenario included 352
 353 here for completeness, and (3) symmetric degradation, 353
 354 which may arise from the use of uncontrolled web-scraped 354
 355 galleries. For each scenario, we apply both blur (BL: $\sigma = 4$) 355
 356 and resolution (RE: 18×18) degradations. This design iso- 356
 357 lates the individual and joint effects of probe and gallery 357
 358 quality on classification accuracy. 358

359 **Results.** Quality degradation substantially impacts clas- 359
 360 sification performance, with different methods showing 360
 361 distinct failure modes. For probe-only degrada- 361
 362 tion, 1-consistency maintains strong ING recall (87.2– 362
 87.9%). Score-thresholding collapses catastrophically 361
 (44.7–47.1%), while gap-thresholding performs moderately 362

Table 3. Best and worst demographic by accuracy for each method. Diff = percentage point difference.

Exp	Degradation		1-Consistency			Score-Threshold			Gap-Threshold		
	Probe	Gal	Best	Worst	Diff	Best	Worst	Diff	Best	Worst	Diff
1	—	—	BLF	WHF	2.1	BLM	WHF	1.0	BLF	WHF	1.6
2	RE	—	BLM	BLF	7.6	BLM	WHF	9.9	BLM	BLF	12.2
	—	RE	BLM	WHF	2.4	BLM	WHF	12.4	BLM	BLF	12.0
	RE	RE	BLM	BLF	6.1	WHM	BLF	24.2	BLM	BLF	9.5
3	—	—	WHM	WHF	1.3	WHM	WHF	1.1	WHM	WHF	1.9
	RE	—	BLM	BLF	8.0	BLM	WHF	14.3	BLM	BLF	13.5

363 (70.6–71.4%), still substantially lower than 1-consistency.
 364 OOG recall remains high for all methods (95.6–99.9%).
 365 Overall accuracy favors 1-consistency (91.4–91.9%) over
 366 both score-thresholding (72.3–73.5%) and gap-thresholding
 367 (85.1–85.5%). Demographic disparities increase substan-
 368 tially for all methods but remain lowest for 1-consistency
 369 (7.6–8.0 pp) compared to score-thresholding (9.9–12.5 pp)
 370 and gap-thresholding (12.9–13.5 pp).

371 For gallery-only degradation, the pattern is similar:
 372 1-consistency achieves ING recall of 91.6–92.1% while
 373 threshold methods again struggle (45.7–49.7% for score,
 374 77.2–78.6% for gap), demonstrating that both probe
 375 and gallery quality independently affect threshold-based
 376 method performance.

377 For symmetric degradation (both probe and gallery de-
 378 graded), score-thresholding achieves the highest ING re-
 379 call (91.2–97.0%) but catastrophically fails on OOG re-
 380 call (31.0–63.6%), resulting in poor overall accuracy (64.0–
 381 77.4%). Gap-thresholding achieves near-perfect OOG re-
 382 call (99.3–99.4%) but lower ING recall (70.8–71.0%).
 383 1-consistency provides the most balanced performance,
 384 achieving the highest overall accuracy (86.8–87.5%) and
 385 lowest demographic disparities (6.2–7.4 pp) despite not
 386 winning either individual recall metric.

387 These results show that threshold-based methods are
 388 highly sensitive to quality degradation, with score-
 389 thresholding primarily failing on ING recall and occasion-
 390 ally on OOG recall under dual degradation. In contrast,
 391 1-consistency maintains more balanced performance across
 392 quality conditions, trading modest OOG recall for substan-
 393 tially higher ING recall and lower demographic disparities.

394 Given that probe-only degradation is the most opera-
 395 tionally relevant scenario and that blur and resolution degra-
 396 dations yield comparable results within each method for
 397 this case, subsequent experiments use resolution-degraded
 398 probes.

399 4.3. Experiment 3: Variable Gallery Composition

400 **To what extent does gallery composition—variable en-**
 401 **rollment size and demographic heterogeneity—impact**
 402 **classification performance?** In this experiment, the gallery

grows in size: it comprises 58,110 total images of all four
 demographics. Each identity has 2–17 enrolled images, and
 the average count differs across demographics: 3 for WHF,
 4 for WHM and BLF, and 6 for BLM. Note that the number
 of represented *identities* is still balanced across demograph-
 ics.

We consider the impact of gallery composition for
 original mugshot-quality probes and resolution-degraded
 probes. For original probes, we compare results to those
 of experiment 1. For degraded probes, we compare results
 to those of experiment 2.

Results. When probes are mugshot-quality, gallery compo-
 sition has negligible impact: across all method/metric pair-
 ings, the average difference from Experiment 1 baseline is
 only 0.1 percentage points. All three methods achieve near-
 perfect performance (96.0–99.4% accuracy).

When probes are resolution-degraded, gallery compo-
 sition introduces modest variation: the average difference
 from Experiment 2 is 2.5 percentage points. The largest
 shifts occurs for ING recall, which *improves* for all meth-
 ods: up 2.1 pp for 1-consistency, 3.3 pp for gap, and 8.8
 for score (though the latter two still remain low overall). 1-
 consistency still maintains the highest ING recall (89.3%)
 and overall accuracy (92.7%), with the lowest demographic
 disparities (8.0 pp vs. 13.7–14.3 pp for threshold methods).

These results confirm that, thus far, image quality re-
 mains the dominant factor determining classification per-
 formance.

4.4. Experiment 4: Larger-Scale Gallery

This experiment evaluates classification under a large-scale,
 uncontrolled gallery structure: 1–108 images per identity,
 single-demographic composition (BLM only), and probe-
 only degradation. The gallery includes 21,106 BLM iden-
 tities—the 3,400 base identities (each contributing all re-
 maining non-probe images, typically 2–17) plus 17,706
 additional BLM identities drawn from the full MORPH
 dataset—yielding 155,715 total images. This configuration
 maximizes both gallery size and within-identity variability
 while maintaining demographic homogeneity, stress-testing
 classification methods under large-scale conditions repre-

443 tentative of operational databases.

444 **Results.** Under mugshot-quality conditions, all three meth- 495
445 ods achieve near-perfect performance, comparable to Ex- 496
446 periment 1 baseline (controlled single-demographic gallery 497
447 with 2 images per identity). Comparing Experiment 4 to
448 Experiment 1: 1-consistency ING recall remains essen-
449 tially identical (100.0% vs. 99.8%), score-thresholding
450 ING recall increases marginally (100.0% vs. 99.7%), and
451 gap-thresholding ING recall increases slightly (99.9% vs.
452 99.4%). Overall accuracy ranges from 97.1–99.8% in Ex-
453 periment 4 versus 95.8–99.4% in Experiment 1. The large
454 gallery size (155,715 images, 21,106 identities) does not
455 impair classification when image quality is high.

456 Under degraded probe conditions, we compare to Exper- 508
457 iment 2 (controlled single-demographic gallery, resolution- 509
458 degraded probes). 1-consistency ING recall remains nearly 510
459 stable (87.9% in Experiment 4 vs. 87.2% in Experi- 511
460 ment 2), demonstrating robustness to gallery scale. Score- 512
461 thresholding ING recall improves substantially from 47.1% 513
462 to 71.9%—a 24.8 pp increase. Gap-thresholding ING 514
463 recall decreases slightly from 70.6% to 67.2%—a 3.4 515
464 pp decline. Notably, this represents the first configura- 516
465 tion where gap-thresholding performs worse than score- 517
466 thresholding on ING recall. OOG recall remains high 518
467 for all methods (96.9–99.9%). Overall accuracy favors 1- 519
468 consistency (92.4%) over score-thresholding (85.8%) and 520
469 gap-thresholding (83.6%). 521

470 These results demonstrate that 1-consistency’s rank- 522
471 stability approach remains largely insensitive to gallery 523
472 scale and within-identity variability, while threshold-based 524
473 methods show moderate sensitivity—score-thresholding 525
474 improving and gap-thresholding declining under large-scale 526
475 conditions. 527

476 4.5. Experiment 5: Simulated Operational Gallery 528

477 This experiment simulates an operational investigative 529
478 scenario with the most realistic conditions tested: un- 530
479 controlled gallery structure (1–108 images per identity), 531
480 cross-demographic composition reflecting actual opera-
481 tional distributions, and both synthetically degraded and
482 real surveillance-quality probes. The gallery includes
483 two additional demographics, Hispanic females and males
484 (HSF, HSM) and 13,092 additional identities, for a total of
485 136,074 total images. The demographic composition ap-
486 proximates the distribution observed in U.S. state and fed-
487 eral prisons [19]: 35% BLM, 33% WHM, 25% HSM, 3%
488 BLF, 3% WHF, 2% HSF. This configuration represents a re-
489 alistic mugshot database used for law enforcement deploy-
490 ments.

491 We evaluate two probe types: (1) resolution-degraded 540
492 MORPH probes (1,500 WHM and 1,500 BLM), and (2) 541
493 real surveillance imagery from the QMUL SurvFace dataset 542
494 [9]. For QMUL, we selected a subset with square aspect ra- 543
544
545

495 tios and at least 18×18 resolution to enable resizing to 496
497 112×112 for matcher ingestion and facilitate comparison 498

498 **Results.** For resolution-degraded MORPH probes, we com- 499
499 pare to Experiment 3 (semi-controlled multi-demographic 500
500 gallery, resolution-degraded probes). 1-consistency ING re- 501
501 call improves slightly from 89.3% to 91.3%, while score- 502
502 thresholding improves from 55.9% to 62.3%, and gap- 503
503 thresholding improves from 73.9% to 75.5%. OOG re- 504
504 call remains high for all methods, with score-thresholding 505
505 achieving perfect classification (100.0%). Overall accu- 506
506 racy favors 1-consistency (94.3%) over score-thresholding 507
507 (81.2%) and gap-thresholding (87.7%). 508

508 The modest improvements across all methods suggest 509
509 that the operational gallery structure—despite its larger 510
510 scale (136K vs. 58K images), demographic imbalance, 511
511 and broader demographic coverage (6 vs. 4 groups)—does 512
512 not substantially impair classification compared to the more 513
513 controlled multi-demographic setting of Experiment 3. This 514
514 reinforces that image quality remains the dominant factor. 515

515 For QMUL surveillance probes, all three methods 516
516 achieved 100% OOG recall, correctly rejecting every 517
517 probe. However, the QMUL images exhibited severe 518
518 degradation across multiple dimensions—extreme occlu- 519
519 sions (hats, glasses), substantial pose variation, and poor 520
520 lighting—representing quality conditions that would ideally 521
521 be flagged during human review and excluded from oper- 522
522 ational searches. This unanimous rejection suggests that 523
523 when probe quality is sufficiently poor, all methods reli- 524
524 ably detect unreliability, though the scenario provides lim- 525
525 ited discriminative insight into method differences. 526

526 These results demonstrate that 1-consistency maintains 527
527 its advantages under the most operationally realistic condi- 528
528 tions tested, achieving the highest accuracy and balanced 529
529 performance across ING and OOG recall even in large- 530
530 scale, demographically heterogeneous, structurally uncon- 531
531 trolled galleries. 532

532 5. Discussion 533

533 Across all experimental configurations, 1-consistency 534
534 demonstrates superior and more robust performance for 535
535 ING/OOG classification compared to threshold-based 536
536 methods. This advantage is most pronounced under the 537
537 conditions that matter most for operational deployment: de- 538
538 graded probe quality matched against high-quality gallery 539
539 databases. The results reveal three key findings. 540

540 **(1) Image quality drives classification performance.** 541
541 Probe degradation causes threshold-based methods to col- 542
542 lapse catastrophically on ING recall (often by 40–50 per- 543
543 centage points), while 1-consistency maintains stable per- 544
544 formance (typically within 10–15 pp of baseline). For 545
545 OOG recall, threshold methods maintain marginal advan-

546	tages (typically 3–5 pp higher than 1-consistency), but this	598
547	slight gain comes at an unacceptable operational cost: dis-	599
548	carding 40–50% of valid in-gallery probes undermines the	600
549	fundamental purpose of investigative identification systems.	601
550	For demographic fairness, probe degradation increases dis-	602
551	parities across all methods, but 1-consistency maintains	
552	the smallest increases (5–6 pp from baseline) compared to	
553	score-thresholding (9–11 pp) and gap-thresholding (11–12	
554	pp), consistently achieving the lowest demographic differ-	
555	entials under challenging conditions.	
556	(2) Classification is relatively stable across gallery struc-	603
557	tures. We evaluated galleries varying in number of images	
558	per identity (fixed at 2 versus variable from 1–108), number	604
559	of identities (from 3,400 to 26,692), and demographic com-	605
560	position (from one to six represented groups). For all three	606
561	methods, these variations produced only minor performance	607
562	effects.	608
563	However, our experiments do not extend to galleries	609
564	scaled by orders of magnitude beyond those tested, e.g.,	610
565	100× or 1000× more identities within the same demo-	611
566	graphic, as may occur in certain state-level deployments.	612
567	While our results suggest that quality effects dominate	613
568	within moderate-to-large gallery settings, performance un-	614
569	der extreme-scale single-demographic galleries remains an	615
570	open question for future study.	616
571	(2) 1-consistency achieves the best balance across com-	617
572	peting objectives. It wins on ING recall in 10 of 15 confi-	618
573	gurations (by 16.3 pp on average), on overall accuracy in 12	619
574	of 15 configurations (by 6.7 pp on average), and on demo-	
575	graphic fairness in 8 of 10 configurations reporting demo-	
576	graphics (by 4.2 pp on average). Threshold-based methods	
577	lead only on OOG recall—and even there, 1-consistency	
578	trails by just 5.1 pp on average. This asymmetry matters	
579	operationally: failing to retain valid in-gallery probes (low	
580	ING recall) directly undermines investigative workflows,	
581	while the modest OOG recall tradeoff represents a far less	
582	costly error mode.	
583	(3) The most operationally relevant quality scenario—	620
584	the comparison of degraded probes to a good-quality	
585	gallery—is where 1-consistency excels. Across the six	
586	configurations testing this scenario, 1-consistency achieves	
587	15–27 pp higher ING recall than threshold methods, 6–12	
588	pp higher overall accuracy, and 2–6 pp lower demographic	
589	disparities. Threshold methods’ OOG recall advantage in	
590	these scenarios is negligible (0.1–0.3 pp).	
591	This performance gap reflects a fundamental distribu-	
592	tional problem, not merely a calibration issue. Distribu-	
593	tional separability analysis (Table 7) reveals how degra-	
594	dation affects each method’s ability to distinguish ING	
595	from OOG probes. Under probe degradation, threshold-	
596	based separability <i>collapses</i> by 60–70% from baseline,	
597	while 1-consistency’s rank-stability signal actually becomes	
	<i>more</i> discriminative, improving by 12–29%. This means	598
	threshold-based methods cannot be improved through recal-	599
	ibration—the underlying score distributions have collapsed	600
	to the point where no single threshold can effectively sepa-	601
	rate ING from OOG probes.	602
	5.1. Operational Implications	603
	The one scenario where threshold-based methods show	604
	competitive performance is the idealized baseline: mug-	605
	shot-quality probes and galleries with controlled	606
	enrollment. Under these conditions, all methods perform	607
	nearly identically (differences under 1 pp). However, such	608
	conditions rarely hold in operational systems, where probe	609
	quality degradation is the norm rather than the exception.	610
	By operating on rank agreement across independently	611
	trained matcher instances rather than on matcher-specific	612
	score distributions, 1-consistency achieves greater invari-	613
	ance to the quality perturbations that characterize real in-	614
	vestigative deployments. This combination of higher accu-	615
	racy, better demographic fairness, and structural robustness	616
	to degradation makes 1-consistency the most operationally	617
	viable approach for distinguishing in-gallery from out-of-	618
	gallery probes in realistic face identification systems.	619
	6. Conclusion	620
	We investigated 1-consistency, a rank-stability approach	621
	for distinguishing in-gallery from out-of-gallery probes in	622
	face identification systems. Across 15 experimental con-	623
	figurations spanning controlled to uncontrolled galleries,	624
	mugshot to degraded image quality, and single- to multi-	625
	demographic compositions, 1-consistency consistently out-	626
	performed threshold-based methods on the metrics that	627
	matter most for operational deployment: in-gallery recall,	628
	overall accuracy, and demographic fairness.	629
	The advantage is most pronounced under realistic op-	630
	erational conditions—degraded probe quality against high-	631
	quality gallery databases—where threshold-based methods	632
	suffer catastrophic performance losses. Distributional anal-	633
	ysis reveals this is not a calibration problem but a structural	634
	collapse: similarity score distributions lose discriminative	635
	power under degradation (60–70% reduction in separabil-	636
	ity), while rank stability actually becomes more discrimina-	637
	tive (12–29% improvement).	638
	Remarkably, even under the most challenging dual-	639
	degradation conditions, 83.6–87% of true in-gallery probes	640
	achieve perfect rank-1 consensus across all 10 independ-	641
	ently trained matcher instances, while 90–91% of out-of-	642
	gallery probes fail to achieve consensus. This demonstrates	643
	that rank stability provides a fundamentally more robust sig-	644
	nal for gallery membership than similarity scores, making	645
	1-consistency the most operationally viable approach for	646
	open-set face identification under realistic deployment con-	647
	ditions.	648

649

References650
651652
653654
655656
657658
659660
661662
663664
665666
667668
669670
671672
673674
675676
677678
679680
681682
683684
685686
687688
689690
691692
693694
695696
697698
699700
701702
703

704

- [1] Clearview AI. Clearview ai 2.0, 2025. Accessed: 2025-02-20. 1
- [2] Abhijit Bendale and Terrance E. Boult. Towards open world recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1893–1902, 2015. 2
- [3] Abhijit Bendale and Terrance E. Boult. Towards open set deep networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1563–1572, 2016. 2
- [4] Aman Bhatta, Gabriella Pangelinan, Michael C King, and Kevin W Bowyer. Impact of blur and resolution on demographic disparities in 1-to-many facial identification. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 412–420, 2024. 1
- [5] Aman Bhatta, Maria Dhakal, Michael C King, and Kevin W Bowyer. Are you in or out (of gallery)? wisdom from the same-identity crowd. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5946–5955, 2025. 2
- [6] Aman Bhatta, Michael C King, and Kevin W Bowyer. Deep cnn face matchers inherently support revocable biometric templates. In *2025 IEEE 19th International Conference on Automatic Face and Gesture Recognition (FG)*, pages 1–10. IEEE, 2025. 3
- [7] Navaneeth Bodla, Jingxiao Zheng, Hongyu Xu, Jun-Cheng Chen, Carlos Castillo, and Rama Chellappa. Deep heterogeneous feature fusion for template-based face recognition. In *Proceedings of the IEEE Winter Conference on Applications of Computer Vision*, pages 586–595, 2017. 2
- [8] Chris Burt. Clearview facial recognition searches double, database reaches 50b images. *Biometric Update*, 2024. Accessed: 2025-02-20. 1
- [9] Zhiyi Cheng, Xi Tian Zhu, and Shaogang Gong. Surveillance face recognition challenge. *arXiv preprint arXiv:1804.09691*, 2018. 7
- [10] Jiankang Deng, Jia Guo, Xiang An, Zheng Zhu, and Stefanos Zafeiriou. Masked face recognition challenge: The insight-face track report. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1437–1444, 2021. 3
- [11] Federal Office for Information Security (BSI). Open Face Image Quality (OFIQ): Short Report. Technical report, Federal Office for Information Security (BSI), 2021. Accessed: 2025-07-08. 3
- [12] Patrick Grother, Mei Ngan, and Kayee Hanaoka. Face Recognition Vendor Test (FRVT) Part 2: Identification. <https://doi.org/10.6028/NIST.IR.8271>, 2019. 2
- [13] Patrick Grother, Mei Ngan, and Kayee Hanaoka. *Face recognition vendor test (FVRT): Part 3, demographic effects*. National Institute of Standards and Technology Gaithersburg, MD, 2019. 2
- [14] Patrick Grother, Mei Ngan, and Elham Tabassi. Face Analysis Technology Evaluation (FATE) Part 11: Face Image

Quality Vector Assessment. <https://doi.org/10.6028/NIST.IR.8485>, 2023. 3 705
706[15] Manuel Gunther, Steve Cruz, Ethan M. Rudd, and Ter- 707
rance E. Boult. Toward open-set face recognition. In *Pro- 708**ceedings of the IEEE Conference on Computer Vision and 709**Pattern Recognition Workshops*, pages 71–80, 2017. 2 710[16] Minchul Kim, Feng Liu, Anil K Jain, and Xiaoming Liu. 711
Cluster and aggregate: Face recognition with large probe 712set. *Advances in Neural Information Processing Systems*, 713
35:36054–36066, 2022. 2 714[17] Davis E King. Dlib-ml: A machine learning toolkit. 715
The Journal of Machine Learning Research, 10:1755–1758, 716

2009. 3 717

[18] Qiang Meng, Shichao Zhao, Zhida Huang, and Feng Zhou. 718
Magface: A universal representation for face recognition and 719quality assessment. In *Proceedings of the IEEE/CVF con- 720**ference on computer vision and pattern recognition*, pages 721
14225–14234, 2021. 3 722[19] Derek Mueller and Rich Kluckow. Prisoners in 2023 – sta- 723
tistical tables. Technical Report NCJ 310197, U.S. Depart- 724ment of Justice, Office of Justice Programs, Bureau of Justice 725
Statistics, 2025. 7 726[20] Gabriella Pangelinan, Aman Bhatta, Haiyu Wu, Michael C 727
King, and Kevin W Bowyer. Analyzing the impact of de- 728

mographic and operational variables on 1-to-many face id 729

search. *IEEE Transactions on Technology and Society*, 2024. 730
1 731[21] Gabriella Pangelinan, Michael C. King, and Kevin W. 732
Bowyer. When probe and gallery are low quality: Decreasing 733accuracy and increasing demographic disparities in 1:n 734
identification. In *Proceedings of the Winter Conference on 735**Applications of Computer Vision*, 2026. To appear. 1, 3 736[22] Walter J. Scheirer, Anderson Rocha, Ross J. Micheals, and 737
Terrance E. Boult. Meta-recognition: The theory and prac- 738tice of recognition score analysis. *IEEE Transactions on Pat- 739**tern Analysis and Machine Intelligence*, 33(8):1689–1695, 740
2011. 2 741[23] Walter J. Scheirer, Anderson de Rezende Rocha, Archana 742
Sapkota, and Terrance E. Boult. Toward open set recogni- 743tion. *IEEE Transactions on Pattern Analysis and Machine 744**Intelligence*, 35(7):1757–1772, 2012. 745[24] Walter J. Scheirer, Lalit P. Jain, and Terrance E. Boult. Prob- 746
ability models for open set recognition. *IEEE Transactions 747**on Pattern Analysis and Machine Intelligence*, 36(11):2317– 748
2324, 2014. 2 749[25] Jiaolong Yang, Peiran Ren, Dongqing Zhang, Dong Chen, 750
Fang Wen, Hongdong Li, and Gang Hua. Neural aggregation 751network for video face recognition. In *Proceedings of the 752**IEEE Conference on Computer Vision and Pattern Recogni- 753**tion*, pages 4362–4371, 2017. 2 754