



Cross-Resolution Face Recognition Adversarial Attacks

Fabio Valerio **Massoli**^{a,**}, Fabrizio **Falchi**^a, Giuseppe **Amato**^a

^a*ISTI-CNR, via G. Moruzzi 1, 56124 Pisa, Italy*

ABSTRACT

Face Recognition is among the best examples of computer vision problems where the supremacy of deep learning techniques compared to standard ones is undeniable. Unfortunately, it has been shown that they are vulnerable to adversarial examples - input images to which a human imperceptible perturbation is added to lead a learning model to output a wrong prediction. Moreover, in applications such as biometric systems and forensics, cross-resolution scenarios are easily met with a non-negligible impact on the recognition performance and adversary's success. Despite the existence of such vulnerabilities set a harsh limit to the spread of deep learning-based face recognition systems to real-world applications, a comprehensive analysis of their behavior when threatened in a cross-resolution setting is missing in the literature. In this context, we posit our study, where we harness several of the strongest adversarial attacks against deep learning-based face recognition systems considering the cross-resolution domain. To craft adversarial instances, we exploit attacks based on three different metrics, i.e., L_1 , L_2 , and L_∞ , and we study the resilience of the models across resolutions. We then evaluate the performance of the systems against the face identification protocol, open- and close-set. In our study, we find that the deep representation attacks represents a much dangerous menace to a face recognition system than the ones based on the classification output independently from the used metric. Furthermore, we notice that the input image's resolution has a non-negligible impact on an adversary's success in deceiving a learning model. Finally, by comparing the performance of the threatened networks under analysis, we show how they can benefit from a cross-resolution training approach in terms of resilience to adversarial attacks.

© 2020 Elsevier Ltd. All rights reserved.

1. Introduction

Face Recognition (Wang and Deng, 2018; Deng et al., 2019) (FR) represents one of the most astonishing applications of Neural Networks (NNs), especially considering Deep Convolutional Neural Networks (DCNNs), that ultimately overcame standard computer vision techniques such as Gabor-Fisher (Liu and Wechsler, 2002) and local binary patterns (Ahonen et al., 2006). The study of such a problem began in the early 90s when Turk and Pentland (1991) proposed the Eigenfaces approach, and it only required two decades for Deep Learning (DL) approaches to start to dominate the field reaching recognition performance up to 99.80% Wang and Deng (2018), thus overcoming human ability. DL-based FR systems do not exploit the output of a classifier directly. Instead, they leverage

the representation power (LeCun et al., 2015) of the learning models to extract face descriptors, i.e., multidimensional vectors, also called deep features or deep representations, to fulfill the recognition task.

Although FR systems obtain very high performance when trained with datasets comprising images acquired under controlled conditions, e.g., high-resolution, they suffer a drastic drop in reliability when tested against cross-resolution (CR) scenarios (Massoli et al., 2019) that naturally arise, for example, in surveillance applications (Zou and Yuen, 2011; Amato et al., 2019; Cheng et al., 2018). To counteract such a weakness, Ekenel and Sankur (2005) and Luo et al. (2019) proposed approaches that were not based on NNs. Instead, only recently such a problem has been tackled in the DL field (Massoli et al., 2020; Zhang et al., 2018).

To make the situation even worse, recently Szegedy et al. (2013); Biggio et al. (2013) showed that DL models are vulner-

**Corresponding author

e-mail: fabio.massoli@isti.cnr.it (Fabio Valerio Massoli)

able to the so-called *adversarial* examples - images to which a specific amount of noise, undetectable to humans, is added to induce a NN to output a wrong prediction. Unfortunately, the ability of an insightful adversary to jeopardize these learning models, considering both the digital (Dong et al., 2019; Song et al., 2018; Qiu et al., 2019; Kakizaki and Yoshida, 2019; Goswami et al., 2018) and physical (Sharif et al., 2016; Kurakin et al., 2016) domains, represents a significant concern in security-related applications such as DL-based biometrics systems (Sundararajan and Woodard, 2018) and forensics (Spaun, 2011). Thus, limiting their adoption in these fields.

In this context, we posit our contribution that we summarize as follows: i) we threaten two DCNNs by exploiting adversarial attacks based on three different metrics, i.e., L_1 , L_2 , and L_∞ ; ii) we generate attacks not only towards a classification objective but also against a similarity one. Indeed, FR systems typically do not exploit a DCNN classification output. Instead, they leverage the ability of NNs to generate discriminative deep representations among which a similarity criterion is evaluated to fulfill the recognition task; iii) we conduct the attacks in a cross-resolution domain, thus emulating a real-world scenario for an FR system; iv) we analyze the success rates of the various attacks across resolutions, studying if a DL model can benefit from a cross-resolution training procedure in terms of robustness to adversarial attacks; v) we analyze the robustness of the models through the face identification protocol (Grother et al., 2019) considering both the open- and close-set settings.

The rest of the paper is structured as follows. In Section 2, we briefly present some related works, while in Section 3, we describe the attacks algorithms we use. Subsequently, in section 4, we explain our experimental procedure and the dataset we use, while in Section 5, we present the results from the experimental campaign. Finally, in Section 6, we report our conclusions.

2. Related Works

To the best of our knowledge, this is the first work that tackles the problem of adversarial attacks against FR systems in a CR scenario. For such a reason, in what follows, we briefly cite a few articles related to the topics of the cross-resolution FR and adversarial attacks against an FR system.

2.1. Cross-Resolution Face Recognition

CR scenarios are met whenever images at different resolutions have to be matched. Such a situation typically happens, for example, in biometric and forensics applications. Super-Resolution (SR) techniques are among the most studied solutions to such a problem, and Singh et al. (2018) proposed to synthesize high-resolution faces from low-resolution ones by employing a multi-level sparse representation of the given inputs. Zangeneh et al. (2020) formulated a mapping of the low- and the high-resolution images to a common space by leveraging a DL architecture made by two distinct branches, one for each image. Luo et al. (2019) exploited the dictionary learning approach based on learning multiple dictionaries, each being associated with a resolution. The most comprehensive study

and widely tested method to improve an FR system’s performance in a CR scenario was recently proposed by Massoli et al. (2020). In their work, the authors formulated a training procedure to fine-tune a state-of-the-art model to the CR domain. They tested their models on several benchmark datasets by showing their superior performance compared to the results available in the literature.

2.2. Face Recognition Adversarial Attacks

As we mentioned at the beginning of this section, we are the first to study adversarial attacks in a cross-resolution domain. Due to the lack of papers that can be directly compared to our study, in what follows we only briefly cite a few articles concerning adversarial attacks against FR systems. Sharif et al. (2016) demonstrated the feasibility and effectiveness of physical attacks by impersonating other identities using eye-glass frames with a malicious texture. Zhong and Deng (2020) observed the superior transferability properties of feature-based attacks compared to label-based ones. Moreover, they proposed a drop-out method for DCNNs to enhance further the transferability of the attacks. Song et al. (2018) proposed a three-player GAN architecture that leveraged a face recognition network as the third player in the competition between generator and discriminator. Dong et al. (2019) successfully performed black-box attacks on FR models and demonstrated their effectiveness in a real-world deployed system.

3. Adversarial Attacks

3.1. Carlini and Wagner

Carlini and Wagner (Carlini and Wagner, 2017) (CW) formulated one of the strongest currently available attacks. The CW- L_2 attack is formalized as:

$$\min_{\mathbf{x}} c \cdot f(\frac{1}{2}\tanh(\mathbf{w}) + 1) + \|\frac{1}{2}(\tanh(\mathbf{w}) + 1) - \mathbf{x}\|_2^2,$$

where $f(\cdot)$ is the objective function, \mathbf{x} is the input image, \mathbf{w} is the adversarial example in the tanh space, and c is a positive constant which value is set by exploiting a binary search procedure.

3.2. Elastic Net Attack to DNNs

The Elastic Net Attack (Chen et al., 2018) (EAD), leverages the elastic-net regularization which is a well known technique in solving high-dimensional feature selection problems (Zou and Hastie, 2005). It is based on the objective proposed in Carlini and Wagner (2017) and it conceives the CW- L_2 attack as a special case. EAD is formulated as:

$$\min_{\mathbf{x}} c \cdot f(\mathbf{x}, t) + \beta \|\mathbf{x} - \mathbf{x}_0\|_1 + \|\mathbf{x} - \mathbf{x}_0\|_2^2,$$

where $f(\cdot)$ is the objective as in the CW- L_2 attack, t is the target class, \mathbf{x}_0 is the input image, t is the target label, \mathbf{x} is the adversarial instance, c is a parameter found by binary search, and β represents the weight of the L_1 penalty term.

132 3.3. Jacobian Saliency Map Attack

133 The Jacobian Saliency Map Attack (Papernot et al., 2016)
 134 (JSMA) exploits an “input-perturbation-to-output” mapping.
 135 Differently from the backpropagation-based attacks, JSMA
 136 leverages the model derivative concerning the classification out-
 137 put rather than the derivative of the loss function. The attack is
 138 formalized as: $\arg \min_{\delta_x} \|\delta_x\|$ s.t. $\mathbf{F}(\mathbf{X} + \delta_x) = \mathbf{Y}^*$, where
 139 \mathbf{F} is the function learned by the DNN, \mathbf{X} and \mathbf{Y}^* are the input
 140 and output of the model, respectively, and δ_x is the adversarial
 141 perturbation defined upon the evaluation of the model input sa-
 142 liency map.

143 3.4. Deep Representations Attacks

144 Differently from the previously mentioned attacks, the Deep
 145 Representations (Sabour et al., 2015) (DR) attack focuses on
 146 the manipulation of image features. It is formulated as an opti-
 147 mization problem which aims at finding the closest perturbed
 148 image, to the original one, whose descriptor is as close as pos-
 149 sible to the one of a target image named the “guide image”.
 150 Specifically, the adversarial crafting procedure is the follow-
 151 ing: $\mathbf{I}_\alpha = \arg \min_{\mathbf{I}} \|\phi_k(\mathbf{I}) - \phi_k(\mathbf{I}_g)\|_2^2$; subject to $\|\mathbf{I} - \mathbf{I}_s\|_\infty <$
 152 δ , where $\phi(\cdot)_k$ is the descriptor extracted at layer k of the
 153 threatened model, \mathbf{I}_s and \mathbf{I}_g are the source and target images, re-
 154 spectively, \mathbf{I}_α is the adversarial example, and δ is he maximum
 155 allowed perturbation in terms of the L_∞ norm.

156 4. Experimental Approach

157 4.1. Dataset and Models

158 In our experiments, we use the ~ 2.9 M images shared among
 159 the 8631 identities contained in the training set of the VGG-
 160 Face2 (Cao et al., 2018) dataset. To construct the gallery and
 161 the queries, we divide the training set into two splits. Concern-
 162 ing the gallery, we evaluate a single template for each identity
 163 as the average features vector among all the corresponding face
 164 images. Regarding the queries, we randomly select 100 identi-
 165 ties, and for each of them, we randomly pick ten correctly
 166 classified images, ending up with 1000 queries.

167 Concerning the learning models, we analyze the performance
 168 of two DCNNs: the face classifier from Cao et al. (2018) and
 169 the CR-trained one from Massoli et al. (2020). They share the
 170 same structure, i.e., a ResNet-50 (He et al., 2016) architecture
 171 equipped with Squeeze-and-Excitation (Hu et al., 2017) blocks.
 172 For both models, we adopt the same preprocessing steps for
 173 the images. First, following the same procedure as in Massoli
 174 et al. (2020), we synthesize different resolution versions of the
 175 input that allow us to evaluate the performance of the models
 176 in a cross-resolution scenario. Specifically, in our analysis, we
 177 consider images at 16, 24, 64, and 256 pixels (shortest side).
 178 Next, each image is resized to have the shortest side of 256
 179 pixels, and then it is cropped to a square picture of size 224×224
 180 pixels. Finally, we subtract the channel mean from each pixel.

4.2. Adversarial Attacks

181 Concerning the generation of the adversarial instances, we 182
 183 exploit the five algorithms we described in Section 3. We 184
 185 use the implementations available in the *foolbox* library (<https://foolbox.readthedocs.io/en/stable/>), with the only exception of 186
 187 the DR one that we build on top of the L-Broyden-Fletcher- 188
 189 Goldfarb-Shanno (L-BFGS) (Szegedy et al., 2013), optimiza- 190
 191 tion procedure. More precisely, the L-BFGS algorithm requires 192
 193 a function to optimize. To our aim, we implement such a func- 194
 195 tion by employing a k-NN algorithm as guidance in the ad- 196
 197 versarial search. We fit the classifier to the gallery templates
 198 we mentioned at the beginning of this section. Then, we start
 199 the crafting procedure and stop it as soon as the k-NN classifies
 200 the malicious image as belonging to the targeted identity. In
 201 Figure 1, we report a schematic view of the procedure we just
 202 described.

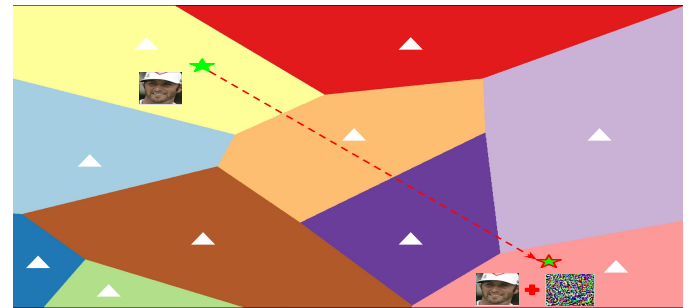


Figure 1: Schematic representation of our approach to crafting DR attacks. The colored regions are the k-NN decision boundaries for ten different identity templates (white triangles). The initial location of the green star represents a correctly classified features vector. The adversarial features vector’s final position is represented by the red encircled star.

4.3. Face Identification Metrics

197 FR systems typically deal with sensitive scenarios such as 198
 199 biometric and forensics applications. Hence, different error 200
 201 types have distinct relevance while evaluating system perform- 202
 203 ance, and a simple accuracy measure is not enough to properly 204
 205 evaluate and compare the performance of FR systems. Instead, 206
 207 as mentioned in Section 1, we focus our study on the face iden- 208
 209 tification protocol. Specifically, we consider both the close- and 210
 211 open-set settings. 212

213 Concerning the close-set setting, we evaluate the Cumulative 214
 215 Match Characteristic (CMC), a metric that represents a sum- 216
 217 marized accuracy evaluated on mated searches only, i.e., con- 218
 219 sidering queries that correspond to identities already available 220
 221 the gallery. The CMC value at rank one is usually named “hit
 222 rate,” and it is the most typical summary indicator of an al-
 223 gorithm’s efficacy. As we mentioned above, we select 100 iden-
 224 tities to construct the queries. Thus, we end up with a gallery
 225 containing 8631 identities that comprise a hundred mated ones
 226 and 8531 un-mated ones acting as “distractors”.

227 In the open-set setting, differently from the close-set one, we
 228 consider both mated and un-mated queries. To this aim, we re-
 229 move half of the queries identities from the gallery, ending up
 230 with 50 mated and 50 un-mated persons and a gallery contain-
 231 ing 8581 templates. With that set, there are two different types

of errors that are usually evaluated, i.e., the False Positive Identification Rate (FPIR) and the False Negative Identification Rate (FNIR) or “miss rate”. Concerning the former, it represents the number of un-mated queries that return a positive match at or above a specific similarity threshold. On the other hand, the FNIR represents the number of mated searches that return candidates with a similarity score below the threshold or outside the top R ranks.

The FNIR and FPIR, parametrized by the similarity threshold, can be combined to construct the Detection Error Tradeoff (DET), which is typically used to report the two types of error trade-off. We use the DET to evaluate the performance of the learning models in the experiments.

5. Experimental Results

We dedicate this section to report the results of our experimental campaigns. As we mentioned in Section 1, we aim to study the behavior of DL-based FR systems when threatened by adversarial attacks in a CR domain. Concerning the FR, as backbone features extractors, we consider the well-known DCNN from Cao et al. (2018) that set the state-of-the-art on the NIST datasets (Klare et al., 2015; Whitelam et al., 2017; Maze et al., 2018) and the CR model from Massoli et al. (2020) that set the state-of-the-art in the cross-resolution domain.

To craft adversarial examples, we harness the algorithms we described in Section 3. Moreover, being interested in the CR scenario, we consider input faces at 16, 24, 64, and 256 pixels (shortest side). Concerning the FR task, we keep the gallery at the original resolution.

As mentioned in Section 2, to our knowledge, we are the first to conduct this type of study. Thus, a direct comparison with previously published works is not possible. Hence, in what follows, we only report our results. We hope that our study will stimulate further researches in this direction. Throughout this section, we refer to the model from Cao et al. (2018) as “Base” model and to the one from Massoli et al. (2020) as “Cross-Resolution” model.

5.1. Threatening the Classification

We report the results from the attacks against the classification in Table 1. Concerning the attacks, we use the following configurations. For JSMA, we consider 1000 iterations, a perturbation per pixel equals to 0.1, 0.3, and 0.5 (percentage over the allowed pixel range), and a maximum number of times each pixel can be modified of 10. For CW- L_2 , we consider 10 binary search steps and 10 and 100 iterations. Concerning EAD, we use the same parameters as for the CW- L_2 attack and a value for the weight of the L_1 penalty term equals to 0.1 and 1. Furthermore, since the DR (Sabour et al., 2015) attack is the least time demanding compared to the others, we enlarge the set of hyperparameters for it. Thus, we dedicate Figure 2 to report their results.

From Table 1, we notice that there is no clear signature for which model is more robust against adversarial attacks. On the other hand, we see that, on average, an adversary’s success rate decreases as the resolution increases while keeping the attack configuration fixed. Let us now turn our attention to a single attack, for example, CW- L_2 . It is interesting to notice the impact of a different choice of hyperparameters. Indeed, even though from the configuration (10-10), the “Base” model seems to be more resilient compared to the “Cross-Resolution” one, this is not true. Indeed, by just increasing the strength of the attack, i.e., (10-100) configuration for which we grow the number of steps, we reach 100% of attack success rate for both models.

From Figure 2 we observe that it is undeniable that the deep features extracted by the “Cross-Resolution” model are much more robust than those extracted from the “Base” NN. Thus, confirming our previous assertion about the benefit of CR training. From the first plot of Figure 2, we see that the success rate of the attack is almost 0% for the “Base” model. Instead, in the second plot, it looks like that both models have the same resilience. This is not in contrast with our previous conclusions. Indeed, as it has been shown in appendix 1 of Massoli et al. (2020), the “Base” model is not able to generate meaningful deep representation at very low resolutions. Thus, it is almost impossible to craft targeted attacks based on deep features. To sustain even more our assertion, we run a test with untargeted DR attacks in which we easily reach a success rate of 100% for

Table 1: Attack success rate against classification for “Base” and “Cross-Resolution” models. The first column reports the specific configuration used for each attack. The four values reported in the second and third main columns represent the success rate at a resolution of 16, 24, 64, and 256 pixels, respectively.

Attack Configuration	Attack Success Rate (%)							
	Base Model				Cross-Resolution Model			
	16	24	64	256	16	24	64	256
JSMA (1000-0.1-1.0)	76.1	61.8	25.5	11.5	65.5	62.8	17.1	6.9
JSMA (1000-0.3-1.0)	96.6	92.5	75.7	61.2	96.0	94.7	70.0	50.1
JSMA (1000-0.5-1.0)	98.5	95.8	86.4	76.6	97.6	97.0	100.	69.6
CW- L_2 (10-10)	82.9	72.9	45.9	32.7	86.4	83.3	52.8	37.4
CW- L_2 (10-100)	100.	100.	100.	100.	100.	100.	100.	100.
EAD (10-0.1-10)	95.7	98.2	94.5	87.0	96.7	99.6	98.8	98.5
EAD (10-0.1-100)	100.	100.	100.	100.	100.	100.	100.	100.
EAD (10-1.0-10)	83.4	85.1	50.2	27.9	72.6	94.4	86.9	73.8
EAD (10-1.0-100)	98.5	99.8	98.7	91.0	97.5	99.8	100.	99.6

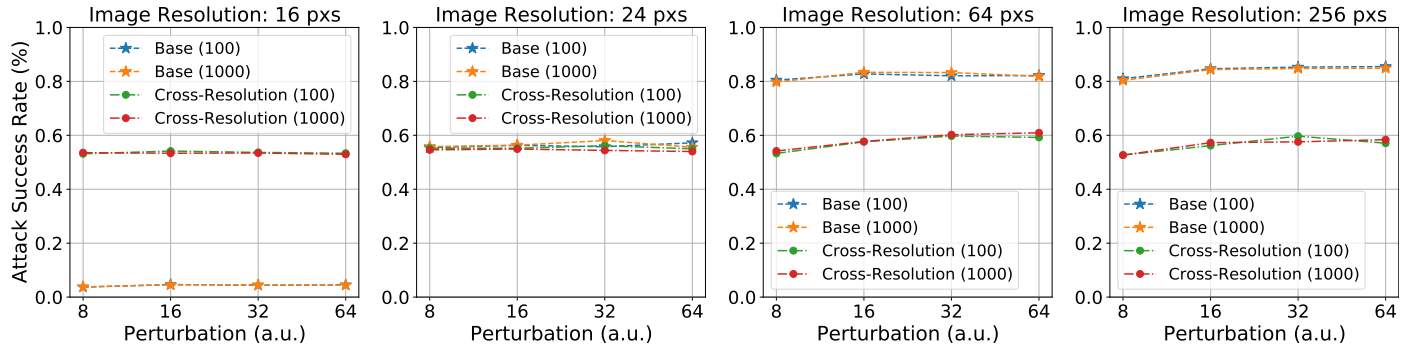


Figure 2: DR (Sabour et al., 2015) attack success rate as function of the maximum allowed perturbation δ considering 100 and 1000 iteration steps. Each plot represents a different input resolution.

297 the “Base” model.

298 Finally, we can notice that from our results, there is no clear
 299 evidence in favor of a specific metric since with the proper hy-
 300 perparameters, we reached high success rates with the L_1 , L_2 ,
 301 and L_∞ .

302 5.2. Threatening the Face Recognition

303 We now turn our attention to DL-based FR systems. We be-
 304 gin our analysis by considering the face identification protocol
 305 in the close-set scenario, and we then move the open-set one.
 306 We refer the reader to Section 4 for a detailed description of the
 307 metrics we use to assess the performance of the systems under
 308 analysis.

309 5.2.1. Close-set

310 As mentioned in Section 4, we use the CMC to evaluate the
 311 performance of the threatened models in the close-set scenario.
 312 Specifically, we summarize our results in Table 2 by reporting
 313 the hit rate, i.e., the CMC value at a rank equals to one, with
 314 the exception of the DR (Sabour et al., 2015) attack to which
 315 we dedicate Figure 3. From a defensive point of view, the more
 316 resilient a model, the lower the hit rate, while from an attacker
 317 perspective, it is the other way round.

318 By looking at Table 2 and Figure 3 we can assert that the DR
 319 attack is much more effective in fooling a DL-based FR sys-

tem than the classification-based ones with respect to any type 320
 of metric. From the attacker’s point of view, this is a funda- 321
 mental result. Indeed, by comparing the results from Table 1 322
 and Table 2, we see that even though the attacks fool the clas- 323
 sification, it is not guaranteed that they can evade a similarity- 324
 based system. Thus, deep representation attacks might be a 325
 better choice to attack an FR system. Moreover, we see how 326
 the “Cross-Resolution”-based system exhibits higher robust- 327
 ness than the one based on the “Base” model. Thus, again, we 328
 find that DCNNs benefit from a CR training approach (Mas- 329
 soli et al., 2020) in terms of resilience to adversarial attacks. 330
 Indeed, it is undeniable that the “Cross-Resolution”-based sys- 331
 tem is much more resilient against adversarial attacks than the 332
 “Base”-based one across all resolutions. 333

5.2.2. Open-set 334

To report the results for the face identification protocol in 335
 the open-set setting, we exploit the DET. Two fundamental as- 336
 pects differentiate the DET from the CMC. Indeed, the former 337
 applies a threshold among the similarity of the features, and it 338
 comprises queries of identities that are not present in the gal- 339
 lery. Instead, the latter does not use any threshold, i.e., it does 340
 not discern among “weak” and “strong” similarity scores, and 341
 it requires queries related to already known identities. 342

As we mentioned in Section 4, the DET represents the er- 343

Table 2: Attacks hit rate. The first column reports the configuration for each attack. The four values reported in the second and third main columns are the results at a resolution of 16, 24, 64, and 256 pixels, respectively. As a reference, we report in the first row the hit rate for the authentic images.

Attack Configuration	Hit Rate (%)							
	Base Model				Cross-Resolution Model			
	16	24	64	256	16	24	64	256
Auth	79.5	95.3	99.8	99.9	96.7	98.8	99.4	99.7
JSMA (1000-0.1-1.0)	12.1	10.7	12.9	12.2	11.9	9.8	9.4	13.0
JSMA (1000-0.3-1.0)	14.0	9.3	10.7	10.6	9.8	10.0	7.4	8.9
JSMA (1000-0.5-1.0)	13.6	10.6	10.0	10.3	10.0	10.2	3.0	6.8
CW- L_2 (10-10)	10.9	6.5	6.1	3.7	10.8	9.3	5.5	5.1
CW- L_2 (10-100)	7.6	4.1	6.1	2.3	9.2	9.3	3.6	4.6
EAD (10-0.1-10)	31.8	32.6	27.8	25.1	19.2	16.8	19.4	19.7
EAD (10-0.1-100)	17.5	9.7	6.3	6.2	13.8	11.6	6.8	5.3
EAD (10-1.0-10)	44.8	38.0	26.7	25.5	20.8	25.7	20.1	21.7
EAD (10-1.0-100)	34.8	30.3	20.7	16.8	17.3	16.5	17.4	17.2

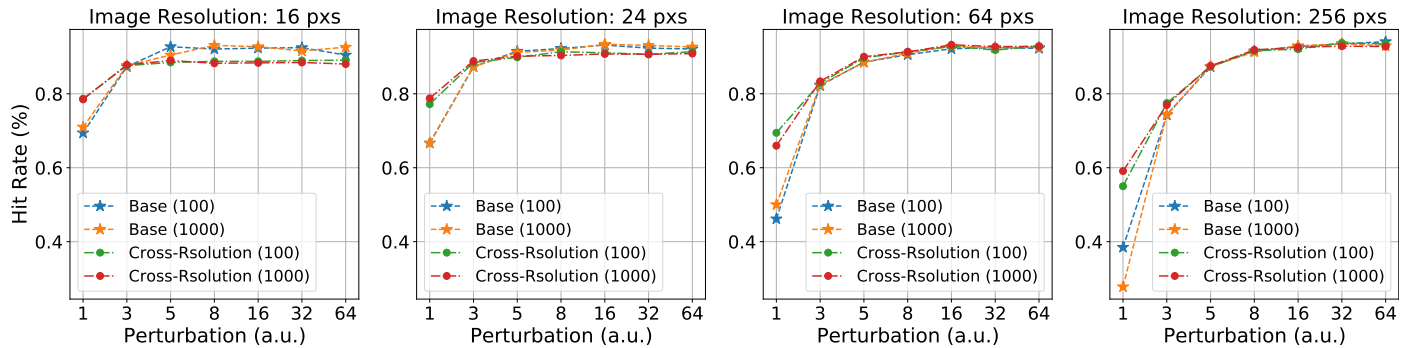


Figure 3: DR (Sabour et al., 2015) hit rate as function of the maximum allowed perturbation δ considering 100 and 1000 attack steps. Each plot represents a different input resolution.

344 ror trade-off between the FNIR and the FPIR. To summarize
 345 the performance of the FR systems, we report the FPIR at a
 346 reference value of the FNIR equals to $1.e^{-2}$. Compared to the
 347 close-set settings, the adversary’s goal is to lower the curve as
 348 much as possible, while from a defensive point of view, a higher
 349 curve represents a more resilient model. The results are reported
 350 in Table 3 with the exception of DR (Sabour et al., 2015) to
 351 which we dedicate Figure 4.

352 Analyzing the results reported in Table 3 and Figure 4 we
 353 obtain the same conclusions we report for the close-set setting.
 354 Specifically, by comparing the results from Table 3 to the ones
 355 in Figure 4 we see that the DR attack is much more effective
 356 in fooling the FR system compared to others and that the
 357 “Cross-Resolution”-based system is much more resilient than
 358 the “Base”-based one against adversarial attacks.

359 6. Conclusions

360 DCNN-based FR systems leverage the representation power
 361 of learning models. Unfortunately, they also share their weak-
 362 nesses. Indeed, it has been recently shown that these systems
 363 suffer a drastic drop in their performance when tested in a cross-
 364 resolution domain. The situation becomes even worse when
 365 an adversary comes into play. Indeed, an FR system can be
 366 deceived by adversarial examples. These weaknesses pose a

367 severe limit to the spread of these systems to sensitive real-
 368 world applications such as biometric systems and forensics.

369 In such a context, we proposed our analysis in which we
 370 compared the resilience to adversarial attacks of FR systems
 371 based on the deep features extracted by NNs in a CR scenario.
 372 We studied two different DCNN models: a former one, trained
 373 only on high-resolution images and a latter one, trained on a
 374 cross-resolution domain. To generate adversarial instances, we
 375 harnessed several algorithms based on different metrics and ob-
 376 jectives, and we craft malicious samples considering input im-
 377 ages at a resolution of 16, 24, 64, and 256 pixels. Concerning
 378 the measures of the performance of the FR systems, we adopted
 379 the face identification protocol. Specifically, we considered the
 380 close- and open-set settings for which we evaluated the CMC
 381 and DET.

382 From our analysis, we notice that, given a specific configura-
 383 tion, the attack success rate is higher at lower resolutions, for
 384 example, at 16 and 24 pixels, than at higher ones, such as 64
 385 and 256 pixels. Such behavior was somehow expected since, at
 386 a very low-resolution part of the face information can be lost,
 387 thus simplifying the effort of an adversary.

388 By looking at the results from the FR systems, it is evident
 389 that a DCNN benefits from a CR training procedure since it
 390 empowers the learning model to extract more robust deep rep-
 391 resentations. Moreover, we observed that DR attacks represent

Table 3: FPIR@FNIR= $1.e^{-2}$. The first column reports the configuration for each attack. The four values reported in the second and third main columns are the results at a resolution of 16, 24, 64, and 256 pixels, respectively. As a reference, we report in the first row the results for the authentic images.

Attack Configuration	FPIR@FNIR= $1.e^{-2}$							
	Base Model				Cross-Resolution Model			
	16	24	64	256	16	24	64	256
Auth	75.0	40.8	0.8	1.0	38.6	20.2	3.6	3.2
JSMA (1000-0.1-1.0)	99.3	99.1	100.	95.1	99.1	98.4	100.	98.1
JSMA (1000-0.3-1.0)	99.0	99.1	97.2	99.7	97.8	98.6	99.0	100.
JSMA (1000-0.5-1.0)	98.0	98.1	98.2	97.0	99.4	98.6	99.0	98.7
CW- L_2 (10-10)	99.5	98.1	99.5	97.4	99.0	98.1	98.9	98.9
CW- L_2 (10-100)	100.	99.0	99.5	99.4	99.6	98.1	99.6	99.2
EAD (10-0.1-10)	95.3	93.2	98.7	99.5	98.4	98.8	96.0	97.6
EAD (10-0.1-100)	98.0	99.4	99.4	99.0	100.	98.8	98.6	99.2
EAD (10-1.0-10)	95.6	96.3	98.3	95.3	96.3	98.1	96.7	96.7
EAD (10-1.0-100)	98.8	97.9	97.1	98.6	98.6	98.1	99.0	97.7

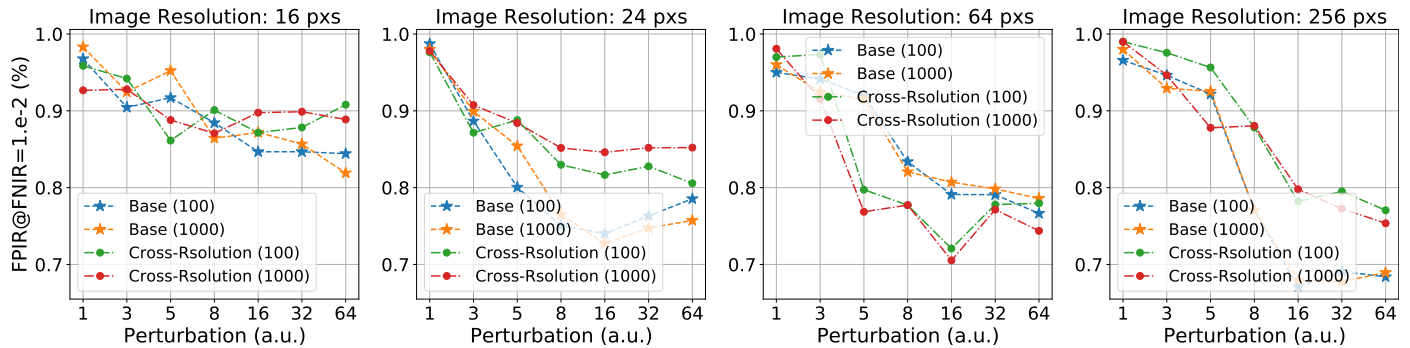


Figure 4: FPIR@FNIR=1.e⁻² for the DR (Sabour et al., 2015) attack as function of the maximum allowed perturbation δ considering 100 and 1000 attack steps. Each plot represents a different input resolution.

392 a much greater menace to an FR system than the ones based
 393 on the classification output of the threatened models for each of
 394 the considered metrics, i.e., L_1 , L_2 and L_∞ . Such a result held
 395 for the close- as well as for the open-set settings.

396 References

397 Ahonen, T., Hadid, A., Pietikainen, M., 2006. Face description with local binary
 398 patterns: Application to face recognition. *IEEE TPAMI*, 2037–2041.
 399 Amato, G., Falchi, F., Gennaro, C., Massoli, F.V., Passalis, N., Tefas, A.,
 400 Trivilini, A., Vairo, C., 2019. Face verification and recognition for digital
 401 forensics and information security, in: *ISDFS, IEEE*. pp. 1–6.
 402 Biggio, B., Corona, I., Maiorca, D., Nelson, B., Šrđić, N., Laskov, P., Giacinto,
 403 G., Roli, F., 2013. Evasion attacks against machine learning at test time, in:
 404 *ECML PKDD, Springer*. pp. 387–402.
 405 Cao, Q., Shen, L., Xie, W., Parkhi, O.M., Zisserman, A., 2018. Vggface2: A
 406 dataset for recognising faces across pose and age, in: *International Confer-*
 407 *ence on Automatic Face & Gesture Recognition, IEEE*. pp. 67–74.
 408 Carlini, N., Wagner, D., 2017. Towards evaluating the robustness of neural
 409 networks, in: *Symposium on security and privacy, IEEE*. pp. 39–57.
 410 Chen, P.Y., Sharma, Y., Zhang, H., Yi, J., Hsieh, C.J., 2018. Ead: elastic-net
 411 attacks to deep neural networks via adversarial examples, in: *Thirty-second*
 412 *AAAI conference on artificial intelligence*.
 413 Cheng, Z., Zhu, X., Gong, S., 2018. Surveillance face recognition challenge.
 414 arXiv preprint arXiv:1804.09691 .
 415 Deng, J., Guo, J., Xue, N., Zafeiriou, S., 2019. Arcface: Additive angular
 416 margin loss for deep face recognition, in: *CVPR, IEEE*. pp. 4690–4699.
 417 Dong, Y., Su, H., Wu, B., Li, Z., Liu, W., Zhang, T., Zhu, J., 2019. Efficient
 418 decision-based black-box adversarial attacks on face recognition, in: *CVPR,*
 419 *IEEE*. pp. 7714–7722.
 420 Ekenel, H.K., Sankur, B., 2005. Multiresolution face recognition. *Image and*
 421 *Vision Computing* 23, 469–477.
 422 Goswami, G., Ratha, N., Agarwal, A., Singh, R., Vatsa, M., 2018. Unravel-
 423 ing robustness of deep learning based face recognition against adversarial
 424 attacks, in: *Thirty-Second AAAI Conference on Artificial Intelligence*.
 425 Grother, P., Grother, P., Ngan, M., Hanaoka, K., 2019. Face Recognition Vendor
 426 Test (FRVT) Part 2: Identification. US Department of Commerce, National
 427 Institute of Standards and Technology.
 428 He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image
 429 recognition, in: *CVPR, IEEE*. pp. 770–778.
 430 Hu, J., Shen, L., Sun, G., 2017. Squeeze-and-excitation networks. arxiv.
 431 Kakizaki, K., Yoshida, K., 2019. Adversarial image translation: Unrestricted
 432 adversarial examples in face recognition systems. [arXiv:1905.03421](https://arxiv.org/abs/1905.03421).
 433 Klare, B.F., Klein, B., Taborsky, E., Blanton, A., Cheney, J., Allen, K., Grother,
 434 P., Mah, A., Jain, A.K., 2015. Pushing the frontiers of unconstrained face
 435 detection and recognition: Iarpa janus benchmark a, in: *CVPR, IEEE*. pp.
 436 1931–1939.
 437 Kurakin, A., Goodfellow, I., Bengio, S., 2016. Adversarial examples in the
 438 physical world. arXiv preprint arXiv:1607.02533 .
 439 LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. *nature* 521, 436–444.
 440 Liu, C., Wechsler, H., 2002. Gabor feature based classification using the en-
 441 hanced fisher linear discriminant model for face recognition. *IEEE Transac-*
 442 *tions on Image processing* 11, 467–476.

Luo, X., Xu, Y., Yang, J., 2019. Multi-resolution dictionary learning for face
 443 recognition. *Pattern Recognition* 93, 283–292. 444
 Massoli, F.V., Amato, G., Falchi, F., 2020. Cross-resolution learning for face
 445 recognition. *Image and Vision Computing*, 103927. 446
 Massoli, F.V., Amato, G., Falchi, F., Gennaro, C., Vairo, C., 2019. Improving
 447 multi-scale face recognition using vggface2, in: *International Conference on*
 448 *Image Analysis and Processing, Springer*. pp. 21–29. 449
 Maze, B., Adams, J., Duncan, J.A., Kalka, N., Miller, T., Otto, C., Jain, A.K.,
 450 Niggel, W.T., Anderson, J., Cheney, J., et al., 2018. Iarpa janus benchmark-
 451 c: Face dataset and protocol, in: *ICB, IEEE*. pp. 158–165. 452
 Papernot, N., McDaniel, P., Jha, S., Fredrikson, M., Celik, Z.B., Swami, A.,
 453 2016. The limitations of deep learning in adversarial settings, in: *2016*
 454 *IEEE European symposium on security and privacy, IEEE*. pp. 372–387. 455
 Qiu, H., Xiao, C., Yang, L., Yan, X., Lee, H., Li, B., 2019. Semanticadv: Gen-
 456 erating adversarial examples via attribute-conditional image editing. arXiv
 457 preprint arXiv:1906.07927 . 458
 Sabour, S., Cao, Y., Faghri, F., Fleet, D.J., 2015. Adversarial manipulation of
 459 deep representations. arXiv preprint arXiv:1511.05122 . 460
 Sharif, M., Bhagavatula, S., Bauer, L., Reiter, M.K., 2016. Accessorize to a
 461 crime: Real and stealthy attacks on state-of-the-art face recognition, in:
 462 *SIGSAC CCS, ACM*. pp. 1528–1540. 463
 Singh, M., Nagpal, S., Singh, R., Vatsa, M., Majumdar, A., 2018. Magnifyme:
 464 Aiding cross resolution face recognition via identity aware synthesis. arXiv
 465 preprint arXiv:1802.08057 . 466
 Song, Q., Wu, Y., Yang, L., 2018. Attacks on state-of-the-art face recogni-
 467 tion using attentional adversarial attack generative network. arXiv preprint
 468 arXiv:1811.12026 . 469
 Spaun, N.A., 2011. Face recognition in forensic science, in: *Handbook of face*
 470 *recognition. Springer*, pp. 655–670. 471
 Sundararajan, K., Woodard, D.L., 2018. Deep learning for biometrics: a survey.
 472 *ACM Computing Surveys (CSUR)* 51, 65. 473
 Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I.,
 474 Fergus, R., 2013. Intriguing properties of neural networks. arXiv preprint
 475 arXiv:1312.6199 . 476
 Turk, M.A., Pentland, A.P., 1991. Face recognition using eigenfaces, in: *Pro-*
 477 *ceedings. 1991 IEEE Computer Society Conference on Computer Vision*
 478 *and Pattern Recognition, IEEE*. pp. 586–591. 479
 Wang, M., Deng, W., 2018. Deep face recognition: A survey. arXiv preprint
 480 arXiv:1804.06655 . 481
 Whitelam, C., Taborsky, E., Blanton, A., Maze, B., Adams, J., Miller, T., Kalka,
 482 N., Jain, A.K., Duncan, J.A., Allen, K., et al., 2017. Iarpa janus bench-
 483 mark b face dataset, in: *CVPR Workshops, IEEE*. pp. 90–98. 484
 Zangeneh, E., Rahmati, M., Mohsenzadeh, Y., 2020. Low resolution face re-
 485 cognition using a two-branch deep convolutional neural network architec-
 486 ture. *Expert Systems with Applications* 139, 112854. 487
 Zhang, K., Zhang, Z., Cheng, C.W., Hsu, W.H., Qiao, Y., Liu, W., Zhang, T.,
 488 2018. Super-identity convolutional neural network for face hallucination, in:
 489 *European Conference on Computer Vision (ECCV)*, pp. 183–198. 490
 Zhong, Y., Deng, W., 2020. Towards transferable adversarial attack against
 491 deep face recognition. arXiv preprint arXiv:2004.05790 . 492
 Zou, H., Hastie, T., 2005. Regularization and variable selection via the elastic
 493 net. *Journal of the royal statistical society: series B* 67, 301–320. 494
 Zou, W.W., Yuen, P.C., 2011. Very low resolution face recognition problem.
 495 *IEEE Transactions on image processing* 21, 327–340. 496