# LOTS about Attacking Deep Features

Andras Rozsa, Manuel Günther, and Terrance E. Boult
Vision and Security Technology (VAST) Lab
University of Colorado, Colorado Springs, USA

{arozsa,mgunther,tboult}@vast.uccs.edu

arXiv:1611.06179v5 [cs.CV] 31 May 2018

## Abstract

*Deep neural networks provide state-of-the-art performance on various tasks and are, therefore, widely used in real world applications. DNNs are becoming frequently utilized in biometrics for extracting deep features, which can be used in recognition systems for enrolling and recognizing new individuals. It was revealed that deep neural networks suffer from a fundamental problem, namely, they can unexpectedly misclassify examples formed by slightly perturbing correctly recognized inputs. Various approaches have been developed for generating these so-called adversarial examples, but they aim at attacking end-to-end networks. For biometrics, it is natural to ask whether systems using deep features are immune to or, at least, more resilient to attacks than end-to-end networks. In this paper, we introduce a general technique called the layerwise origin-target synthesis (LOTS) that can be efficiently used to form adversarial examples that mimic the deep features of the target. We analyze and compare the adversarial robustness of the end-to-end VGG Face network with systems that use Euclidean or cosine distance between gallery templates and extracted deep features. We demonstrate that iterative LOTS is very effective and show that systems utilizing deep features are easier to attack than the end-to-end network.*

Figure 1: WHO ARE THEY? Although we might not be able to recognize the presented celebrities, we can still differentiate them from one another. In fact, eight of these images are manipulated in such a way that their VGG Face descriptors mimic Kate McKinnon's (bottom-left) and cause systems that apply Euclidean or cosine distance to classify each image incorrectly as her.

## 1. Introduction

In the last few years, the most advanced deep neural networks (DNNs) have managed to reach or even surpass human level performance on a wide range of challenging machine learning tasks [13, 19, 23, 7], including face recognition. DNNs trained to perform specific tasks are able to learn representations that generalize well to other datasets [25, 13], and the extracted generic descriptors can be utilized to tackle other diverse problems [15, 3]. For visual recognition tasks in biometrics – e.g., facial attribute classification [11] and face recognition [25, 22, 13] – features obtained from DNNs are widely used in the literature.

Although we are capable of designing and training

DNNs that perform well, our understanding of these complex networks is still incomplete. This was highlighted by the intriguing properties of machine learning models discovered by Szegedy *et al.* [24]. Namely, machine learning models – including the state-of-the-art DNNs – suffer from an unexpected instability as they misclassify adversarial examples formed by adding imperceptibly small perturbations to otherwise correctly recognized inputs. Due to their excellent generalization capabilities, DNNs are expected to be robust to such small perturbations to their inputs, therefore the existence of adversarial examples challenges our understanding of DNNs and raises questions about the applications of such vulnerable learning models.

Considering the revealed adversarial instability of the end-to-end machine learning models, it is natural to ask whether systems utilizing extracted features from DNNs are also vulnerable to such perturbations. In case they are susceptible to adversarial examples, are they more or less robust than end-to-end DNNs? To be able to answer these questions, first we need to design a novel adversarial example generation technique that is capable of efficiently attacking those systems.

In this paper, we introduce the layerwise origin-target synthesis (LOTS) technique designed to perturb samples in such ways that their deep feature representations mimic any selected target activations. We experimentally demonstrate the effectiveness of LOTS in terms of forming high quality adversarial examples. We analyze and compare the robustness of the end-to-end VGG Face network [13] to adversarial perturbations with other face recognition systems that utilize deep features extracted from the same network using Euclidean or cosine distance. Our results show that LOTS is capable of successfully attacking each system, and that face recognition systems using the extracted deep features are less robust than the end-to-end network.

## 2. Related Work

Automatic face recognition has a long history and many different approaches have been proposed in the literature [26, 8, 20, 6]. While these traditional face recognition algorithms perform well on facial images with decent quality [12], they are not able to handle pose variations [6]. Only the development of deep neural networks [25, 21, 13] has overcome this issue, and nowadays these methods are the quasi standard for face recognition in uncontrolled scenarios. For example, the DNNs used by Chen *et al.* [1] provide the current state-of-the-art results on the IJB-A benchmark [9], which is outperformed by the (unpublished) DNN of Ranjan *et al.* [14].

In biometric recognition, training and evaluation can use different identities not just different images, which means that identities cannot be directly classified by an end-to-end network. Instead, the last layer of the network is removed, and deep features extracted from the penultimate layer of the DNN are used as a representation of the face [25, 13]. To form a more robust representation of an identity, deep features of several images are averaged [1]. Finally, the comparison of deep features is obtained via simple distance measures in the deep feature space such as Euclidean [13] or cosine distance [1]. Our experiments use deep features extracted with the publicly available VGG Face network [13].

Since Szegedy *et al.* [24] presented the problem posed by adversarial examples and introduced the first method capable of reliably finding such perturbations, various approaches were proposed in the literature. Compared to the computationally expensive box-constrained optimiza-

tion technique (L-BFGS) that Szegedy *et al.* [24] used, a more lightweight, still effective technique was introduced by Goodfellow *et al.* [5]. Their fast gradient sign (FGS) method relies on using the sign of the gradient of loss with respect to the input, which needs to be calculated only once per adversarial example generation. The authors demonstrated that using an enhanced objective function that implicitly incorporates FGS examples, the overall performance and the adversarial robustness of the trained models can be improved. Later, Rozsa *et al.* [16] showed that by not using the sign, the formalized fast gradient value (FGV) approach forms different adversarial samples than FGS and those yield a greater improvement when used for training.

The aforementioned two adversarial example generation techniques – FGS and FGV – rely on simply ascending the gradient of loss used for training the network. Namely, the formed perturbation causes misclassification by increasing the loss until the particular original class does not have the highest probability. In their recent paper focusing on adversarial training, Kurakin *et al.* [10] proposed extensions over the FGS method to be able to target a specific class or by calculating and applying gradients iteratively compared to a single one for conducting a line-search via FGS.

A few approaches that do not rely on using the gradient of training loss were also proposed by researchers. Rozsa *et al.* [16] introduced the hot/cold approach producing adversarial examples by both reducing the prediction probability of the original class of the input as well as increasing the probability of a specified target class. To do so, the hot/cold approach defines a Euclidean loss with varying target classes on the pre-Softmax layer and uses its gradients as directions for forming adversarial perturbations. This approach is capable of producing multiple adversarial examples per input, but still targets training classes, so cannot be directly applied to deep features.

Finally, the approach introduced by Sabour *et al.* [17] produces adversarial examples that not only cause misclassifications but also mimic the internal representations of the targeted inputs. However, their technique relies on using the computationally expensive L-BFGS technique, which limits its application.

Since, in general, biometric systems operate on a dataset different than the end-to-end network was trained on, such systems cannot be attacked by end-to-end adversarial generation techniques. Our novel LOTS method can be considered an extension of the hot/cold approach to deeper layers, and it also shows similarities to the technique of Sabour *et al.* [17] in terms of directly adjusting internal feature representations – without relying on the L-BFGS algorithm.

## 3. Approach

This section describes the targeted face recognition systems, introduces our approach to form adversarial perturba-
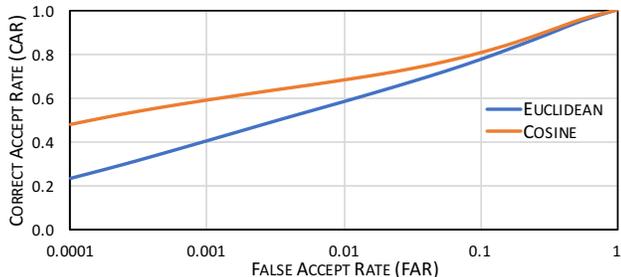
Figure 2: ROC - EUCLIDEAN AND COSINE DISTANCES.
Receiver operating characteristic (ROC) curves are shown for systems using Euclidean or cosine distance between gallery templates and VGG Face descriptors of probe images.

tions on those systems and, finally, presents the metric that we use for quantifying the quality of adversarial examples.

### 3.1. Face Recognition Systems

For our systems, we use the publicly available VGG Face dataset[1] [13] which contains 2,604,175 images of 2,622 identities. We chose this dataset because of its quality and size, and due to the fact that there is a publicly available end-to-end classification network, called the VGG Face network[2] [13], that was trained on this dataset.

The VGG Face network is intended to be used for extracting deep features, so-called VGG Face descriptors – the authors successfully utilized the captured representations of the FC7 layer as face descriptors on the labeled face in the wild (LFW) dataset [13] – or for being fine-tuned on other datasets. The network was trained on the VGG Face dataset, which was divided into three subsets containing 1,822,894, 520,835, and 260,446 images for training, validation, and test purposes, respectively. Splitting the dataset happened proportionately with respect to images per identity: each identity in the dataset has $\leq$1000 images, 70%, 20%, and 10% of those are training, validation, or test images.

To be able to directly compare the robustness of the end-to-end VGG Face network with systems using the extracted VGG Face descriptors – which are extracted at the FC7 layer before ReLU – in latter systems we need to have the same identities as the VGG Face dataset has. Therefore, we utilize the test set. We form a gallery template for each identity by calculating the mean VGG Face descriptor of the first half of the test images ($\leq$50 per identity). The VGG Face descriptors from the other half of the test images serve as probes, where each probe is compared to each gallery template, yielding 130,233 positive (same identity) and 341,340,693 negative (different identity) comparisons.

Using the positive and negative comparisons, we calculate Euclidean or cosine distance among them to compute

---

ROC curves and, finally, we identify distance thresholds for attacking deep features in Sec. 4.1. These ROC curves are displayed in Fig. 2. Since we would like to compare the adversarial robustness of the end-to-end network with systems having characteristics like real-world applications, we define them to have a low false accept rate (FAR) of 0.001 – 1 successful zero-effort impostor attack out of 1000 attempts – which translates into thresholds of 630.46 for Euclidean and 0.1544 for cosine distance.

### 3.2. Attacking Deep Features with LOTS

Let us consider a network $f$ with weights $w$ in a layered structure, i.e., having layers $y^{(l)}, l = \{1, \ldots, L\}$, with their respective weights $w^{(l)}$. For a given input $x$, the output of the network can be formalized as:

$$f(x) = y^{(L)}\left(y^{(L-1)}\left(\ldots\left(y^{(1)}(x)\right)\ldots\right)\right), \quad (1)$$

while the internal representation (the deep feature) of the given input $x$ captured at layer $l$ is:

$$f^{(l)}(x) = y^{(l)}\left(y^{(l-1)}\left(\ldots\left(y^{(1)}(x)\right)\ldots\right)\right). \quad (2)$$

Our layerwise origin-target synthesis (LOTS) approach adjusts the internal representation of an input $x_o$, the *origin*, to get closer to the *target* internal representation $t$. In order to do so, we use a Euclidean loss defined on the internal representation $f^{(l)}(x_o)$ of the origin at layer $l$ and the target $t$, and apply its gradient with respect to the origin to manipulate the internal features of the origin, formally:

$$\eta^{(l)}(x_o, t) = \nabla_{x_o}\left(\frac{1}{2}\left\|t - f^{(l)}(x_o)\right\|^2\right). \quad (3)$$

The target $t$ can be chosen without any constraints. We can manipulate origin's features at layer $l$ to get closer to the feature representation of a specific targeted input $x_t$ using $t = f^{(l)}(x_t)$ or specify any arbitrary feature representation that the origin should mimic.

We can use the direction defined by the gradient of the Euclidean loss and form adversarial perturbations using a line-search – similar to the fast gradient sign (FGS) method [5] or the hot/cold approach [16]. Compared to those previous techniques, LOTS has the potential to form dramatically greater quantities of diverse perturbations for each input due to the billions of possible targets and the number of layers it can be applied on. To form higher quality adversarial examples – with less perceptible perturbations – we can use iterative LOTS as detailed in Alg. 1. This "step-and-adjust" algorithm perturbs $x_p$ initialized with the origin $x_o$ to get closer to the target $t$ step by step until the origin mimics the target, i.e., the Euclidean or cosine distance between $f^{(l)}(x_o)$ and $t$ is smaller than a predefined threshold value, or $x_p$ is classified by the end-to-end network as desired. The perturbed image $x_p$ mimicking the

**Algorithm 1** DEEP FEATURE MIMICKING VIA LOTS. Iterative LOTS is a generic algorithm that perturbs origin $x_o$ in order to have deep features mimicking the specified target representation $t$. The function `mimicked` depends on the targeted system.

```
 1: procedure MIMIC(x_o, t)          ▷ Origin x_o mimics target t
 2:     x_p ← x_o
 3:     x'_p ← x_o
 4:     while not mimicked(x_p, t) do
 5:         grad ← η^(l)(x'_p, t)                      ▷ Eq. (3)
 6:         peak ← max(abs(grad))
 7:         grad_s ← grad / peak      ▷ Elementwise division
 8:         x'_p ← clip(x'_p − grad_s)
 9:         x_p ← round(x'_p)
10:     end while
11:     return x_p              ▷ Image with features mimicking t
12: end procedure
```

target $t$ has discrete pixel values in $[0, 255]$, however, while taking steps towards the target, the algorithm is designed to temporarily utilize non-discrete pixel values within $x'_p$ in order to obtain better adversarial quality. Finally, note that we apply a scaled gradient (line 7 in Alg. 1) with $L_\infty = 1$ to move faster towards the specified target.

### 3.3. Quantifying Adversarial Quality

In order to analyze and compare the various attacks described in Sec. 4.1, we need to assess the quality of adversarial images that iterative LOTS can generate. While $L_p$ norms are commonly used to quantify perturbations, some researchers [17, 16] concluded that those measures are not matched well to human perception. To address the problem, Rozsa et al. [16] proposed the psychometric called the perceptual adversarial similarity score (PASS) to better measure the quality of adversarial images. While $L_2$ and $L_\infty$ norms focus strictly on the perturbation – regardless of how visible it is on the distorted image – PASS is designed to better quantify the distinguishability or similarity of the original image $x_o$ and the perturbed image $x_p$ with respect to human perception.

The calculation of PASS takes two steps: alignment by maximizing the enhanced correlation coefficient (ECC) [4] of the image pair with homography transform $\Psi(x_p, x_o)$, followed by quantifying the similarity between the aligned original and perturbed images using the structural similarity (SSIM) index [27]. By design, the alignment via ECC takes place before SSIM calculation as small translations or rotations can remain imperceptible to the human eye, thus, PASS eliminates those before determining the structural similarity of the image pair. Consequently, PASS can be formalized as:

$$\text{PASS}(x_p, x_o) = \text{SSIM}(\Psi(x_p, x_o), x_o), \qquad (4)$$

where $\text{PASS}(x_p, x_o) = 1$ indicates perfect similarity.

As the structural similarity via SSIM can be calculated only on grayscale images, we align the converted grayscale images using OpenCV's ECC with termination criteria of 100 iterations or $\epsilon = 0.01$, then we calculate the structural similarity of the aligned images using SSIM.[3]

## 4. Experiments

The primary goal of this paper is to answer the question whether systems relying on extracted deep features of DNNs are vulnerable to adversarial perturbations, and if they are, how their adversarial robustness compares to end-to-end classification networks'. To be able to conduct a fair comparison, we need to design our experiments carefully.

### 4.1. Adversaries and Attack Scenarios

For analyzing the capabilities of LOTS and studying the adversarial robustness of various face recognition systems, we use a dozen adversaries – 6 identities hand-picked from the VGG Face dataset, along with 6 manually chosen external identities not contained in the VGG Face dataset. We manually selected them in order to obtain a diverse set of adversaries as shown in Fig. 3. As each adversary is represented by a single image, internal adversaries need to be chosen carefully. Therefore, from the validation set of the VGG Face dataset we selected an image for each internal adversary that is correctly classified by the end-to-end VGG Face network, and by both systems using Euclidean or cosine distance between the gallery templates and the extracted VGG Face descriptors of probe images, cf. Sec. 3.1. Images representing the external adversaries were hand-picked and manually cropped.

With having both internal and external adversaries, our goal is to analyze whether VGG Face descriptors generalize well to novel identities or if they are more specific to the VGG Face dataset in terms of better representing those identities present in the dataset. In the latter case, attacks conducted with external adversaries would outperform attacks by internal adversaries.

We conduct four sets of experiments utilizing the iterative layerwise origin-target synthesis (LOTS) approach, as detailed by Alg. 1. We perturb images of adversaries such that they mimic the VGG Face dataset identities yielding misclassifications on different face recognition systems.

*First*, on the end-to-end VGG Face network, we use LOTS on representations extracted from the Softmax layer. While origins are external and internal adversaries, we specify the targeted identity by using a particular one-hot vector on the Softmax layer as target $t$. This can be considered a traditional or more conventional approach for forming adversarial perturbations on end-to-end classification networks. In fact, this utilization of LOTS can be interpreted as

---

[3]Python implementation of SSIM by Antoine Vacavant:
http://isit.u-clermont1.fr/~anvacava/codes/ssim.py

| (a) Daniel Craig | (b) Hugh Laurie | (c) Idris Elba | (d) Kate Beckinsale | (e) Kristen Bell | (f) Thandie Newton |

| (g) Denzel Washington | (h) Ewan McGregor | (i) Halle Berry | (j) Naomi Watts | (k) Penelope Cruz | (l) Pierce Brosnan |

Figure 3: ADVERSARIES - INTERNAL AND EXTERNAL. These are the adversaries that we use throughout our experiments. The internal adversaries shown in the top row are images from the VGG Face dataset that are correctly classified by each of our systems. The external adversaries displayed in the bottom are not contained in the VGG Face dataset.

a slightly adjusted, iterative variant of the hot/cold approach introduced by Rozsa *et al.* [16]. *Second*, we aim to generate adversarial perturbations on the end-to-end network using iterative LOTS on VGG Face descriptors of adversaries to mimic gallery templates that we formed by using the mean face descriptors of VGG identities (cf. Sec. 3.1). This scenario can be viewed as attackers computing mean face descriptors from several images of targeted identities – e.g., taken from the Internet – and using them as target $t$. We conduct these two experiments to assess the effectiveness of LOTS on the end-to-end VGG Face network. The results also allow comparison of the more traditional approach of manipulating representations of the Softmax layer with the novel approach of mimicking VGG Face descriptors.

*Third* and *fourth*, we conduct experiments to generate adversarial examples on face recognition systems that use Euclidean or cosine distance between the gallery templates and the extracted VGG Face descriptors of probe images. Using iterative LOTS, our goal is to get face descriptors of adversaries closer to templates than Euclidean or cosine distance thresholds of systems having FAR = 0.001, as defined in Sec. 3.1.

Throughout our experiments, we attempt to target every possible identity of the VGG Face dataset with each adversary. For internal adversaries, this yields 2,621 subjects, while external adversaries can aim at impersonating all 2,622 identities. To limit the computational costs, we constrain iterative LOTS to 500 steps. In case the algorithm exceeds the limit, the particular attempt is considered a failure. As we will see, this constraint has little effect on our experiments. Furthermore, based on our experience, iterative LOTS taking more than 500 steps produces perturba-

tions that are highly visible, in other words, those examples are not adversarial at all.

## 4.2. Results

The results obtained by conducting the four sets of experiments using the selected adversaries are presented in Tab. 1. Comparing the collected metrics on the two types of attacks on the end-to-end VGG Face network, we can conclude that, in general, iterative LOTS operating on VGG Face descriptors produces examples with better adversarial quality than the traditional attack working on the Softmax layer. Considering all internal and external adversaries, there is only one exception: for external adversary Denzel Washington, the formed examples using Softmax features contain less perceptible perturbations in average, as indicated by the higher PASS. Furthermore, we can note with respect to the attacks on the end-to-end face recognition network that there is a small proportion of targeted identities for each adversary – varying between 41 and 45 – where iterative LOTS limited to 500 steps failed. By analyzing these unsuccessful attempts, interestingly, we find that using the diverse set of adversaries our algorithm failed to form perturbations more or less for the same targeted identities. We conjecture that those subjects are hard to reach via iterative LOTS because they are simply more difficult to be recognized by the end-to-end network – Doddington *et al.* [2] dubbed them as "goats." Consequently, those directions provided by the calculated gradients via iterative LOTS simply cannot find a way to those identities.

We can see in Tab. 1 that iterative LOTS performs better on face recognition systems that use Euclidean or cosine distance on extracted VGG Face descriptors with FAR =

Table 1: ADVERSARIAL EXAMPLE GENERATION VIA ITERATIVE LOTS. These results are obtained using iterative LOTS with the listed internal and external adversaries. With each adversary, we attacked every possible subject by mimicking their gallery templates to cause misclassifications on the end-to-end VGG Face network (End-To-End FD), and on systems using Euclidean or cosine distance between gallery templates and the extracted VGG Face descriptors. Furthermore, we attacked each identity on the end-to-end network by manipulating their representations at the Softmax layer (End-To-End SM) targeting the appropriate one-hot vector. We list the mean and standard-deviation of PASS, followed by the percentage of successful attacks, i.e., when the perturbed images were classified as the target.

| | ADVERSARY | END-TO-END SM | END-TO-END FD | EUCLIDEAN DISTANCE | COSINE DISTANCE |
|---|---|---|---|---|---|
| INTERNAL | Daniel Craig | $0.9833 \pm 0.0076$ (98.44%) | $0.9846 \pm 0.0075$ (98.32%) | $0.9873 \pm 0.0083$ (100.00%) | $0.9900 \pm 0.0055$ (100.00%) |
| | Hugh Laurie | $0.9606 \pm 0.0186$ (98.44%) | $0.9697 \pm 0.0123$ (98.32%) | $0.9805 \pm 0.0116$ (100.00%) | $0.9850 \pm 0.0081$ (100.00%) |
| | Idris Elba | $0.9643 \pm 0.0206$ (98.36%) | $0.9686 \pm 0.0147$ (98.32%) | $0.9844 \pm 0.0130$ (100.00%) | $0.9894 \pm 0.0075$ (100.00%) |
| | Kate Beckinsale | $0.9804 \pm 0.0113$ (98.44%) | $0.9840 \pm 0.0107$ (98.32%) | $0.9900 \pm 0.0066$ (100.00%) | $0.9921 \pm 0.0049$ (100.00%) |
| | Kristen Bell | $0.9704 \pm 0.0156$ (98.44%) | $0.9821 \pm 0.0115$ (98.32%) | $0.9883 \pm 0.0062$ (100.00%) | $0.9905 \pm 0.0049$ (100.00%) |
| | Thandie Newton | $0.9792 \pm 0.0099$ (98.44%) | $0.9849 \pm 0.0085$ (98.28%) | $0.9881 \pm 0.0068$ (100.00%) | $0.9904 \pm 0.0055$ (100.00%) |
| EXTERNAL | Denzel Washington | $0.9866 \pm 0.0107$ (98.44%) | $0.9839 \pm 0.0093$ (98.32%) | $0.9869 \pm 0.0095$ (100.00%) | $0.9900 \pm 0.0060$ (100.00%) |
| | Ewan McGregor | $0.9925 \pm 0.0063$ (98.44%) | $0.9936 \pm 0.0042$ (98.32%) | $0.9944 \pm 0.0041$ (100.00%) | $0.9957 \pm 0.0027$ (100.00%) |
| | Halle Berry | $0.9913 \pm 0.0066$ (98.44%) | $0.9918 \pm 0.0057$ (98.32%) | $0.9943 \pm 0.0040$ (100.00%) | $0.9955 \pm 0.0027$ (100.00%) |
| | Naomi Watts | $0.9721 \pm 0.0150$ (98.44%) | $0.9823 \pm 0.0100$ (98.32%) | $0.9869 \pm 0.0071$ (100.00%) | $0.9891 \pm 0.0065$ (100.00%) |
| | Penelope Cruz | $0.9745 \pm 0.0151$ (98.44%) | $0.9867 \pm 0.0084$ (98.32%) | $0.9906 \pm 0.0052$ (100.00%) | $0.9923 \pm 0.0048$ (100.00%) |
| | Pierce Brosnan | $0.9835 \pm 0.0106$ (98.44%) | $0.9808 \pm 0.0095$ (98.28%) | $0.9877 \pm 0.0085$ (100.00%) | $0.9907 \pm 0.0056$ (100.00%) |

0.001 than it does on the end-to-end system. The better performance is highlighted by both the higher percentage of successful attacks and the overall adversarial quality shown by generally higher PASS. Statistical testing – two-sided heteroscedastic t-tests with Bonferroni correction – show very significant ($p < 0.00001$) difference between each pair of the four attack scenarios. While the differences in PASS may seem to be small, the large sample size of over 2,600 identities results in the strong rejection of the hypothesis that the four attacks provide similar results. For the methods shown in Tab. 1, the quality of the generated adversarial images statistically significantly increases from left to right, supporting the conclusion: *Systems utilizing deep features are easier to attack and admit less perceptible perturbations than the end-to-end network.*

While iterative LOTS forms perturbations for adversaries to reach nearly all targeted identities, the manipulated images also maintain high adversarial quality. To demonstrate the effectiveness of iterative LOTS in terms of reaching the targeted subjects via small distortions, in Fig. 4 we show examples with VGG Face descriptors that are closer to gallery templates than Euclidean or cosine thresholds. Most of these examples are indeed adversarial images as those perturbations are imperceptible. The displayed examples are external adversaries targeting identities of the selected images used as internal adversaries. We show these particular examples simply because we already introduced those identities by displaying an image for each internal adversary in Fig. 3, thus we can associate a face to those subjects.

As indicated by the collected metrics, the formed perturbations that we obtained on the system using cosine distance yield even slightly better adversarial quality than collected on the Euclidean system. This means that the system with the higher recognition accuracy (cf. Fig. 2) is also easier to attack. To highlight the differences among systems with

respect to adversarial vulnerability, we visualize perturbations for some distorted examples to show what it takes to cause misclassifications. These can be seen in Fig. 5, where we display perturbed examples causing misclassifications on the three systems by manipulating VGG Face descriptors via iterative LOTS. To be able to directly compare the various distortions, we show PASS as well as $L_2$ and $L_\infty$ norms of perturbations in sub-captions.

Finally, we can observe that the collected metrics on the produced examples generated via iterative LOTS vary among adversaries. While we cannot see a trend differentiating internal and external adversaries, the distorted examples produced using the various adversaries have significantly different adversarial qualities. We believe this is normal – a face close to the average is naturally closer to others, contrarily, a very characteristic face is farther away and, thus, needs stronger perturbations to be turned to others. For example, as the internal adversary of Hugh Laurie has a unique and very characteristic face among adversaries, it is not surprising that the distorted images of that adversary have one of the worst overall adversarial qualities. On the other hand, we have two external adversaries – Halle Berry and Ewan McGregor – that can be easily turned to other subjects with smaller, less perceptible perturbations relative to other adversaries. Doddington *et al.* [2] referred to such identities as "wolves."

## 5. Conclusion

Since researchers mainly focus on adversarial example generation techniques and, in general, adversarial robustness on end-en-end classification networks, the primary goal of this paper was to extend research to systems that utilize deep features extracted from deep neural networks (DNNs), which is common in biometrics. In this paper, we

Figure 4: ITERATIVE LOTS ON VGG FACE DESCRIPTORS WITH EXTERNAL ADVERSARIES. These perturbed images of external adversaries mimic targeted gallery templates. The VGG Face descriptors of all examples are incorrectly verified to match the gallery templates with Euclidean or cosine distances below FAR = 0.001 thresholds: from top row to bottom, images match gallery templates of Daniel Craig, Hugh Laurie, Idris Elba, Kate Beckinsale, Kristen Bell, and Thandie Newton (cf. Fig. 3).

have introduced our novel layerwise origin-target synthesis (LOTS) algorithm. LOTS is generic and can be efficiently used iteratively to form adversarial examples both on end-to-end classification networks and on systems that use extracted deep features of DNNs.

We have experimentally demonstrated the capabilities of iterative LOTS by generating high quality adversarial exam-

ples on different systems. We have conducted large-scale experiments to compare the adversarial robustness of three face recognition approaches using a dozen adversaries targeting all possible identities. We have generated adversarial examples on the end-to-end VGG Face network via iterative LOTS working on the extracted VGG Face descriptors, and, more traditionally, on features of the Softmax layer. Fur-
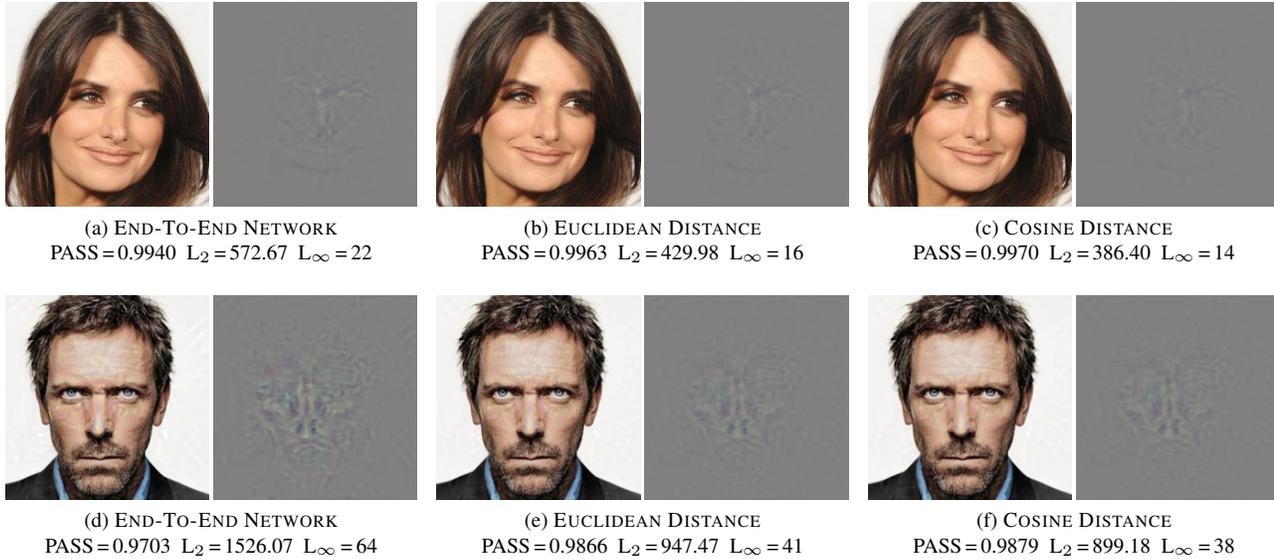
(a) END-TO-END NETWORK
PASS $= 0.9940$ $L_2 = 572.67$ $L_\infty = 22$

(b) EUCLIDEAN DISTANCE
PASS $= 0.9963$ $L_2 = 429.98$ $L_\infty = 16$

(c) COSINE DISTANCE
PASS $= 0.9970$ $L_2 = 386.40$ $L_\infty = 14$

(d) END-TO-END NETWORK
PASS $= 0.9703$ $L_2 = 1526.07$ $L_\infty = 64$

(e) EUCLIDEAN DISTANCE
PASS $= 0.9866$ $L_2 = 947.47$ $L_\infty = 41$

(f) COSINE DISTANCE
PASS $= 0.9879$ $L_2 = 899.18$ $L_\infty = 38$

Figure 5: ADVERSARIAL EXAMPLES VIA ITERATIVE LOTS ON VGG FACE DESCRIPTORS TARGETING KRISTEN BELL. This figure shows adversarial examples paired with their corresponding perturbations that yield incorrect classifications on the end-to-end VGG Face network, and on systems using Euclidean or cosine distance between the extracted VGG Face descriptors and the gallery template of Kristen Bell. The sub-captions show the targeted system, the PASS between the origin and the perturbed image, and the $L_2$ and $L_\infty$ norms of the perturbation.

thermore, using iterative LOTS, we have formed adversarial perturbations on systems that use VGG Face descriptors with Euclidean or cosine distance that are closer to the targeted gallery templates than the FAR $= 0.001$ thresholds.

To assess the robustness of the targeted systems, we have quantified the quality of the produced adversarial examples using the perceptual adversarial similarity score (PASS), and we have measured the percentage of successful attempts where the perturbed images are classified as the targeted identities. A less vulnerable system allows adversaries to impersonate fewer of their targeted identities and/or requires adversaries to form stronger, thus more visible perturbations in order to achieve the targeted misclassifications. Based on the collected metrics, we have concluded that the end-to-end system is more robust to adversarial perturbations formed by iterative LOTS, and the system utilizing cosine distance is the most vulnerable among all. While adversaries could not reach all their targeted identities on the end-to-end VGG Face network, they could achieve that on the other systems utilizing the extracted face descriptors – along with better adversarial qualities. Unfortunately, the system most vulnerable to iterative LOTS is preferred in biometrics due to the fact that, in general, cosine distance provides better performing systems than those that utilize Euclidean distance.

Finally, although we have performed our experiments only using VGG Face descriptors to form adversarial examples, we assume that our results will be portable to other network architectures. We were only targeting "raw" deep features, while deep features are often processed by triplet-loss embedding [13, 18] before the applicable distances are calculated. As these projections are external to the DNN, and the triplet-loss projection matrix from Parkhi *et al.* [13] is not available, we cannot attack deep features after triplet-loss embedding. We conjecture that iterative LOTS is capable of forming examples causing incorrect recognition on systems applying triplet-loss embedding or lower FAR thresholds, the only question is whether the produced examples would be adversarial in terms of human perception. In future work, we will consider attacking such face recognition systems.

## Acknowledgments

# References

[1] J.-C. Chen, V. M. Patel, and R. Chellappa. Unconstrained face verification using deep CNN features. In *Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2016.

[2] G. R. Doddington, W. Liggett, A. F. Martin, M. A. Przybocki, and D. A. Reynolds. Sheep, goats, lambs and wolves: A statistical analysis of speaker performance in the NIST 1998 speaker recognition evaluation. In *International Conference on Spoken Language Processing (ICSPL)*, 1998.

[3] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. DeCAF: A deep convolutional activation feature for generic visual recognition. In *International Conference on Machine Learning (ICML)*, 2014.

[4] G. D. Evangelidis and E. Z. Psarakis. Parametric image alignment using enhanced correlation coefficient maximization. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 30(10), 2008.

[5] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representation (ICLR)*, 2015.

[6] M. Günther, L. El Shafey, and S. Marcel. *Face Recognition Across the Imaging Spectrum*, chapter Face Recognition in Challenging Environments: An Experimental and Reproducible Research Survey. Springer, 1 edition, 2016.

[7] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2016.

[8] R. Jafri and H. R. Arabnia. A survey of face recognition techniques. *Journal of Information Processing Systems (JIPS)*, 5(2), 2009.

[9] B. F. Klare, B. Klein, E. Taborsky, A. Blanton, J. Cheney, K. Allen, P. Grother, A. Mah, and A. K. Jain. Pushing the frontiers of unconstrained face detection and recognition: IARPA Janus benchmark A. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2015.

[10] A. Kurakin, I. J. Goodfellow, and S. Bengio. Adversarial machine learning at scale. In *International Conference on Learning Representation (ICLR)*, 2017.

[11] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *International Conference on Computer Vision (ICCV)*. IEEE, 2015.

[12] A. J. O'Toole, P. J. Phillips, F. Jiang, J. Ayyad, N. Penard, and H. Abdi. Face recognition algorithms surpass humans matching faces over changes in illumination. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 29, 2007.

[13] O. M. Parkhi, A. Vedaldi, and A. Zisserman. Deep face recognition. In *British Machine Vision Conference (BMVC)*, 2015.

[14] R. Ranjan, C. D. Castillo, and R. Chellappa. L2-constrained softmax loss for discriminative face verification, 2017. under review; arXiV preprint: 1703.09507.

[15] A. S. Razavian, H. Azizpour, J. Sullivan, and S. Carlsson. CNN features off-the-shelf: an astounding baseline for recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*. IEEE, 2014.

[16] A. Rozsa, E. M. Rudd, and T. E. Boult. Adversarial diversity and hard positive generation. In *Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*. IEEE, 2016.

[17] S. Sabour, Y. Cao, F. Faghri, and D. J. Fleet. Adversarial manipulation of deep representations. In *International Conference on Learning Representation (ICLR)*, 2016.

[18] S. Sankaranarayanan, A. Alavi, C. D. Castillo, and R. Chellappa. Triplet probabilistic embedding for face verification and clustering. In *Biometrics Theory, Applications and Systems (BTAS)*. IEEE, 2016.

[19] F. Schroff, D. Kalenichenko, and J. Philbin. FaceNet: A unified embedding for face recognition and clustering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2015.

[20] Á. Serrano, I. Martín de Diego, C. Conde, and E. Cabello. Recent advances in face biometrics with Gabor wavelets: A review. *Pattern Recognition Letters*, 31(5), 2010.

[21] Y. Sun, Y. Chen, X. Wang, and X. Tang. Deep learning face representation by joint identification-verification. In *Neural Information Processing Systems (NIPS)*, 2014.

[22] Y. Sun, X. Wang, and X. Tang. Deep learning face representation from predicting 10,000 classes. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2014.

[23] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2015.

[24] C. J. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. In *International Conference on Learning Representation (ICLR)*, 2014.

[25] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2014.

[26] X. Tan, S. Chen, Z.-H. Zhou, and F. Zhang. Face recognition from a single image per person: A survey. *Pattern Recognition*, 39, 2006.

[27] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: From error visibility to structural similarity. *Transactions on Image Processing (TIP)*, 13(4), 2004.