

Query Attack via Opposite-Direction Feature: Towards Robust Image Retrieval

Zhedong Zheng, Liang Zheng, Yi Yang, *Senior Member, IEEE* and Fei Wu, *Senior Member, IEEE*

Abstract—Most existing works of adversarial samples focus on attacking image recognition models, while little attention is paid to the image retrieval task. In this paper, we identify two inherent challenges in applying prevailing image recognition attack methods to image retrieval. First, image retrieval demands discriminative visual features, which is significantly different from the one-hot class prediction in image recognition. Second, due to the disjoint and potentially unrelated classes between the training and test set in image retrieval, predicting the query category from predefined training classes is not accurate and leads to a sub-optimal adversarial gradient. To address these limitations, we propose a new white-box attack approach, Opposite-Direction Feature Attack (ODFA), to generate adversarial queries. Opposite-Direction Feature Attack (ODFA) effectively exploits feature-level adversarial gradients and takes advantage of feature distance in the representation space. To our knowledge, we are among the early attempts to design an attack method specifically for image retrieval. When we deploy an attacked image as the query, the true matches are prone to receive low ranks. We demonstrate through extensive experiments that (1) only crafting adversarial queries is sufficient to fool the state-of-the-art retrieval systems; (2) the proposed attack method, ODFA, leads to a higher attack success rate than classification attack methods, validating the necessity of leveraging characteristics of image retrieval; (3) the adversarial queries generated by our method have good transferability to other retrieval models without accessing their parameters, *i.e.*, the black-box setting.

Index Terms—Image Retrieval, Adversarial Samples, Convolutional Neural Network, Deep Learning.

1 INTRODUCTION













HUMAN can quickly find similar images from limited candidate images, but fail to scale to millions of images. The image retrieval technique is to help user efficiently and effectively find images of interest from a large scale of digital images, which has been applied to many real-world tasks, such as online shopping [28], [23], [26], and tourism recommendation [48], [59], [5]. Recent advances in this field are due to two factors, *i.e.*, the availability of large-scale datasets and the discriminative features extracted from deeply-learned models. The state-of-the-art methods even achieve over 90% Recall@1 accuracy and surpass the human-level performance [73], [43], [47], [30], [23]. However, one problem remains whether these deeply-learned models are accurate and robust enough to various images in the realistic scenario. The robustness evaluation of the image retrieval system has not been well studied. Inspired by the strong ability of human towards small variants, we investigate the adversarial examples to evaluate the retrieval system via adding small, human-imperceptible noise to the original image. The crafted adversarial examples are to imitate the extreme case in the real applications and cheat the retrieval model of predicting a totally different ranking result. In this way, the adversarial attack serves as an indicator to assess the robustness of the target model and understand the weakness of current image retrieval methods.

For attacking image retrieval models, there is no well-studied method. Most existing methods for generating adversarial examples focus on the classification setting, in which the source and target sets share exactly the same classes [54], [11], [20], [32], [9], [61], [62]. However, in image retrieval, the target set has limited overlap or even no overlap classes with the source set (see Figure 1 (a)) [10], [25], [52], [41], [77], [58]. Instead of predicting the class of the input, image retrieval is to compute the similarity between the query and database images, and find relevant images according to the similarity score. In this context, we consider the task of generating adversarial examples from the feature space to fool the retrieval system.

There are two main challenges when deploying existing classification attack methods in the retrieval scenario. First, classification attack methods target class predictions to generate adversarial examples, but this strategy is inconsistent with the testing procedure of image retrieval (see Figure 1 (b)). In image retrieval, we aim to retrieve relevant images from the database images by matching visual features. Therefore, attacking the class prediction does not directly affect the retrieval task, which relies on the intermediate deep features. Second, classification methods attack the class prediction, which usually does not contain the query class in image retrieval. Given a query image of an unseen class, traditional attack methods lead to the inferior adversarial gradient, which compromises the effectiveness of the attack.

To address the potential problems of classification attack approaches, this work focuses on generating adversarial examples tailored for image retrieval. Image retrieval target catching the true matches in the top-K of the ranking list. To achieve a successful adversarial attack on the retrieval system, all true matches should be ranked as low as possible in the ranking list. In this setting, an alternative solution to attacking the query image is to attack the database (candidate image pool). However, the

- Z. Zheng and Y. Yang are with the Australian Artificial Intelligence Institute (AAIL), University of Technology Sydney, Australia, 2007. E-mail: Zhedong.Zheng@student.uts.edu.au, Yi.Yang@uts.edu.au
- L. Zheng is with the Research School of Computer Science, Australian National University, Australia, 0200. E-mail: liang.zheng@anu.edu.au
- F. Wu is with the College of Computer Science and Technology, Zhejiang University, China, 310027. E-mail: wufei@cs.zju.edu.cn

	Source			Target		
Image Recognition	class 1 	class 2 	class 3 	class 1 	class 2 	class 3 
Image Retrieval	class 1 	class 2 	class 3 	class ? 	class ? 	class ? 

(a)

	Evaluation	Adversary
Image Recognition	Class Prediction	Attack Class Prediction
Image Retrieval	Ranking List	Attack Class Prediction?

(b)

Fig. 1. Comparison between image recognition and image retrieval. (a) Problem definition. Image recognition usually recognizes the same classes in the source set and target set. The target set in image retrieval has very few or even no overlapping classes with the source set. (b) The existing classification attack methods employ the classification adversary, which is consistent with its target testing procedure but is inconsistent with the testing procedure in image retrieval. **In this work, we directly attack the visual feature and explore feature-level adversarial gradients, yielding significant attack success rate on the image retrieval task.**

database is sometimes inaccessible and contains thousands of images. Attacking a large number of database images is prohibitive in terms of time cost. In this paper, therefore, we focus on crafting adversarial query images. Without knowledge of the database, we report that adversarial queries alone are sufficient to fool the retrieval system and the cost of generating an adversarial query is relatively low.

We propose a white-box attack method for adversarial example generation in the retrieval context, named Opposite-Direction Feature Attack (ODFA). ODFA exploits the gradient in the feature space, and is based on the target testing procedure, *i.e.*, similarity computation between the query and database images using their respective features. Our key idea underpinning the attack is to explicitly push the feature of the adversarial example away from its original representation. In particular, we first define the *opposite-direction feature*, which, as its name implies, faces the opposite direction to the feature of the original query. During the adversarial attack, we then force the query feature to move towards the *opposite-direction feature*. Due to the revised direction of the feature vector of the adversarial query, the similarity between the database true matches and the adversarial query will decrease. The retrieval model is therefore prone to treating all the true matches as outliers when an adversarial query is used. In the experiment, we show that the proposed ODFA method leads to a large drop in ranking accuracy on five image retrieval datasets, *i.e.*, Food-256 [18], CUB-200-2011 [56], Market-1501 [75], Oxford5k [38] and Paris6k [39]. Under various levels of image perturbation, ODFA outperforms the conventional attack methods, *i.e.*, fast-gradient sign method [11], basic iterative method [20] and iterative least-likely class method [20]. Moreover, when ODFA is adopted on image recognition systems such as ResNet [13] on Cifar-10 [19], the effect of its attack does not show clear superiority over this same set of methods [11], [20]. This indicates that the specificity of our method on image retrieval problems. Additionally, we observe that ODFA has good transferability in the black-box setting; that is, the adversarial queries crafted for one white-box model remain adversarial for another black-box model in the retrieval scenario. In summary, our contributions are itemized below.

- To our knowledge, we are among the early attempts toward attacking image retrieval. While deeply-learned retrieval systems have achieved impressive retrieval performance, they still suffer from the lack of robustness against trivial

noise. We find that crafting adversarial queries alone is sufficient to fool the state-of-the-art retrieval systems.

- We propose a new adversarial attack method, Opposite-Direction Feature Attack (ODFA), tailored for attacking image retrieval. Specifically, our method works on the feature level (feature matters most for retrieval) instead of the classifier level (classifier matters most for classification). We can thus effectively and efficiently generate the adversarial queries. In addition, we extend ODFA to ODFA-MS for multi-scale inputs.
- We conduct experiments on five image retrieval datasets. The experimental results show that the adversarial queries with human-imperceptible noise can successfully cause large performance drops on not only the baseline victim model but also prevailing retrieval approaches.
- Without the knowledge of the target model, the proposed adversarial queries have good transferability and could be applied to the black-box setting.

2 RELATED WORK

Image retrieval. Recently the progress in image retrieval has been due to two factors: the availability of large-scale datasets and the learned representation using the deep neural network. Some early works directly apply the off-the-shelf Convolutional Neural Network (CNN) pretrained on the large-scale dataset like ImageNet [45] to extract visual representations and sort the candidate images according to the feature similarity [70], [55], [23], [7], [64], [17], [25], [65]. The bias between the source set ImageNet and the target datasets, *e.g.*, Oxford5k, compromises the ranking performance. Most works, therefore, collect the related large-scale datasets like Landmarks dataset [1], in which the data shares similar distribution with the target set, and fine-tune the CNN-based model. The visual representation is tuned in an end-to-end manner and shows a stronger discriminative ability [42], [55], [66]. The demon is usually in the details, and another line of works is to explore the detailed information [68], [30], [2], [57]. Some researchers design the objectives to force the CNN model to learn discriminative features and better similarity metric. For instance, triplet loss with hard sampling policy is widely-applied [44], [50], [76], [69]. Other works pay attentions to the local patterns, and explicitly involve the image parts into training [51], [52], [3]. For instance, Radenović *et al.* [43] propose a trainable pooling layer to deal with the

scale problem, and arrive the state-of-the-art performance. Zhang *et al.* [73] propose a dynamic part-matching method and achieve the result surpassing the human-level performance.

Despite the impressive performance, no prior works have explicitly explored the robustness of the retrieval system. One related practice is to simply collect more distractor images and add them to the testing database, such as 100,000 images, to validate the system robustness [38], [39], [75], [12]. This line of the practice considers to increase the complexity of the database and is orthogonal to our method focusing on the query variants. Compared with our method, increasing the database is limited in two aspects. 1) Most collected images are not “strong negatives”, which could cheat the retrieval system of changing the ranking list. The true matches with high similarity score are still in the top-K of the ranking results. 2) Increasing the database also leads to larger time cost. When testing, one needs to extract the feature for a large number of candidate images. In contrast, the proposed adversarial queries do not increase the evaluation time cost and efficiently affect the final ranking result.

Adversarial attack. Adversarial Attack is to craft the sample from the real data to fool the learned model and helps to evaluate the robustness of the target models [4], [72], [21]. Szegedy *et al.* first show that the adversarial images, while looking pretty much the same as the original ones, can mislead the CNN model to classify them into a specific class [54]. It raises the security problem of the current state-of-the-art models [46], [9] and also provides us more insights into the CNN mechanism [11]. The adversarial attack literature can be broadly divided into two classes, *i.e.*, gradient-based attack and score-based attack. Given an input image, gradient-based methods need to know the gradient of the applied model. One of the earliest works is the fast-gradient sign method [11], which generates adversarial examples in one step. Some works extend [11] to iteratively updating the adversarial images with small step sizes, *i.e.*, basic iterative method [20], deep fool [32] and momentum iterative method [8]. Compared with the fast-gradient sign method, the perturbation generated by iterative methods is more smooth, and the adversarial samples are, therefore, more imperceptible to the human. On the other hand, another line of methods relies on searching the input space. Jacobian-based saliency map attack greedily modifies the input instance [37]. In [34], Narodytska *et al.* further shows that single pixel perturbation, which is out of the valid image range, can successfully lead to misclassification on small-scale images. They also extend the method to large-scale images by local greedy searching. Except for pixel modification, sample generation via spatial transform also can result in the adversarial examples [62]. The closest inspiring work is the iterative least-likely class method [20], which makes the classification model output interesting mistakes, *e.g.*, classifying an image of the class *vehicle* into the class *cat*. They achieve this effect by constraining to increase the predicted probability of the least-likely class. This work adopts a similar spirit. In order to fool the retrieval model into assigning the true matches with possibly low ranks, we constrain to increase the similarity of the query feature vector with a vector of an opposite direction in the feature space. Here we emphasize that our work is different from [20] in two aspects. First, Kurakin *et al.* [20] focus on image recognition and rely on class predictions to obtain the least-likely class. In the retrieval setting, the classification model faces images from unseen classes. The inaccurate class prediction compromise the iterative least-likely class method. In this respect, the proposed method directly works on the intermediate feature level and alleviates this problem. Second, Kurakin *et al.* [20]

increase the probability of the least-likely class but do not decrease the probability of the most-likely class. So the true match images / classes are still in the top-K prediction. In comparison, our method explicitly constrains to decrease the similarity of the adversarial image and its original image in the feature space, so that the similarity between the adversarial image and original true-matches also drops. The model is prone to ranking all the true-matches out of the top-K.

3 NOTATION

3.1 Problem Setup

We denote the original query image and its adversarial example as X and X' , respectively. Given a query image X , the image retrieval model extracts the visual feature $f_X = F(X)$, and then ranks the database images according to the similarity score in the feature space. $F(\cdot)$ denotes a nonlinear mapping function. In this work, we study the widely-used image retrieval models, *i.e.*, convolutional neural network (CNN) as the mapping function F . For two images X_m and X_n , the similarity score in the feature space can be formulated as the cosine similarity: $D(X_m, X_n) = \frac{f_{X_m} \cdot f_{X_n}}{\|f_{X_m}\|_2 \cdot \|f_{X_n}\|_2}$, where $\|\cdot\|_2$ is the L2-norm. A high D score indicates that the two images are very similar. To successfully attack the retrieval result, we intend to find the adversarial query X' to lower the similarity score between the query and all true matches. In the mean time, we demand that the difference between the adversarial query and the original query could be as small as possible, which ensures that the adversarial perturbation is imperceptible to the human. In particular, we follow the practice in [20] to keep pixel difference within a valid value range. We clip the pixels whose values fall out of the valid range, and remove the distortions which are larger than the hyper-parameter perturbation rate ϵ : $\text{Clip}_{X,\epsilon}\{X'\} = \min\{255, X + \epsilon, \max\{0, X - \epsilon, X'\}\}$. It ensures that X' with $\|X' - X\|_\infty \leq \epsilon$. Since a large ϵ will make the perturbation perceptible to the human, we set the $\epsilon \leq 16$ in this work.

3.2 Victim Model

We call the model to be attacked as the victim model and adopt the white-box assumption as the conventional gradient-based methods [54], [11], [20], [8]; that is, the parameters of the victim model are accessible. Under the assumption, we can obtain the gradient to the inputs. To verify the effectiveness of the proposed method, we mainly adopt two kinds of widely-used retrieval models as victim model, *i.e.*, *Classification-based retrieval model* and *Ranking-based retrieval model*. The main difference between these two models is whether we can access the category prediction.

Classification-based retrieval model. Classification-based retrieval model exploits category recognition as the proxy task to learn the projection function from data to the feature space. Following the common practice in [1], [60], [40], [22], we train the CNN-based model mapping the training data to the class label, and adopt the feature before the final classification layer as the retrieval feature f . To compare with the traditional classification attack methods, which rely on the class prediction, we preserve the final classification prediction p . The final linear classifier can be formulated as: $p = Wf + b$, which maps the feature f to the class probability p . W and b are learned weights of the classification layer.

Ranking-based retrieval model. Ranking-based retrieval model is trained with the distance-based objectives, *e.g.*, ranking loss [50],

[14], [74] and contrastive loss [43], which pull examples with different class labels apart from each other and push examples from the same classes closer to each other. This line of model does not involve the category prediction part. The conventional classification attack method, therefore, could not work on this kind of models, when our attack method is still feasible.

To compare with the traditional attacking methods and illustrate our intuition, we assume that the category prediction is available when attacking. We adopt the *Classification-based retrieval model* as the victim model in the following Section 4. However, we note that the proposed method does not depend on the category prediction, and also could work on the *Ranking-based retrieval model*. In Experiment, we show the result of the proposed model attacking both kinds of retrieval models.

4 METHOD

In this section, we first extend the traditional adversarial methods to the retrieval scenario and discuss the limitation of these methods. We next introduce the proposed Opposite-Direction Feature Attack (ODFA) method, which addresses the weakness of the traditional methods. Furthermore, we extend ODFA to ODFA-MS, attacking the common evaluation trick, *i.e.*, the feature fusion of multi-scale inputs.

4.1 Adoption of Classification Attack in Image Retrieval

Previous works in the adversarial example generation are designed for image recognition and aim to attack the class prediction [11], [20]. When the prediction is changed, the intermediate activation in the network is also implicitly impacted. Although these methods are not designed for the retrieval problem, they still work for the retrieval scenario with a minor modification. We assume that we could acquire the label prediction of the victim model, and extend these existing classification attack methods to generate the adversarial queries.

Specifically, for the fast-gradient sign method [11] and basic iterative method [20], we deploy the label predicted by the victim model as the pseudo label $y_{max} = \arg \max_y \{p(y|X)\}$. To fool the model, the objective is to decrease the probability $p(y_{max})$ that the adversarial query X' is classified into the pseudo class. The objective is written as,

$$\arg \min_{X'} J(X') = \log(p(y_{max}|X')). \quad (1)$$

For the iterative least-likely class method [20], we calculate the least-likely class $y_{min} = \arg \min_y \{p(y|X)\}$. The attack objective is to increase the probability $p(y_{min})$ so that the input is classified as the least-likely class. The objective is,

$$\arg \max_{X'} J(X') = \log(p(y_{min}|X')). \quad (2)$$

When generating adversarial samples, the weight of the victim model is fixed and we only update the input. For the fast-gradient sign method, $X' = X + \epsilon \text{sign}(\nabla J(X))$. For the iterative methods, *i.e.*, basic iterative method and iterative least-likely class method, we initialize X' with X : $X'_0 = X$, and then update the adversarial samples T times: $X'_{t+1} = X'_t + \alpha \text{sign}(\nabla J(X'_t))$, where α is a relatively small hyper-parameter. Following the practice [20], we set $\alpha = 1$ and the number of the iterations $T = \min(\epsilon + 4, 1.25 \times \epsilon)$. The clip function $\text{Clip}_{X,\epsilon}\{X'\}$ is also added in every iteration to keep pixels of the adversarial query in the valid range.

Algorithm 1 Opposite-Direction Feature Attack (ODFA)

Require: The victim model F ; a real query image X ;

Require: The perturbation rate ϵ .

Ensure: An adversarial example X' with $\|X' - X\|_\infty \leq \epsilon$.

- 1: $X'_0 = X$;
 - 2: $T = \min(\epsilon + 4, 1.25 \times \epsilon)$;
 - 3: **for** $t = 0$ to $T - 1$ **do**
 - 4: Input X'_t to F , extract feature f , calculate the objective $J(X'_t)$;

$$J(X'_t) = \left(\frac{f_{X'_t}}{\|f_{X'_t}\|_2} + \frac{f_X}{\|f_X\|_2} \right)^2. \quad (3)$$
 - 5: Update X'_{t+1} by applying the sign gradient as

$$X'_{t+1} = X'_t + \alpha \text{sign}(\nabla J(X'_t)). \quad (4)$$
 - 6: Keep pixels of the adversarial query in the valid range

$$X'_{t+1} = \text{Clip}_{X,\epsilon}\{X'_{t+1}\}. \quad (5)$$
 - 7: **end for**
 - 8: **return** $X' = X'_T$.
-

Discussion. *Why the conventional classification attack methods can work for retrieval?* The retrieval system learns a semantic projection function, mapping input images to the feature space, which is highly relevant to the semantics category. Although classification attack methods do not target changing the representation, they make changes to the category prediction p of the query. In particular, according to the prediction function $p = Wf + b$ (note that W and b are fixed), the intermediate feature f is also implicitly affected when optimizing the adversarial objective. Due to the changes to the feature, the semantic similarity between the adversarial example and the original image implicitly decreases. The traditional attack methods, therefore, could generate effective adversarial samples for the retrieval system.

What are the disadvantages of the classification attack for retrieval? There are two main disadvantages. First, the source set and the query usually do not contain the same set of classes. The predefined training classes in the source set cannot well represent the semantics of the unseen query. So the predicted most-likely label may not really be the most-likely one, and the predicted least-likely label may not really be the least-likely one, either. Second, the above-mentioned three classification attack methods [20], [11] work on the prediction score and do not explicitly change the similarity in the feature space. The classification attack, therefore, usually obtains a sub-optimal gradient and is limited in their adversarial performance on the retrieval system.

4.2 Opposite-Direction Feature Attack

To overcome the above disadvantages of the classification attack methods, we propose a new attack method named Opposite-Direction Feature Attack (ODFA), which directly works on the intermediate feature without the prerequisite to the category predictions. Specifically, given a query image X , the retrieval model extracts the original feature f_X . We argue that the similarity score $D(X, X_{gt})$ between query X and its true match X_{gt} is relatively high in a well-learned retrieval system. To attack the retrieval model, our target is to minimize the similarity score $D(X', X_{gt})$ between the adversarial query X' and its true match image X_{gt} . To achieve this goal, we define the loss objective as,

$$\arg \min_{X'} J(X') = \left(\frac{f_{X'}}{\|f_{X'}\|_2} + \frac{f_X}{\|f_X\|_2} \right)^2. \quad (6)$$

This loss function aims to push the feature $f_{X'}$ of the adversarial image to the opposite side of the original query feature f_X . We

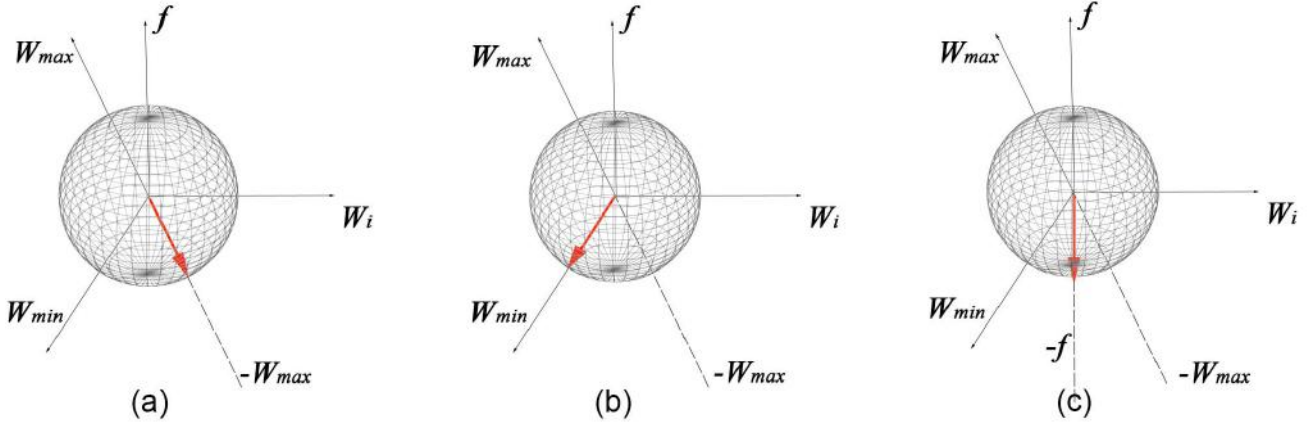


Fig. 2. Geometric interpretation of (a) the fast-gradient sign method [11] and the basic iterative method [20], (b) the iterative least-likely class method [20], and (c) the proposed ODFA. The red arrows represent the direction of the gradient on the original feature f . W_{max} denotes the weight of most-likely class y_{max} and W_{min} denotes the weight of the least-likely class y_{min} . The proposed method does not rely on the class prediction scores and deploys a straightforward opposite gradient direction $-f$ in attacking the retrieval features.

name $-f_X$ as the *opposite-direction feature*. When the objective $J(X') \rightarrow 0$, $\frac{f_{X'}}{\|f_{X'}\|_2}$ will be close to $-\frac{f_X}{\|f_X\|_2}$, $D(X, X') \rightarrow -1$. The similarity score between the adversarial query and the true match images is,

$$D(X', X_{gt}) = \frac{f_{X'}}{\|f_{X'}\|_2} \times \frac{f_{X_{gt}}}{\|f_{X_{gt}}\|_2} \rightarrow -\frac{f_X}{\|f_X\|_2} \times \frac{f_{X_{gt}}}{\|f_{X_{gt}}\|_2} = -D(X, X_{gt}). \quad (7)$$

Since $D(X, X_{gt})$ is usually high in the original retrieval model, we can deduce that the similarity score $D(X', X_{gt})$ is low. To generate an adversarial query X' , we adopt an iterative method to update X' : $X'_0 = X$, $X'_{t+1} = X_t + \alpha \text{sign}(\nabla J(X'_t))$. The clip function is also added to keep pixels in the adversarial sample within the valid range. The overall process of crafting the adversarial query is present in Algorithm 1.

Discussion. We provide a 2D geometric interpretation to illustrate the difference of the gradient direction between the proposed method and traditional attack methods (see Figure 2). The classification attacks use the class prediction $p = Wf + b$, where W is the learned weight and b is the bias term. The weight $W = \{W_1, W_2, \dots, W_K\}$ contains K weights for the K classes in the source set. We use W_{max} to denote the weight of most-likely class y_{max} and W_{min} to denote the weight of the least-likely class y_{min} . For the fast-gradient sign method and the basic iterative method, the gradient on feature f equals to,

$$\frac{\partial J(X')}{\partial f_{X'}} = -W_{max} \times \frac{\partial J(X')}{\partial p(y_{max})}. \quad (8)$$

Note that $\frac{\partial J(X')}{\partial p(y_{max})}$ is a positive constant. So the direction of the gradient is the direction of $-W_{max}$ (see Figure 2 (a)). For the iterative least-likely class method, the gradient equals to,

$$\frac{\partial J(X')}{\partial f_{X'}} = W_{min} \times \frac{\partial J(X')}{\partial p(y_{min})}. \quad (9)$$

The gradient has the same direction with W_{min} (see Figure 2 (b)). For the unseen images of new classes, i.e., query images, $-W_{max}$ and W_{min} are not accurate to describe the adversary of the original query, so the adversarial attack effect is limited. In this

Algorithm 2 Opposite-Direction Feature Attack with Multiple-Scale Inputs (ODFA-MS)

Require: The victim model F ; a real query image X ;

Require: The multiple-scale factors S ; The perturbation rate ϵ .

Ensure: An adversarial example X' with $\|X' - X\|_\infty \leq \epsilon$.

- 1: $X'_0 = X$;
 - 2: $T = \min(\epsilon + 4, 1.25 \times \epsilon)$;
 - 3: **for** $t = 0$ to $T - 1$ **do**
 - 4: Resize the X'_t to $X_t^{s'}$ for s in S
 - 5: Input $X_t^{s'}$ to F , extract features of different scales $f_{X_t^{s'}} = F(X_t^{s'})$.
 - 6: Calculate the objectives $J(X_t^{s'})$ as Eq. 3 and gradients of different-scale inputs;
 - 7: Resize the gradients to the original scale and average the gradients
 - $\nabla J(X'_t) = \frac{1}{n_s} \sum \nabla \tilde{J}(X_t^{s'}). \quad (11)$
 - 8: Update X'_{t+1} by applying the sign gradient as Eq. 4
 - 9: Keep pixels of the adversarial query in the valid range as Eq. 5
 - 10: **end for**
 - 11: **return** $X' = X'_T$.
-

paper, instead of using class predictions, we directly attack the representation in the feature space. According to the Eq. 6, the adversarial gradient of the proposed method is written as,

$$\frac{\partial J(X')}{\partial f_{X'}} = -2 \times \left(\frac{f_{X'}}{\|f_{X'}\|_2} + \frac{f_X}{\|f_X\|_2} \right), \quad (10)$$

where f_X is the feature of the original query image. In Figure 2 (c), we draw the gradient direction of the first iteration. In the first iteration, $f_{X'_0} = f_X$, $\frac{\partial J(X'_0)}{\partial f_{X'_0}} = -4 \frac{f_X}{\|f_X\|_2}$. Our method leads the feature to the opposite direction of the original feature, so the similarity of true matches drops more quickly. The observation in the experiment, as shown in Figure 3 and Figure 4, also verifies that the proposed method is more efficient than the conventional methods.

4.3 Opposite-Direction Feature Attack with Multiple-Scale Inputs

Fusing the features of multiple-scale inputs is a common practice in many image retrieval systems, such as landmark instance retrieval [42], [43]. In particular, when testing, the input image is resized with multiple scale factors $S = \{s_1, s_2, \dots, s_{n_s}\}$, and then the model extracts the feature from inputs of different scales. The mean of the normalized features is used as the final retrieval representation. Since the final representation fuses the feature of the multi-scale inputs, the retrieval system is more robust in terms of the scale variants. In the experiment, we observe that only calculating the adversarial gradient upon the input of the original scale is less effective to fool the image retrieval system. It is due to that the designed imperceptible perturbation is deprecated when resizing images.

To successfully attack the multiple-scale inputs, we further extend the proposed ODFA to ODFA-MS. The whole pipeline is summarized in Algorithm 2. We first view the inputs of different scales as independent inputs X_t^s . Similar to the single scale setting, we calculate the adversarial gradient based on each scale $\nabla J(X_t^s)$. To generate the adversarial gradient towards the original input, we resize all gradients to the original scale $\nabla \tilde{J}(X_t^{s'})$ and average the multi-scale adversarial gradients as follow,

$$\nabla J(X_t') = \frac{1}{n_s} \sum \nabla \tilde{J}(X_t^{s'}), \quad (12)$$

in which n_s is the number of scale factors. Similar to ODFA, we add the sign gradient to the original input and iteratively update the input to obtain the adversarial samples. Since we explicitly consider multi-scale adversarial gradients, the ODFA-MS significantly outperform the regular ODFA in terms of multiple-scale evaluation. More details can be found in Section 5.3.

5 EXPERIMENT

5.1 Datasets and Settings

We evaluate the attack performance on five image retrieval datasets, *i.e.*, Food-256, CUB-200-2011, Market-1501, Oxford5k, Paris6k, and one image recognition dataset, *i.e.*, Cifar-10.

Food-256 is a large-scale food retrieval dataset [18]. The author collects images of 256-kind cuisines. Most of the cuisine categories in this dataset are popular foods in Japan and other countries. There are 31,395 images of 256 cuisines. Following the train / test split in [27], we use 27,849 images of 224 cuisines as the source set and the rest 3,546 images of another 32 cuisines as the target set. In the target set, we select 512 images as query and the rest 3,034 are used as the database images. There is no overlapping class (food category) between the source and target sets.

CUB-200-2011 consists of 11,788 images of 200 bird species [56]. Following [50], we use the CUB-200-2011 dataset for fine-grained image retrieval. The first 100 classes (5,864 images) are used as source set and we evaluate the model on the rest 100 classes (5,924 images).

Market-1501 is a large-scale public pedestrian retrieval dataset [75]. This type of retrieval task is also known as person re-identification (re-ID), which aims at spotting a person of interest across the camera network. For instance, the technique can help quickly find the lost child in a large park, campus, or airport. The author collects images under six different cameras at a university campus. There are 32,668 detected images of 1,501 identities. Following the standard train / test split, we use 12,936 images

of 751 identities as the source set and the rest 19,732 images of another 750 identities as the target set. There is no overlap class (identity) between the source and target sets.

Oxford5k & Paris6k are two widely-used landmark retrieval datasets. Oxford5k contains 5,062 images of 11 particular Oxford buildings [38], and Paris6k contains 6,412 images of 12 particular Paris landmarks [39], respectively. Both datasets are only used as the target set for evaluation. Following the practice in [43], we deploy the non-overlapping building images collected on Flickr as source set to train the model. The source set contains 133k images. **Cifar-10** is a widely-used image recognition dataset, containing 60,000 images of 10 classes [19]. There are 50,000 training images and 10,000 test images. We conduct the image recognition evaluation on this dataset.

Evaluation Metric. With the limited image perturbation, we compare the methods by the drop of the accuracy. The lower accuracy indicates that the adversarial examples make more true matches receive low ranks. For image retrieval, we use two evaluation metrics, *i.e.*, Recall@K and mean average precision (mAP). **Recall@K** is the probability that the right match appears in the top-K of the ranking list. Given a ranking list, the average precision (AP) calculates the space under the recall-precision curve. **mAP** is the mean of the average precision of all queries. For image recognition, we use Top-1 and Top-5 accuracy. **Top-K** is the probability that the right class appears in the top-K predicted classes.

Implementation Details of the Victim Model. For the classification-based retrieval victim model, we follow the common practice in [14], [6] to fine-tune the ResNet-50 [13] by class prediction on Food-256, CUB-200-2011 and Market-1501. During training, the cuisine images in Food-256 are resized to 256×256 , while the pedestrian images of Market-1501 is resized to 256×128 following the previous practices [53], [77]. The images in CUB-200-2011 are first resized with its shorter side to 256, and we then apply a 256×256 random crop to the images. The learning rate is 0.01 for the first 40 epochs and decays to 0.001 for the last 20 epochs. For the ranking-based retrieval victim model, we follow the setting in [43] to train the ResNet-101 [13] on the collected building dataset [43] with contrastive loss. For image recognition, our implementation employs ResNet with 20 layers for the Cifar-10 dataset [13]. The size of the input image is 32×32 . The training policy follows the practice in [13]. The learning rate starts from 0.1 and is divided by 10 after the 150th and 225th epoch. We stop training after 300 epochs.

5.2 Effectiveness of ODFA in the Classification-based Retrieval Model

We first demonstrate the superior attack performance of ODFA to other conventional attack methods. The quantitative results, *i.e.*, Recall@1, Recall@10 and mAP, on Food-256 using clean and adversarial queries are summarized in Figure 3. The victim model using clean queries arrives at a relatively high performance: Recall@1 = 66.41% and mAP = 34.56%. As mentioned, the classification attack method changes the semantic prediction, which *implicitly* changes the retrieval features. When the perturbation rate $\epsilon = 8$, the adversarial images generated by the three classification attacks lead to more than 50% rank-1 error. When $\epsilon = 16$, the iterative least-likely class method even yields a Recall@1 = 10.55%. Nevertheless, these methods are not very effective to move true matches out of the top-10 rank. Although Recall@10

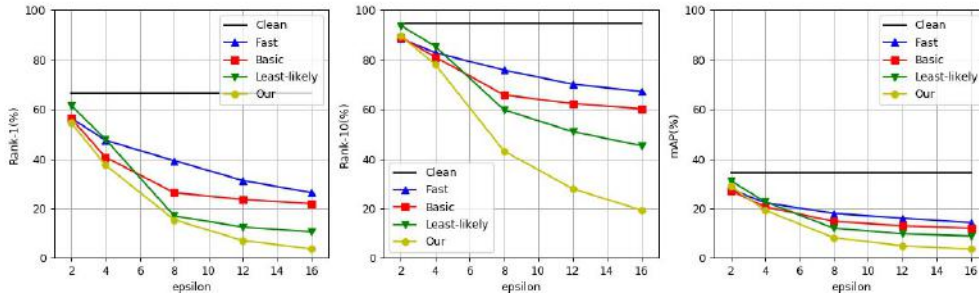


Fig. 3. Recall@1 (%), Recall@10 (%) and mAP (%) of the victim model on Food-256 under the attack by different methods and different perturbation rates ϵ . “Clean” denotes the result obtained by using the original query without any attack. The victim model using clean queries arrives at Recall@1 = 66.41%, Recall@10 = 94.53% and mAP = 34.56%.

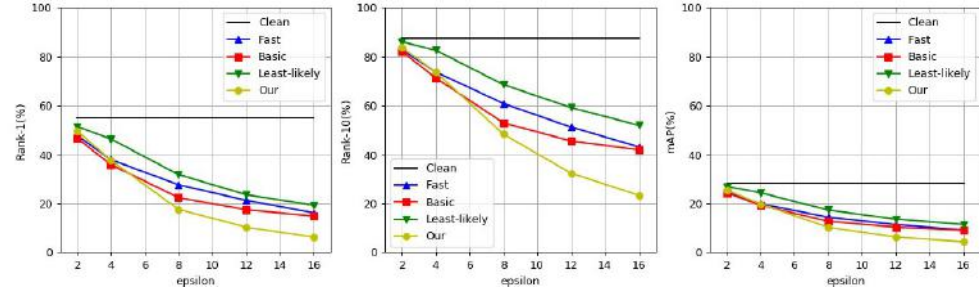


Fig. 4. Recall@1 (%), Recall@10 (%) and mAP (%) of the victim model on CUB-200-2011 under the attack by different methods and different perturbation rates ϵ . “Clean” denotes the result by inputting the original query without any attack. The victim model using clean queries arrives at Recall@1 = 54.86%, Recall@10 = 87.51% and mAP = 28.29%.

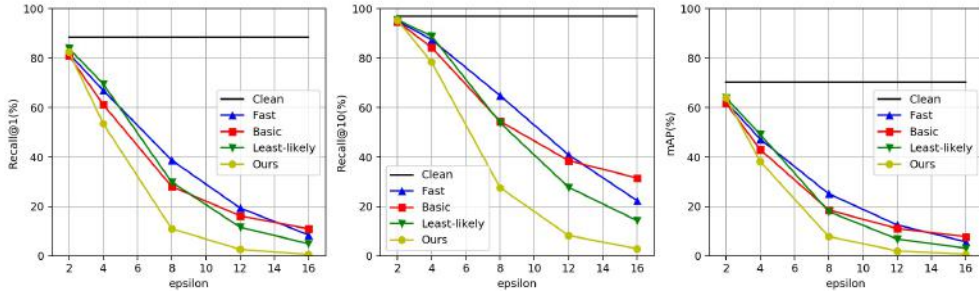


Fig. 5. Recall@1 (%), Recall@10 (%) and mAP (%) of the victim model on Market under the attack by different methods and different perturbation rates ϵ . “Clean” denotes the result by inputting the original query without any attack. The victim model using clean queries arrives at Recall@1 = 88.56%, Recall@10 = 97.03% and mAP = 70.28%.

TABLE 1

Attack the classification-based retrieval model. We mainly compare the proposed method with other attack methods on three datasets, *i.e.*, Food-256, CUB-200-2011 and Market-1501. Here we show the results in % (Lower is better). The perturbation rate is fixed to $\epsilon = 16$. We compare the three classification attack methods, *i.e.*, Fast [11], Basic [20], Least-likely [20].

Methods	Food-256		CUB-200-2011		Market-1501	
	Recall@1	mAP	Recall@1	mAP	Recall@1	mAP
Victim	66.41	34.56	54.86	28.29	88.56	70.28
Fast [11]	26.37	14.25	8.74	4.88	8.49	5.74
Basic [20]	21.88	12.09	8.88	5.57	10.87	7.77
Least-likely [20]	10.55	8.87	9.60	4.78	4.87	3.09
ODFA	3.71	3.59	1.81	1.72	0.68	0.72

TABLE 2

Attack the ranking-based retrieval model. We mainly evaluate the attack performance on two datasets, *i.e.*, Oxford5k and Paris6k, with and without multiple-scale (MS) evaluation. Here we show the results in % (Lower is better). The perturbation rate is fixed to $\epsilon = 16$.

Methods	Oxford5k		Paris6k	
	Recall@1	mAP	Recall@1	mAP
Victim	100.00	86.24	100.00	90.66
ODFA	0.00	0.77	3.69	2.86
Victim (MS)	100.00	88.17	100.00	92.52
ODFA	92.73	73.80	98.18	87.98
ODFA-MS	1.82	2.24	3.64	4.78

TABLE 3

Attack the image recognition on Cifar-10. Here we show the results in % (Lower is better). The perturbation rate is fixed to $\epsilon = 16$. We compare the three classification attack methods, *i.e.*, Fast [11], Basic [20], Least-likely [20].

Methods	Cifar-10	
	Top-1	Top-5
Victim	93.14	99.76
Fast	14.95	67.55
Basic	4.74	21.47
Least-likely	0.03	45.58
ODFA	0.06	0.76

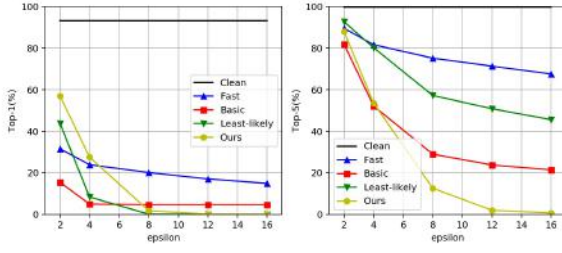


Fig. 6. Top-1 (%) and Top-5 (%) accuracy of the victim model on Cifar-10 under the attack by different methods and different perturbation rates ϵ . “Clean” denotes the result by using original query images.

continues to decrease with increasing ϵ , the best method, *i.e.*, the iterative least-likely class method, only achieves a Recall@10 of 45.31%. In comparison, the proposed ODFA achieves a lower Recall@1 and Recall@10 when $\epsilon = 8$. This can be attributed to the opposite gradient direction attack mechanism. Since the distance between the feature of the adversarial query and that of the original query is much larger, the true matches, which are close to the original query, are thus far from the adversarial query in the feature space. As we increase the perturbation rate ϵ to 16, the victim model yields Recall@1 = 3.71%, Recall@10=19.34%, mAP = 3.59%, which is lower than all the traditional classification attack methods.

The experiments on the fine-grained image retrieval dataset, *i.e.*, CUB-200-2011, and pedestrian retrieval dataset, *i.e.*, Market-1501, indicate similar observation (see Figure 4 and Figure 5). First, due to the subtle differences among the fine-grained classes, the baseline victim model does not arrive a relatively high performance: Recall@1 = 54.86%, Recall@10 = 87.51% and mAP = 28.29% using clean queries. We, however, still can use the proposed method to make the retrieval accuracy even worse. When $\epsilon = 16$, we arrive at Recall@1 = 1.81%, Recall@10 = 8.76% and mAP = 1.72%. Second, compared with the three classification attack methods, our method achieves a larger accuracy drop. Since there are no overlapping bird classes in the source and target sets, the impact of the classification attack is limited. When $\epsilon = 16$, the best classification attack, *i.e.*, fast-gradient sign method arrives at Recall@1=8.74%, Recall@10 = 31.79% and mAP = 4.88%. This accuracy drop is smaller than the drop of the proposed feature-based method. It shows that the proposed ODFA can more effectively and efficiently fool the target retrieval model with the small, human-imperceptible noise. Similarly, on the Market-1501 dataset, the proposed method successfully fools the victim model of predicting worse ranking predictions. The mAP accuracy drops from 70.28% to 0.72%. More quantitative results are shown in Table 1.

5.3 Effectiveness of ODFA in the Ranking-based Retrieval Model

Ranking-based retrieval models use the distance metric, and do not contain the classification prediction part. The traditional classification attack methods, which depend on category prediction, could not work on this line of retrieval models. We, therefore, only evaluate the proposed method to attack the victim model in Table 2. The victim model using clean queries arrives at a high performance: Recall@1 = 100.00%, mAP = 86.24% on Oxford5k and Recall@1 = 100.00%, mAP = 90.66% on Paris6k. When $\epsilon = 16$, the proposed method successfully fools the victim model,

and the accuracy drops to Recall@1 = 0.00%, mAP = 0.77% on Oxford5k and Recall@1 = 3.69%, mAP = 2.86% on Paris6k, respectively.

Furthermore, we evaluate ODFA on attacking the multiple-scale inputs. Following the practice in [43], we extract and fuse the features of multiple-scale inputs. The fusion of multiple-scale features leads to a robust representation towards scale variants, and slightly improves the victim retrieval performance. The victim model arrives at Recall@1 = 100.00%, mAP = 88.17% on Oxford5k and Recall@1 = 100.00%, mAP = 92.52% on Paris6k. We observe that the imperceptible noise generated by ODFA are somehow deprecated after resizing the image, and the attack performance is limited. As shown in Table 2, ODFA works on the original-scale inputs and only achieves a small performance drop. In contrast, the extended ODFA-MS, benefiting of considering the multiple-scale adversarial gradients, successfully fools the victim model. The victim model with multi-scale inputs also drops to a low precision at Recall@1 = 1.82%, mAP = 2.24% on Oxford5k and Recall@1 = 3.64%, mAP = 4.78% on Paris6k.

5.4 Performance of ODFA in Image Recognition

We further test ODFA in image recognition. Results are shown in Figure 6. We observe that our attack does not achieve the largest drop of top-1 accuracy when ϵ is small. This can be explained by the adversarial target. The iterative least-likely class method aims to make the model mis-classify the adversarial example into the least-likely class. In comparison, our method does not increase the probability of a specific class. Although the confidence score of the correct class decreases, there are no competitors to replace the correct top-1 class which already has a high confidence score. Nevertheless, as for top-5 misclassification, the proposed method converges to a lower point than other methods. Since the value of the bias term b for 10 classes is close, we ignore the impact of b . When our method converges, the original top-1 prediction $p = Wf$ becomes the lowest probability $p' = -Wf$. So the correct class is moved out of the top-5 classes quickly. When $\epsilon = 16$, the adversarial images generated by our method compromise the top-5 accuracy from 99.76% to 0.76%. The attacked top-1 accuracy 0.06% is also competitive to the result of iterative least-likely class method 0.03%. In summary, the proposed ODFA method reports competitive performance and is not evidently superior to the competing methods as the case in image retrieval (see Table 3).

5.5 Effectiveness of ODFA in the Black-box Setting

As shown in previous works [54], [37], [35], [36], [31], [33], adversarial examples have good transferability that can successfully attack other black-box models in the recognition scenario, because the models learn a similar decision boundary in the classification space. In this section, we study the transferability of the adversarial queries in terms of the retrieval scenario.

For the classification-based retrieval model, we train a stronger victim with *DenseNet-121* [15] as the black-box model, which arrives at Recall@1 = 89.96% and mAP = 73.39% using “clean” images on Market-1501. The adversarial queries are independently generated by the white-box *ResNet-50* ($\epsilon = 16$). The experiment shows that adversarial samples generated by *ResNet-50* also compromise the performance of *DenseNet-121*: Recall@1 = 10.24% and mAP = 7.88%.

We observe a similar phenomenon on attacking the ranking-based retrieval model. We train the white-box model with *ResNet-101* and use the *ResNet-101* generated adversarial query to attack






























































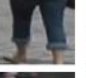
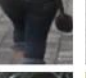

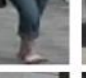
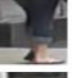
Datasets	Queries	Ranking Results: Rank1→Rank10										
Food-256	Original Query											
	Adversarial Query											
CUB-200-2011	Original Query											
	Adversarial Query											
Market-1501	Original Query											
	Adversarial Query											

Fig. 7. Ranking results of the original queries and the adversarial queries generated by our method. The proposed approach introduces trivial noise on original queries to fool the retrieval system, while the human is robust to such noise. Three original queries are from Food-256 [18], CUB-200-2011 [56] and Market-1501 [75], respectively. The corresponding top-10 retrieval results are also provided. The proposed adversarial queries successfully fool the retrieval model to predict irrelevant ranking results. (Best viewed when zoomed in)

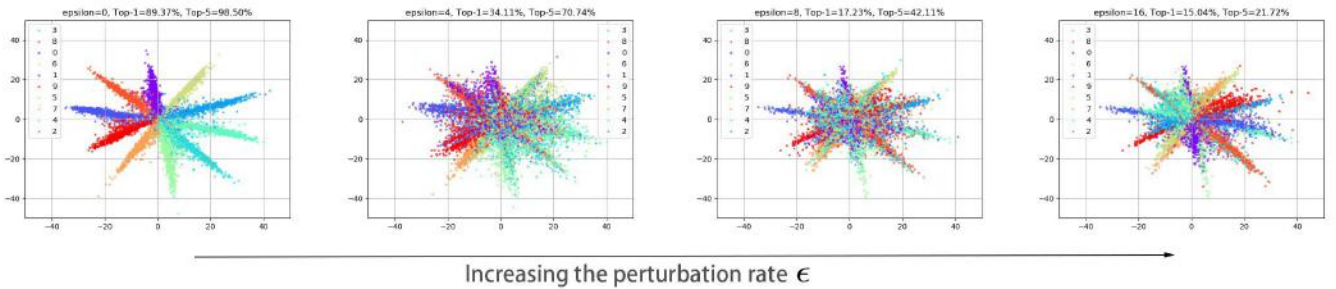


Fig. 8. Feature visualization on Cifar-10 (best viewed in color). As the perturbation rate ϵ increases, the feature gradually moves to the opposite side of the original direction. Top-1 (%) and Top-5 (%) accuracy of the victim model (in the title of each subfigure) also decrease.

the black-box model based on *VGG-16* [49]. The generation process of the adversarial queries are totally independent with the black-box model. The accuracy of the black-box model also drops from 100.00% to 0.00% Recall@1 and 85.24% to 0.79% mAP on Oxford5k. It verifies that the adversarial queries have good transferability and could also be applied to the black-box setting.

5.6 Further Analysis and Discussions

Attack against State-of-the-art Methods. Furthermore, we evaluate our method on some state-of-the-art models, which achieve competitive accuracy on benchmarks. Specifically, for person retrieval (image retrieval), we attack a recent ECCV'18 model called *PCB* [53]. On Market-1501, our re-implementation arrives Recall@1 = 92.70%, mAP = 77.14% using clean queries for the victim model. As shown in Figure 9 (a,b), Recall@1 and

mAP drops to 34.00% and 21.52% respectively by the proposed ODFA. The second best method, Fast-gradient sign method, also arrives a relatively low accuracy 37.11% and 24.40%, but is still smaller than the accuracy drop of the proposed method. For image recognition, we evaluate our method on the prevailing WideResNet-28 [71]. Our re-implementation arrives Top-1 accuracy 96.14% and Top-5 accuracy 99.91% using clean queries, respectively. As shown in Figure 9 (c,d), we have consistent observations with the baseline victim models, *i.e.*, competitive top-1 accuracy drop and largest top-5 accuracy drop. Our method arrives Top-1 accuracy of 0.34% and Top-5 accuracy of 1.29%.

Visualization of Retrieval Results. We provide one qualitative comparison on the retrieval results with original queries and adversarial queries in Figure 7. Since we employ an iterative policy with small steps, the adversarial queries generated by our

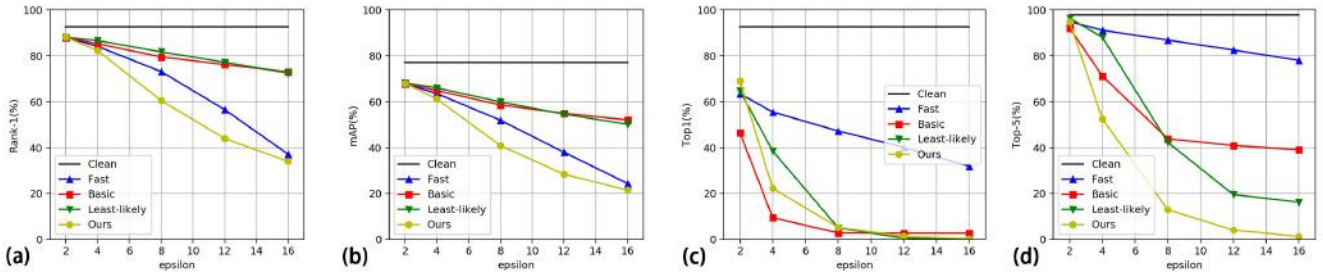


Fig. 9. Performance of attacking state-of-the-art models. (a) and (b): Recall@1 (%) and mAP(%) on Market-1501 when attacking the victim model PCB [53]. (c) and (d): Top-1 and Top-5 accuracy (%) on Cifar-10 when attacking WideResNet-28 [71].

method are visually close to the original query, which simulates extreme retrieval cases to evaluate the model robustness. In these examples, the ranking results obtained by the original queries are good. However, when using the adversarial queries, the top-10 ranked images are all false matches with a significantly different appearance to the query. The adversarial query successfully makes the victim model predict low ranks to the true matches.

Visualization of Attacked Features. Following the visualization trick in [29], we insert an additional 2-dim fully-connected layer into the CNN model to visualize the feature. We train the victim model with an extra fully-connected layer on Cifar-10 and then extract the 2-dim feature of every test image to plot maps. Due to applying the visualization trick (using the 2-dim feature to classify 10 classes), the accuracy of the new victim model is a little bit lower than the baseline result in Table 3, but still arrives at a relatively high accuracy, Top-1=89.37%, and Top-5=98.50%. It is good enough to verify our intuition in the feature space. As shown in Figure 8, the points in the same color belong to the same class. We plot four maps with different perturbation rates $\epsilon = 0, 4, 8, 16$ to see the feature movement. $\epsilon = 0$ is the output of the victim model on clean test images. The features gradually move to the opposite side of the original direction, when ϵ increases. The observation verifies the effectiveness of our objective, *i.e.*, moving to the opposite direction. Comparing the figure of $\epsilon = 0$ with the figure of $\epsilon = 16$, the feature of most adversarial examples successfully move to the opposite side of the original feature. Due to the change of the immediate features, the classification accuracy, as shown in the title of every subfigure, also gradually drops. The observation validates the mechanism of the proposed method.

6 CONCLUSION

In this paper, we consider the adversarial attack in a new realistic setting, *i.e.*, image retrieval, and propose a new attack method named Opposite-Direction Feature Attack (ODFA) tailored for the retrieval scenario. Different from previous works, the proposed attack method does not depend on the category prediction. Instead, ODFA takes the advantage of the intermediate feature and explicitly considers the feature distance in the representation space. On five image retrieval datasets, *i.e.*, Food-256 [18], CUB-200-2011 [56], Market-1501 [75], Oxford5k [38] and Paris6k [39], we validate the effectiveness of the proposed method on two kinds of retrieval victims, *i.e.*, classification-based retrieval model and ranking-based retrieval model. The proposed ODFA leads to a large performance drop in ranking accuracy with human imperceptible perturbation. We also extend the ODFA to adapt the multi-scale evaluation and

verify the effectiveness of ODFA on black-box models. Moreover, we visualize the change of the feature direction to further support our intuition in the feature space. In the future, we will investigate into applying the proposed attack to shallow layers of neural networks and study its effect on other tasks, such as semantic segmentation, action recognition and object detection [63], [24], [16], [67].

REFERENCES

- [1] A. Babenko, A. Slesarev, A. Chigorin, and V. Lempitsky. Neural codes for image retrieval. In *ECCV*, 2014. 2, 3
- [2] S. Bai, X. Bai, and Q. Tian. Scalable person re-identification on supervised smoothed manifold. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 2
- [3] X. Bai, M. Yang, T. Huang, Z. Dou, R. Yu, and Y. Xu. Deep-person: Learning discriminative deep features for person re-identification. *Pattern Recognition*, 98:107036, 2020. 2
- [4] A. Chakraborty, M. Alam, V. Dey, A. Chattopadhyay, and D. Mukhopadhyay. Adversarial attacks and defences: A survey. *arXiv:1810.00069*, 2018. 3
- [5] J. Chen and C.-W. Ngo. Deep-based ingredient recognition for cooking recipe retrieval. In *ACM Multimedia*, 2016. 1
- [6] Y. Chen, X. Zhu, and S. Gong. Person re-identification by deep learning multi-scale representations. In *ICCV*, 2017. 6
- [7] C. Deng, X. Yang, F. Nie, and D. Tao. Saliency detection via a multiple self-weighted graph-based manifold ranking. *IEEE Transactions on Multimedia*, 2019. 2
- [8] Y. Dong, F. Liao, T. Pang, H. Su, J. Zhu, X. Hu, and J. Li. Boosting adversarial attacks with momentum. *CVPR*, 2018. 3
- [9] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno, and D. Song. Robust physical-world attacks on deep learning visual classification. In *CVPR*, 2018. 1, 3
- [10] Y. Fu, T. M. Hospedales, X. Tao, and S. Gong. Transductive multi-view zero-shot learning. *IEEE transactions on pattern analysis and machine intelligence*, 37(11):2332–2345, 2015. 1
- [11] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. *ICLR*, 2015. 1, 2, 3, 4, 5, 7
- [12] H. Guo, C. Zhao, Z. Liu, W. Jinqiao, and L. Hanqing. Learning coarse-to-fine structured feature embedding for vehicle re-identification. *aaai*, 2018. In *AAAI*, 2018. 3
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 2, 6
- [14] A. Hermans, L. Beyer, and B. Leibe. In defense of the triplet loss for person re-identification. *arXiv:1703.07737*, 2017. 4, 6
- [15] G. Huang, Z. Liu, L. V. D. Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In *CVPR*, 2017. 8
- [16] Y.-G. Jiang, Z. Wu, J. Wang, X. Xue, and S.-F. Chang. Exploiting feature and class relationships in video categorization with regularized deep neural networks. *IEEE transactions on pattern analysis and machine intelligence*, 40(2):352–364, 2017. 10
- [17] L. Jin, K. Li, Z. Li, F. Xiao, G.-J. Qi, and J. Tang. Deep semantic-preserving ordinal hashing for cross-modal similarity search. *IEEE transactions on neural networks and learning systems*, 30(5):1429–1440, 2018. 2
- [18] Y. Kawano and K. Yanai. Automatic expansion of a food image dataset leveraging existing categories with domain adaptation. In *ECCV Workshop on Transferring and Adapting Source Knowledge in Computer Vision (TASK-CV)*, 2014. 2, 6, 9, 10

- [19] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. 2009. [2](#), [6](#)
- [20] A. Kurakin, I. Goodfellow, and S. Bengio. Adversarial examples in the physical world. *ICLR Workshop*, 2017. [1](#), [2](#), [3](#), [4](#), [5](#), [7](#)
- [21] J. Li, R. Ji, H. Liu, X. Hong, Y. Gao, and Q. Tian. Universal perturbation attack against image retrieval. In *ICCV*, 2019. [3](#)
- [22] K. Li, G.-J. Qi, and K. A. Hua. Learning label preserving binary codes for multimedia retrieval: A general approach. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 14(1):2, 2018. [3](#)
- [23] Y. Li, T. Yao, Y. Pan, H. Chao, and T. Mei. Deep metric learning with density adaptivity. *IEEE Transactions on Multimedia*, 2019. [1](#), [2](#)
- [24] X. Liang, L. Lin, Y. Wei, X. Shen, J. Yang, and S. Yan. Proposal-free network for instance-level object segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2978–2991, 2017. [10](#)
- [25] K. Lin, J. Lu, C.-S. Chen, J. Zhou, and M.-T. Sun. Unsupervised deep learning of compact binary descriptors. *IEEE transactions on pattern analysis and machine intelligence*, 41(6):1501–1514, 2018. [1](#), [2](#)
- [26] K. Lin, H.-F. Yang, J.-H. Hsiao, and C.-S. Chen. Deep learning of binary hash codes for fast image retrieval. In *CVPRW*, 2015. [1](#)
- [27] M.-Y. Liu, X. Huang, A. Mallya, T. Karras, T. Aila, J. Lehtinen, and J. Kautz. Few-shot unsupervised image-to-image translation. In *CVPR*, 2019. [6](#)
- [28] S. Liu, Z. Song, G. Liu, C. Xu, H. Lu, and S. Yan. Street-to-shop: Cross-scenario clothing retrieval via parts alignment and auxiliary set. In *CVPR*, 2012. [1](#)
- [29] W. Liu, Y. Wen, Z. Yu, and M. Yang. Large-margin softmax loss for convolutional neural networks. In *ICML*, 2016. [10](#)
- [30] X. Liu, W. Liu, T. Mei, and H. Ma. Provid: Progressive and multi-modal vehicle reidentification for large-scale urban surveillance. *IEEE Transactions on Multimedia*, 20(3):645–658, 2017. [1](#), [2](#)
- [31] Y. Liu, X. Chen, C. Liu, and D. Song. Delving into transferable adversarial examples and black-box attacks. In *ICLR*, 2017. [8](#)
- [32] S. M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard. Deepfool: a simple and accurate method to fool deep neural networks. In *CVPR*, 2016. [1](#), [3](#)
- [33] S. M. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard. Universal adversarial perturbations. In *CVPR*, 2017. [8](#)
- [34] N. Narodytska and S. P. Kasiviswanathan. Simple black-box adversarial perturbations for deep networks. *CVPR Workshop*, 2017. [3](#)
- [35] N. Papernot, P. McDaniel, and I. Goodfellow. Transferability in machine learning: from phenomena to black-box attacks using adversarial samples. *arXiv:1605.07277*, 2016. [8](#)
- [36] N. Papernot, P. McDaniel, I. Goodfellow, S. Jha, Z. B. Celik, and A. Swami. Practical black-box attacks against machine learning. In *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security*, pages 506–519. ACM, 2017. [8](#)
- [37] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, and A. Swami. The limitations of deep learning in adversarial settings. *European Symposium on Security & Privacy*, 2016. [3](#), [8](#)
- [38] J. Philbin, O. Chum, M. Isard, J. Sivic, and A. Zisserman. Object retrieval with large vocabularies and fast spatial matching. In *CVPR*, 2007. [2](#), [3](#), [6](#), [10](#)
- [39] J. Philbin, O. Chum, M. Isard, J. Sivic, and A. Zisserman. Lost in quantization: Improving particular object retrieval in large scale image databases. In *CVPR*, 2008. [2](#), [3](#), [6](#), [10](#)
- [40] X. Qian, Y. Fu, Y.-G. Jiang, T. Xiang, and X. Xue. Multi-scale deep learning architectures for person re-identification. In *ICCV*, 2017. [3](#)
- [41] X. Qian, Y. Fu, T. Xiang, Y.-G. Jiang, and X. Xue. Leader-based multi-scale attention deep architecture for person re-identification. *IEEE transactions on pattern analysis and machine intelligence*, 42(2):371–385, 2019. [1](#)
- [42] F. Radenović, G. Tolias, and O. Chum. Cnn image retrieval learns from bow: Unsupervised fine-tuning with hard examples. In *ECCV*, 2016. [2](#), [6](#)
- [43] F. Radenović, G. Tolias, and O. Chum. Fine-tuning cnn image retrieval with no human annotation. *IEEE transactions on pattern analysis and machine intelligence*, 2018. [1](#), [2](#), [4](#), [6](#), [8](#)
- [44] E. Ristani and C. Tomasi. Features for multi-target multi-camera tracking and re-identification. In *CVPR*, 2018. [2](#)
- [45] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015. [2](#)
- [46] M. Sharif, S. Bhagavatula, L. Bauer, and M. K. Reiter. Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. ACM, 2016. [3](#)
- [47] C. Shen, Z. Jin, W. Chu, R. Jiang, Y. Chen, G.-J. Qi, and X.-S. Hua. Multi-level similarity perception network for person re-identification. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 15(2):32, 2019. [1](#)
- [48] B. Sigurbjörnsson and R. Van Zwol. Flickr tag recommendation based on collective knowledge. In *WWW*, 2008. [1](#)
- [49] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv:1409.1556*, 2014. [9](#)
- [50] H. O. Song, Y. Xiang, S. Jegelka, and S. Savarese. Deep metric learning via lifted structured feature embedding. In *CVPR*, 2016. [2](#), [4](#), [6](#)
- [51] Y. Suh, J. Wang, S. Tang, T. Mei, and K. Mu Lee. Part-aligned bilinear representations for person re-identification. In *ECCV*, 2018. [2](#)
- [52] Y. Sun, L. Zheng, Y. Li, Y. Yang, Q. Tian, and S. Wang. Learning part-based convolutional features for person re-identification. *IEEE transactions on pattern analysis and machine intelligence*, 2019. [1](#), [2](#)
- [53] Y. Sun, L. Zheng, Y. Yang, Q. Tian, and S. Wang. Beyond part models: Person retrieval with refined part pooling. *ECCV*, 2018. [6](#), [9](#), [10](#)
- [54] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. *ICLR*, 2014. [1](#), [3](#), [8](#)
- [55] G. Tolias, R. Sivic, and H. Jégou. Particular object retrieval with integral max-pooling of cnn activations. *ICLR*, 2015. [2](#)
- [56] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD Birds-200-2011 Dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011. [2](#), [6](#), [9](#), [10](#)
- [57] J. Wang, K. Sun, T. Cheng, B. Jiang, C. Deng, Y. Zhao, D. Liu, Y. Mu, M. Tan, X. Wang, et al. Deep high-resolution representation learning for visual recognition. *IEEE transactions on pattern analysis and machine intelligence*, 2020. [2](#)
- [58] J. Wang, T. Zhang, N. Sebe, H. T. Shen, et al. A survey on learning to hash. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):769–790, 2017. [1](#)
- [59] Y. Wang, X. Lin, L. Wu, and W. Zhang. Effective multi-query expansions: Collaborative deep networks for robust landmark retrieval. *IEEE Transactions on Image Processing*, 26(3):1393–1404, 2017. [1](#)
- [60] Y. Wei, Y. Zhao, C. Lu, S. Wei, L. Liu, Z. Zhu, and S. Yan. Cross-modal retrieval with cnn visual features: A new baseline. *IEEE transactions on cybernetics*, 47(2):449–460, 2016. [3](#)
- [61] C. Xiao, B. Li, J.-Y. Zhu, W. He, M. Liu, and D. Song. Generating adversarial examples with adversarial networks. *IJCAI*, 2018. [1](#)
- [62] C. Xiao, J.-Y. Zhu, B. Li, W. He, M. Liu, and D. Song. Spatially transformed adversarial examples. *ICLR*, 2018. [1](#), [3](#)
- [63] C. Xie, J. Wang, Z. Zhang, Y. Zhou, L. Xie, and A. Yuille. Adversarial examples for semantic segmentation and object detection. In *ICCV*, 2017. [10](#)
- [64] E. Yang, C. Deng, C. Li, W. Liu, J. Li, and D. Tao. Shared predictive cross-modal deep quantization. *IEEE transactions on neural networks and learning systems*, 29(11):5292–5303, 2018. [2](#)
- [65] H.-F. Yang, K. Lin, and C.-S. Chen. Supervised learning of semantics-preserving hash via deep convolutional neural networks. *IEEE transactions on pattern analysis and machine intelligence*, 40(2):437–451, 2017. [2](#)
- [66] X. Yang, P. Zhou, and M. Wang. Person reidentification via structural deep metric learning. *IEEE Transactions on Neural Networks and Learning Systems*, (99):1–12, 2018. [2](#)
- [67] Y. Yang, C. Deng, S. Gao, W. Liu, D. Tao, and X. Gao. Discriminative multi-instance multitask learning for 3d action recognition. *IEEE Transactions on Multimedia*, 19(3):519–529, 2016. [10](#)
- [68] Q. Yu, X. Chang, Y.-Z. Song, T. Xiang, and T. M. Hospedales. The devil is in the middle: Exploiting mid-level representations for cross-domain instance matching. *arXiv preprint arXiv:1711.08106*, 2017. [2](#)
- [69] R. Yu, Z. Dou, S. Bai, Z. Zhang, Y. Xu, and X. Bai. Hard-aware point-to-set deep metric for person re-identification. In *ECCV*, 2018. [2](#)
- [70] J. Yue-Hei Ng, F. Yang, and L. S. Davis. Exploiting local features from deep networks for image retrieval. In *CVPRW*, 2015. [2](#)
- [71] S. Zagoruyko and N. Komodakis. Wide residual networks. *arXiv:1605.07146*, 2016. [9](#), [10](#)
- [72] S. Zhang, R. Ji, J. Hu, X. Lu, and X. Li. Face sketch synthesis by multidomain adversarial learning. *IEEE transactions on neural networks and learning systems*, 30(5):1419–1428, 2018. [3](#)
- [73] X. Zhang, H. Luo, X. Fan, W. Xiang, Y. Sun, Q. Xiao, W. Jiang, C. Zhang, and J. Sun. Alignedreid: Surpassing human-level performance in person re-identification. *arXiv:1711.08184*, 2017. [1](#), [3](#)
- [74] X. Zhang, R. Zhang, J. Cao, D. Gong, M. You, and C. Shen. Part-guided attention learning for vehicle re-identification. *arXiv:1909.06023*, 2019. [4](#)
- [75] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian. Scalable person re-identification: A benchmark. In *ICCV*, 2015. [2](#), [3](#), [6](#), [9](#), [10](#)
- [76] Z. Zheng, L. Zheng, and Y. Yang. A discriminatively learned cnn embedding for person reidentification. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 14(1):13, 2018. [2](#)
- [77] Z. Zhong, L. Zheng, Z. Luo, S. Li, and Y. Yang. Learning to adapt invariance in memory for person re-identification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. [1](#), [6](#)



Zhedong Zheng received the B.S. degree in computer science from Fudan University, China, in 2016. He is currently a Ph.D. student with the School of Computer Science at University of Technology Sydney, Australia. His research interests include robust learning for image retrieval, generative learning for data augmentation, and unsupervised domain adaptation.



Liang Zheng is a Lecturer and a Computer Science Futures Fellow in the Research School of Computer Science, Australian National University. He received the Ph.D degree in Electronic Engineering from Tsinghua University, China, in 2015, and the B.E. degree in Life Science from Tsinghua University, China, in 2010. He was a postdoc researcher in the Centre for Artificial Intelligence, University of Technology Sydney, Australia. His research interests include image retrieval, classification.



Yi Yang received the Ph.D. degree in computer science from Zhejiang University, Hangzhou, China, in 2010. He is currently a professor with University of Technology Sydney, Australia. He was a Post-Doctoral Research with the School of Computer Science, Carnegie Mellon University, Pittsburgh, PA, USA. His current research interest includes machine learning and its applications to multimedia content analysis and computer vision, such as multimedia analysis and video semantics understanding.



Fei Wu received his PhD degree in Computer Science from Zhejiang University in 2002. He is currently a full professor with the College of Computer Science and Technology, Zhejiang University. His current research interests include artificial intelligence, multimedia retrieval and machine learning.