TOWARDS MAKING A TROJAN-HORSE ATTACK ON TEXT-TO-IMAGE RETRIEVAL

Fan Hu Aozhu Chen Xirong Li*

AIMC Group, MoE Key Lab of DEKE, Renmin University of China

ABSTRACT

While deep learning based image retrieval is reported to be vulnerable to adversarial attacks, existing works are mainly on image-to-image retrieval with their attacks performed at the front end via query modification. By contrast, we present in this paper the first study about a threat that occurs at the back end of a text-to-image retrieval (T2IR) system. Our study is motivated by the fact that the image collection indexed by the system will be regularly updated due to the arrival of new images from various sources such as web crawlers and advertisers. With malicious images indexed, it is possible for an attacker to indirectly interfere with the retrieval process, letting users see certain images that are completely irrelevant w.r.t. their queries. We put this thought into practice by proposing a novel Trojan-horse attack (THA). In particular, we construct a set of Trojan-horse images by first embedding word-specific adversarial information into a QR code and then putting the code on benign advertising images. A proof-of-concept evaluation, conducted on two popular T2IR datasets (Flickr30k and MS-COCO), shows the effectiveness of the proposed THA in a white-box mode.

Index Terms— Text-to-image retrieval, Trojan-horse attack, adversarial patch generation

1. INTRODUCTION

Text-to-image retrieval (T2IR), with its cross-modal matching ability, allows us to retrieve *unlabeled* images by naturallanguage queries. Given the increasing amounts of unlabeled or subjectively labeled images in both public and private domains, T2IR is crucial for next-generation image search [2, 3]. The state-of-the-art of T2IR relies on big vision-language models, with CLIP [1] as a pronounced manifestation. In this paper, we present a novel *Trojan-horse*-style attack on a CLIP-based T2IR system, see Fig. 1.

Existing attacks on deep learning based image retrieval are conducted mostly in the context of image-to-image retrieval (I2IR), where one uses a specific image as query to find visually similar images in a given collection [4, 5, 6, 7]. In order to fool a targeted I2IR system, the given query image has to be modified. PIRE [4] and DAIR [6], for instance, modify the query image to make it adversarial to the underlying I2IR model so that visually similar images cannot be top-ranked. Such an attack might be used for anti-plagiarism detection or anti-geo-localization. TMAA [5] attempts to hide a user's real query by embedding the query information in an invisible manner into a benign or carrier image, which is then submitted to the I2IR system. All these good efforts discuss threats to (image-to-)image retrieval at the front end. By contrast, we study a potential threat that occurs at the back end.

Our motivation is as follows. The data collection indexed by an image search engine is not fixed. Rather, it has to be regularly updated or expanded as new images are gathered from various sources, *e.g.*, web crawlers and advertisers. With malicious images indexed, an attacker may indirectly interfere with the retrieval process, making users see certain images that are not supposed to be retrieved in a normal condition.

Notice that QR codes are commonly seen in advertising images. We therefore use advertising images randomly collected from the Internet, see Fig. 1(a), as our choice of benign images, and use QR codes as adversarial patches. As illustrated in Fig. 1(b), we construct malicious images by overlaying modified QR codes on the benign images. In particular, we learn to generate the adversarial patches in a word-specific manner so that the malicious images will be only activated by queries containing a specific word. Such a behavior is in a way similar to that of a Trojan-horse malware. We thus term such sorts of malicious images as Trojan-horse (TH) images. Once the TH images are indexed by the image search engine, a Trojan-horse attack (THA) may occur, see Fig. 1(c).

To sum up, our main contributions are as follows. To the best of our knowledge, this paper is the first study that discusses a back-end threat to a (big-model driven) T2IR system. Such a threat can occur if the system unconsciously indexes certain malicious images. We implement the threat by proposing a novel Trojan-horse attack (THA) on CLIP-based T2IR. We provide a proof-of-concept evaluation on two T2IR datasets (Flickr30k and MS-COCO), showing the effectiveness of THA in a white-box mode. Source code is released¹.

2. PROPOSED METHOD

We formalize a Trojan-horse attack (THA) on a given text-toimage retrieval (T2IR) system as follows. Suppose the sys-

^{*}Corresponding author: Xirong Li (xirong@ruc.edu.cn)

¹https://github.com/fly-dragon211/tth

tem, driven by a deep cross-modal matching network \mathcal{N} , has indexed a set of n_0 images X_0 . Each image $x \in X_0$ has been represented by a cross-modal feature vector denoted by e(x). The system answers a natural-language query s, by first encoding the query into a cross-modal feature e(s) that shares the same feature space as e(x). The relevance of each image *w.r.t.* the query is computed in terms of certain (dis)similarity between the corresponding features. The top k most relevant images are returned as the search result. A THA is to con-



(b) Converting benign images to Trojan-horse (TH) images Query: The dog is running around the **cow**.



(c) Top-5 image retrieval results, without or with TH images indexed

Fig. 1: Illustration of our proposed Trojan-horse attack (THA) on CLIP-based [1] text-to-image retrieval (T2IR). Given a specific word w, say cow, and a set of benign images X_b (advertising images from the Internet), we construct a set of TH images $X_{h,w}$ by overlaying a word-specific adversarial patch δ on each image in X_b . The patch is derived by iteratively making an attack on a known CLIP model to enforce the TH images to be more close to the given word in the cross-modal feature space. Consequently, the TH images may appear in the search results of a query containing the word.

struct a set of n_h TH images X_h such that once the indexed collection is expanded as $X_0 \cup X_h$, the top k images will contain items from X_h . Consequently, users are shown with images the attacker wants them to see, even though the images can be completely irrelevant to their information needs.

2.1. Adversarial Patch based TH Image Generation

We start with a set of n_h benign images X_b . To simulate a common procedure that expands the database of an image search engine by adding advertising images, we instantiate X_b with such types of images randomly collected from the Internet, see Fig. 1(a). The TH image set X_h is generated by modifying certain amount of pixels of X_b to embed the attack.

In order to let $x_h \in X_h$ be ranked higher w.r.t. the given query s, the similarity between $e(x_h)$ and e(s) shall be larger. However, due to the *ad-hoc* nature of queries in T2IR, s is not known a priori. Directly targeting the query is thus difficult. Alternatively, we aim to construct X_h for a specific word w so that the THA remains effective for a given query s_w that contains w. A word-specific X_h is denoted as $X_{h,w}$. In order to let $X_{h,w}$ be more close to w in the cross-modal feature space, we introduce a loss as follows

$$\ell(X_h, w) = \frac{1}{n_h} \sum_{x_h \in X_h} (1 - \cos(e(w), e(x_h))), \quad (1)$$

where *cos* indicates the cosine similarity as commonly used for cross-modal matching [8]. We generate $X_{h,w}$ by minimizing $\ell(X_h, w)$.

THA embedding via adversarial patches. As putting a QR code on an advertising image is common, we propose a patch-based THA attack where the adversarial information is embedded into the QR code yet without affecting its usability. Specifically, we use δ to indicate an adversarial patch. Such a patch is practically obtained in an iterative manner, so we use δ_i to denote the patch after the *i*-th iteration, i = 1, 2, ..., t, where *t* is a pre-specified maximum number of iterations. Accordingly, a TH image derived from a specific benign image x_b can be formally expressed as

$$x_{h,i} = (1 - M) \odot x_b + M \odot (\text{zero-padding}(\delta_i)), \quad (2)$$

where M is a pre-specified binary mask that determines where x_b is overlaid with δ_i and \odot represents pixel-wise multiplication. In this work, the patch is placed at the top-right corner to prevent occlusion of the main part of the image, see Fig. 1(b). To make the patch less significant in $x_{h,i}$, it is downsized, with the ratio of its size to the benign images being 0.1, unless stated otherwise. Hence, zero-padding on δ_i is needed in Eq. 2. The initial state of the patch², denoted by δ_o , is fixed to be the QR code of the Trojan-horse wiki page³.

²The patch can also be initialized with a region of a benign image. ³https://en.wikipedia.org/wiki/Trojan_horse

Minimizing Eq. 1 alone will introduce distortion to the patch that makes the QR code not scannable. To preserve the code's usability, we add an *l*2 distance-based constrain, obtaining a combined loss as

$$\underbrace{\ell(X_h, w)}_{\text{Attack effectiveness}} + \lambda \underbrace{\|\delta - \delta_o\|^2}_{\text{QR-code usability}}, \quad (3)$$

where λ is a positive hyper-parameter that strikes a balance between the attack effectiveness and the QR-code usability. Per iteration, given ∇_i as the back-propagated gradient *w.r.t.* Eq. 3, the adversarial patch is updated as

$$\delta_i = \max(0, \min(255, \delta_{i-1} + \eta \cdot \nabla_i)), \tag{4}$$

with η as the learning rate. Note the max-min operation is used to ensure the validity of the pixel values.

Concerning the word embedding e(w), a straightforward choice is to compute e(w) by feeding w into the textual encoder of the network \mathcal{N} . However, the meaning of a word is context-dependent, subject to the sentence that uses the word. In order to obtain a contextualized embedding of a given word, we gather m sentences having w, denoted as S_w , from a training corpus, and subsequently perform mean pooling over the sentence embeddings, *i.e.*,

$$e(w) = \frac{1}{m} \sum_{s \in S_w} e(s).$$
(5)

As e(w) is fixed, maximizing the cosine similarity between $e(x_h)$ and e(w) means performing an iterative images-toword move in the cross-modal feature space, see Fig. 1(b). The entire procedure is summarized as Algorithm 1.

Algorithm 1: Trojan-horse image set generation				
input : A given word w ;				
A benign-image set X_b ;				
A normal QR code δ_o ;				
A cross-modal matching network \mathcal{N} ;				
output: A Trojan-horse image set $X_{h,w}$				
1 Compute word embedding $e(w)$ by Eq. 5;				
2 for $i = 1,, t$ do				
3 Generate Trojan-horse images				
$X_i = (1 - M) \odot X_b + M \odot (\text{zero-padding}(\delta_{i-1}));$				
Compute image embeddings $e(X_i)$ using \mathcal{N} ;				
Compute the combined loss by Eq. 3;				
6 Update the patch δ_i by Eq. 4;				
7 $X_{h,w} \leftarrow X_t$				

2.2. Deep Cross-Modal Matching Network

As our proposed method is generic, any cross-modal matching network that produces e(w) and e(x) in an end-to-end manner can in principle be used. We instantiate \mathcal{N} with CLIP (ViT-B/32) [1], an up-to-date open-source model for image-text matching⁴. CLIP consists of a GPT for text embedding and a ViT for image embedding. Both e(w) and e(x) are 512-dimensional. The model has been pre-trained on web-scale image-text corpora by contrastive learning.

3. EXPERIMENTS

3.1. Experimental Setup

Datasets. We use two popular T2IR datasets: Flickr30k [9] and MS-COCO [10]. Flickr30k has 30k images for training, 1k images for validation and 1k images for test, while the data split of COCO is 83k / 5k / 5k for training / validation / test. Per dataset, we finetune the original CLIP, terming the resultant models CLIP-*flickr* and CLIP-*coco*, respectively.

Implementation. Recall that our TH images are keywordspecific. We build a diverse set of 24 keywords by randomly selecting 8 nouns (*jacket, dress, floor, female, motorcycle, policeman, cow, waiter*), 8 verbs (*smiling, climbing, swimming, reading, run, dancing, floating, feeding*) and 8 adjectives (*blue, front, little, green, yellow, pink, navy, maroon*) from Flickr30k captions. We gathered at maximum m = 500sentences per word. For patch generation, we use an initial learning rate of 0.01 and perform 300 iterations. The hyper-parameter λ in Eq. 3 is empirically set to 0.3.

Performance metric. We report Recall at 10 (R10), *i.e.*, the percentage of test queries that have relevant images included in top-10 search results. A THA shall increase R10 of the TH images and decrease R10 of truly relevant images.

3.2. Experiment 1. THA in a Fully White-box Setup

We first try THA in a fully white-box setup, where both the T2IR model and the data source upon which T2IR is performed are known. Tab. 1 shows R10 scores of specific T2IR models with or without THA. Consider T2IR on Flickr30k for instance. When THA is applied, R10 of truly relevant images by CLIP decreases from 94.9 to 52.1. Meanwhile, R10 of the TH images is 97.2, meaning that for 97 out of 100 queries, there will be at least one TH image shown in their top-10 search results. Similar results are observed on COCO, suggesting that the origial CLIP model is vulnerable to THA.

The finetuned CLIPs (CLIP-*flickr* and CLIP-*coco*) seem to be less vulnerable to THA, which obtain lower R10 of 77.1 and 44.9 *w.r.t.* the TH images, respectively. Our interpretation is as follows. As the TH images are to compete with the relevant images for the top-10 positions, it is more difficult to make a THA given a better T2IR model which tends to place more relevant images at the top. Nonetheless, we believe that showing TH images with R10 of 44.9 will make a noticeably negative impact on a user's search experience.

⁴https://github.com/openai/CLIP

Table 1: Evaluating THA in a fully white-box setup. Metric: Recall at 10 (R10). Lower recall scores of relevant images and higher recall scores of TH images are better.

T2IR model	THA	Benign / TH (†)	Rel. images (\downarrow)	
Dataset: Flickr30k				
CLIP	-	0.0	94.9	
CLIP	+	97.2	52.1	
CLIP-flickr	-	0.8	98.6	
CLIP-flickr	+	77.1	95.3	
Dataset: COCO				
CLIP	-	0.0	80.5	
CLIP	+	82.3	51.2	
CLIP-coco	-	1.8	90.1	
CLIP-coco	+	44.9	83.9	

The effect of two main hyper-parameters in THA, *i.e.*, λ and patch size, is shown in Fig. 2. Choosing λ between 0.3 and 1 allows us to generate adversarial patches effective for THA and scannable. As for the path size, using a small ratio of 0.05 (to the benign images) is sufficient to affect over 70% of the test queries.



Fig. 2: The effect of hyper-parameters in THA, *i.e.*, λ and patch size. Data source: Flickr30k. Model: CLIP-*flickr*.

3.3. Experiment 2. THA in a Cross-dataset Setup

We now try THA in a cross-dataset setup, where the T2IR model is known yet the data source upon which T2IR is performed is unknown. More specifically, when T2IR is performed on Flickr30k (COCO), we use COCO (Flickr30k) as a surrogate dataset to generate TH images. The results are summarized in Tab. 2. Although the R10 scores of TH images are relatively lower than their counterparts in Tab. 1, *e.g.*, 58.6 *vs* 77.1 for CLIP-*flickr* and 44.5 *vs* 44.9 for CLIP-*coco*, the proposed THA has successfully affected a substantial part of the test queries.

3.4. Experiment 3. THA in a Cross-weights Setup

Lastly, we evaluate THA in a cross-weights setup, where the weights of the targeted T2IR model, *i.e.*, CLIP-*flickr* on

Table 2: Evaluating THA in a cross-dataset setup.

T2IR model	TH images (†)	Rel. images (↓)
T2IR on Flick	r30k, with TH gener	ration using COCO:
CLIP	83.1	70.5
CLIP-flickr	58.6	97.7
T2IR on COC	O, with TH generat	ion using Flickr30k:
CLIP	77.8	53.5
CLIP-coco	44.5	84.0

Flickr30k and CLIP-*coco* on COCO, is unknown to the attacker. The original CLIP is therefore used for TH image generation. The performance of THA in the cross-weights setup is shown in Tab. 3. Due to the natural discrepancy between the weights of CLIP and its finetuned conterparts, the percentage of test queries affected by THA is clearly reduced (18.5 on Flickr30k and 7.8 on COCO). How to make a successful THA in a black-box mode is challenging and deserves future investigation.

Table 3: Evaluating THA in a cross-weights setup, wherethe original CLIP is used for TH image generation.

T2IR model	TH images (\uparrow)	Rel. images (\downarrow)
CLIP-flickr	18.5	97.8
CLIP-coco	7.8	90.1

4. CONCLUSIONS AND REMARKS

We propose Trojan-horse attack (THA), a new form of threat to CLIP-based text-to-image retrieval (T2IR). Our pilot study with T2IR experiments on Flickr30k and MS-COCO allows us to draw conclusions as follows. In a white-box mode where the targeted T2IR model is known to the attacker, a substantial amount of test queries (58.6% on Flickr30k and 44.5% on MS-COCO) will be effectively affected by the THA. For these queries, Trojan-horse images are ranked in the top-10 search results, although the images are completely irrelevant *w.r.t.* to the queries. Our experiments also indicate a clear performance gap between THA in the white-box mode and that in a black-box mode. We believe that by enhancing THA with proper black-box attack techniques, the gap can be much reduced in the near future.

Acknowledgments. This work was supported by NSFC (62172420), Tencent Marketing Solution Rhino-Bird Focused Research Program, the Outstanding Innovative Talents Cultivation Funded Programs 2022 of Renmin University of China, and Public Computing Cloud, Renmin University of China.

5. REFERENCES

- [1] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning transferable visual models from natural language supervision," in *ICML*, 2021.
- [2] C. Liu, Z. Mao, W. Zang, and B. Wang, "A neighboraware approach for image-text matching," in *ICASSP*, 2019.
- [3] Z. Ma, F. Liu, J. Dong, X. Qu, Y. He, and S. Ji, "Hierarchical similarity learning for language-based product image retrieval," in *ICASSP*, 2021.
- [4] Z. Liu, Z. Zhao, and M. Larson, "Who's afraid of adversarial queries?: The impact of image modifications on content-based image retrieval," in *ICMR*, 2019.
- [5] G. Tolias, F. Radenovic, and O. Chum, "Targeted mismatch adversarial attack: Query with a flower to retrieve the tower," in *ICCV*, 2019.

- [6] M. Chen, J. Lu, Y. Wang, J. Qin, and W. Wang, "DAIR: A query-efficient decision-based attack on image retrieval systems," in *SIGIR*, 2021.
- [7] S. Hu, Y. Zhang, X. Liu, L. Y. Zhang, M. Li, and H. Jin, "AdvHash: Set-to-set targeted attack on deep hashing with one single adversarial patch," in ACMMM, 2021.
- [8] J. Dong, X. Li, C. Xu, X. Yang, G. Yang, and X. Wang, "Dual encoding for video retrieval by text," *TPAMI*, vol. 44, no. 8, pp. 4065–4080, 2022.
- [9] B. A. Plummer, L. Wang, C. M. Cervantes, J. C. Caicedo, J. Hockenmaier, and S. Lazebnik, "Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models," in *ICCV*, 2015.
- [10] T. Lin, M. Maire, S. J. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common objects in context," in *ECCV*, 2014.