
Protect, Show, Attend and Tell: Empowering Image Captioning Models with Ownership Protection

Jian Han Lim¹ Chee Seng Chan¹ Kam Woh Ng² Lixin Fan² Qiang Yang³

¹Universiti Malaya, Kuala Lumpur, Malaysia

²WeBank AI Lab, Shenzhen, China

³Hong Kong University Science and Technology, Hong Kong

Abstract

By and large, existing Intellectual Property (IP) protection on deep neural networks typically i) focus on image classification task only, and ii) follow a standard digital watermarking framework that was conventionally used to protect the ownership of multimedia and video content. This paper demonstrates that the current digital watermarking framework is insufficient to protect image captioning tasks that are often regarded as one of the frontier AI problems. As a remedy, this paper studies and proposes two different embedding schemes in the hidden memory state of a recurrent neural network to protect the image captioning model. From empirical points, we prove that a forged key will yield an unusable image captioning model, defeating the purpose of infringement. To the best of our knowledge, this work is the first to propose ownership protection on image captioning task. Also, extensive experiments show that the proposed method does not compromise the original image captioning performance on all common captioning metrics on Flickr30k and MS-COCO datasets, and at the same time it is able to withstand both removal and ambiguity attacks. Code is available at <https://github.com/jianhanlim/ipr-imagecaptioning>

1 Introduction

Recent advances in deep neural networks (DNN) had significantly improved the overall model performance in multiple artificial intelligence (AI) domains, for example, natural language processing, computer vision, gaming, etc. As a result of this, it has enabled a growing number of AI start-ups and companies to offer their DNN solutions in terms of Software as a Service (SaaS). As such, the protection of the Intellectual Property (IP) of DNN has become a necessity in order to protect the model against IP infringement to preserve the owner's competitive advantage in an open market.

For the past few years, IP protection on DNN [1–8] has been a significant research area. Ideally, the goal is the IP protection solution should not degrade the performance of the original model, and at the same time, it must also be resilient against both ambiguity and removal attacks. Although all these existing solutions have achieved this goal, it is unsatisfactory in our view as we found out that all existing DNN watermarking methods have been i) following a standard digital watermarking framework that was conventionally used to protect the ownership of multimedia and video content, and ii) focusing on DNNs for classification tasks that map images to labels and DNNs for other tasks are forgotten such as image captioning that map images to texts.

A natural question is then why not directly apply existing watermarking methods [1, 8] designed for the classification DNNs to watermark the DNNs in image captioning. Unfortunately, it is not the case for the white-box watermarking methods. The obstacles lie in several fundamental differences between these two kinds of DNNs. First, DNNs for classification output a label. In contrast, DNNs

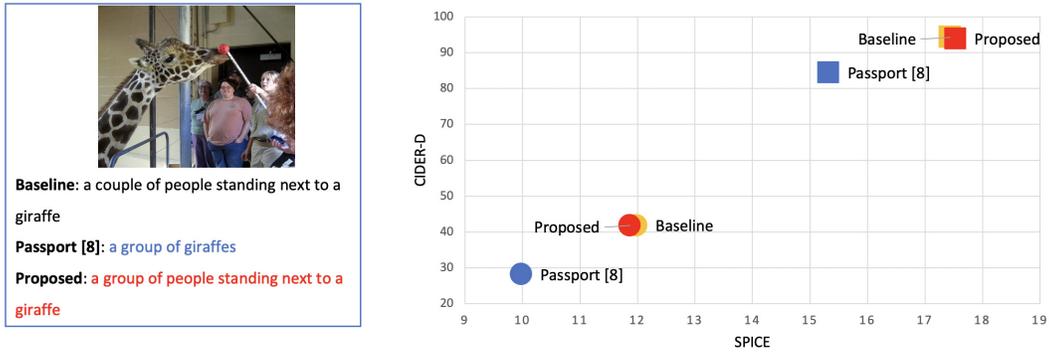


Figure 1: Comparison of the baseline model (in yellow), passport model [8] (in blue), and proposed model (in red) on two public datasets. Best view in colour.

for image captioning output a sentence. Second, classification is about finding the decision boundaries among different classes, whereas image captioning is not only to understand the image content in depth beyond category or attribute levels but also to connect its interpretation with a language model to create a natural sentence [9]. We demonstrated in Figure 1, Figure 4, and Sect. 4.2 that a recent digital watermarking framework [8] that used to protect deep-based classification model is insufficient to protect image captioning models. Figure 1 shows that [8] is insufficient to be deployed to protect the IP of the image captioning task as compared to our proposed model, against the baseline model. For instance, it can be noticed that the caption generated by [8] is incomplete and incorrect, while the caption generated from our model is very near to the baseline model. In addition, the CIDEr-D and SPICE score of our model is very near to the baseline as compared to [8] on MS-COCO and Flickr30k datasets.

As a remedy, this paper proposes a novel embedding framework that consists of two different embedding schemes to embed a unique secret key into the recurrent neural network (RNN) cell [10] to protect the image captioning model against various attacks. Specifically, we show that embed a secret key into the hidden memory state of an RNN is the best choice for the image captioning task such that a forged key will immediately yield an unusable image captioning model in terms of poor quality outputs, defeating the purpose of infringement.

On the one hand, our solution bears a similarity to digital watermarking - they both embed certain digital entities into models during training sessions. In terms of IP protection, however, embedded watermarks only enable the verification of the ownership of models. One has to rely on government investigation and enforcement actions to discourage IP infringement. Whether this kind of approach can provide *reliable*, *timely*, and *cost-effective* juridical protection remains questionable. On the other hand, our key-protected models will not function normally unless the valid key is provided, thus immediately preventing the unlawful usages of the models with no extra costs. Indeed, we regard this proactive protection as the most prominent advantage of our solution over digital watermarking. For instance, in Fig. 6, the protected model with a valid key demonstrated almost identical performance as that of the original model, in contrast the same model presented with a forged key has a huge performance drop in all metric scores.

The contributions are threefold: i) We renovate the paradigm of digital watermarking based neural network IP protection, by proposing a key-based strategy that provides *reliable*, *preventive* and *timely* IP protection (Sect. 3.1) at virtually *no extra cost* (Sect. 4.6) for image captioning task. ii) This paper formulates the problem and proposes a generic solution as well as concrete implementation schemes that embed a unique key into RNN models through the hidden memory state (Section 3; Fig. 2b). We prove that a forged key will yield a useless image model. Also, we empirically show the effectiveness of our approach against various attacks and prove the ownership of the model (Sect. 4); and iii) To the best of our knowledge, we are the first to propose IP protection on image captioning model and we demonstrated that the proposed method does not compromise the original image captioning performance on all common captioning metrics on Flickr30k and MS-COCO datasets (Table 1).

2 Related Work

Conventionally, digital watermarks were extensively used in protecting the ownership of multimedia contents, including images, videos, audio, or functional designs. It is a process of embedding a marker into the content and subsequently using it to verify the ownership. In deep learning, the IP protection on the models can be categorized into i) white-box based solution [1, 2], ii) black-box based solution [3–6, 11] or iii) a combinatorial of both white and black based solution [7, 8].

The first work that introduced digital watermarks for DNN was proposed by Uchida et al. [1], where the authors embedded a watermark into the weights parameters via parameter regularizer during the training as white-box protection. For verification, owners are required to access the model parameters to extract the watermark. To remedy this issue, [3–6, 11] proposed digital watermarks in a black-box setting. In this setting, a set of trigger set images is generated as random image and label pairs. During training, the feature distributions of those images are distant from the labeled training samples. During verification, the trigger set watermark can be extracted remotely without the need to access the model weights. For example, Zhang et al. [3] introduced three different key generations which are content-based, noise-based, and unrelated-based images respectively. Adi et al. [4] proposed a watermarking method similar to [3] but their main contribution is the model verification. While Merrer et al. [6] proposed to use adversarial examples as the watermark key set to modify the model decision boundary. Quan et al. [11] aimed to develop a black-box watermarking method in images to images tasks such as image denoising and super-resolution by exploiting the overparameterization of the model.

Recently, [7, 8] presented a watermarking framework that works in both white-box and black-box settings. Rouhani et al. [7] embedded watermark in activation of selected layers of the DNN by integrating two additional regularization loss terms, binary cross-entropy loss, and Gaussian Mixed Model (GMM) agent loss. It is robust against pruning, fine-tuning, and overwriting attack but require more computation. The work that most closer to us is Fan et al. [8] added special “passport” layers into the DNN model to enable ownership verification. With a forged passport, the performance of the model will significantly deteriorate. This design relies on the secrecy of passport layer weights that requires the owner to keep the passport layer weights secret from the attacker. However, empirically, we demonstrated that [8] does not able to protect the image captioning model effectively.

3 Approach

Our image captioning framework of interest is a simplified variant of the *Show, Attend and Tell* model [9]. It is a popular framework that forms the basis for subsequent state-of-the-art works on image captioning [12–20]. It follows the encoder-decoder framework, where a convolutional neural network (CNN) is used to encode an image into a fixed-size representation, and the long short-term memory (LSTM) is employed to generate the captions.

Given an image I , it is encoded into a fixed-size feature vector using CNN as followed:

$$X = f_c(I) \tag{1}$$

where $f_c(\cdot)$ represents CNN encoder in our models, X is the image feature vector.

In the decoder, the image feature vector X is fed into LSTM at each time step t to output the probability of next word S_t as:

$$h_t = LSTM(X, h_{t-1}, m_{t-1}) \tag{2}$$

$$p(S_t | S_0, \dots, S_{t-1}, I) = F_1(h_t) \tag{3}$$

where h_{t-1} is the previous LSTM’s hidden state, m_{t-1} is the previous memory cell, $F_1(\cdot)$ is a nonlinear function that outputs the probability of S_t , p is the probability of next word S_t with image I and previous words S_0, \dots, S_{t-1} .

Unless otherwise stated, our models are trained under the maximum likelihood estimation (MLE) framework, where the probability of generating a correct caption of length T with tokens $\{S_0, \dots, S_{T-1}\}$ for an image I is directly maximized:

$$\log p(S | I) = \sum_{t=0}^T \log p(S_t | I, S_{0:t-1}, c_t) \tag{4}$$

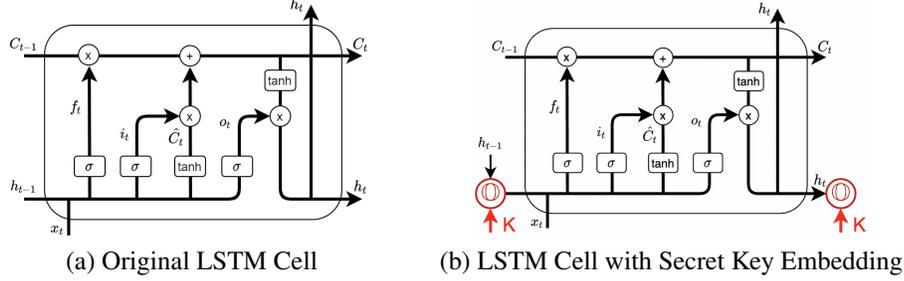


Figure 2: An overview of our approach. (a) The original LSTM Cell and (b) LSTM Cell with key embedding operation $\odot \in \{M_{\oplus}, M_{\otimes}\}$ (see Section 3.2)

where t is the time step, $p(S_t | I, S_{0:t-1}, c_t)$ is the probability of generating a word given an image I , previous words $S_{0:t-1}$, and context vector c_t .

3.1 Problem formulation

Let \mathcal{N} denote an image captioning model to be protected by a *secret* key k , after a training process, the image captioning model embedded with the key is denoted by $\mathcal{N}[K]$ as shown in Figure 2. The inference of such a protected model can be characterized as a process M that *modifies* the model behavior according to the running-time key l :

$$M(\mathcal{N}[K], l) = \begin{cases} \mathcal{M}_K, & \text{if } l = K, \\ \mathcal{M}_{\bar{K}}, & \text{otherwise,} \end{cases} \quad (5)$$

in which \mathcal{M}_K is the network performance with key correctly verified, and $\mathcal{M}_{\bar{K}}$ is the performance with the incorrect key i.e., $\bar{K} \neq K$.

The properties of $M(\mathcal{N}[K], l)$ defined below are desired for the sake of IP protection:

Definition 1. If $l = K$, the performance \mathcal{M}_K should be as close as possible to that of the original network \mathcal{N} . Specifically, if the performance *inconsistency* between \mathcal{M}_K and that of \mathcal{N} is smaller than a desired threshold, then the protected network is called *functionality-preserving*.

Definition 2. If $l \neq K$, on the other hand, the performance $\mathcal{M}_{\bar{K}}$ should be as far as possible to that of \mathcal{M}_K . The discrepancy between \mathcal{M}_K and $\mathcal{M}_{\bar{K}}$ therefore can be defined as the *protection-strength*.

3.2 Embedding operation

Figure 3a shows the overview of the embedding process. Our *embedding* process can be represented as $E_{\odot}(\mathbf{D}, \mathbf{g}, \mathcal{N}[\cdot], L) = \mathcal{N}[\mathbf{W}, \mathbf{g}]$, is a RNN learning process. It takes inputs *training data* $\mathbf{D} = \{I, S\}$, and optionally signature \mathbf{g} , and optimizes the model $\mathcal{N}[\mathbf{W}, \mathbf{g}]$ by minimizing the given loss L . In this paper, we introduce two different key embedding operations \odot which are i) element-wise addition model (M_{\oplus}) or ii) element-wise multiplication model (M_{\otimes}):

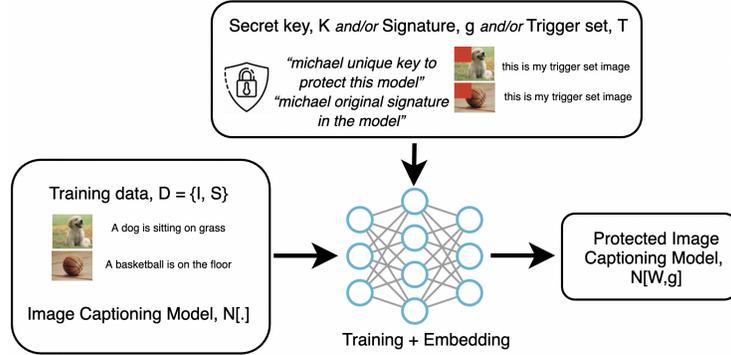
$$\odot(K, h_{t-1}, e) = \begin{cases} K \oplus h_{t-1}, & \text{if } e = \oplus, \\ K \otimes h_{t-1}, & \text{else.} \end{cases} \quad (6)$$

where $k_f = \{k_{f,i}\}_{i=1}^N$ with N is the size of the hidden state, $k_{f,i} \in \mathbb{R} : -1 \leq k_{f,i} \leq 1$ and $k_b = \{k_{b,i}\}_{i=1}^N$ with $k_{b,i} \in \{-1, 1\}$.

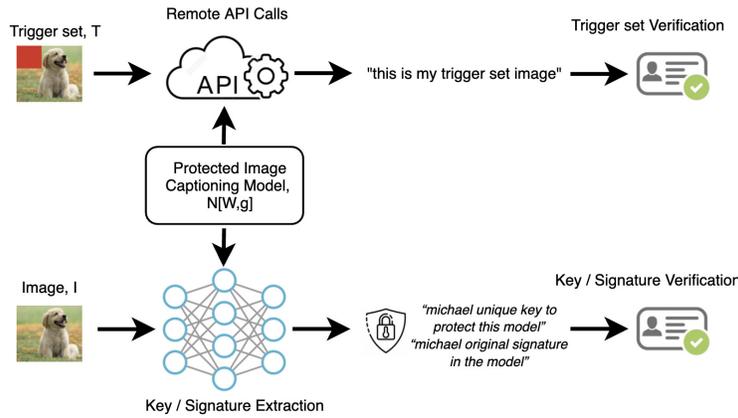
Then, the embedded key K is represented in terms k_b is generated by converting the string provided by owner to a binary vector BE . However, we found that the binary vector for very near alphanumeric, e.g., string A and C, has only a 1-bit difference. Therefore, we proposed a new transformation function \mathbb{T} :

$$\mathbb{T}(BE, BC) = BE \otimes BC = k_b \quad (7)$$

where BC is a binary vector sampled from value of -1 or 1 according to the seed provided by user, to alleviate this issue.



(a) An embedding process E_O , takes inputs training data $D = \{I, S\}$, secret key K and/or signature g and/or trigger set T , model $N[.]$ to produce protected model $N[W, g]$.



(b) A verification process V takes as inputs, either an image I or a trigger set T , and outputs the result to verify the ownership.

Figure 3: Visual explanation for embedding and verification processes.

3.3 Sign of hidden state as signature

In order to further strengthen our model, we follow [8] to add the sign loss regularization term into the loss function as:

$$L_g(h, G, \gamma) = \sum_{i=1}^N \max(\gamma - h_i g_i, 0) \quad (8)$$

where $G = \{g_i\}_{i=1}^N$ with $g_i \in \{-1, 1\}$ consists of the designated binary bits for hidden state h . To enforce the hidden state to have a magnitude greater than 0, a hyperparameter γ is introduced into the sign loss. However, one of the main differences of our approach compared to [8] is our signature is not embedded in the model weights, but it is embedded in the hidden state which is the output of the LSTM cell. This is because we found out that embedding signature in the model weights can be easily attacked with a channel permutation, i.e., change the signature but remains the output of the model.

3.4 Ownership verification

Figure 3b shows the overview of the verification process. Suppose an owner tries to verify the ownership of a target model, three verification methods are proposed: 1) V_1 : Key-based verification; 2) V_2 : Signature-based verification; and 3) V_3 : Trigger set verification.

V_1 : **Secret key-based verification** - In this verification scheme, there are two different approaches, depending on the secret key is either public or private. Formerly, the trained model and the public key will be provided to the clients. For model inferences, the public key will be required as part of the input to the model to ensure the model performance is preserved. The ownership of the model

can be verified directly by the provided key. Latter, a private key is directly embedded into the model. For inference, only image is required as the model input. However, for ownership verification, one has to have access to the model and extract the key from the LSTM cell.

V_2 : Signature verification - In this verification scheme, a unique signature is embedded in the sign of the hidden state during the training process via sign loss regularization. To verify the signature, the owner is required to access the trained model. Then, an image will be sent to the model to generate the caption. During inference time, the sign of the hidden state of the LSTM cell will be extracted and compared with our signature to verify the ownership. This binary bits signature can be transformed back to a human-readable string for example the name of the owner.

V_3 : Trigger set verification - V_1 and V_2 are considered as white-box verification, where the owner is required to have access to the model physically in order to verify the ownership. Hence, we introduce trigger set verification that can be conducted remotely via API calls. First, a set of trigger set image-caption pairs¹ are generated, and then they are used together with the original training samples to train an image captioning model. In this paper, the trigger set images are generated by adding noise (e.g., red color patch) to an original image so that the model is trained to generate the trigger set caption (e.g., this is my trigger set image). For verification, the owner will send the trigger set images to the model and test whether the model returns the trigger set caption.

4 Experiments

This section presents the experiment results of our approaches in terms of resilience to ambiguity and robustness to removal attacks. Qualitative analysis is also carried out to compare the quality of the caption generated by different approaches. We compare the following models: (i) **Baseline** is implemented based on the soft attention model as to [9], it is an unprotected model. (ii) **Passport** [8] is the work that most closer to us that added "passport" layers into the DNN model to enable ownership verification. (iii) M_{\oplus} is our element-wise addition model presented in Section 3.2. (iv) M_{\otimes} is our element-wise multiplication model presented in Section 3.2.

We used ResNet-50 [21] pre-trained on the ImageNet dataset as the encoder. The image features are extracted using ResNet-50 without fully connected layers, resulting in $7 \times 7 \times 2048$ dimensional outputs. For the decoder part, we use LSTM with a dropout rate of 30%. Both the word embedding and hidden state are set to 512. We set the attention loss factor to 0.01. The LSTM decoder is trained using a learning rate of 1e-4 for 8 epochs and finetune the CNN with a learning rate of 1e-5 up to 20 epochs. The model is trained by cross-entropy loss with a mini-batch size of 32 using Adam [22] optimizer with β_1 set to 0.9, β_2 set to 0.999, and ϵ set to 1e-6. We apply gradient clipping to prevent any gradient to have norm greater than 5.0 to prevent exploding gradients. We repeated all experiments 3 times to get the average performance. The beam size is set as 3 in the inference stage.

4.1 Dataset and metrics

We train and evaluate our approaches on the MS-COCO [23] and Flickr30k [24] datasets, which are widely used for the image captioning task. We followed the widely used split in [15] for both datasets. MS-COCO contains 113,287 training images with 5 human-annotated captions for each image. The validation and test sets contain 5,000 images each. Flickr30k contains 1,000 images for validation, 1,000 for test, and the rest for training. We truncated captions longer than 20 words and converted all the words into lower case. Fixed vocabulary size of 10,000 is used for both datasets.

We evaluate our approaches using all common metrics in the image captioning task: CIDEr-D [25], SPICE [26], BLEU [27], METEOR [28], and ROUGE-L [29]. However, CIDEr-D and SPICE have been shown to have a higher correlation with human judgments compared to BLEU and ROUGE [25, 26], but it is common practice to report all the aforementioned metrics in the image captioning task.

4.2 Comparison with CNN-based watermarking framework

For comparison with the existing digital watermarking framework, we re-implement [8] using the official repository and refer this model as Passport. We choose [8] because the work is somehow

¹It is actually a list of data wrongly labeled by purpose

Methods	MS-COCO								Flickr30k							
	B-1	B-2	B-3	B-4	M	R	C	S	B-1	B-2	B-3	B-4	M	R	C	S
Baseline	72.14	55.70	41.86	31.14	24.18	52.92	94.30	17.44	63.40	45.18	31.68	21.90	18.04	44.30	41.80	11.98
Passport [8]	68.50	53.30	38.41	29.12	21.03	48.80	84.45	15.32	48.30	38.23	26.21	17.88	15.02	32.25	28.22	9.98
M_{\oplus}	72.53	56.07	42.03	30.97	24.00	52.90	*91.40	*17.13	62.43	44.40	30.90	21.13	*17.53	43.63	*40.07	*11.57
M_{\otimes}	*72.47	*56.03	*41.97	*30.90	*23.97	52.90	91.60	17.17	*62.30	*44.07	*30.73	*21.10	17.63	*43.53	40.17	11.67

Table 1: Comparison between our approaches (M_{\oplus}, M_{\otimes}) with baseline and Passport [8] on MS-COCO and Flickr30k datasets, across 5 common metrics where B-N, M, R, C, and S are BLEU-N, METEOR, ROUGE-L, CIDEr-D, and SPICE scores. **BOLD** is the best result and * is the second best result.

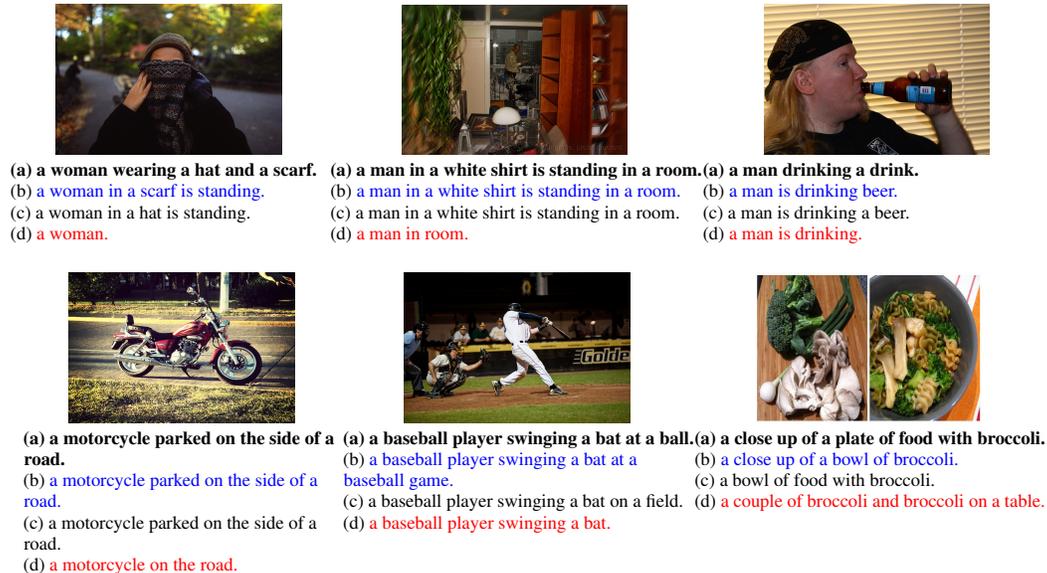


Figure 4: Comparison of captions generated by (a) Baseline, (b) M_{\oplus} , (c) M_{\otimes} , and (d) Passport [8]. The first row is the images from Flickr30k. The second row is the images from the MS-COCO dataset. It is noticed that the quality of the captions generated by our models is very close to the baseline.

similar to ours in terms of technical implementation. In Table 1, we can notice that the overall performance of the Passport model on both MS-COCO and Flickr30k is very poor compared with the baseline and our proposed methods. For example, the CIDEr-D score dropped 10.45% (MS-COCO) and 32.49% (Flickr30K), respectively when compared to the baseline. In contrast, both of our proposed methods only dropped 3-4% on both MS-COCO and Flickr30K when compared to the baseline.

In terms of qualitative comparison, Figure 4 shows the sample captions generated by the Passport model are relatively brief when compared to the baseline and our approaches. For instance, the first image in Figure 4, our proposed model generated *a woman in a scarf is standing*, it matches with the ground truth provided by the baseline, but the Passport model only generated *a woman*, missing the rest of the rich context altogether. a similar observation is found for the rest of the images.

Furthermore, we conducted an experiment to attack the Passport model with a forged passport and found out the Passport model still can have a relatively high CIDEr-D score. In Table 2, we show the quantitative results of the Passport model with the correct passport vs. Passport model with a forged passport. We found out that the Passport model with forged passport still can achieve very similar results on MS-COCO and Flickr30k dataset as to the Passport model with correct passport. For example, it has a CIDEr-D score of 83.0 (84.45) and 26.5 (28.22) on both datasets. In terms of qualitative comparison, Figure 5 shows the captions generated by the Passport model with (a) correct and (b) forged passport. It can be noticed that both models generate almost similar captions in



Figure 5: Comparison of captions generated by (a) Passport [8] with correct passport, (b) Passport [8] with forged passport. The first row is the images from Flickr30k. The second row is the images from the MS-COCO dataset.

Methods	MS-COCO								Flickr30k							
	B-1	B-2	B-3	B-4	M	R	C	S	B-1	B-2	B-3	B-4	M	R	C	S
Passport	68.50	53.30	38.41	29.12	21.03	48.80	84.45	15.32	48.30	38.23	26.21	17.88	15.02	32.25	28.22	9.98
<i>Passport</i> (forged)	67.50	52.65	37.15	29.01	20.95	47.90	83.00	15.00	47.30	37.87	26.01	17.10	14.82	31.88	26.50	9.90

Table 2: Comparison between Passport [8] with (top) correct passport and (bottom) forged passport on MS-COCO and Flickr30k datasets, across 5 common metrics where B-N, M, R, C, and S are BLEU-N, METEOR, ROUGE-L, CIDEr-D, and SPICE scores. It can be clearly seen that Passport model [8] is clearly insufficient to protect the image captioning model as the performance with correct and incorrect passport across all 5 metrics are almost similar.

terms of word selection and caption length. As a conclusion, we deduce that the conventional digital watermarking framework is insufficient to protect the image captioning model.

4.3 Fidelity Evaluation

Fidelity is defined as matching the performance of the original model. In this section, we show that our proposed embedding schemes do not degrade the overall model performance in terms of metrics, as well as the quality of the generated sentences. According to Table 1, it shows that the overall performance of our approaches and baseline model on MS-COCO and Flickr30k in all 5 image captioning metrics. Specifically, we can observe that M_{\oplus} performed the best as it out-performed baseline in BLEU1-3 score on MS-COCO dataset, and BLEU-1 in Flickr30K dataset, respectively. For the rest of the metric score, we can also observe that M_{\oplus} came as 2nd best score. In contrast, [8] performed poorly with at least a 10% drop in all metrics.

Subsequently, Table 4 shows the comparison of the uniqueness of generated caption from our approaches and baseline model. A caption is considered unique if the generated caption does not exist in the training dataset. On both datasets, it shows that our approaches have very similar uniqueness and average caption length compared to baseline. This is consistent with the caption generated shown in Figure 4. For example, both of our models have an exact caption generated as to baseline on the first image. And subsequently, in the rest of the images, the choice of words generated (i.e., shirt, room, road, swinging) are also very consistent with the baseline.

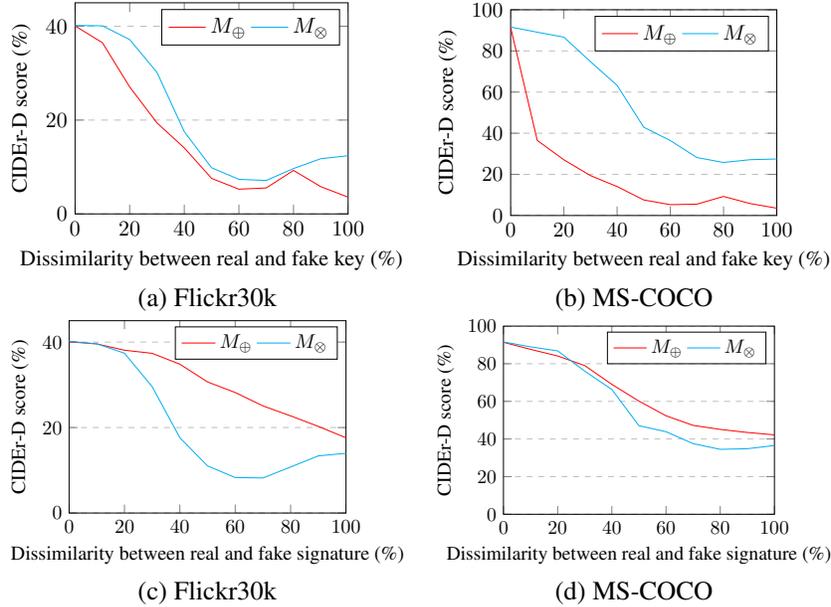


Figure 6: CIDEr-D on Flickr30k and MS-COCO under ambiguity attack on (a-b) key; (c-d) signature.

Methods	MS-COCO		Flickr30k	
	MS-COCO	Flickr30k	Flickr30k	MS-COCO
Baseline	- (94.30)	- (37.70)	- (41.80)	- (88.50)
M_{\oplus}	100 (91.40)	70.40 (37.50)	100 (40.07)	72.50 (87.30)
M_{\otimes}	99.99 (91.60)	71.50 (37.8)	99.99 (40.17)	71.35 (86.50)

Table 3: Fine-tuning attack: CIDEr-D (in-bracket) of baseline and proposed models (Left: MS-COCO fine-tune on Flickr30k. Right: vice-versa.) Accuracy (%) outside bracket is the signature detection rate.

4.4 Resilience against ambiguity attacks

4.4.1 Protection against forged key

In this case, we assume the attacker somehow has the access to the model but does not have the correct secret key and so tries to attack the model with a random forged key. Figure 6(a-b) show the CIDEr-D score of proposed models under ambiguity attack on the secret key in Flickr30k and MS-COCO. Accordingly, we can observe that in general, the model performance will drop when a forged key is deployed. In particular, we would like to highlight that the CIDEr-D score on MS-COCO drops significantly (almost 50% difference) in M_{\otimes} even a forged key that has a 75% similarity to the real key is being deployed. This shows that our proposed method is resilient against this forged key attack.

Figure 7 shows another six sample images and the respective captions generated by our proposal and baseline. From the first image, it shows that given the correct secret key, our proposed method is able to generate a caption that consists of object, scene, and attributes that are very similar to baseline. When a forged key is used, in this example, we show in Figure 7(c) a forged key that has a 75% similarity to the correct secret key and in Figure 7(d) another forged key that has a 50% matching to the correct secret key, the generated caption is either not meaningful at all with repetitive words (i.e., a man and a man) or a very brief caption (i.e., two dogs). According to Table 4, we can also observe similar patterns. For instance, \widehat{M}_{\oplus} has almost 100% uniqueness with the longest average caption length on both datasets. From Figure 7, we can understand that this is due to repetitive words. Meanwhile, \widehat{M}_{\otimes} has 88.44% uniqueness with the shortest average caption length on the Flickr30K dataset. Yet again, this phenomenon is observed from the generated caption in Figure 7.

Methods	MS-COCO		Flickr30k	
	Unique	Avg-L	Unique	Avg-L
Baseline	62.93%	8.86	88.80%	9.50
M_{\oplus}	70.96%	8.81	88.00%	9.30
M_{\otimes}	70.26%	8.91	87.10%	9.28
\widehat{M}_{\oplus}	100.00%	19.97	97.40%	18.56
\widehat{M}_{\otimes}	88.44%	12.71	53.40%	7.69

Table 4: Comparison of the uniqueness of caption generated by our approaches and baseline model. \widehat{M}_{\oplus} and \widehat{M}_{\otimes} are M_{\oplus} and M_{\otimes} , respectively but with forged secret key. Avg-L stands for average length.



Figure 7: Comparison of captions generated from (a) Baseline, (b) M_{\otimes} , (c) M_{\otimes} with the forged key has 75% similarity as to real key and (d) M_{\otimes} with the forged key has 50% similarity as to real key. The first row is the images from Flickr30k. The second row is the images from the MS-COCO dataset.

4.4.2 Protection against fake signature

In this case, we assume the secret key is exposed to the attacker and one can use the model with original performance. However, the signature is able to use as proof of ownership. As such, the attacker will try to attack the signature by attempting to change the sign of the signature. Figure 6(c-d) show the overall performance of our proposed models (CIDEr-D score) will decrease when the signature is being compromised on both Flickr30k and MS-COCO. For instance, even very small changes (only 10% of the sign are toggle), we can observe at least a 10-15% drop of performance in terms of CIDEr-D score; and when half of the sign are toggle, it is seen that the model performance is almost useless. In Figure 8, when 10% of the sign are modified, it can be seen that the captions generated by the model are relatively brief and shorter than the original model. For example in the second image, it shows that the generated caption is without “a rock” when 10% of the sign are modified. When 50% of the sign are toggle, the captions are repetitive words.

4.5 Robustness against removal attacks

4.5.1 Model pruning

Generally, model pruning is used to reduce the weights and computation overhead of a DNN model. However, the attacker might leverage it to remove the signature in the model. In order to test our approaches is robust to this attack, we implemented class-blind pruning method [30]. Figure 9 shows

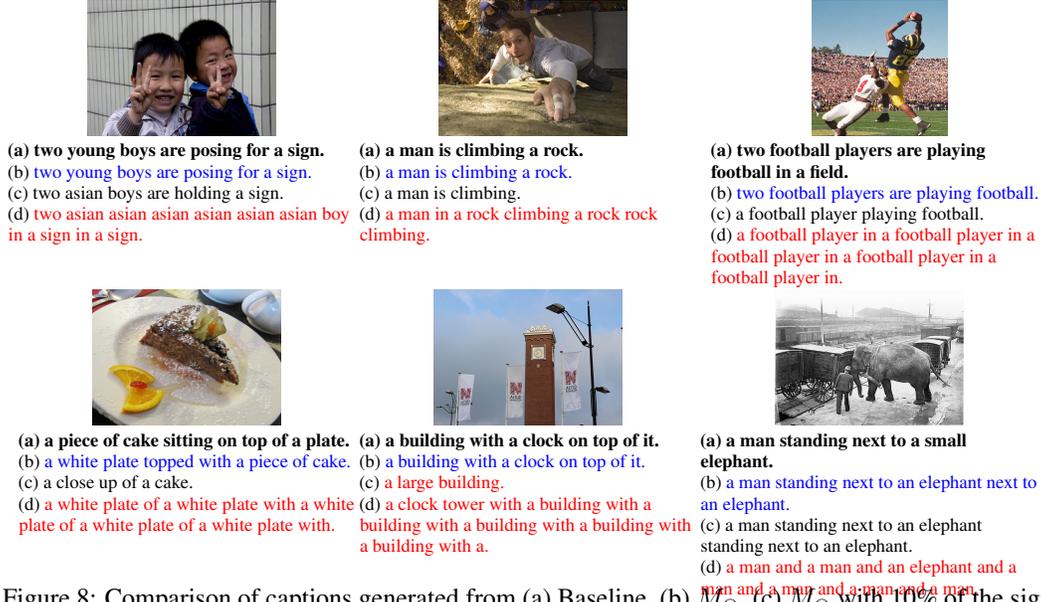


Figure 8: Comparison of captions generated from (a) Baseline, (b) M_{\otimes} , (c) M_{\otimes} with 10% of the sign are toggle and (d) M_{\otimes} with 50% of the sign are toggle. The first row is the images from Flickr30k. The second row is the images from the MS-COCO dataset.

the CIDEr-D score and signature detection rate on M_{\oplus} and M_{\otimes} against different pruning rates. We show that even 60% of the network parameters are pruned, the signature detection rate is still intact at more than 84% and 91% on both Flickr30k and MS-COCO datasets, respectively. Figure 10 shows that the pruned model can still generate a meaningful caption for the images from both Flickr30k and MS-COCO datasets. As a summary, we show that even 60% of the network parameters are pruned, the signature detection rate and the quality of the caption are still intact.

4.5.2 Fine-tuning

Here, we simulate an attacker who fine-tunes the stolen model with a new dataset to obtain a new model that inherits the performance of the stolen model while attempting to remove the embedded signature. Table 3 shows the signature detection rate and CIDEr-D score of the proposed model after perform fine-tuning. The signature can be detected at almost 100% accuracy for our approaches in the original task. After fine-tuning the model (e.g., from MS-COCO to Flickr30k or Flickr30k to MS-COCO), we show that our approaches achieve comparable CIDEr-D score as to the baseline, however, we observe the signature detection rate decreased to around 70%. This is one of the limitations of the proposed method but overall it does not compromise the IP protection of the model as we still have the secret key to act as proof of ownership. Therefore, the proposed secret key working together with the signature in this paper can act as complete protection for ownership verification.

4.5.3 Key pruning

If an attacker knows a key is in place, instead of pruning the model, we simulate an attacker to prune the key rather than model weights. Figure 11 shows the CIDEr-D score and signature detection rate on M_{\oplus} and M_{\otimes} against different key pruning rates. We show that even the attacker prunes 100% on the key, the signature detection rate still remains at more than 98% and 95% on both MS-COCO and Flickr30k datasets. Since key pruning has changed the original key, it degrades the performance of the proposed model as well. As shown in Figure 11, the CIDEr-d score continues to drop when the key pruning rate increases. Hence, the proposed secret key and signature can protect the ownership of the model against the key pruning attack.

4.5.4 Fine-tuning key and signature

This experiment is different from Section 4.5.2, here, we simulate an attacker who knows everything about the model, i.e., the training procedure, the training parameters, the dataset used, and the

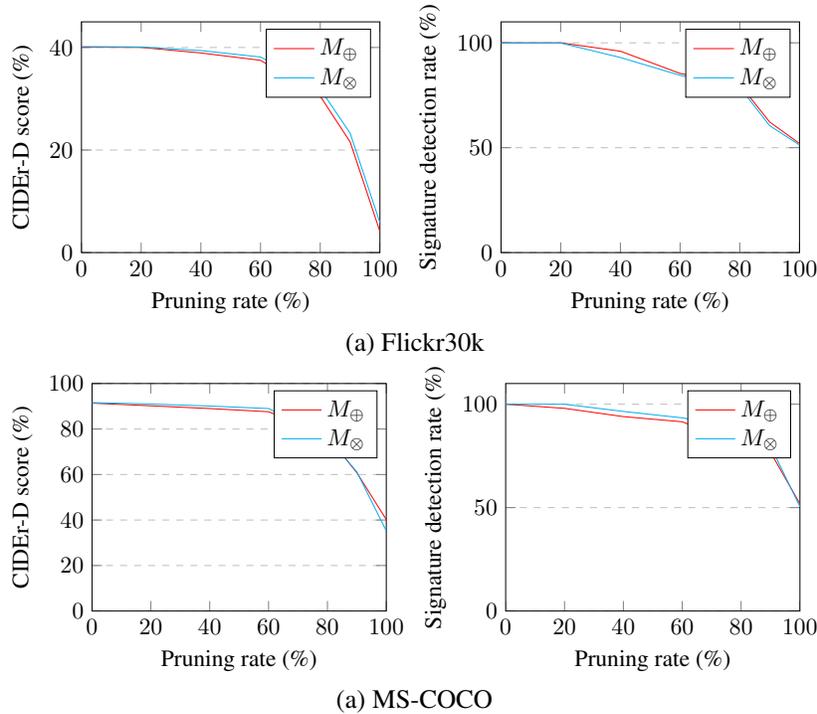


Figure 9: Removal attack (Pruning): CIDEr-D score and signature detection rate of our approaches on both MS-COCO and Flickr30k against different pruning rates.

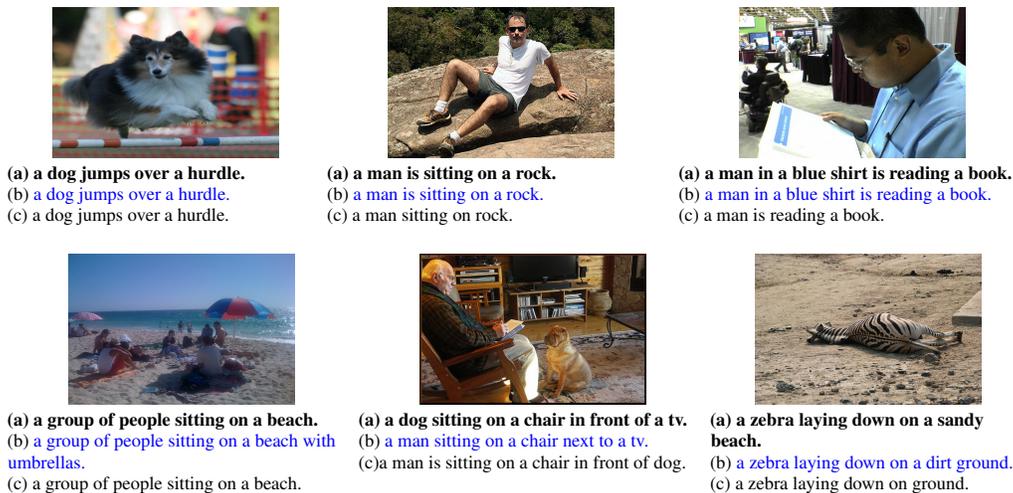


Figure 10: Comparison of captions generated from (a) Baseline, (b) M_{\otimes} , (c) M_{\oplus} with 60% pruning rate. The first row is the images from Flickr30k. The second row is the images from the MS-COCO dataset.

key/signature. The attacker fine-tunes the model with a different key/signature following the same training steps. Table 5 shows the CIDEr-D score and the signature detection rate of the proposed model after perform key and signature fine-tuning. After fine-tuning the model to a new key and signature, we show that our approaches achieve a slightly lower CIDEr-D score (-1% to -5%) as compared to the protected model. However, the signature detection rate dropped from almost 100% to around 68%. This is the worst-case scenario where the model is difficult to protect against the attacker who knows everything about the model including the training steps.

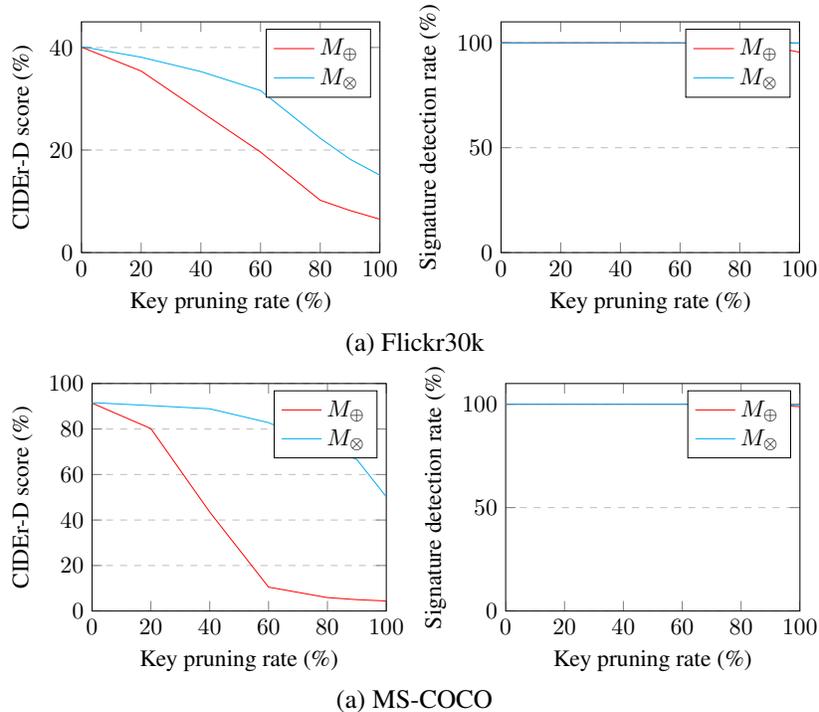


Figure 11: Removal attack (Key Pruning): CIDer-D score and signature detection rate of our approaches on both MS-COCO and Flickr30k against different key pruning rates.

Methods	MS-COCO		Flickr30k	
	Protect	Attack	Protect	Attack
M_{\oplus}	100 (91.40)	68.08 (89.60)	100 (40.07)	69.14 (39.70)
M_{\otimes}	99.99 (91.60)	68.16 (89.6)	99.99 (40.17)	67.96 (38.1)

Table 5: Fine-tuning key and signature attack: CIDer-D (in-bracket) of proposed models (Left: MS-COCO. Right: Flickr30k.) Accuracy (%) outside bracket is the signature detection rate.

4.6 Network Complexity

Methods	MSCOCO		Flickr30k	
	Training time	Inference time	Training time	Inference time
Baseline	20h 10m	10.60s	4h 25m	10.64s
M_{\oplus}	20h 31m	10.58s	4h 30m	10.60s
M_{\otimes}	20h 31m	10.45s	4h 30m	10.44s

Table 6: Comparison of network complexity of our approaches and baseline model. Training time is calculated over the entire training set for 20 epochs. Inference time is calculated over a single iteration, where h, m, s are hour, minute, and second.

Our approach with the key and signature embedding in the image captioning model does not cause the extra cost. We conducted an experiment to compare the training and inference time between baseline, M_{\oplus} and M_{\otimes} . All the experiments used TITAN V GPU with the same setting and hyperparameter as stated in Section 4. On both Flickr30k and MS-COCO datasets, the complexity of our approach (with the key and signature embedded) compared to the baseline model is almost negligible. Herein, we show the complete results in Table 6 to compare the training time and inference time. It is noticed that our approaches only have an incremental of 1.89% or below in the training time on Flickr30k and MS-COCO datasets. For the inference time, it is almost similar to the baseline model for a single

iteration. As a summary, we show that our proposed model provides reliable, preventive, and timely IP protection at virtually no extra cost for image captioning tasks.

4.7 Comparison with different types of image captioning frameworks

Despite showing our approach in the popular image captioning framework *Show, Attend and Tell* model [9]. We further applied our approach in the different types of image captioning frameworks, which are Up-Down [12] and SCST [14]. Up-Down model used the bottom-up attention techniques to find the most relevant regions based on bounding boxes and two LSTM layers are used to selectively attend to the image features to generate the caption. While SCST model applied the reinforcement learning method to image captioning by optimizing the model directly on those objective evaluation metrics like CIDEr score. Both Up-Down and SCST are having different frameworks compared to our baseline. Therefore, we need to adapt our proposed key embedding approach to the frameworks.

Up-Down model consists of two LSTM layers, the first LSTM layer as a top-down visual attention model and the second LSTM layer as a language model. It has an attention module in between the first and second LSTM layers. We apply our proposed key embedding approach in the second LSTM layer and the attention module. The key embedding process is defined as:

$$h_{t-1}^2 = C \cdot (h_{t-1}^2 \oplus K) + (1 - C) \cdot (h_{t-1}^2) \quad (9)$$

$$\hat{v}_t = C \cdot (\hat{v}_t) + (1 - C) \cdot (\hat{v}_t \oplus K) \quad (10)$$

where h_{t-1}^2 is the previous hidden state of the second LSTM layer, \hat{v}_t is the attended image features and $C \in \{0, 1\}$, 0 indicates forged key, 1 indicates real key.

SCST model modified the architecture of the attention model for captioning given in [9], and input the attention-derived image feature only to the cell node of the LSTM. We apply our proposed key embedding approach in the attention-derived image feature and the hidden state. The key embedding process is defined as:

$$h_{t-1} = C \cdot (h_{t-1} \oplus K) + (1 - C) \cdot (h_{t-1}) \quad (11)$$

$$I_t = C \cdot (I_t) + (1 - C) \cdot (I_t \oplus K) \quad (12)$$

where h_{t-1} is the previous hidden state of the LSTM layer, I_t is the attention-derived image feature.

Table 7 shows the performance of our approach with the Up-Down model on the MS-COCO dataset in all 5 image captioning metrics. UD- M_{\oplus} refers to our proposed key embedding approach in the Up-Down model. While UD- \widehat{M}_{\oplus} is similar to UD- M_{\oplus} but with the forged secret key. We can observe that UD- M_{\oplus} is having a lower score compared to the original Up-Down model. With the forged secret key, the performance of the model degrades significantly which showing our approach can protect the Up-Down captioning model. For example, the CIDEr-D score drops 17.27% from 101.93 to 84.33.

According to Table 8, it shows the performance of our approach with the SCST model on the MS-COCO dataset. We follow the experiment in the original paper [14] to evaluate the model in 4 image captioning metrics: BLEU-4, METEOR, ROUGE-L, and CIDEr-D. SCST- M_{\oplus} refers to our proposed key embedding approach in the SCST model. While SCST- \widehat{M}_{\oplus} is similar to SCST- M_{\oplus} but with the forged secret key. Our proposed key embedding approach SCST- M_{\oplus} does perform poorly compared to the original SCST model but is able to protect the model against the forged secret key. With the forged secret key, the CIDEr-D score drops from 101.87 to 90.53.

5 Limitations

We show that our proposed secret key together with signature can protect the model against unauthorized usage as simulated in Section 4. However, the embedding-based methods have some limitations that are inevitable. For example, the secret key can be removed and the signature detection rate is dropped in the worst-case scenario in which the attacker knows everything about the model as shown in Section 4.5.4. This hinders open-sourcing the model as we need to prevent others from knowing the secret key, the training procedure, and the exact training parameters.

Methods	Evaluation Metric					
	B-1	B-4	M	R	C	S
Up-Down [12]	76.97	36.03	26.67	56.03	111.13	19.90
UD- M_{\oplus}	71.57	33.83	25.33	52.43	101.93	18.60
UD- \widehat{M}_{\oplus}	65.20	29.50	20.33	48.60	84.33	16.53

Table 7: Comparison between our approach (UD- M_{\oplus}) with Up-Down [12] model on MS-COCO dataset, across 5 common metrics where B-N, M, R, C and S are BLEU-N, METEOR, ROUGE-L, CIDEr-D and SPICE scores. UD- \widehat{M}_{\oplus} is similar to UD- M_{\oplus} but with the forged secret key.

Methods	Evaluation Metric			
	B-4	M	R	C
SCST [14]	33.87	26.27	55.23	111.33
SCST- M_{\oplus}	31.60	24.97	52.43	101.87
SCST- \widehat{M}_{\oplus}	29.83	22.63	50.17	90.53

Table 8: Comparison between our approach (SCST- M_{\oplus}) with SCST [14] model on MS-COCO dataset. SCST- \widehat{M}_{\oplus} is similar to SCST- M_{\oplus} but with the forged secret key.

6 Conclusion

IP protection on DNN has been a significant research area and we take the first step to implement the ownership protection on the image captioning task. The protection is achieved in two different embedding schemes, using the hidden memory state of RNN so that the image captioning functionalities are paralyzed for unauthorized usage. We demonstrated with extensive experiments that our proposed, on the one hand, the image captioning functionalities are well-preserved in the presence of valid secret key and well-protected for unauthorized usages on the other hand. The proposed key-based protection is, therefore, more cost-effective, proactive, and timely, as compared with watermarking-based protections which have to rely on government investigation and juridical enforcing actions. However, the proposed key-based protection also has some weaknesses, where the overall protection will be compromised when the attacker knows everything about the model. This is the direction in our future work to solve these weaknesses and ensure the model can be fully protected against different types of attackers.

6.1 Broader Impact

Our work is mainly focused on protecting the image captioning model with ownership verification. The engineer or researcher of the image captioning model might benefit from this research to protect their model against IP infringement. This is crucial as the development and training of an image captioning model is expensive especially when it involves a very large dataset. IP protection in image captioning is critical to fostering innovation in this field to achieve top-level performance that can benefit the society. We believe that no one in genuine may be put at disadvantage from this work. In case of the failure of our work in protecting the model, the worst scenario is an attacker can access the model without owner acknowledgement. In short, our work is bringing benefit to the society especially to AI start-ups to secure their advantage in the open market. We will also make the source code of this work publicly available for people to reproduce and follow up.

References

- [1] Yusuke Uchida, Yuki Nagai, Shigeyuki Sakazawa, and Shin’ichi Satoh, “Embedding watermarks into deep neural networks,” in *Proceedings of the 2017 ACM on International Conference on Multimedia Retrieval*, New York, NY, USA, 2017, ICMR ’17, p. 269–277, Association for Computing Machinery.
- [2] Huili Chen, Bitar Darvish Rouhani, Cheng Fu, Jishen Zhao, and Farinaz Koushanfar, “Deepmarks: A secure fingerprinting framework for digital rights management of deep learning

- models,” in *Proceedings of the 2019 on International Conference on Multimedia Retrieval*, 2019, pp. 105–113.
- [3] Jialong Zhang, Zhongshu Gu, Jiyong Jang, Hui Wu, Marc Ph Stoecklin, Heqing Huang, and Ian Molloy, “Protecting intellectual property of deep neural networks with watermarking,” in *Proceedings of the 2018 on Asia Conference on Computer and Communications Security*, 2018, pp. 159–172.
 - [4] Yossi Adi, Carsten Baum, Moustapha Cisse, Benny Pinkas, and Joseph Keshet, “Turning your weakness into a strength: Watermarking deep neural networks by backdooring,” in *USENIX*, 2018, pp. 1615–1631.
 - [5] Jia Guo and Miodrag Potkonjak, “Watermarking deep neural networks for embedded systems,” in *2018 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*. IEEE, 2018, pp. 1–8.
 - [6] Erwan Le Merrer, Patrick Perez, and Gilles Trédan, “Adversarial frontier stitching for remote neural network watermarking,” *Neural Computing and Applications*, pp. 1–12, 2019.
 - [7] Bitar Darvish Rouhani, Huili Chen, and Farinaz Koushanfar, “Deepsigns: An end-to-end watermarking framework for ownership protection of deep neural networks,” in *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*, 2019, pp. 485–497.
 - [8] Lixin Fan, Kam Woh Ng, and Chee Seng Chan, “Rethinking deep neural network ownership verification: Embedding passports to defeat ambiguity attacks,” in *Advances in Neural Information Processing Systems*, 2019, pp. 4716–4725.
 - [9] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio, “Show, attend and tell: Neural image caption generation with visual attention,” in *International Conference on Machine Learning*, 2015, pp. 2048–2057.
 - [10] Sepp Hochreiter and Jürgen Schmidhuber, “Long short-term memory,” *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
 - [11] Yuhui Quan, Huan Teng, Yixin Chen, and Hui Ji, “Watermarking deep neural networks in image processing,” *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
 - [12] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang, “Bottom-up and top-down attention for image captioning and visual question answering,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6077–6086.
 - [13] Raffaella Bernardi, Ruket Cakici, Desmond Elliott, Aykut Erdem, Erkut Erdem, Nazli Ikiçler-Cinbis, Frank Keller, Adrian Muscat, and Barbara Plank, “Automatic description generation from images: A survey of models, datasets, and evaluation measures,” *Journal of Artificial Intelligence Research*, vol. 55, pp. 409–442, 2016.
 - [14] Steven J Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel, “Self-critical sequence training for image captioning,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7008–7024.
 - [15] Andrej Karpathy and Li Fei-Fei, “Deep visual-semantic alignments for generating image descriptions,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3128–3137.
 - [16] Songtao Ding, Shiru Qu, Yuling Xi, Arun Kumar Sangaiah, and Shaohua Wan, “Image caption generation with high-level image features,” *Pattern Recognition Letters*, vol. 123, pp. 89–95, 2019.
 - [17] Xinwei He, Baoguang Shi, Xiang Bai, Gui-Song Xia, Zhaoxiang Zhang, and Weisheng Dong, “Image caption generation with part of speech guidance,” *Pattern Recognition Letters*, vol. 119, pp. 229–237, 2019.
 - [18] Xinyu Xiao, Lingfeng Wang, Kun Ding, Shiming Xiang, and Chunhong Pan, “Dense semantic embedding network for image captioning,” *Pattern Recognition*, vol. 90, pp. 285–296, 2019.
 - [19] Junbo Wang, Wei Wang, Liang Wang, Zhiyong Wang, David Dagan Feng, and Tieniu Tan, “Learning visual relationship and context-aware attention for image captioning,” *Pattern Recognition*, vol. 98, pp. 107075, 2020.

- [20] Junzhong Ji, Zhuoran Du, and Xiaodan Zhang, “Divergent-convergent attention for image captioning,” *Pattern Recognition*, p. 107928, 2021.
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [22] D Kinga and J Ba Adam, “A method for stochastic optimization,” in *ICLR*, 2015, vol. 5.
- [23] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick, “Microsoft coco: Common objects in context,” in *European Conference on Computer Vision*. Springer, 2014, pp. 740–755.
- [24] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier, “From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions,” *T-ACL*, vol. 2, pp. 67–78, 2014.
- [25] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh, “Cider: Consensus-based image description evaluation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 4566–4575.
- [26] Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould, “Spice: Semantic propositional image caption evaluation,” in *European Conference on Computer Vision*. Springer, 2016, pp. 382–398.
- [27] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu, “Bleu: a method for automatic evaluation of machine translation,” in *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*. Association for Computational Linguistics, 2002, pp. 311–318.
- [28] Satanjeev Banerjee and Alon Lavie, “Meteor: An automatic metric for mt evaluation with improved correlation with human judgments,” in *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, 2005, pp. 65–72.
- [29] Chin-Yew Lin, “ROUGE: A package for automatic evaluation of summaries,” in *Text Summarization Branches Out*, Barcelona, Spain, July 2004, pp. 74–81, Association for Computational Linguistics.
- [30] Abigail See, Minh-Thang Luong, and Christopher D Manning, “Compression of neural machine translation models via pruning,” *arXiv preprint arXiv:1606.09274*, 2016.