

# A Novel Actor Dual-Critic Model for Remote Sensing Image Captioning

Ruchika Chavhan<sup>\*</sup>, Biplab Banerjee<sup>\*</sup>, Xiao Xiang Zhu<sup>†</sup>, Subhasis Chaudhuri<sup>\*</sup>

<sup>\*</sup> Indian Institute of Technology Bombay, India

<sup>†</sup>Signal Processing in Earth Observation, Technical University of Munich, Germany

**Abstract**—We deal with the problem of generating textual captions from optical remote sensing (RS) images using the notion of deep reinforcement learning. Due to the high inter-class similarity in reference sentences describing remote sensing data, jointly encoding the sentences and images encourages prediction of captions that are semantically more precise than the ground truth in many cases. To this end, we introduce an Actor Dual-Critic training strategy where a second critic model is deployed in the form of an encoder-decoder RNN to encode the latent information corresponding to the original and generated captions. While all actor-critic methods use an actor to predict sentences for an image and a critic to provide rewards, our proposed encoder-decoder RNN guarantees high-level comprehension of images by sentence-to-image translation. We observe that the proposed model generates sentences on the test data highly similar to the ground truth and is successful in generating even better captions in many critical cases. Extensive experiments on the benchmark Remote Sensing Image Captioning Dataset (RSICD) and the UCM-captions dataset confirm the superiority of the proposed approach in comparison to the previous state-of-the-art where we obtain a gain of sharp increments in both the ROUGE-L and CIDEr measures.

## I. INTRODUCTION

Image Captioning is the task of generating a natural language description for a given image. Describing the images in properly framed sentences is a task that humans can perform with utmost ease [1], but was an inconceivable task for computers before the advent of deep learning. The primary requirements to automatically generate a caption are capturing the essence of the image and describing the correlations between objects to localize them in the image. Image captioning is more difficult than other tasks [2] that the computer vision community has concentrated on, as it involves the integration of semantic interpretation of an image into textual explanations.

Image captioning data is derived from different modalities: images that are represented by pixel intensities and textual descriptions represented as discrete word count vectors. The center of image captioning is recognizing the collective interpretation of these various modalities. The problem setting requires studying an abstract interpretation of the contents of the image and the semantic associations between the objects that can be constructed as a natural language description. Therefore, the problem of image captioning was defined as: Given an image  $I$ , a model is trained to maximize the likelihood  $p(W|I)$  where  $W = w_1, w_2, w_3, \dots, w_n$  where all  $w_i$  are words from a pre-defined vocabulary. One of the first

methods to describe images in human interpretable language [3] employed a Convolutional Neural Network (CNN) as an encoder to extract meaningful features from an image in the form of a fixed vector and a recurrent neural network (RNN) as a decoder to generate sentences. The primary inspiration of this model arises from the intuition that images can be translated to sentences, similar to the problem of machine translation [4] which follows an encoder-decoder architecture to translate the source sentence into target format through a bottleneck latent space. The goal of subsequent works in image captioning [5], [6] has been to produce more diverse captions that also accurately represent the image content. The aim of most image captioning works is to produce multiple contextual explanations that novelly manifest the local associations of objects in the image.

Reinforcement Learning [7] is a domain of Machine Learning that enables an agent to explore an environment by performing an action determined by a policy. While Deep Learning offers the best set of models for learning representations of multi-modal data, Reinforcement Learning is a framework for learning sequential decision-making tasks. Therefore, RL is the mainstream algorithm used to solve complex environments in different games to achieve super-human performance [8], robotics [9], and recommendation systems. Most RL agents acquire a stochastic mapping of states and actions, called the policy and pursue a trajectory by carrying out actions determined by this policy to maximize the total expected return in a given time phase. A model-free RL agent does not pre-specify a structural model of the environment, instead, it gradually learns the best policies based on trial and error and adequate exploration of the environment [10]. Most environments consist of multi-modal data and a single joint representation is learnt by an RL agent as a decision-making framework.

All reinforcement learning-based approaches are an exploration-based approach in which the agent first collects knowledge about the entire environment and state-space to anticipate reward-maximizing behavior. The agent is able to make highly optimized decisions in that way. Therefore, addressing the task of image captioning in the Reinforcement learning domain enables the generation of semantically more coherent captions. All supervised learning methods aim to generate sentences that are exactly identical to ground truth, while policies trained in an RL setup allow one to predict sentences much better than ground truth. Policy-generated

natural language descriptions can be dramatically enhanced by increasing the degree of environmental analysis consisting of paired images and captions.

Recently, acquisition of an unprecedented volume of satellite images from different sensors is observed. The major product derived from the satellite images has been the land-cover maps which assigns semantic class-labels at the pixel locations. However, the task of pixel-labeling is redundant for some emergency applications where an overall description of the scene under consideration is encouraged. Under this premise, the task of generating automatic captions for satellite scenes holds much potential. However, there does not exist extensive literature in this front specifically for remote sensing data.

In this paper, we implement an Actor Dual-Critical (ADC) training setup to address the issue of high inter-class image and caption similarity in satellite data. We have performed our proposed experiment on the Remote Sensing Image Captioning Dataset (RSICD) [11] and the UCM-captions dataset [12], [13]. The key problem in the datasets for many images is the strong inter-class similarity and the identical reference sentences. The contributions and results of the paper are summarized as follows:

- To the best of our knowledge, our approach is the first to use Reinforcement Learning to produce captions on remote sensing images. Unlike all existing methods, we employ an additional encoder-decoder RNN as a critic for Actor Dual-Critic (ADC) training setup. This critic plays a key role in creating a variety of different sentences that represent identical visual perception.
- The captions predicted by our proposed model are more semantically related to the objects in the image, explicitly describe object localization, and specifically focus on the existence of the entities in the image.
- We also perform cross-dataset captioning and obtain superior results on the RSICD dataset by models trained on the UCM-captions data set to demonstrate the generalization capability of a policy trained using the ADC setup as compared to previous state-of-the-art methods.

## II. RELATED WORK

Actor-critic methods in the Reinforcement Learning aim to train the participant in choosing actions taking into account the critic's reward. In the context of image captioning implemented in an actor-critic training strategy [14], given an image  $I$  and partially generated sentence  $S = (w_1, w_2, \dots, w_t)$ ,  $w_{t+1}$  is viewed as an action that the policy predicts. An actor's job is to learn a policy  $\pi(a_t|s_{t-1})$ , where  $a_t$  is the action performed and  $s_t$  is the state at time  $t$ . The critic provides the value function  $v_\theta^\pi$ , where  $\theta$  denotes the parameters of the value network. The reward metric for the generated caption is the score evaluated by ROUGE-L and  $v_\theta^\pi$  is used as an advantage baseline for the advantage factor  $A^\pi(s_t|a_{t+1}) = (\gamma^{T-t-1}r_T - v_\theta^\pi)$ . Subsequently, the actor is trained to optimize expected value of reward over the trajectory using a cross-entropy loss and the advantage factor

using the REINFORCE Algorithm where gradient are given by  $\mathbf{E}[\sum_{t=0}^T A^\pi(s_t|a_{t+1}) \nabla \log \pi(a_t|s_{t-1})]$ . Hence, this training method is named Advantage-Actor Critic (A2C). Instead of calculating a value function for baseline, the test-time reward metric can also be used a baseline to normalize the rewards experienced during training, thus eliminating the need for a critic, making the training setup self-critical [15], [16].

Recent work on remote sensing image captioning [12], [17] involves using deep multi-modal networks and analyzing the quality of generated captions by experimenting with various types of CNN, RNN, and LSTM combinations. Similar experiments were performed on the RSICD dataset [11] to produce accurate captions. Because of the wide area of the Earth's surface covered by remote sensing images, the main criteria for caption generation are to recognize semantic uncertainty in remote sensing images by analyzing key instances [18] and performing image context and landscape analysis.

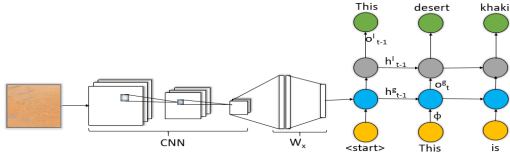
## III. PROPOSED METHODOLOGY

We first discuss the working principles of our proposed Actor Dual-Critic (ADC) training and elaborate on the three components employed in the process. We formulate the problem statement as follows: Given an image  $I$ , the model is expected to generate a natural language description  $S = (w_1, w_2, \dots, w_n)$  where  $w_i$  is a word from a pre-defined vocabulary. The actor generates a sentence given the image and the critics provide rewards based on the quality and relevance of this sentence. We utilize two critics: an RNN critic and an encoder-decoder RNN critic to provide two rewards to update parameters of the actor using the REINFORCE Algorithm. The training algorithm is mentioned in Algorithm 1.

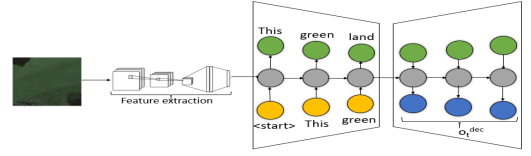
**Learning the Policy:** We propose a three-component model for the actor to learn a policy. The actor provides a measure of confidence  $q_\pi(a_t|s_t)$  to predict the next action  $\mathbf{a}_t = \mathbf{w}_{t+1} \in \mathbb{R}^d$  according to the current state. For the extraction of features from an image, we utilize a pre-trained CNN. The extracted features  $\mathbf{f} \in \mathbb{R}^n$  are passed as an input to a Gated Recurrent Unit (GRU). The hidden state of GRU  $\mathbf{h}_t^g \in \mathbb{R}^n$  evolves with time as the predicted words  $w_{t-1}$  are fed into it. The GRU acts indirectly as the generator of a context vector [19] to assist the decoder to surmount difficulties in learning due longer sentences. We observed a significant increase in performance by replacing the normal LSTM with a Layer Norm LSTM [20] with dropout. The following equations explain the functionality of the actor:

$$\begin{aligned} f &= W_x(\text{CNN}(I)) \\ \phi_0 &= f \\ o_t^g, h_t^g &= \text{GRU}(\phi_{t-1}, h_{t-1}^g) \\ o_t^l, h_t^l &= \text{LSTM}(o_t^g, h_{t-1}^l) \\ q_\pi(a_t|s_t) &= \psi(o_t^l) \\ \phi_t &= \zeta(w_{t-1}) \end{aligned} \tag{1}$$

Here,  $W_x$  is the weight of the linear embedding model of the CNN,  $o_t^g$  and  $o_t^l$  are the outputs of the GRU and



(a) Working of the actor. The words are converted into the embedding space (not shown) before being fed into the GRU (denoted by blue) and LSTM (denoted by grey).



(b) Working of the proposed Encoder-Decoder RNN critic. The outputs and the hidden states of the encoder RNN are transformed by  $\psi_1$  and  $\psi_2$  (not shown) and serve as input to the decoder RNN.

Fig. 1: Working Principle of the proposed ADC setup

LSTM respectively at time step  $t$ .  $\psi: \mathbb{R}^n \mapsto \mathbb{R}^d$  is a non-linear function that transforms the output of the LSTM to a space whose dimension is equal to the dimension of the vocabulary.  $\zeta: \mathbb{R}^d \mapsto \mathbb{R}^n$  denotes the embedding model to represent words in a common embedding space. We denote the policy network by  $\pi(a_t|s_{t-1})$ . Please refer to Figure 1a for complete description and pipeline of the model.

The model is trained to optimize the objective:

$$\min_{\pi} \sum_{t=0}^T \log(q_{\pi}(a_t|s_t)) \quad (2)$$

**Value Network:** This critic consists of an RNN which outputs a value function  $v_{\theta}^{\pi}(s_t)$  given the words predicted by the current policy and features extracted by the CNN in Figure 1a. We initialize the hidden state of the RNN by these extracted features. We will denote this critic as  $V(s_t)$ . The output of this network is directed to be the expected value of future rewards  $\mathbb{E}[\sum_{l=0}^{T-t-1} \gamma^l r_{t+l+1} | a_{t+1}, \dots, a_T \sim \pi, I]$  for choosing a state  $s_t$  given the current policy  $\pi$ . We set  $\gamma = 1$  similar to [14] and calculate the reward  $r_T$  after the prediction of the entire sequence implying  $r_t = 0, \forall t < T$ .

The reward  $r_T$  for the entire generated sentence is obtained by using the evaluation scores of ROUGE-L or BLEU. We observed more stable training using the Huber Loss instead of the norm of the difference between the reward and  $v_{\theta}$ :

$$L = \begin{cases} \|v_{\theta}^{\pi} - r_T\|^2 & \|v_{\theta}^{\pi} - r_T\| \leq \delta = 0.5 \\ \delta \|v_{\theta}^{\pi} - r_T\| - \frac{1}{2} \delta^2 & \text{otherwise} \end{cases} \quad (3)$$

**Encoder-Decoder LSTM Critic:** The addition of this critic is the main contribution of this paper. The intuition behind this critic is as follows:

- Theoretically, image captioning is defined as translation of images into sentences that aptly describe the images. We hypothesize that the sentences translated back into images should generate a quantity which closely resembles the features extracted from images.
- This procedure ensures maximum accumulation of information about the image through a textual description relatively similar to the information captured by ground truth sentences.

Please refer to Figure 1b for the entire pipeline of this model. The working principle of this critic denoted by  $D(S)$  is governed by the following equations:

$$\begin{aligned} h_0^{enc} &= W_x(\text{CNN}(I)) \\ \eta_t &= \zeta(S) \\ o_t^{enc}, h_t^{enc} &= \text{RNN}_{enc}(\eta_t, h_{t-1}^{enc}) \\ h_0^{dec} &= \psi_2(h_T^{enc}) \\ i_1^{dec} &= \psi_1(o_T^{enc}) \\ o_t^{dec}, h_t^{dec} &= \text{RNN}_{dec}(i_t^{enc}, h_{t-1}^{dec}) \end{aligned} \quad (4)$$

Here,  $\text{RNN}_{enc}$  and  $\text{RNN}_{dec}$  are the encoder-decoder RNN respectively.  $S = (w_1, w_2, \dots, w_T)$  denotes a natural language description of the image. We have used  $\psi_1, \psi_2: \mathbb{R}^n \mapsto \mathbb{R}^n$  as a linear function with dimension equal to that of the embedding space of sentences, along with the ReLU activation function.

This critic is trained by optimizing the mean squared error (MSE) between the output of the decoder and the features:

$$L = \left( \frac{\sum_{t=0}^T o_t^{dec}}{|S|} - f \right)^2 \quad (5)$$

Accuracy for the decoder output is given by the cosine similarity between the output of the decoder and features:

$$A_{gen} = \frac{\frac{\sum_{t=0}^T o_t^{dec}}{|S|} f}{\|f\| \left\| \frac{\sum_{t=0}^T o_t^{dec}}{|S|} \right\|} \quad (6)$$

$A_{gen}$  and  $A_{orig}$  are the accuracies of the network when captions generated by the actor and ground truth captions are fed into the encoder respectively. We defined an advantage factor for this critic to be:

$$A_{ed} = A_{gen} - \delta_t A_{orig} \quad (7)$$

The encoder-decoder critic is pre-trained on features extracted and corresponding original sentences to learn a ground truth latent representation. We observed that  $\delta_t = 1$  initially leads to exploding gradients and a non-converging policy. Therefore, we begin with  $\delta_t = 0.01$  and increase it linearly to 1.0 over the epochs. For generating captions for a test image, we pass the image  $I_{test}$  to the actor to generate sentences till the <end> token is encountered.

#### IV. EXPERIMENTS AND RESULTS

In this section, we study the performance of our model and its results on two remote sensing image captioning datasets: RSICD and UCM-captions. We also visualise the output of our novel critic and analyze the validity of its working principle.

---

##### Algorithm 1 Training Algorithm

---

**Input:** Pre-trained models  $\pi(a_t|s_{t-1})$ ,  $V(s_t)$  and  $D(S)$  using the objectives given by the equations 2, 3 and 5 respectively as done in [14].

- 1: **for**  $episode = 1$  to total episodes **do**
- 2:   Given an Image  $I$  sample action  $(a_1, a_2, \dots, a_T)$  from the current policy using a multinomial distribution given by  $q_\pi(s_t|a_t)$ ;
- 3:   Calculate advantage factor  $A^\pi$  using the reward  $r_T$  for the value network;
- 4:   Update the parameters of the policy using  $A^\pi$  by the REINFORCE Algorithm;
- 5:   Update parameters of the critic by optimising Eq. 3;
- 6:   Calculate advantage factor  $A_{ed}$  using the encoder-decoder critic;
- 7:   Update the parameters of the policy using  $A_{ed}$  by the REINFORCE Algorithm;
- 8:   Update parameters of the critic using  $A_{orig}$ .
- 9: **end for**

---

**Datasets:** The vocabulary of words obtained from the RSICD dataset contains a collection of 1653 words including the <start>, <unk> and <end> tokens. The dataset contains image-caption pairs of 30 classes namely airport, bare-land, baseball-field, beach, church, commercial, dense-residential, meadow, river, bridge etc. The dataset contains a total of 10921 images, split into 8734 training, 1094 validation, and 1093 testing images. The UCM-captions dataset contains identical captions for images from the same class and thus spans a vocabulary of only 210 words. It contains 21 classes land use image, including agricultural, airplane, buildings, chaparral, overpass, parking lot, river, runway etc and with 100 images for each class. Because of the broad variety of image classes they provide, the above two datasets are widely used in remote sensing image captioning research.

**Experiment Details:** We observed a substantial improvement in performance for extracting features from images by employing AlexNet [23]. A fully connected layer with dimension of 256 and batch normalization along with the ReLU activation function is applied to the output of the feature extractor. We use an embedding module of dimension 256 to encode each word into an embedding space. For the value Network, a fully connected layer of dimension  $256 \times 1$ , with hyperbolic tan has been applied as the activation function to obtain a value function in the range  $[-1, 1]$ . We also observed that by employing a GRU in the encoder-decoder critic instead of an RNN results in a sharp increment in training speed. The functions  $\psi_1$  and  $\psi_2$  in Equation 4 is a fully connected layer with dimension 256. We used the Adam Optimiser [24] with

a learning rate of  $5e-4$  and is decreased by factor of 0.9 after every 10 iterations. The networks are trained for a total of 100 epochs.

**Metrics for Comparison of results:** We compare our results qualitatively and quantitatively with the [11] who trained a deep multi-modal neural network (referred to as MM) with different types of CNNs, RNNs, and LSTMs for semantic understanding of high-resolution remote sensing images in the RSICD and UCM-captions dataset. The authors also performed experiments on their dataset using the "hard" and "soft" attention mechanisms proposed by [6] denoted by HA and SA respectively. These three methods are not methods based on reinforcement learning, except for the mechanism of hard attention which uses the REINFORCE algorithm for attention but not for the prediction of the caption. Table I and Table II quantitatively compares the three methods (denoted by MM, SA, HA respectively) and the results of the A2C training setup with our proposed method. In Figure 4 and Figure 5, we compare the results of the A2C training setup replacing their LSTM by a Layer Norm LSTM like in our ADC training setup. Figure 2b shows the values of the reward metrics ROUGE-L, BLEU-1, BLEU-2, BLEU-3 and BLEU-4 during training on the RSICD dataset.

**Demonstrating validity of proposed critic:** The aim of this experiment is to visualise the relevance of the output of the decoder with respect to the output of the feature extractor CNN for different image-sentence pairs. To verify if the critic can distinguish between correct and incorrect captions, we also input a test image of the same class with a different sentence from the reference sentence of the ground truth (Figure 2a). The figures in this section are normalized vector representation resized for the sake of visualization. Figure 3c represents the difference between the output of the decoder for a reference sentence and the features extracted by the CNN. The difference between the output of the decoder for the ground truth and predicted sentence is shown in Figure 3e. As expected, this difference is negligible (faint grey lines). This means that a generalized high-level correspondence between images and sentences that defines it semantically has been learned by this critic. Figure 3g is the difference between the representations learnt by the encoder-decoder critic for an image from the same class on a different ground truth caption. As expected, this difference is high implying that a sentence not describing an image does not get translated into abstract features corresponding to that image. We deduce that this critic does not trivially learn the identity function with respect to input features, due to the presence of non-invertible functions  $\psi_1$  and  $\psi_2$  in Equation 4. It means that the critic takes into account the correspondences of both image and sentence for two sentences with entirely different word2vec representations. We may infer from the above experiment that the critic effectively learns the correlations between images and reference sentences and encourages the generation of different sentences that encapsulate the same correlations.

**Qualitative Analysis of Results:** As observed from the result of the sentence generation, our proposed method can generate

TABLE I: Results of ADC setup on the RSICD dataset

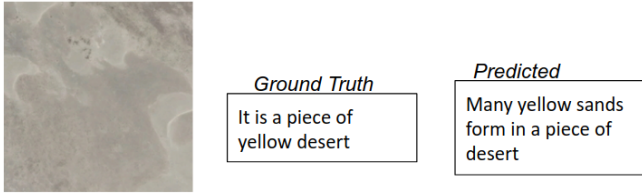
Metric	B-1	B-2	B-3	B-4	METEOR [21]	ROUGE-L	CIDEr [22]
MM [11]	0.57905	0.41871	0.32628	0.26552	0.26103	0.51913	2.05261
SA [11]	0.65638	0.51489	0.41764	0.34464	0.32924	0.61039	1.87415
HA [11]	0.68968	0.5446	0.44396	0.36895	<b>0.33521</b>	0.62673	1.98312
A2C [14]	0.60157	0.41991	0.364516	0.28788	0.19382	0.63185	2.098
Ours	<b>0.73973</b>	<b>0.55259</b>	<b>0.46353</b>	<b>0.41016</b>	0.22126	<b>0.71311</b>	<b>2.243</b>

TABLE II: Results of ADC setup on the UCM dataset

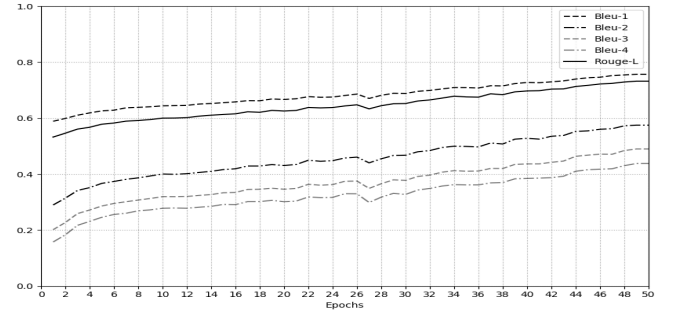
Metric	B-1	B-2	B-3	B-4	METEOR [21]	ROUGE-L	CIDEr [22]
MM [11]	0.37066	0.32344	0.32346	0.23259	0.40476	0.4236	1.708
SA [11]	0.79693	0.71345	0.6514	0.59895	0.74952	0.41676	2.12846
HA [11]	0.78498	0.70929	0.65182	0.60167	0.77357	0.43058	2.19594
A2C [14]	0.373089	0.23776	0.15857	0.12222	0.39645	0.35989	2.381
Ours	<b>0.85330</b>	<b>0.75679</b>	<b>0.67854</b>	<b>0.61165</b>	<b>0.83242</b>	<b>0.80872</b>	<b>4.865</b>

TABLE III: Results of cross dataset captioning on the RSICD dataset

Metric	B-1	B-2	B-3	B-4	METEOR [21]	ROUGE-L	CIDEr [22]
MM [11]	0.19618	0.01481	0.00721	0.00445	0.07416	0.2457	0.08015
A2C [14]	0.19405	0.04137	0.00714	0.00175	0.18855	0.1846	0.961
Ours	<b>0.38810</b>	<b>0.08643</b>	<b>0.019065</b>	<b>0.00608</b>	<b>0.23964</b>	<b>0.2888</b>	<b>2.013</b>



(a) Original Image with ground truth and predicted sentences. The image is from the class 'desert'. A test image from the same class but with the same captions as the test input image are passed to the encoder for this experiment. The test image used for this case belongs to the class 'desert'.



(b) Reward during Training Process.

Fig. 2



(a) Features extracted from image



(b) Output of the decoder for reference sentence



(c) Difference between (a) and (b)



(d) Output of the decoder for predicted sentence



(e) Difference between (b) and (d)

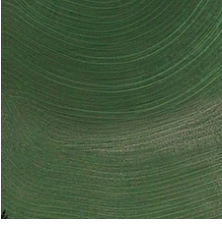


(f) Output of the decoder for different image from same class



(g) Difference between (b) and (f)

Fig. 3: Qualitative results of the experiment comparing output of the decoder for different image-sentence pairs



(a) **A2C:** It is a piece of green meadow.  
**Ours:** A dirt lines are in this meadow.



(b) **A2C:** It is a piece of yellow desert.  
**Ours:** It is a rather flat desert stained with several black stains



(c) **A2C:** Many rectangular buildings and green trees are in a dense area.  
**Ours:** Houses with red roofs on both sides of the road.

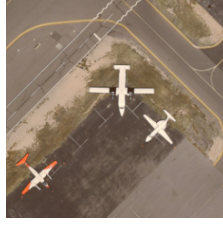


(d) **A2C:** Many white boats are in the port.  
**Ours:** Two rows of white boats are in port

Fig. 4: Qualitative results of image captioning of the ADC setup on the RSICD dataset



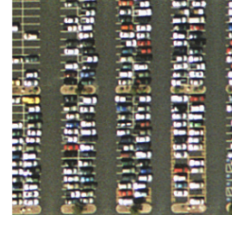
(a) **A2C:** There are many airports at the airport.  
**Ours:** There is a red air-plane with lots of cars.



(b) **A2C:** There are many airports at the airport.  
**Ours:** There is a red air-plane in the airport.



(c) **A2C:** There are many buildings.  
**Ours:** There is one road next to many buildings.

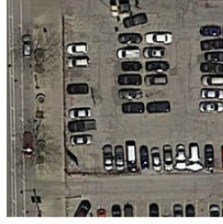


(d) **A2C:** There are lots of cars with some buildings.  
**Ours:** Lots of cars are rectangular and close to each other in the parking lot.

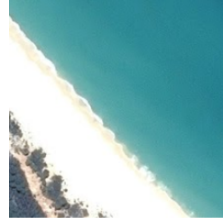
Fig. 5: Qualitative results of image captioning of the ADC setup on the UCM-captions dataset



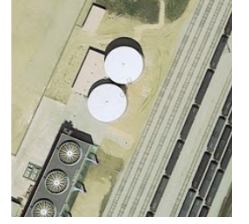
(a) **Original:** Many tall buildings are in a commercial area.  
**Ours:** There is one road next to many buildings.



(b) **Original:** Many cars are parked in the parking lot.  
**Ours:** Lots of cars are parked neatly in a parking lot.



(c) **Original:** Yellow beach is between green ocean and green trees.  
**Ours:** This is a beach with blue sea and white sands.



(d) **Original:** On the ground, there are two spherical storage tank.  
**Ours:** Two small storage tanks are on the ground.

Fig. 6: Qualitative results of cross dataset image captioning of the ADC setup on the RSICD dataset

more complicated yet explainable sentences having longer lengths, more words with low frequency in caption labels. From the examples in Figure 7, it is evident that the generated captions contain phrases that provide a highly accurate semantic explanation of the nature and localization of objects in the scene. We note, however, that these phrases are absent in the caption of ground truth sentences. We analyzed ground-truth captions of other images of the same class containing such phrases to explain the occurrence of these phrases, and compared the test image with the corresponding images in the data set. From Figure 7, it can be deduced that the addition

of an encoder-decoder RNN critic has significantly improved the quality of the policy’s performance on remote sensing images as compared to the baseline method (A2C). This demonstrates that the policy has successfully investigated the environment consisting of images and captions and has gained more knowledge compared to the baseline approach due to this critic’s extra upgrade step in the optimization of policy objective. As observed in Table I, our proposed approach provides a rapid increase in six out of seven scores used for comparison for the RSICD dataset. However, our approach proves superior for the UCM-captions dataset than previous







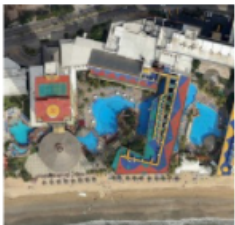


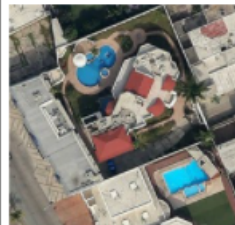





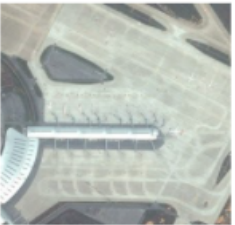


 <p><b>Ground truth:</b> A town is located on the bank of the river with a <i>island</i>.</p>	<p><b>A2C:</b> Green trees are in two sides of a curved river.</p> <p><b>Ours:</b> This <i>s shaped green river</i> with an <i>island</i> in it goes through this land divided into blocks of farms.</p>	 <p><b>Ground truth:</b> The <i>s shaped river</i> surrounds the triangular grey green land</p>	 <p><b>Ground truth:</b> An <i>s shaped river</i> with forest belt alongside is lying on a meadow</p>	 <p><b>Ground truth:</b> A mid lake <i>island</i> is in the <i>s shaped river</i> which is between two living region.</p>
 <p><b>Ground truth:</b> Next to the <i>beach</i> is a large resort</p>	<p><b>A2C:</b> Some buildings and green trees are in a resort.</p> <p><b>Ours:</b> Several <i>buildings</i> with <i>swimming pools</i> and some <i>green plants</i> are near a <i>beach</i>.</p>	 <p><b>Ground truth:</b> Many buildings and several <i>swimming pools</i> are near a beach.</p>	 <p><b>Ground truth:</b> Many <i>green trees</i> and buildings with swimming pools are near a beach.</p>	 <p><b>Ground truth:</b> Some <i>buildings</i> with two swimming pools are near a road in a resort</p>
 <p><b>Ground truth:</b> There is a road passing by a residential area where rows of high rise housings are standing.</p>	<p><b>A2C:</b> Many tall buildings are in a commercial area.</p> <p><b>Ours:</b> <i>Three rows of skyscrapers</i> stands at this <i>prosperous commercial area</i>.</p>	 <p><b>Ground truth:</b> Plenty of <i>skyscrapers</i> stands at this <i>prosperous commercial area</i>.</p>	 <p><b>Ground truth:</b> <i>Three rows of tall buildings</i> are in a commercial area.</p>	 <p><b>Ground truth:</b> There are several blocks in this <i>prosperous commercial area</i>.</p>
 <p><b>Ground truth:</b> Many planes are parked around a large building in an airport with runways</p>	<p><b>A2C:</b> Some planes are parked near an airport with a parking lot.</p> <p><b>Ours:</b> Several <i>white planes</i> are around a <i>circle building</i> with a parking lot.</p>	 <p><b>Ground truth:</b> Semi <i>circle building</i> connected with a long narrow buildings with planes parked along.</p>	 <p><b>Ground truth:</b> Many <i>white planes</i> are in an airport near some buildings.</p>	 <p><b>Ground truth:</b> The terminal buildings including three <i>circle buildings</i> tarmac and runways are built in the field of sparse grass.</p>

Fig. 7: Qualitative comparison of results of image captioning of the ADC setup on the RSICD dataset.

methods for all the seven scores. CIDEr captures consensus-based human judgment better than established metrics through sentences created from different sources. The implementation of an encoder-decoder LSTM critic has resulted in strong increments in CIDEr for both the datasets.

**Cross Dataset Captioning:** Testing trained models across datasets with similar domains gives an understanding of the model's ability to generalize and utilise it for real time predictions. We have therefore tested a model trained on the UCM-captions dataset on the RSICD dataset and made qualitative as well as quantitative comparisons with previous methods [11]. The models trained on another dataset experience a rapid decline in metrics compared with the outcome of model trained on the corresponding dataset. From Table III, it is noted that the inclusion of our proposed critic has resulted in a substantial improvement in all 7 scores for the captioning of the cross datasets. We also note from Figure 6 that our policy generates the captions which convey the same meaning as the sentences of ground truth. This ensures that the policy trained using the ADC system can be used in Remote Sensing applications in real life.

## V. CONCLUSION

In this paper, we proposed an Actor Dual-critic (ADC) method for Image Captioning for the Remote Sensing Image Captioning Dataset. We are introducing another critic to the A2C training setup to encourage the prediction of sentences capturing relevant details along with sentence diversity. In the sense of converting sentences back to original images we suggest an encoder-decoder model for this critic. We also proposed a training strategy with an advantage factor based on a weighted difference of cosine similarities to update the policy parameters. We prove the superiority of our method quantitatively using the metrics BLEU-1, BLEU-2, BLEU-3, BLEU-4, ROUGE-L, and CIDEr. Ultimately, we use a model trained over the UCM-captions dataset and validate its superiority over other approaches to generate textual descriptions of images in the RSICD dataset.

## REFERENCES

- [1] L. Fei-fei, A. Iyer, C. Koch, and P. Perona, "What do we perceive in a glance of a real-world scene?" *Journal of vision*, vol. 7 1, p. 10, 2007.
- [2] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge," *International Journal of Computer Vision (IJCV)*, vol. 115, no. 3, pp. 211–252, 2015.
- [3] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, "Show and tell: A neural image caption generator," in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015, pp. 3156–3164.
- [4] K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using rnn encoder-decoder for statistical machine translation," in *EMNLP, A. Moschitti, B. Pang, and W. Daelemans, Eds. ACL*, 2014, pp. 1724–1734.
- [5] A. Karpathy and L. Fei-Fei, "Deep visual-semantic alignments for generating image descriptions," in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015, pp. 3128–3137.
- [6] K. Xu, J. L. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhutdinov, R. S. Zemel, and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention," in *Proceedings of the 32nd International Conference on Machine Learning - Volume 37, ser. ICML'15*. JMLR.org, 2015, p. 2048–2057.
- [7] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. Cambridge, MA, USA: A Bradford Book, 2018.
- [8] S. D. Holcomb, W. K. Porter, S. V. Ault, G. Mao, and J. Wang, "Overview on deepmind and its alphago zero ai," in *Proceedings of the 2018 International Conference on Big Data and Education*, ser. ICBDE '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 67–71. [Online]. Available: <https://doi.org/10.1145/3206157.3206174>
- [9] S. Levine, C. Finn, T. Darrell, and P. Abbeel, "End-to-end training of deep visuomotor policies," *J. Mach. Learn. Res.*, vol. 17, no. 1, p. 1334–1373, Jan. 2016.
- [10] B. D. Ziebart, A. Maas, J. A. Bagnell, and A. K. Dey, "Maximum entropy inverse reinforcement learning," in *Proceedings of the 23rd National Conference on Artificial Intelligence - Volume 3, ser. AAAI'08*. AAAI Press, 2008, p. 1433–1438.
- [11] X. Lu, B. Wang, X. Zheng, and X. Li, "Exploring models and data for remote sensing image caption generation," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 4, pp. 2183–2195.
- [12] B. Qu, X. Li, D. Tao, and X. Lu, "Deep semantic understanding of high resolution remote sensing image," in *2016 International Conference on Computer, Information and Telecommunication Systems (CITS)*, 2016, pp. 1–5.
- [13] Y. Yang and S. Newsam, "Bag-of-visual-words and spatial extensions for land-use classification," 01 2010, pp. 270–279.
- [14] L. Zhang, F. Sung, F. Liu, T. Xiang, S. Gong, Y. Yang, and T. Hospedales, "Actor-critic sequence training for image captioning," 06 2017.
- [15] S. Rennie, E. Marcheret, Y. Mroueh, J. Ross, and V. Goel, "Self-critical sequence training for image captioning," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1179–1195, 2017.
- [16] J. Gao, S. Wang, S. Wang, S. Ma, and W. Gao, "Self-critical n-step training for image captioning," 04 2019.
- [17] X. Zhang, Q. Wang, S. Chen, and X. Li, "Multi-scale cropping mechanism for remote sensing image captioning," in *IGARSS 2019 - 2019 IEEE International Geoscience and Remote Sensing Symposium*, 2019, pp. 10 039–10 042.
- [18] Z. Shi and Z. Zou, "Can a machine generate humanlike language descriptions for a remote sensing image?" *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 6, pp. 3623–3634, June 2017.
- [19] Y. Bengio and Y. LeCun, Eds., *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. [Online]. Available: <https://iclr.cc/archive/www/doku.php%3Fid=iclr2015:accepted-main.html>
- [20] J. Ba, J. Kiros, and G. Hinton, "Layer normalization," 07 2016.
- [21] A. Lavie and A. Agarwal, "Meteor: An automatic metric for mt evaluation with high levels of correlation with human judgments," in *Proceedings of the Second Workshop on Statistical Machine Translation*, ser. StatMT '07. USA: Association for Computational Linguistics, 2007, p. 228–231.
- [22] R. Vedantam, C. L. Zitnick, and D. Parikh, "Cider: Consensus-based image description evaluation," *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4566–4575, 2014.
- [23] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems* 25, F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2012, pp. 1097–1105. [Online]. Available: <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>
- [24] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *International Conference on Learning Representations*, 12 2014.