

# aiTPR: Attribute Interaction-Tensor Product Representation for Image Caption

Chiranjib Sur

Computer & Information Science & Engineering Department, University of Florida.

Email: chiranjib@ufl.edu

**Abstract**—Region visual features enhance the generative capability of the machines based on features, however they lack proper interaction attentional perceptions and thus ends up with biased or uncorrelated sentences or pieces of misinformation. In this work, we propose Attribute Interaction-Tensor Product Representation (aiTPR) which is a convenient way of gathering more information through orthogonal combination and learning the interactions as physical entities (tensors) and improving the captions. Compared to previous works, where features are added up to undefined feature spaces, TPR helps in maintaining sanity in combinations and orthogonality helps in defining familiar spaces. We have introduced a new concept layer that defines the objects and also their interactions that can play a crucial role in determination of different descriptions. The interaction portions have contributed heavily for better caption quality and has out-performed different previous works on this domain and MSCOCO dataset. We introduced, for the first time, the notion of combining regional image features and abstracted interaction likelihood embedding for image captioning.

**Index Terms**—language modeling, representation learning, tensor product representation, image description, sequence generation, image understanding, automated textual feature extraction

## I. INTRODUCTION

OBJECT recognition and segmentation has provided ample scope to understand the semantic relationship in images among the different objects [83] as a spatio-topological property. This will help in understanding the contexts [69] and events of the images than mere detection of the objects. The gradual demand and rising interest in industry and AI related frameworks, the requirement is generation and synthesis of reply in structured form. There was a sudden rise in application related to reverse synthesis, caption generation [85], dialogues instead of just prediction of likelihood of decisions. Most of the application of synthesis is based on the requirement to perceive and development of topological dependence of contexts and proceed for generation and synthesis. The basis of our image captioning [86] model is based on reviving the underlying understanding of the representational aspects from both the context and the role prospective for languages. Here, we define Attribute Interaction-Tensor Product Representation (aiTPR), where the context is more about understanding and the roles play the role of creating the dependency or the interaction [85]. Here, we defined a layer that samples the combination, transforms them into a understandable space and help generate the captions. While "English" language is based on structural principles, these are roles and contexts

are pseudo-principles (as a word can have different roles like noun, verb etc and also different contexts or meaning) and this creates a large spectrum of operation and valid possibilities, which can be narrowed only through understanding the interactions. While, we can only deal with the objects and the role transformation as our aiTPR, we provided evidence that the interaction terms are far more important and can help bridging the gap created by the object subspace, like in [83], where the object space is mere structural from images or the features space is randomly defined to fit.

Previous works in image captioning focused on visual features from pre-trained network capable to detect objects with high precision like Vgg [47], ResNet ([83], [50]), Inception [43] etc, later gradually going multimodal ([47], [49], [84]) through different stages of feature, expecting that different layers will capture different aspects. Later on, gradually attention [45], [43], [56], [50], [34], [41], [11]) based features became more popular due to the limitations of LSTM to keep its memory intact and limitations of the weighted transformations to capture all the relevant features required for diversified and correct attributed captions. Tensor Product Representation (TPR) concept is widely studied mathematical model with several mathematical properties, popular among people who wants to deal with orthogonality. However, the computational TPR space is yet to be explored. Recently, TPR [88] is used to solve the question answering problem and the image captioning problem, where three models were proposed: dual generation model [85] for capturing both the context and roles for next word prediction, attention based model [85] for capturing the interdependence of the different words and by-passing the role of semantic features, and lastly the graphical tensor product representation model [85] combining both semantics and structural perceptions of triplet graphs for local features extraction and establishing the spatial temporal relationship for sentence generation. In this work, we will mostly concentrate on some novel structures for caption generation where the main focus is on categorically understand the individual components and produce most favorable representation for best reverse generation. The reverse generation is mostly done through language and some cases through images, like finding the most relevant images or generate images through GAN architectures. Reverse generation training has the problem of gathering biasness for the models and restricts generation of diverse sentences and out of box thinking or combination possibilities. Image captioning has always provided scope of further possibilities for video narration and understanding of

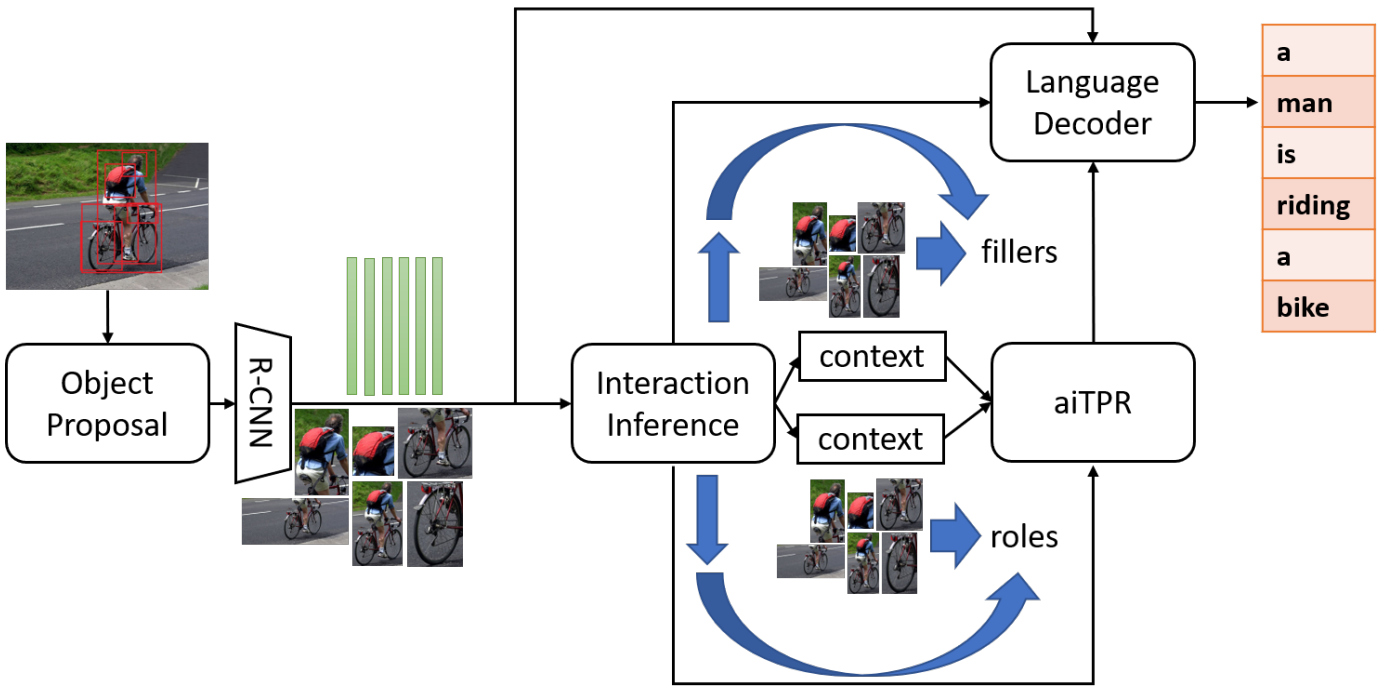


Fig. 1. Overview Architecture of Attribute Interaction-Tensor Product Representation (aiTPR) as a Combination of Tensors Derived from Image Level Attributes and Inferred Interactions.

sequentially related image spaces and story-telling.

The rest of the document is arranged as Section ?? with the revisit of the existing works, descriptions of the problem in Section II, methodology and architectural details in Section III, experimental details, numerical and qualitative results and analysis in Section IV and conclusion and discussion of future works in Section V.

The main contribution of this paper are the followings: 1) defining a new architecture and concept of representation generation and fusion known as Attribute Interaction-Tensor Product Representation (aiTPR) 2) interactions bind different features together along with the orthogonality of the computational TPR architecture 3) our model contributed to the fact that orthogonality decides feature spaces and restricts mixing and mapping to undefined feature spaces 4) our model outperformed many previous works done on this domain and on the MSCOCO data.

## II. DESCRIPTION

The existing problem in language generation problem is the disability of the machines to differentiate between similar quite of situations, which arises because of the similar kind of combination of the aspects. In this work, our proposal architecture handles this problem in a different way, where we derive different levels of fine details of an image as likelihood, then the interactions are derived and then we derive different context and role generations for the sentence to be completed. However, compared to many previous works like in [85], this approach is much more exploratory and performed better in caption generation.

### A. Problem Description

Generative application is gradually gained potential as we need to make the machine talk to humans. However, biasness in sentence generation held the topological exploration for a longer time. In this work, we explored the possibility of creating interaction criteria as an estimation and proceed them with the regional image features for derivation of the representations for context and roles as we can call them as tensor product when we consider them as a combined representation ( $T_t$  in Equation 32). We deliberately create this tensor product so that we can multiply them as tensors and can define the semantics for the languages.

### B. Difficulties Faced So Far

Normal network learning does through these five difficulties apart from vanishing and exploding gradient problems, which is either bypassed or overlooked. These are as follows and we described here, how we can overcome this, mainly focusing on the model we have define in this work.

a) *Object and Combination Problem*: When similar kinds of object combination exists in images, it becomes very difficult for the CNN network to differentiate them and provide an unique solution. In that regard, aiTPR can provide better solutions.

b) *Weight Learning Problem*: While the weights are required to learn large part of the transformation for large amount of data, aiTPR can be a better alternative as it segregates the information in the form of a graphical structure instead of transformation.

c) *Lower and Higher Level Understanding Problem*: While the lower level features are the image features, several combinations created by aiTPR can help in identification of the

unique opportunities for new attributes to crop up in sentences along with attributes.

d) *Summation to Undefined Spaces Problem*: Summation problem is another problem that should be dealt with, mainly with the presence of large number of objects and their forms. Several such summation may converge them to similar spaces or to null spaces.

e) *Representation Approximation*: Due to representation approximation, large part of the information are suppressed or gets mixed into unidentified states. To counter that, aiTPR keeps up the states through orthogonal transformation and are expected to regain the subspace, when required. All these when combined can provide a much better alternative and structurally sound neural network, that can be more explained.

### III. OUR METHODOLOGY

Many organizations used an ensemble approach to avoid these problems instead of solving it naturally. In this work, we have proposed series of solutions for these above problems through use of other learned networks and intermediate inferences, which provides ample scope to neutralize the discrepancies in comprehension of the representations. Here, we have used extensively the region based object features and extracted different layer of information and fused them for representation generation in our network architecture and for each iteration, these representations took a new form so that the captions did not have to go through the biasness of non-linear approximations.

Here, we have provided the concept of Attribute Interaction-Tensor Product Representation (aiTPR) which operates on the regional features and their interaction criteria. While, the interactions are transformed replications of the objects features and are generated through inference, we can still consider two separate strategies to engage the attributes and the interaction segments. Figure 2 provided Late Attribute Interaction-Tensor Product Representation (aiTPR) while Figure 3 provided Early Attribute Interaction-Tensor Product Representation (aiTPR) architectures. The main difference is the way then mixture of the objects representations takes place as a weighted sum of the lower level features of regional images. This kind of strategies are already in the literature like in MIMO antenna, where the aim is to estimate the best possible entropy from series of similar signals. The novelty of our procedure is that we used the whole image and series of spatial relationships as a semantic composition denoting objects and activity-relationships, which provides opportunities for new objects and their interactions get more attention than before.

#### A. Early Attribute Interaction-Tensor Product Representation

Early Attribute Interaction-Tensor Product Representation uses the strategy of linking the attribute and interaction tensors much early in the network, considering the fact that there will be correlations establishment among them much earlier. Mathematically, we represent Early Attribute Interaction-Tensor Product Representation with the following set of equations. If we define  $v \in \mathbb{R}^{2048}$  as the visual feature for image  $\mathbf{I}$  and  $\mathbf{v} = \{v_1, \dots, v_n\}$  as the RCNN based region based

object feature representation where each  $v_i \in \mathbb{R}^{2048}$  for  $i \in \{1, \dots, n\}$  and  $n$  is the number of regional objects detected in the image  $\mathbf{I}$  making the overall  $\mathbf{v} \in \mathbb{R}^{n \times 2048}$ . The initial parameters for the Assembled Selector Layer are initialized as the followings, considering that the biasness of the network must be neutralized and it will also help in estimation of the content of the images.

$$\bar{v} = \frac{1}{k} \sum_{i=1}^{i=k} v_i \quad (1)$$

$$\bar{v} = \mathbf{v} \quad (2)$$

The initial parameters are initialized as the followings.

$$\mathbf{h}_0, \mathbf{c}_0 = \mathbf{W}_{h_0} \bar{v}, \mathbf{W}_{c_0} \bar{v} \quad (3)$$

$\mathbf{W}_{h_0} \in \mathbb{R}^{2048 \times d}$ ,  $\mathbf{W}_{c_0} \in \mathbb{R}^{2048 \times d}$ . The Intermediate Transfer Layer helps in transferring the

$$\mathbf{a}_t = \mathbf{W}_a \tanh(\mathbf{W}_h \mathbf{h}_{t-1}) \quad (4)$$

where  $\mathbf{W}_a \in \mathbb{R}^{b \times d}$ ,  $\mathbf{W}_h \in \mathbb{R}^{d \times k}$ .

$$\alpha_t = \text{softmax}(\mathbf{a}_t) \quad (5)$$

$$\mathbf{a}_t \in \mathbb{R}^k \in \{a_{1,t}, \dots, a_{k,t}\}$$

$$\hat{v}_t = \frac{\sum_{i=1}^{i=k_1} v_i \alpha_{i,t} + \sum_{i=1}^{i=k_2} v'_i \alpha'_{i,t}}{2} \quad (6)$$

with  $k = k_1 + k_2$ ,  $\sum \alpha_i = 1$  and  $\sum \alpha'_i = 1$   $v$  is the regional CNN and  $v'$  is the representation of the objects detected through the regional CNN model.  $\hat{v}_t \in \mathbb{R}^{b \times d}$  where  $b$  is the batch size and  $d$  is the hidden layer dimension. The Assembled Selector Layer can be denoted as the following equations,

$$\mathbf{q}_t = \hat{v}_t \quad (7)$$

$$\mathbf{p}_t = \mathbf{W}_e \mathbf{x}_{t-1} \quad (8)$$

$$\begin{aligned} \mathbf{T}_t &= \mathbf{W}_{s_{12}} \sigma(\mathbf{W}_{s_{11}} \mathbf{h}_{t-1} + \mathbf{b}_1) \\ &\otimes \tanh(\mathbf{W}_{s_{22}} (\mathbf{v}_x \sigma(\mathbf{W}_{s_{21}} \mathbf{h}_{t-1} + \mathbf{b}_2)) + \mathbf{b}_3) \end{aligned} \quad (9)$$

$$\begin{aligned} \mathbf{T}_t &= \mathbf{W}_{s_{12}} \sigma(\mathbf{W}_{s_{11}} \mathbf{h}_{t-1} + \mathbf{W}_{w_1} \sum_{i=0}^{t-1} \mathbf{W}_e \mathbf{x}_i + \mathbf{b}_1) \otimes \\ &\tanh(\mathbf{W}_{s_{22}} (\mathbf{v}_x \sigma(\mathbf{W}_{s_{21}} \mathbf{h}_{t-1} + \mathbf{W}_{w_2} \sum_{i=0}^{t-1} \mathbf{W}_e \mathbf{x}_i + \mathbf{b}_2)) + \mathbf{b}_3) \end{aligned} \quad (10)$$

$$\mathbf{v}_x = f_x(\{v_{a_1}, v_{a_2}, \dots, v_{b_1}, v_{b_2}, \dots\}) \quad (11)$$

where we have  $f_x(\cdot)$  as a function and  $v_{a_i} \in \mathbb{R}^{2048}$  and  $v_{b_i} \in \mathbb{R}^{2048}$  are the attribute and interaction components.  $\otimes$  is an algebraic operation. Here, we considered  $\otimes = \odot$  as we try to rectify one context with the other context.

$$\mathbf{q}_t = \mathbf{W}_{h,m} S \odot \mathbf{W}_{h,n} \mathbf{q}_t \quad (12)$$

$$\mathbf{p}_t = \mathbf{W}_{h,m} S \odot \mathbf{W}_{h,n} \mathbf{p}_t \quad (13)$$

$$\mathbf{i}_t = \sigma(\mathbf{W}_{pi} \mathbf{p}_t + \mathbf{W}_{qi} \mathbf{q}_t + \mathbf{W}_{Ti} \mathbf{T}_t + \mathbf{b}_i) \quad (14)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_{pf} \mathbf{p}_t + \mathbf{W}_{qf} \mathbf{q}_t + \mathbf{W}_{Tf} \mathbf{T}_t + \mathbf{b}_f) \quad (15)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_{po}\mathbf{p}_t + \mathbf{W}_{qo}\mathbf{q}_t + \mathbf{W}_{To}\mathbf{T}_t + \mathbf{b}_o) \quad (16)$$

$$\mathbf{g}_t = \tanh(\mathbf{W}_{pg}\mathbf{p}_t + \mathbf{W}_{qg}\mathbf{q}_t + \mathbf{W}_{Tg}\mathbf{T}_t + \mathbf{b}_g) \quad (17)$$

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{g}_t \quad (18)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t) \quad (19)$$

$$\mathbf{x}_t = \max \arg \text{softmax}(\mathbf{W}_{hx}\mathbf{h}_t) \quad (20)$$

Mathematically, Early Attribute Interaction-Tensor Product Representation (aiTPR), denoted as  $f_{aiTPRE}(\cdot)$ , can be described as the followings probability distribution estimation.

$$\begin{aligned} f_{aiTPRE}(\mathbf{v}) &= \prod_k \Pr(w_k | \mathbf{T}_i, \mathbf{v}, \mathbf{W}_{L_1}) \\ &\prod_i \Pr(\mathbf{T}_i | \mathbf{v}, \mathbf{W}_1) \\ &= \prod_k \Pr(w_k | \mathbf{T}_i, \left(\frac{1}{K} \sum_{m=1}^K v_m\right), \frac{\left(\sum_{m=1}^{N_1} a_m v_m\right)}{2} + \\ &\quad \frac{\left(\sum_{m=1}^{N_2} a'_m v'_m\right)}{2}, \mathbf{W}_{L_1}) \prod_i \Pr(\mathbf{T}_i | \mathbf{v}, \mathbf{W}_1) \\ &= \prod_k Q_{IC}(w_k | \mathbf{T}_i, \left(\frac{1}{K} \sum_{m=1}^K v_m\right), \frac{\left(\sum_{m=1}^{N_1} a_m v_m\right)}{2} + \\ &\quad \frac{\left(\sum_{m=1}^{N_2} a'_m v'_m\right)}{2}) \prod_i Q(\mathbf{T}_i | \mathbf{v}) \end{aligned} \quad (21)$$

using the weights of the LSTM in the architecture is denoted as  $\mathbf{W}_{L_1}$ ,  $w_i$  as words of sentences,  $v_i$  as regional image features,  $a_m s_m$  as intermediate learnt parameters,  $Q_{IC}(\cdot)$  and  $Q(\cdot)$  are the Image Caption and Scene-Graph generator function respectively.  $Q(\cdot)$  derives  $\mathbf{x}$  (Scene-Graph information) from  $\mathbf{v}$  of  $\mathbf{I}$ . These set of equations worked best for this kinds of situation and experimented it for understanding the effects of our new concept of Attribute Interaction-Tensor Product Representation (aiTPR). Apart from that, this network architecture helps in fusion of different set of feature tensors without the scope of influence or suppression and compared to architectures like [2], it is much lighter in the number of weights.

### B. Late Attribute Interaction-Tensor Product Representation

Late Attribute Interaction-Tensor Product Representation (aiTPR) kept the attribute and interaction separated so that the both the network learns equivalently before they fused for the language decoder. The Late aiTPR model operates without any expectation of correlation between the attribute set and the interaction set. While, most of the works in attention are much concentrated on automatic segregation of useful information, we propose that the features must segregate in pure form before getting into the series of linear and non-linear approximations. The contexts from regional object features are coupled along, just like tensor product, but instead of word

embedding, some of the attributes are still the image features from a RCNN network. Before, we discuss the Late aiTPR, we revisit the concept of TPR as a general.

1) *TPR*: For a TPR, we have Here, the context is represented as  $f_i \in \mathbb{R}^{d_n}$  and the topological information vector as  $r_i \in \mathbb{R}^{d_m}$  creating the summation of orthogonal vectors as  $\sum f_i r_i$  and using the same topological information vector, we get back  $f_j$  as  $f_j = r_j^T \sum f_i r_i = f_1 r_1 r_j^T + f_2 r_2 r_j^T + \dots = f_j r_j r_j^T = f_j$ . here  $f$  is object features and  $r$  is orthogonal vectors where both represent an output in caption. The novelty lies in a direct relationship between object to word generation without explicitly knowing the identity of the words. It is aided by the other learnable weights. TPR can be represented as  $\mathbf{s}(\mathbf{w})$  as,

$$\mathbf{s}(\mathbf{w}) = \sum f_i \otimes r_i^T \quad (22)$$

where  $\mathbf{w}$  is the feature vector, and  $\{\mathbf{w} \rightarrow \mathbf{f} : \mathbf{w} \in \mathbf{W}_e\}$  is the transformation,  $\mathbf{W}_e$  are the raw features or the embedding vectors for features which minimizes the context function

Mathematically, the equations that will describe the architecture can be introduced as the followings,

$$\bar{v} = \frac{1}{k} \sum_{i=1}^{i=k} v_i \quad (23)$$

$$\bar{v} = \mathbf{v} \quad (24)$$

The initial parameters are initialized as the followings.

$$\mathbf{h}_0, \mathbf{c}_0 = \mathbf{W}_{h_0} \bar{v}, \mathbf{W}_{c_0} \bar{v} \quad (25)$$

where we have  $\mathbf{W}_{h_0} \in \mathbb{R}^{2048 \times d}$ ,  $\mathbf{W}_{c_0} \in \mathbb{R}^{2048 \times d}$ . The next segment equations can be denoted as the following,

$$\mathbf{a}_t = \mathbf{W}_a \tanh(\mathbf{W}_h \mathbf{h}_{t-1}) \quad (26)$$

where  $\mathbf{W}_a \in \mathbb{R}^{b \times d}$ ,  $\mathbf{W}_h \in \mathbb{R}^{d \times k}$ .

$$\alpha_t = \text{softmax}(\mathbf{a}_t) \quad (27)$$

$$\mathbf{a}_t \in \mathbb{R}^k \in \{a_{1,t}, \dots, a_{k,t}\}$$

$$\hat{v}_t = \left[ \sum_{i=1}^{i=k_1} v_i \alpha_{i,t} + \sum_{i=1}^{i=k_2} v'_i \alpha'_{i,t} \right] \quad (28)$$

$\sum \alpha_i = 1$ ,  $\hat{v}_t \in \mathbb{R}^{b \times d}$  where  $b$  is the batch size and  $d$  is the hidden layer dimension.

$$\mathbf{q}_t = \hat{v}_t \quad (29)$$

$$\mathbf{p}_t = \mathbf{W}_e \mathbf{x}_{t-1} \quad (30)$$

$$\begin{aligned} \mathbf{T}_t &= \mathbf{W}_{s_{12}} \sigma(\mathbf{W}_{s_{11}} \mathbf{h}_{t-1} + \mathbf{b}_1) \\ &\otimes \tanh(\mathbf{W}_{s_{22}} (\mathbf{v} \sigma(\mathbf{W}_{s_{21}} \mathbf{h}_{t-1} + \mathbf{b}_2)) + \mathbf{b}_3) \end{aligned} \quad (31)$$

$$\begin{aligned} \mathbf{T}_t &= \mathbf{W}_{s_{12}} \sigma(\mathbf{W}_{s_{11}} \mathbf{h}_{t-1} + \mathbf{W}_{w_1} \sum_{i=0}^{t-1} \mathbf{W}_e \mathbf{x}_i + \mathbf{b}_1) \otimes \\ &\tanh(\mathbf{W}_{s_{22}} (\mathbf{v} \sigma(\mathbf{W}_{s_{21}} \mathbf{h}_{t-1} + \mathbf{W}_{w_2} \sum_{i=0}^{t-1} \mathbf{W}_e \mathbf{x}_i + \mathbf{b}_2)) + \mathbf{b}_3) \end{aligned} \quad (32)$$

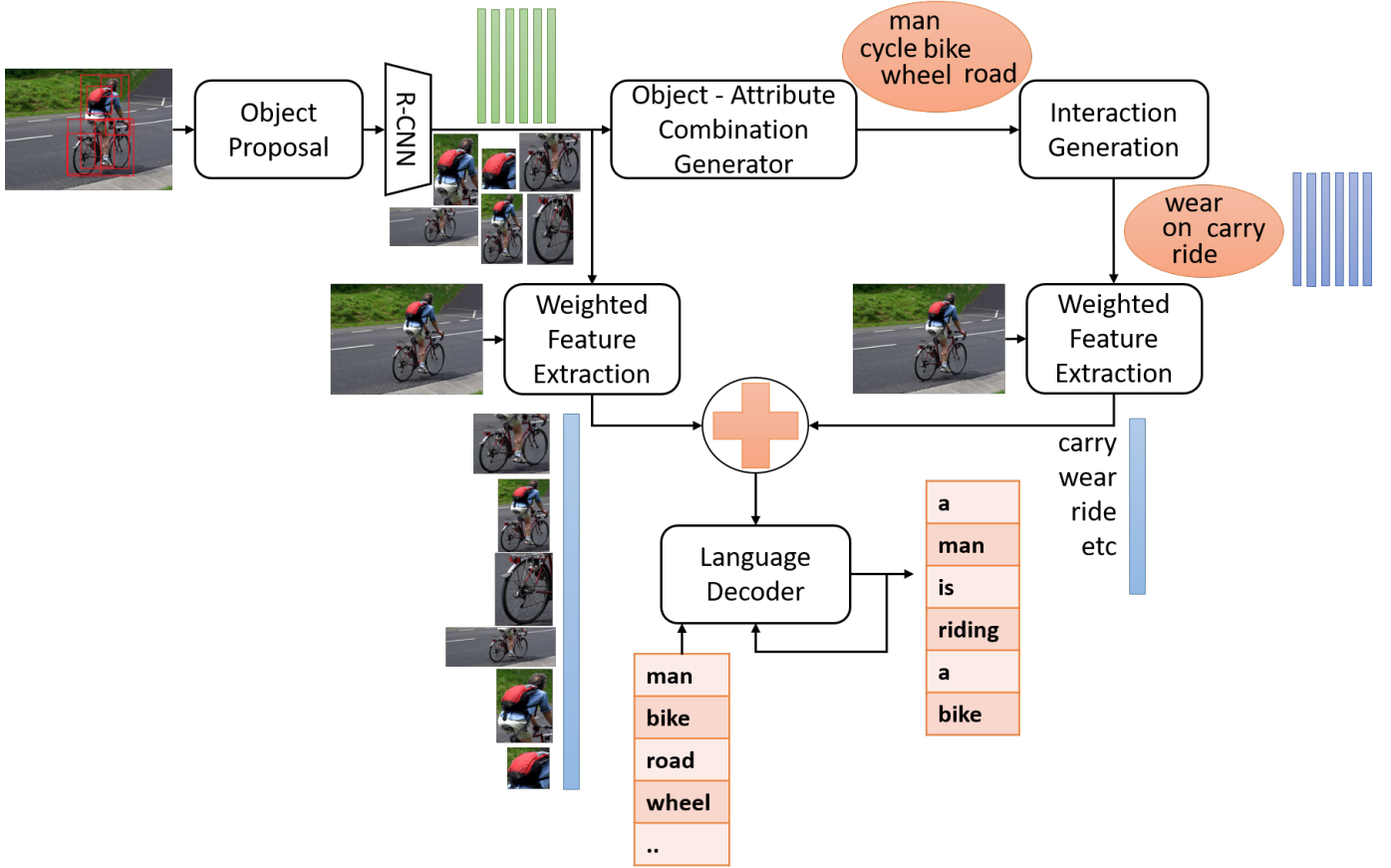


Fig. 2. Early Attribute Interaction-Tensor Product Representation (aiTPR).

$\otimes$  is an algebraic operation. Here, we considered  $\otimes = \odot$  as we try to rectify one context with the other context.

$$\mathbf{q}_t = \mathbf{W}_{h,m} S \odot \mathbf{W}_{h,n} \mathbf{q}_t \quad (33)$$

$$\mathbf{p}_t = \mathbf{W}_{h,m} S \odot \mathbf{W}_{h,n} \mathbf{p}_t \quad (34)$$

$$\mathbf{i}_t = \sigma(\mathbf{W}_{pi} \mathbf{p}_t + \mathbf{W}_{qi} \mathbf{q}_t + \mathbf{W}_{Ti} \mathbf{T}_t + \mathbf{b}_i) \quad (35)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_{pf} \mathbf{p}_t + \mathbf{W}_{qf} \mathbf{q}_t + \mathbf{W}_{Tf} \mathbf{T}_t + \mathbf{b}_f) \quad (36)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_{po} \mathbf{p}_t + \mathbf{W}_{qo} \mathbf{q}_t + \mathbf{W}_{To} \mathbf{T}_t + \mathbf{b}_o) \quad (37)$$

$$\mathbf{g}_t = \tanh(\mathbf{W}_{pg} \mathbf{p}_t + \mathbf{W}_{qg} \mathbf{q}_t + \mathbf{W}_{Tg} \mathbf{T}_t + \mathbf{b}_g) \quad (38)$$

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{g}_t \quad (39)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t) \quad (40)$$

$$\mathbf{x}_t = \max \arg \text{softmax}(\mathbf{W}_{hx} \mathbf{h}_t) \quad (41)$$

Mathematically, Late Attribute Interaction-Tensor Product Representation (aiTPR), denoted as  $f_{aiTPR_L}(\cdot)$ , can be described as the followings probability distribution estimation.

$$\begin{aligned} f_{aiTPR_L}(\mathbf{v}) &= \prod_k \Pr(w_k | \mathbf{T}_i, \mathbf{v}, \mathbf{W}_{L_1}) \\ &\prod_i \Pr(\mathbf{T}_i | \mathbf{v}, \mathbf{W}_1) \\ &= \prod_k \Pr(w_k | \mathbf{T}_i, \left( \frac{1}{K} \sum_{m=1}^K v_m \right), \left( \sum_{m=1}^{N_1} a_m v_m \right) + \\ &\quad \left( \sum_{m=1}^{N_2} a'_m v'_m \right), \mathbf{W}_{L_1}) \prod_i \Pr(\mathbf{T}_i | \mathbf{v}, \mathbf{W}_1) \\ &= \prod_k Q_{IC}(w_k | \mathbf{T}_i, \left( \frac{1}{K} \sum_{m=1}^K v_m \right), \left( \sum_{m=1}^{N_1} a_m v_m \right) + \\ &\quad \left( \sum_{m=1}^{N_2} a'_m v'_m \right)) \prod_i Q(\mathbf{T}_i | \mathbf{v}) \end{aligned} \quad (42)$$

using the weights of the LSTM in the architecture is denoted as  $\mathbf{W}_{L_1}$ ,  $w_i$  as words of sentences,  $v_i$  as regional image features,  $a_m s_m$  as intermediate learnt parameters,  $Q_{IC}(\cdot)$  and  $Q(\cdot)$  are the Image Caption and TPR generator function respectively.

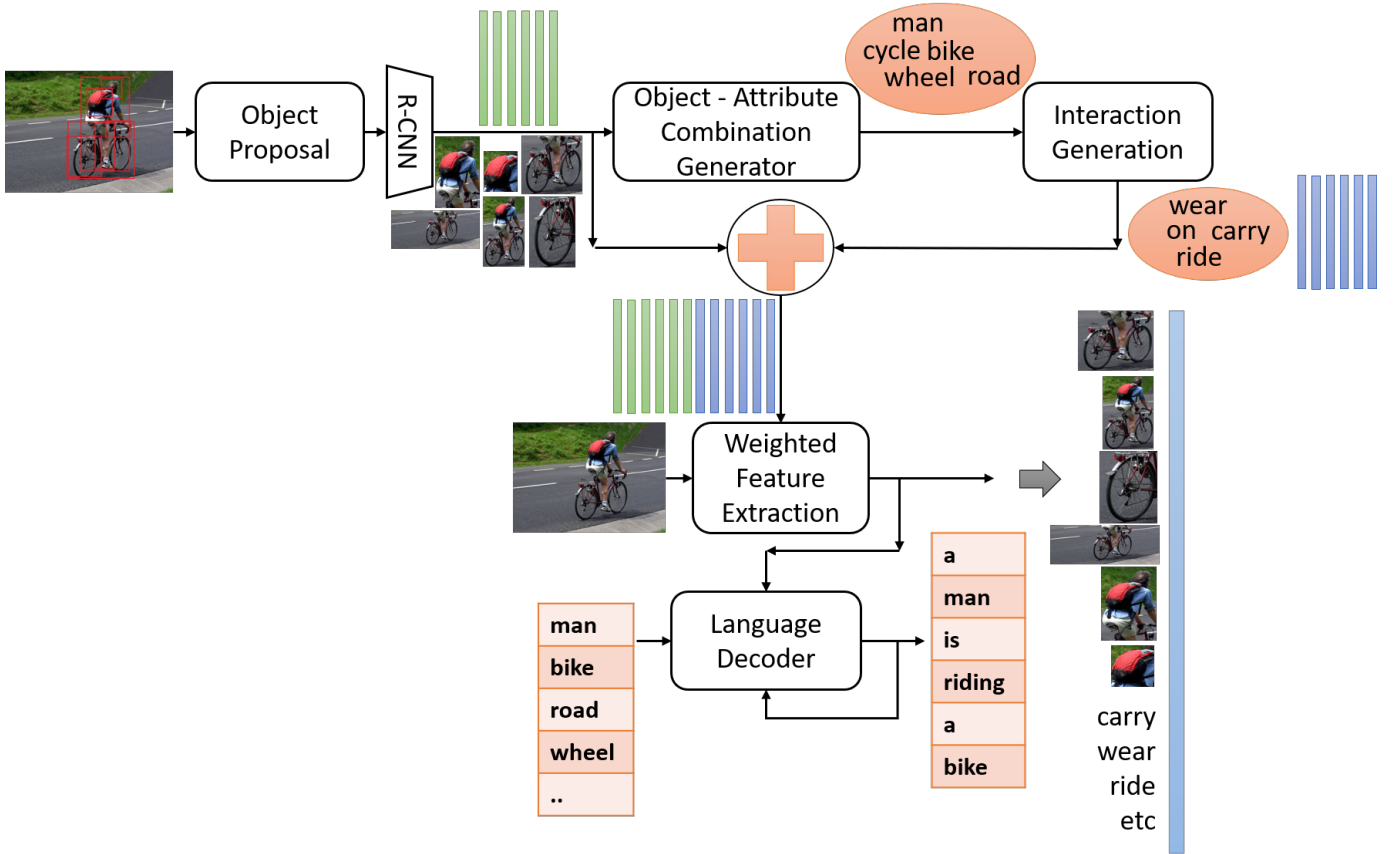


Fig. 3. Late Attribute Interaction-Tensor Product Representation (aiTPR).

### C. Uniqueness of aiTPR

The most components and uniqueness of our architecture is the introduction and fusion of regional image features and several other abstracted likelihood embedding of interaction terms. Unlike other works, where they have used image feature, a different kind of constructed features can bring stability in the variance of the features and can restrict the null spaces. While, many works have demonstrated regional image features based networks, they have left out the scope of introducing the fusion points of these regional influences. In this work, we have introduced such a scheme called interactions and have presented in the form of TPR, where we construct two different contexts from the same feature space and use them as product for maximum influence.

## IV. EXPERIMENTS, RESULTS & ANALYSIS

We have done wide range of experiments to show the behavioral influence of object based attention through this architectural network, where we have defined different levels of information for the generation of captions. Before we analyze the results, short description of the dataset is also provided along with the training session description as we explore the joint distribution of the learning state space.

### A. Data Description

We have used the MSCOCO and the Visual Genome dataset for our analysis. MSCOCO consists of 123287 train+validation

images and 566747 train+validation sentence, where each image is associated with at least five sentences from a vocabulary of 8791 words. There are 5000 images (with 25010 sentences) for validation and 5000 images (with 25010 sentences) for testing. We used the same data split as described in Karpathy's paper [47]. Visual Genome dataset is used for other language semantic information for the MSCOCO datasets and a model is trained to derive the annotations for those images. Roughly, 38% of the MSCOCO data has attribute level annotations in the Visual Genome dataset.

### B. Quantitative Evaluation

Several evaluation metrics like CIDEr-D, Bleu\_4, Bleu\_3, Bleu\_2, Bleu\_1, ROUGE\_L, METEOR and SPICE is used for our experiments. Table I provide a quantitative evaluation of our experiments, mainly focusing on the different architectures, related to the context  $\mathbf{h}_{t-1}$  and previous word embedding  $\mathbf{x}_{t-1}$  semantics, which is often found to enhance the performance for the captions. We found that the performance of our model out-performed many previous models and it was also found that the late fusion model provide better performance, inducing more experiments based on the semantic correction. Clearly, Late aiTPR [aiTPR (3)] with Equation 33 was the winner interms of most of the evaluation metrics, but the other schemes like [aiTPR (1)] and [aiTPR (2)] performed in a competitive way. The main reason, [aiTPR (3)] worked better was because of the influence it produced on the combinations

of regional feature set and their interactions, which ultimately helped in better captions, while Equation 34 helped in the topological continuity of the word for the generated sentence.

### C. Discussion

The main improvement that our models had put into the architecture is the introduction of aiTPR, which is characterized by understanding the interaction level information based representation. This was never tried before and introduction of regional influences attended refined feature levels and provided ample scope of a fitted structure with improvement of the representations with iterations. While, previous works were concentrated on defining better image feature quality, we have paid more importance on the inference level information that generalizes the representations, but create combination level enhancements. While, most of the works were just what the model has learned, we paid more importance what we can create and feed into the network like shown in Equation 21 and Equation 42. With this approach, we have establish a new performance level, which has surpassed other previous works in all the metrics. We have used a RCNN network to find the regional level details along with the coordinates of the regions and used them for interaction level inference without much concerning about the correctness as we are more interested in the defined representation than the inference level likelihood tensors.

### D. Qualitative Evaluation

Normally statistical formulas are the best evaluation of the significance levels of any model, except language ones. This is because the statistical models most concentrate on the content and the numerals associated with it. To evaluate language structures, we need perception, which is also diverse and subject to high variations. Also, whether a model is better than the other cannot be judged through average numerals. Whether overall improved captions are generated or not is also difficult to judge from numerals in Table I, Hence, we have considered some of the generated sentences as our qualitative analysis and will reflect the supremacy of our novel architecture. Figure 4 and Figure 5 provided some of the instances that were derived from the models. The examples provided very good representations of the generated sentences from the corresponding images.

## V. DISCUSSION

While the previous works mainly concentrated on the features and their combinations in an un-thoughtful way, this work produces a technique where you can derive useful interactions for the attributes and generate the most useful tensors and their products, without the requirement for non-linear approximation. We introduced, for the first time, the notion of combining regional image features and abstracted interaction likelihood embedding for image captioning. We call this model as Attribute Interaction-Tensor Product Representation (aiTPR) as this is an good AI technique to consider the attribute-interaction (ai) where the attributes are structured

ones as the lower level features transferred from any pre-trained model, while the interactions are the composed derived from them and then approximated through a language representation that is well fitted with the image feature model. With this work, we have derived a new state-of-the-art result and also a novel work.

### ACKNOWLEDGMENT

The author has used University of Florida HiperGator, equipped with NVIDIA Tesla K80 GPU, extensively for the experiments. The author acknowledges University of Florida Research Computing for providing computational resources and support that have contributed to the research results reported in this publication. URL: <http://researchcomputing.ufl.edu>

### CONFLICT OF INTEREST

The author declares that he has no conflict of interest.

### REFERENCES

- [1] Lu, D., Whitehead, S., Huang, L., Ji, H., & Chang, S. F. (2018). Entity-aware Image Caption Generation. arXiv preprint arXiv:1804.07889.
- [2] Lu, J., Yang, J., Batra, D., & Parikh, D. (2018, March). Neural Baby Talk. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 7219-7228).
- [3] You, Q., Jin, H., & Luo, J. (2018). Image Captioning at Will: A Versatile Scheme for Effectively Injecting Sentiments into Image Descriptions. arXiv preprint arXiv:1801.10121.
- [4] Melnyk, I., Sercu, T., Dognin, P. L., Ross, J., & Mroueh, Y. (2018). Improved Image Captioning with Adversarial Semantic Alignment. arXiv preprint arXiv:1805.00063.
- [5] Wu, J., Hu, Z., & Mooney, R. J. (2018). Joint Image Captioning and Question Answering. arXiv preprint arXiv:1805.08389.
- [6] Kilickaya, M., Akkus, B. K., Cakici, R., Erdem, A., Erdem, E., & Ikizler-Cinbis, N. (2017). Data-driven image captioning via salient region discovery. IET Computer Vision, 11(6), 398-406.
- [7] Chen, F., Ji, R., Su, J., Wu, Y., & Wu, Y. (2017, October). Structcap: Structured semantic embedding for image captioning. In Proceedings of the 2017 ACM on Multimedia Conference (pp. 46-54). ACM.
- [8] Jiang, W., Ma, L., Chen, X., Zhang, H., & Liu, W. (2018). Learning to guide decoding for image captioning. arXiv preprint arXiv:1804.00887.
- [9] Wu, C., Wei, Y., Chu, X., Su, F., & Wang, L. (2018). Modeling visual and word-conditional semantic attention for image captioning. Signal Processing: Image Communication.
- [10] Fu, K., Li, J., Jin, J., & Zhang, C. (2018). Image-Text Surgery: Efficient Concept Learning in Image Captioning by Generating Pseudopairs. IEEE Transactions on Neural Networks and Learning Systems, (99), 1-12.
- [11] Chen, H., Ding, G., Lin, Z., Zhao, S., & Han, J. (2018). Show, Observe and Tell: Attribute-driven Attention Model for Image Captioning. In IJCAI (pp. 606-612).
- [12] Cornia, M., Baraldi, L., Serra, G., & Cucchiara, R. (2018). Paying More Attention to Saliency: Image Captioning with Saliency and Context Attention. ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM), 14(2), 48.
- [13] Zhao, W., Wang, B., Ye, J., Yang, M., Zhao, Z., Luo, R., & Qiao, Y. (2018). A Multi-task Learning Approach for Image Captioning. In IJCAI (pp. 1205-1211).
- [14] Li, X., Wang, X., Xu, C., Lan, W., Wei, Q., Yang, G., & Xu, J. (2018). COCO-CN for Cross-Lingual Image Tagging, Captioning and Retrieval. arXiv preprint arXiv:1805.08661.
- [15] Chen, M., Ding, G., Zhao, S., Chen, H., Liu, Q., & Han, J. (2017, February). Reference Based LSTM for Image Captioning. In AAAI (pp. 3981-3987).
- [16] Tavakoliy, H. R., Shetty, R., Borji, A., & Laaksonen, J. (2017, October). Paying attention to descriptions generated by image captioning models. In Computer Vision (ICCV), 2017 IEEE International Conference on (pp. 2506-2515). IEEE.



TABLE I  
PERFORMANCE EVALUATION FOR DIFFERENT ARCHITECTURES WITHOUT REINFORCEMENT LEARNING

Algorithm	CIDEr-D	Bleu_4	Bleu_3	Bleu_2	Bleu_1	ROUGE_L	METEOR	SPICE
Human [42]	0.85	0.22	0.32	0.47	0.66	0.48	0.2	–
Neural Talk [47]	0.66	0.23	0.32	0.45	0.63	–	0.20	–
MindsEye [49]	–	0.19	–	–	–	–	0.20	–
Google [43]	0.94	0.31	0.41	0.54	0.71	0.53	0.25	–
LRCN [51]	0.87	0.28	0.38	0.53	0.70	0.52	0.24	–
Montreal [45]	0.87	0.28	0.38	0.53	0.71	0.52	0.24	–
m-RNN [56]	0.79	0.27	0.37	0.51	0.68	0.50	0.23	–
[72]	0.81	0.26	0.36	0.49	0.67	–	0.23	–
MSR [46]	0.91	0.29	0.39	0.53	0.70	0.52	0.25	–
[53]	0.84	0.28	0.38	0.52	0.70	–	0.24	–
bi-LSTM [40]	–	0.244	0.352	0.492	0.672	–	–	–
MSR Captivator [50]	0.93	0.31	0.41	0.54	0.72	0.53	0.25	–
Nearest Neighbor [84]	0.89	0.28	0.38	0.52	0.70	0.51	0.24	–
ATT [66]	0.94	0.32	0.42	0.57	0.73	0.54	0.25	–
[42]	0.92	0.31	0.41	0.56	0.73	0.53	0.25	–
Adaptive [35]	1.085	0.332	0.439	0.580	0.742	–	0.266	–
MSM [34]	0.986	0.325	0.429	0.565	0.730	–	0.251	–
ERD [?]	0.895	0.298	–	–	–	–	0.240	–
Att2in [41]	1.01	0.313	–	–	–	–	0.260	–
Top-Down† [37]	1.054	0.334	–	–	0.745	–	0.261	0.192
Attribute-Attention [11]	1.044	0.338	0.443	0.579	0.743	0.549	–	–
LSTM [69]	0.889	0.292	0.390	0.525	0.698	–	0.238	–
SCN [69]	1.012	0.330	0.433	0.566	0.728	–	0.257	–
NBT† [2]	1.07	0.347	–	–	0.755	–	0.271	0.201
Top-Down† [37]	1.135	0.362	–	–	0.772	0.564	0.27	0.203
Early aiTPR with Equation 33	1.387	0.476	0.588	0.718	0.850	0.632	0.318	0.233
Late aiTPR with Equation 33	1.401	0.484	0.595	0.721	0.850	0.637	0.320	0.233
Late aiTPR [aiTPR (1)] without Equation 33 and Equation 34	1.386	0.476	0.586	0.714	0.846	0.631	0.318	0.232
Late aiTPR [aiTPR (2)] with Equation 33 and Equation 34	1.398	0.484	0.593	0.721	0.852	0.635	0.320	0.233
Late aiTPR [aiTPR (3)] with Equation 33	<b>1.401</b>	<b>0.484</b>	<b>0.595</b>	<b>0.721</b>	<b>0.850</b>	<b>0.637</b>	<b>0.320</b>	<b>0.233</b>

†Ensemble & Reinforcement Learning Used with RCNN features and [37] Used 10K Vocabulary Dataset

- [17] Chen, H., Zhang, H., Chen, P. Y., Yi, J., & Hsieh, C. J. (2017). Show-and-fool: Crafting adversarial examples for neural image captioning. arXiv preprint arXiv:1712.02051.
- [18] Ye, S., Liu, N., & Han, J. (2018). Attentive Linear Transformation for Image Captioning. IEEE Transactions on Image Processing.
- [19] Wang, Y., Lin, Z., Shen, X., Cohen, S., & Cottrell, G. W. (2017). Skeleton key: Image captioning by skeleton-attribute decomposition. arXiv preprint arXiv:1704.06972.
- [20] Chen, T., Zhang, Z., You, Q., Fang, C., Wang, Z., Jin, H., & Luo, J. (2018). "Factual" or "Emotional": Stylized Image Captioning with Adaptive Learning and Attention. arXiv preprint arXiv:1807.03871.
- [21] Chen, F., Ji, R., Sun, X., Wu, Y., & Su, J. (2018). GroupCap: Group-Based Image Captioning With Structured Relevance and Diversity Constraints. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 1345-1353).
- [22] Liu, C., Sun, F., Wang, C., Wang, F., & Yuille, A. (2017). MAT: A multimodal attentive translator for image captioning. arXiv preprint arXiv:1702.05658.
- [23] Harzig, P., Brehm, S., Lienhart, R., Kaiser, C., & Schallner, R. (2018). Multimodal Image Captioning for Marketing Analysis. arXiv preprint arXiv:1802.01958.
- [24] Liu, X., Li, H., Shao, J., Chen, D., & Wang, X. (2018). Show, Tell and Discriminate: Image Captioning by Self-retrieval with Partially Labeled Data. arXiv preprint arXiv:1803.08314.
- [25] Chunseong Park, C., Kim, B., & Kim, G. (2017). Attend to you: Personalized image captioning with context sequence memory networks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 895-903).
- [26] Sharma, P., Ding, N., Goodman, S., & Soricut, R. (2018). Conceptual Captions: A Cleaned, Hypernymed, Image Alt-text Dataset For Automatic Image Captioning. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers) (Vol. 1, pp. 2556-2565).
- [27] Yao, T., Pan, Y., Li, Y., & Mei, T. (2017, July). Incorporating copying mechanism in image captioning for learning novel objects. In 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (pp. 5263-5271). IEEE.
- [28] Zhang, L., Sung, F., Liu, F., Xiang, T., Gong, S., Yang, Y., & Hospedales, T. M. (2017). Actor-critic sequence training for image captioning. arXiv preprint arXiv:1706.09601.
- [29] Fu, K., Jin, J., Cui, R., Sha, F., & Zhang, C. (2017). Aligning where to see and what to tell: Image captioning with region-based attention and scene-specific contexts. IEEE transactions on pattern analysis and machine intelligence, 39(12), 2321-2334.
- [30] Ren, Z., Wang, X., Zhang, N., Lv, X., & Li, L. J. (2017). Deep reinforcement learning-based image captioning with embedding reward. arXiv preprint arXiv:1704.03899.
- [31] Liu, S., Zhu, Z., Ye, N., Guadarrama, S., & Murphy, K. (2017, October). Improved image captioning via policy gradient optimization of spider. In Proc. IEEE Int. Conf. Comp. Vis (Vol. 3, p. 3).
- [32] Cohn-Gordon, R., Goodman, N., & Potts, C. (2018). Pragmatically Informative Image Captioning with Character-Level Reference. arXiv preprint arXiv:1804.05417.
- [33] Liu, C., Mao, J., Sha, F., & Yuille, A. L. (2017, February). Attention Correctness in Neural Image Captioning. In AAAI (pp. 4176-4182).
- [34] Yao, T., Pan, Y., Li, Y., Qiu, Z., & Mei, T. (2017, October). Boosting image captioning with attributes. In IEEE International Conference on Computer Vision, ICCV (pp. 22-29).
- [35] Lu, J., Xiong, C., Parikh, D., & Socher, R. (2017, July). Knowing when to look: Adaptive attention via a visual sentinel for image captioning. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (Vol. 6, p. 2).
- [36] Vinyals, O., Toshev, A., Bengio, S., & Erhan, D. (2017). Show and tell: Lessons learned from the 2015 mscoco image captioning challenge. IEEE transactions on pattern analysis and machine intelligence, 39(4), 652-663.
- [37] Anderson, P., He, X., Buehler, C., Teney, D., Johnson, M., Gould, S., & Zhang, L. (2018). Bottom-up and top-down attention for image captioning and visual question answering. In CVPR (Vol. 3, No. 5, p. 6).
- [38] Zhang, M., Yang, Y., Zhang, H., Ji, Y., Shen, H. T., & Chua, T. S. (2018). More is Better: Precise and Detailed Image Captioning using



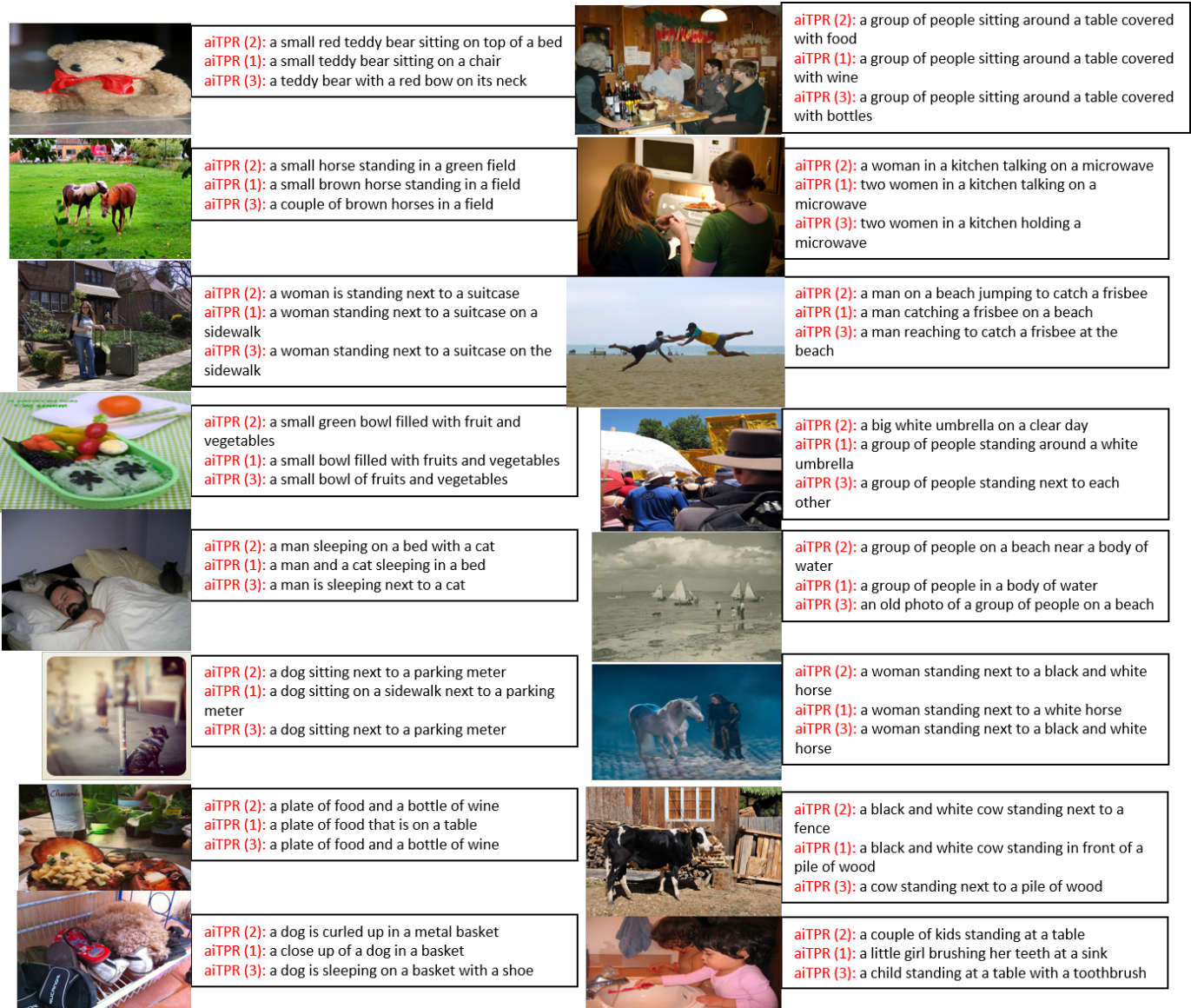


Fig. 4. Qualitative Analysis. Part 1. [aiTPR (1)], [aiTPR (2)] and [aiTPR (3)] are defined Table I.

- Online Positive Recall and Missing Concepts Mining. IEEE Transactions on Image Processing.
- [39] Park, C. C., Kim, B., & Kim, G. (2018). Towards Personalized Image Captioning via Multimodal Memory Networks. IEEE Transactions on Pattern Analysis and Machine Intelligence.
- [40] Wang, Cheng, Haojin Yang, and Christoph Meinel. "Image Captioning with Deep Bidirectional LSTMs and Multi-Task Learning." ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM) 14.2s (2018): 40.
- [41] Rennie, S. J., Marcheret, E., Mroueh, Y., Ross, J., & Goel, V. (2017, July). Self-critical sequence training for image captioning. In CVPR (Vol. 1, No. 2, p. 3).
- [42] Wu, Q., Shen, C., Wang, P., Dick, A., & van den Hengel, A. (2017). Image captioning and visual question answering based on attributes and external knowledge. IEEE transactions on pattern analysis and machine intelligence.
- [43] Vinyals, Oriol, et al. "Show and tell: A neural image caption generator." Proceedings of the IEEE conference on computer vision and pattern recognition. 2015.
- [44] Karpathy, Andrej, Armand Joulin, and Fei Fei Li. "Deep fragment embeddings for bidirectional image sentence mapping." Advances in neural information processing systems. 2014.
- [45] Xu, Kelvin, et al. "Show, attend and tell: Neural image caption generation with visual attention." International Conference on Machine Learning. 2015.
- [46] Fang, Hao, et al. "From captions to visual concepts and back." Proceedings of the IEEE conference on computer vision and pattern recognition. 2015.
- [47] Karpathy, Andrej, and Li Fei-Fei. "Deep visual-semantic alignments for generating image descriptions." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2015.
- [48] Anne Hendricks, Lisa, et al. "Deep compositional captioning: Describing novel object categories without paired training data." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2016.
- [49] Chen, Xinlei, and C. Lawrence Zitnick. "Mind's eye: A recurrent visual representation for image caption generation." Proceedings of the IEEE conference on computer vision and pattern recognition. 2015.
- [50] Devlin, Jacob, et al. "Language models for image captioning: The quirks and what works." arXiv preprint arXiv:1505.01809 (2015).
- [51] Donahue, Jeffrey, et al. "Long-term recurrent convolutional networks for visual recognition and description." Proceedings of the IEEE conference on computer vision and pattern recognition. 2015.
- [52] Gan, Chuang, et al. "StyleNet: Generating attractive visual captions with styles." CVPR, 2017.
- [53] Jin, Junqi, et al. "Aligning where to see and what to tell: image caption with region-based attention and scene factorization." arXiv preprint

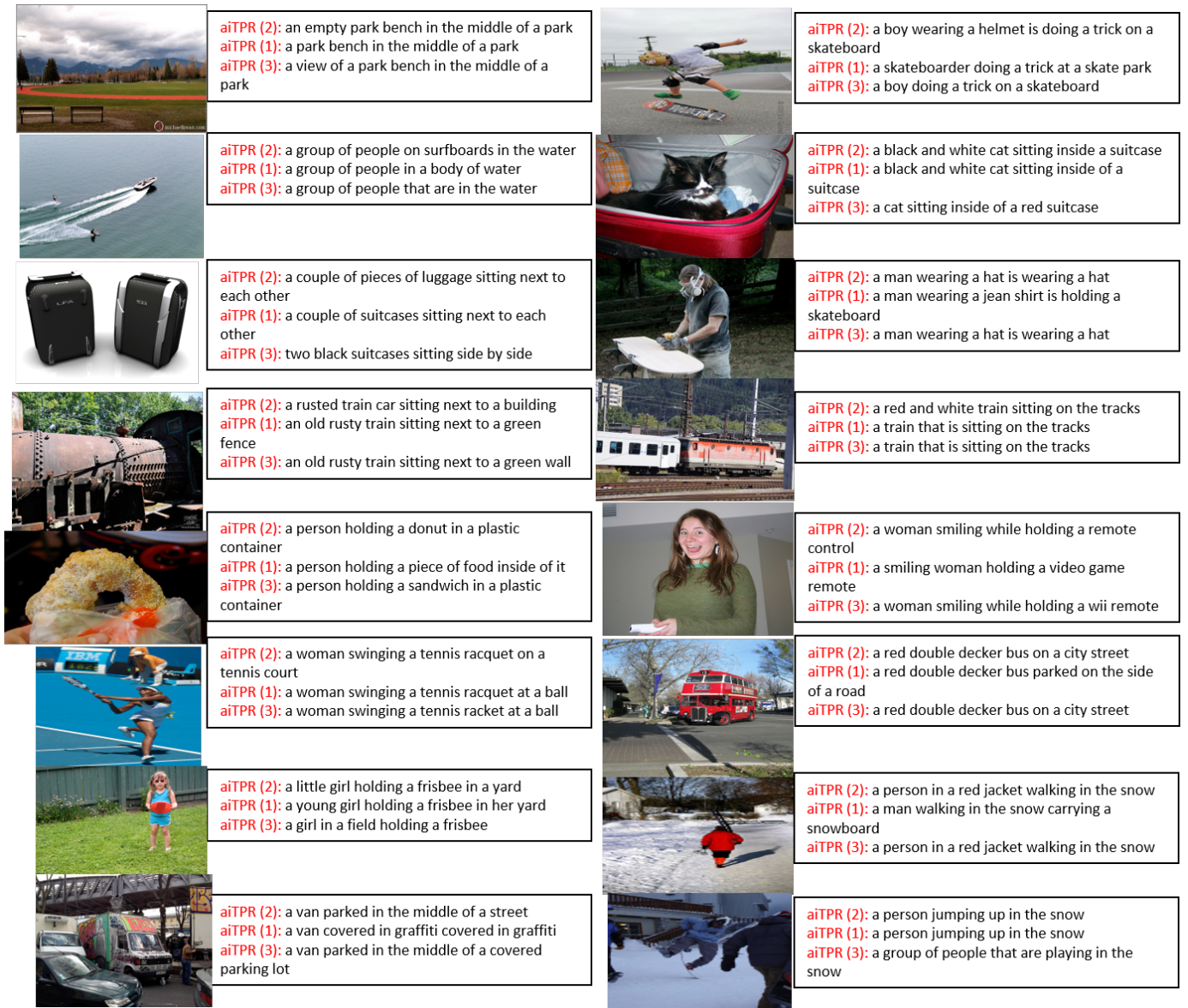


Fig. 5. Qualitative Analysis. Part 2. [aiTPR (1)], [aiTPR (2)] and [aiTPR (3)] are defined Table I.

- arXiv:1506.06272 (2015).
- [54] Kiros, Ryan, Ruslan Salakhutdinov, and Richard S. Zemel. "Unifying visual-semantic embeddings with multimodal neural language models." arXiv preprint arXiv:1411.2539 (2014).
  - [55] Kiros, Ryan, Richard Zemel, and Ruslan R. Salakhutdinov. "A multiplicative model for learning distributed text-based attribute representations." Advances in neural information processing systems. 2014.
  - [56] Mao, Junhua, et al. "Deep captioning with multimodal recurrent neural networks (m-rnn)." arXiv preprint arXiv:1412.6632 (2014).
  - [57] Memisevic, Roland, and Geoffrey Hinton. "Unsupervised learning of image transformations." Computer Vision and Pattern Recognition, 2007. CVPR'07. IEEE Conference on. IEEE, 2007.
  - [58] Pu, Yunchen, et al. "Variational autoencoder for deep learning of images, labels and captions." Advances in Neural Information Processing Systems. 2016.
  - [59] Socher, Richard, et al. "Grounded compositional semantics for finding and describing images with sentences." Transactions of the Association for Computational Linguistics 2 (2014): 207-218.
  - [60] Sutskever, Ilya, James Martens, and Geoffrey E. Hinton. "Generating text with recurrent neural networks." Proceedings of the 28th International Conference on Machine Learning (ICML-11). 2011.
  - [61] Sutskever, Ilya, Oriol Vinyals, and Quoc V. Le. "Sequence to sequence learning with neural networks." Advances in neural information processing systems. 2014.
  - [62] LTran, Du, et al. "Learning spatiotemporal features with 3d convolutional networks." Proceedings of the IEEE international conference on computer vision. 2015.
  - [63] Tran, Kenneth, et al. "Rich image captioning in the wild." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops. 2016.
  - [64] Wu, Qi, et al. "What value do explicit high level concepts have in vision to language problems?." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2016.
  - [65] Yang, Zhilin, et al. "Review networks for caption generation." Advances in Neural Information Processing Systems. 2016.
  - [66] You, Qunzeng, et al. "Image captioning with semantic attention." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2016.
  - [67] Young, Peter, et al. "From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions." Transactions of the Association for Computational Linguistics 2 (2014): 67-78.
  - [68] Farhadi, Ali, et al. "Every picture tells a story: Generating sentences from images." European conference on computer vision. Springer, Berlin, Heidelberg, 2010.

- [69] Gan, Zhe, et al. "Semantic compositional networks for visual captioning." arXiv preprint arXiv:1611.08002 (2016).
- [70] Girshick, Ross, et al. "Rich feature hierarchies for accurate object detection and semantic segmentation." Proceedings of the IEEE conference on computer vision and pattern recognition. 2014.
- [71] Hodosh, Micah, Peter Young, and Julia Hockenmaier. "Framing image description as a ranking task: Data, models and evaluation metrics." Journal of Artificial Intelligence Research 47 (2013): 853-899.
- [72] Jia, Xu, et al. "Guiding the long-short term memory model for image caption generation." Proceedings of the IEEE International Conference on Computer Vision. 2015.
- [73] Krishna, Ranjay, et al. "Visual genome: Connecting language and vision using crowdsourced dense image annotations." International Journal of Computer Vision 123.1 (2017): 32-73.
- [74] Kulkarni, Girish, et al. "Babytalk: Understanding and generating simple image descriptions." IEEE Transactions on Pattern Analysis and Machine Intelligence 35.12 (2013): 2891-2903.
- [75] Li, Siming, et al. "Composing simple image descriptions using web-scale n-grams." Proceedings of the Fifteenth Conference on Computational Natural Language Learning. Association for Computational Linguistics, 2011.
- [76] Kuznetsova, Polina, et al. "TREETALK: Composition and Compression of Trees for Image Descriptions." TACL 2.10 (2014): 351-362.
- [77] Mao, Junhua, et al. "Learning like a child: Fast novel visual concept learning from sentence descriptions of images." Proceedings of the IEEE International Conference on Computer Vision. 2015.
- [78] Mathews, Alexander Patrick, Lexing Xie, and Xuming He. "SentiCap: Generating Image Descriptions with Sentiments." AAAI. 2016.
- [79] Mitchell, Margaret, et al. "Midge: Generating image descriptions from computer vision detections." Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics. Association for Computational Linguistics, 2012.
- [80] Ordonez, Vicente, Girish Kulkarni, and Tamara L. Berg. "Im2text: Describing images using 1 million captioned photographs." Advances in Neural Information Processing Systems. 2011.
- [81] Yang, Yezhou, et al. "Corpus-guided sentence generation of natural images." Proceedings of the Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, 2011.
- [82] Sur, C. (2019). Survey of deep learning and architectures for visual captioningtransitioning between media and natural languages. Multimedia Tools and Applications, 1-51.
- [83] Ren, S., He, K., Girshick, R., & Sun, J. (2015). Faster r-cnn: Towards real-time object detection with region proposal networks. In Advances in neural information processing systems (pp. 91-99).
- [84] Devlin, J., Gupta, S., Girshick, R., Mitchell, M., Zitnick, C. L. (2015). Exploring nearest neighbor approaches for image captioning. arXiv preprint arXiv:1505.04467.
- [85] Sur, Chiranjib. "Representations For Vision Language Intelligence Using Tensor Product Representation." Ph.D dissertation, University of Florida, 2019.
- [86] Sur, C. (2019). Survey of deep learning and architectures for visual captioningtransitioning between media and natural languages. Multimedia Tools and Applications, 1-51.
- [87] Sur, Chiranjib. "UCRLF: unified constrained reinforcement learning framework for phase-aware architectures for autonomous vehicle signaling and trajectory optimization." Evolutionary Intelligence (2019): 1-24.
- [88] Smolensky, Paul. "Tensor product variable binding and the representation of symbolic structures in connectionist systems." Artificial intelligence 46.1-2 (1990): 159-216.