

# Towards Coherent Visual Storytelling with Ordered Image Attention

Tom Braude  
 The Interdisciplinary Center  
 Herzliya, Israel  
 tom.braude@idc.ac.il

Idan Schwartz  
 Technion  
 Haifa, Israel  
 idansc@technion.ac.il

Alexander Schwing  
 University of Illinois at Urbana-Champaign  
 Champaign, IL, US  
 aschwing@illinois.edu

Ariel Shamir  
 The Interdisciplinary Center  
 Herzliya, Israel  
 arik@idc.ac.il

## Abstract

We address the problem of visual storytelling, i.e., generating a story for a given sequence of images. While each sentence of the story should describe a corresponding image, a coherent story also needs to be consistent and relate to both future and past images. To achieve this we develop ordered image attention (OIA). OIA models interactions between the sentence-corresponding image and important regions in other images of the sequence. To highlight the important objects, a message-passing-like algorithm collects representations of those objects in an order-aware manner. To generate the story’s sentences, we then highlight important image attention vectors with an Image-Sentence Attention (ISA). Further, to alleviate common linguistic mistakes like repetitiveness, we introduce an adaptive prior. The obtained results improve the METEOR score on the VIST dataset by 1%. In addition, an extensive human study verifies coherency improvements and shows that OIA and ISA generated stories are more focused, shareable, and image-grounded.

## 1. Introduction

Visual Storytelling (VST) [24, 14] – the task of generating a story based on a sequence of images – goes beyond a basic understanding of visual scenes and can be applied in many real-world scenarios, e.g., to support the visually impaired. Moreover, VST reflects on the creative ability of intelligent systems. Although similar in concept to other cognitive tasks such as image captioning and visual question answering, VST differs as it requires to reason over a sequence of images while simultaneously ensuring coherence across multiple generated sentences. To achieve this, VST methods need to address two major challenges: the first is visual and relates to grounding the story’s text to the images. The second is linguistic and relates to the quality of the story. Both challenges can be described in terms of coherency: the story should be coherent by itself, and coherent with the images.

Prior research on VST started to address the aforementioned challenges. Early works expand captioning [33, 38, 6], focusing sentence generation mainly on the current image [9, 34]. This limits the ability to incorporate complex semantic information, which is necessary for visual reasoning. Prior work also makes limited use of temporal dependence and history, e.g., sentences that have already been generated are not used. Consequently, the output lacks narrative consistency and is prone to linguistic errors such as repetitiveness and incoherence [22]. To mitigate these issues, later works strive to generate more meaningful stories via adversarial and reinforcement learning [36, 13], which remain delicate to train.

Importantly, images are not independent. For example, if the first image in a sequence shows a protest, the model may want to focus on signs in later images. Conversely, if the last image shows a ring on a finger, then the model should pay attention to wedding-related objects and activities in the preceding images. This is important for VST because sentences are created per image but are part of a story. Hence, objects that the model is focusing on in one image should be conditioned on the selection in other images.

To do this we develop a novel model which (1) implicitly reasons over objects, activities, and their temporal dependencies in each image; and which (2) improves the coherency of the narrative. To reason over objects and activities in each image, i.e., to understand their dependencies

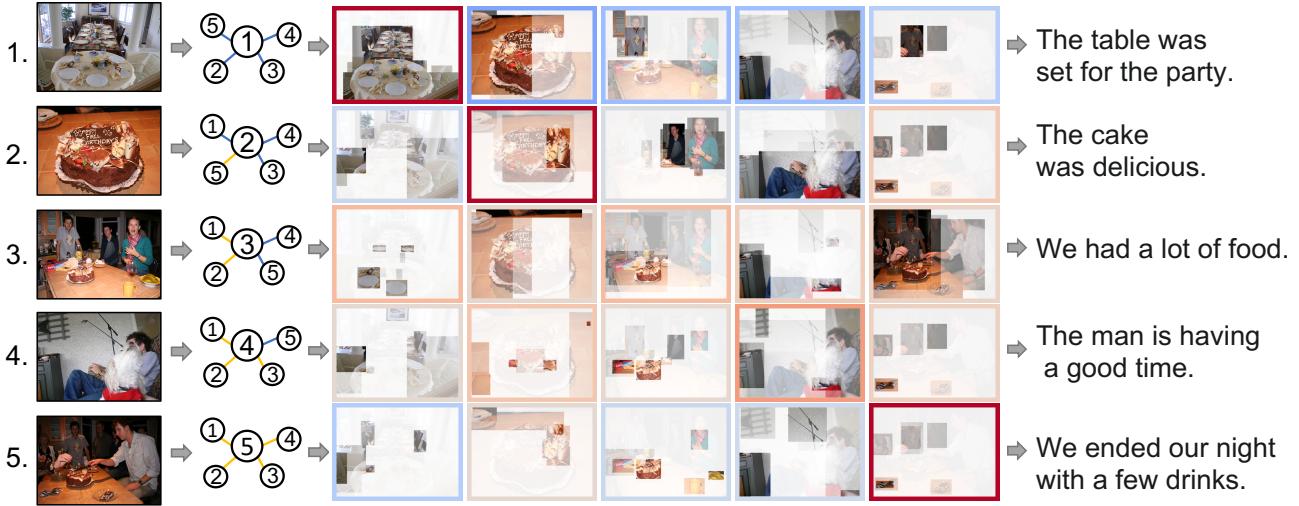


Figure 1: We propose Ordered Image Attention (OIA) to form the structure of a sentence and to encourage coherency. Each row shows the spatial attention of the five images created when generating a specific sentence. We find important objects by collecting directional interactions. The relative order to the sentence-corresponding image determines the connection type, illustrated as the blue and orange edges for preceding and proceeding connections. The attended images’ border indicates the image attention importance formed by the Image-Sentence Attention (ISA). *E.g.*, red indicates a high attention score, meaning the image is essential for generating that sentence. Our model performs this step for all five images in parallel, creating a total of 25 spatial attention maps, that are fed into the decoder to create the sentences in order.

and their temporal ordering, we introduce *ordered image attention* (OIA). As illustrated in Fig. 1, for each image, OIA accumulates representation information from objects detected within the corresponding image into an attended image representation. Importantly, accumulation factors depend on whether the image precedes or succeeds the image for which we are currently generating the sentence, which permits to establish an order. The attended image representations are subsequently summarized into a context embedding via an Image-Sentence Attention (ISA) unit, before being used for sentence decoding.

In addition, to alleviate common linguistic mistakes like repetitiveness and to promote coherence in the story, we incorporate information from the story generated up to the current sentence into the sentence generation decoder. Specifically, the decoding strategy decays the probability of a word if it has already been used in the story. The decoder also maintains a separate prior over the output probability distribution, independent from the language generation unit. This prior is based on counts of the words that were already predicted in the story. Both the prior, and the Recurrent Neural Net (RNN) decoder output are combined to predict the next word in the sentence.

Empirical results on the challenging VIST dataset [14] demonstrate that the proposed method generates stories with an improved narrative quality. The method outperforms prior state-of-the-art by 1% on the METEOR score. Examples of stories generated by the approach are shown in Fig. 1. We also present a user study demonstrating the advantage of the model in terms of coherency.

## 2. Related Work

Vision+Language has been an active area of research for many years, addressing tasks such as image/video captioning, paragraph generation, and visual question answering. We briefly review those related areas in the following.

### 2.1. Image Captioning

Bernard *et al.* [3] first explored annotating images with text. Since then, image/video captioning has seen a surge of research activity. Initial work utilized pre-trained image embeddings from a CNN network. The success of attention mechanisms for language translation quickly transferred to image captioning as well [38]. Later work leveraged advances in object detection and proposed a bottom-up/top-down attention approach to attend to specific objects in the image instead of fixed spatial regions [1]. Different from image captioning, for visual storytelling, both story coherency and visual grounding are important.

### 2.2. Multimodal Attention

Multimodal problems are characterized by input data that comes from different domains, *e.g.*, visual and linguistic. This raises two challenges: 1) how to model interactions between different domains, and 2) how to manage the large input data. Considering those challenges, attention has been a prominent tool as it models interactions to select the important elements. In early work, Xu *et al.* [38] used interaction-based attention with the image at each caption generation step. This idea was later extended to visual question answering [37]. To imitate multi-step reasoning, Yang

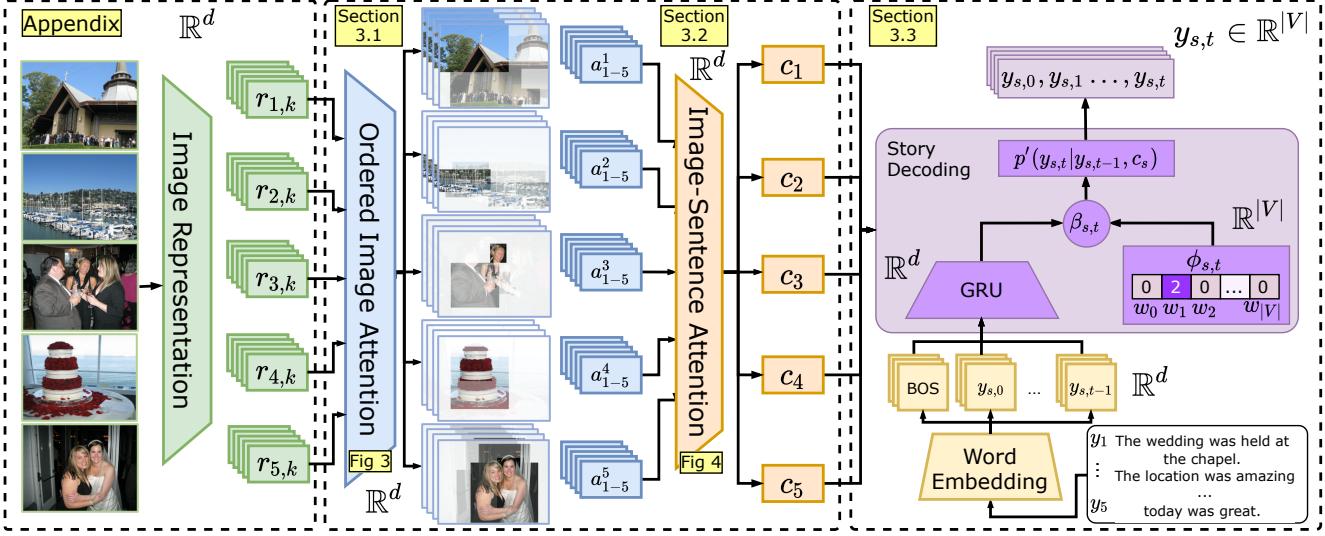


Figure 2: Our architecture for Visual Storytelling synthesis.

*et al.* [40] stacked attention modules sequentially. Later, many works concentrated on better vector-fusion modeling [7, 16, 4, 42]. Importantly, Lu *et al.* [21] suggested attending to the visual and textual modalities separately. Afterward, Kim *et al.* [15] proposed a bilinear module that efficiently generates attention for every pair. Following Lu *et al.* [21], Schwartz *et al.* [27, 28] suggested a general framework that extends attention to any number of utilities via local and interaction-based factors. We improve upon those ideas by suggesting an ordered attention. This ensures that interaction modeling is affected by the image position in a sequence.

### 2.3. Visual Storytelling

Huang *et al.* [14] introduced the Visual Storytelling task. Initially, Gonzalez *et al.* [9] adapted work by Vinyals *et al.* [33] used for captioning. Kim *et al.* [17] presented a Seq2Seq [29] approach with a decoding sampling strategy aimed to reduce the amount of repetition based on a word list. We improve their strategy by using a data-driven approach, penalizing each word differently based on its average counts. Wang *et al.* [36] employ adversarial learning to improve output stories. Huang *et al.* [13] utilize a reinforcement learning (RL) approach based on inter-image relations. Later works by Li *et al.* [19] and Zhang *et al.* [43] rely on preprocessing the data to better ground visual elements to the text while Yang *et al.* [39] and Hsu *et al.* [12] enrich the data with an external word common-sense knowledge graph. Our approach captures inter-image relations via ordered attention and is trained in an end-to-end manner alleviating the computational drawbacks of preprocessing or RL. Recently, state-of-the-art results were obtained by generating scene graphs for each image in the sequence [35]. Conversely, our image representations are dependant on all the images in the sequence.

## 3. Method

The goal of visual storytelling is to generate a story, composed of  $N$  ordered sentences  $\{y_s | 1 \leq s \leq N\}$ , given an ordered sequence of images  $I = \{I_s | 1 \leq s \leq N\}$ . Each sentence  $y_s = (y_{s,0}, \dots, y_{s,t}, \dots)$  is composed of words  $y_{s,t} \in \mathcal{Y}$  from vocabulary  $\mathcal{V}$ .

The order in which the images are given is essential as it defines the plot line of the story. The story should be focused, *i.e.*, each sentence should be related to the remainder of the story. Importantly, the sentences should form a coherent body of text describing the set of images, and not only a set of related information. For instance, the story “*The church was beautiful. The bride and groom walk down the aisle. The cake was amazing.*” is less coherent than: “*We went to the church for the wedding today. The bride and groom were excited for the day. Both cut the cake together.*”

**Overview:** To address this challenge, we develop the model illustrated in Fig. 2. It infers conditional probabilities  $p'(y_{s,t}|y_{s,t-1}, c_s)$  for the  $t$ -th word  $y_{s,t} \in \mathcal{Y}$  in sentence  $y_s$  given the previous word  $y_{s,t-1}$  and the context embedding  $c_s$  for sentence  $s$ . The context embedding  $c_s$  summarizes region representations  $r_{i,k}$  of all  $K$  object regions across all  $N$  images  $I_i$  ( $i \in [1, N]$ ,  $k \in [1, K]$ ) via Ordered Image Attention (OIA) (Sec. 3.1) and Image-Sentence Attention (ISA) (Sec. 3.2). Specifically, when generating sentence  $s$ , OIA computes an attended image representation  $a_i^s$  for every image  $I_i$  by attending to the  $K$  region representations  $r_{i,k}$  (Sec. 3.1). These attended image representations  $a_i^s$  are subsequently summarized into the context embedding  $c_s$  via an image-sentence attention (Sec. 3.2).

Below we first discuss computation of the attended image representation  $a_i^s$  (Sec. 3.1), before detailing computation of the context embedding  $c_s$  (Sec. 3.2) and computation of the conditional probabilities  $p'(y_{s,t}|y_{s,t-1}, c_s)$  (Sec. 3.3).

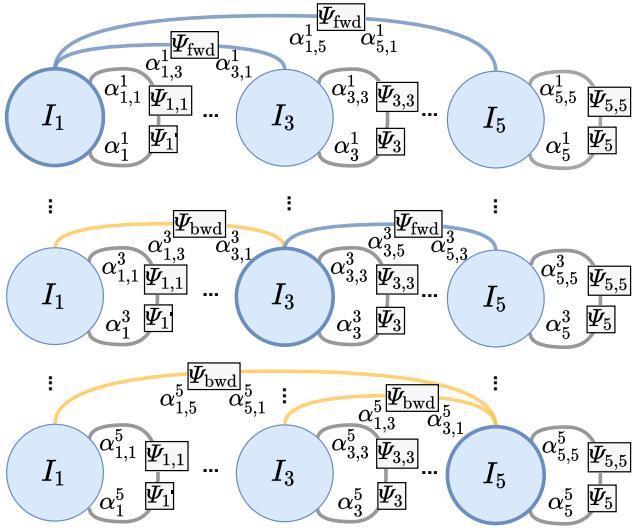


Figure 3: Illustration of Ordered Image Attention. Each node represents an image attention belief. For each sentence, we connect all the images with the sentence-corresponding image. The relative position to this image determines whether the connection is modeled with the  $\Psi_{\text{bwd}}$  factor (for preceding images) or the  $\Psi_{\text{fwd}}$  factor (for subsequent images). We infer the attention belief by collecting interactions and local object information within the image. We use scalars to calibrate the importance of each factor. In total, we generate 25 attention maps, one per image for every sentence.

### 3.1. Ordered Image Attention (OIA)

Ordered Image Attention (OIA) is designed to 1) form a structure across ordered images and to 2) select the relevant objects per image. For this we model preceding and proceeding interactions separately using different attention factors. We calibrate each factor’s importance with trainable scalars, which forms a graph of dependencies between the images. For each sequence of  $N$  images, the model infers a total of  $N^2$  attention maps, one per image for each sentence. We detail this module next.

#### 3.1.1 Attention Belief

For each image  $I_i = \{r_{i,1}, \dots, r_{i,K}\}$  we consider a set of  $K$  regions, represented by their feature vectors  $r_{i,k} \in \mathbb{R}^d$ , where  $d$  is the objects’ embedding dimension. Suppose we are currently generating sentence  $y_s$  ( $1 \leq s \leq N$ ). To do this we first compute an attended image representation  $a_i^s$  as follows

$$a_i^s = \sum_{k=1}^K b_{i,k}^s r_{i,k}, \quad (1)$$

where  $b_{i,k}^s \geq 0$  is the attention belief highlighting the importance of the  $k$ -th object in the  $i$ -th image when generating the  $s$ -th sentence. Importantly, for every image  $I_i$  we require  $b_{i,k}^s$  to be a valid probability distribution, *i.e.*, we also enforce  $\sum_{k=1}^K b_{i,k}^s = 1 \forall s, i$ .

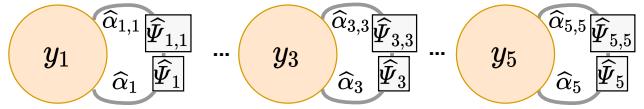


Figure 4: Illustration of ISA. The attention selects the attended image representation per sentence. We model interactions between attended images of the same sentence to compute each image’s importance. Note, each node represents a sentence attention belief over the attended images.

The object attention belief  $b_{i,k}^s$  is dependent on all the input data, *i.e.*, other objects and images. To avoid complex computation, we factorize the belief  $b_{i,k}^s$  into two pairwise dependencies that preserve the order, and a local term. For the pairwise terms we use  $\mu_{j \rightarrow i}^{\text{bwd}}$ , which is a message from a preceding image  $I_j$ , or  $\mu_{j \rightarrow i}^{\text{fwd}}$ , which is a message from a subsequent image  $I_j$ . We also use  $\mu_{i \rightarrow i}$  for self-messages. Additionally, we include a local factor  $\Psi_i(r_{i,k})$  that considers the object representation. Unlike the messages mentioned before, the local factor does not rely on interactions with other objects. We aggregate all the messages along with the local factor as illustrated in Fig. 3. For normalization we employ a softmax.

Formally we compute the attention belief  $b_{i,k}^s$  by distinguishing three cases. If  $i = s$  we have

$$\begin{aligned} b_{i,k}^s &\propto \exp(\alpha_i^s \Psi_i(r_{i,k}) + \alpha_{i,i}^s \mu_{i \rightarrow i}(r_{i,k}) + \\ &\quad \sum_{j < i} \alpha_{i,j}^s \mu_{j \rightarrow i}^{\text{bwd}}(r_{i,k}) + \sum_{j > i} \alpha_{i,j}^s \mu_{j \rightarrow i}^{\text{fwd}}(r_{i,k})). \end{aligned} \quad (2)$$

If  $i < s$  we use

$$\begin{aligned} b_{i,k}^s &\propto \exp(\alpha_i^s \Psi_i(r_{i,k}) + \\ &\quad \alpha_{i,i}^s \mu_{i \rightarrow i}(r_{i,k}) + \alpha_{i,s}^s \mu_{s \rightarrow i}^{\text{bwd}}(r_{i,k})). \end{aligned} \quad (3)$$

If  $i > s$  we obtain

$$\begin{aligned} b_{i,k}^s &\propto \exp(\alpha_i^s \Psi_i(r_{i,k}) + \\ &\quad \alpha_{i,i}^s \mu_{i \rightarrow i}(r_{i,k}) + \alpha_{i,s}^s \mu_{s \rightarrow i}^{\text{fwd}}(r_{i,k})). \end{aligned} \quad (4)$$

In all three cases  $\alpha_i^s, \alpha_{i,i}^s, \alpha_{i,s}^s \in \mathbb{R}$  are scalars used to calibrate the importance of different messages for a given sentence. These scalars form a dependency structure between images for each of the generated sentence indices. Intuitively, when we generate the first sentence, the attention belief might depend more on subsequent images, to correctly identify the story event, *e.g.*, a wedding, a parade, *etc.* Thus, the scalars will promote interaction with later images. An analysis of these scalars is provided in the appendix. Next, we define the different types of messages.

#### 3.1.2 Pairwise Messages and Factors

A message aggregates interaction scores from an image to an object. The three messages  $\mu_{j \rightarrow i}^{\text{bwd}}$ ,  $\mu_{j \rightarrow i}^{\text{fwd}}$  and  $\mu_{i \rightarrow i}(r_{i,k})$

are computed as follows:

$$\mu_{j \rightarrow i}^{\text{bwd}}(r_{i,k}) = \sum_{k'=1}^K \Psi_{\text{bwd}}(r_{i,k}, r_{j,k'}), \quad (5)$$

$$\mu_{j \rightarrow i}^{\text{fwd}}(r_{i,k}) = \sum_{k'=1}^K \Psi_{\text{fwd}}(r_{i,k}, r_{j,k'}), \text{ and} \quad (6)$$

$$\mu_{i \rightarrow i}(r_{i,k}) = \sum_{k'=1}^K \Psi_{i,i}(r_{i,k}, r_{i,k'}). \quad (7)$$

Importantly, these messages collect three different types of order-dependent interaction factors: (1) A backward image interaction, namely  $\Psi_{\text{bwd}}(r_{i,k}, r_{j,k'})$ . This interaction models relations to the preceding  $j$ -th image in the sequence. (2) A forward image interaction, namely  $\Psi_{\text{fwd}}(r_{i,k}, r_{j,k'})$ . This interaction models relations to the subsequent  $j$ -th image in the sequence. (3) The self interaction factor, namely  $\Psi_{i,i}(r_{i,k}, r_{i,k'})$ , which takes into account interactions between objects within the image. We formally define the different factors next.

**Interaction factors:** A commonly used practice to capture interactions across attention mechanisms is to first embed the elements into a joint Euclidean space followed by a dot-product [31, 27, 8, 28]. While we follow the same practice, we define three types of interaction factors to preserve the order. Consider two objects,  $r_{i,k} \in I_i$  from the sentence-corresponding image and  $r_{j,k'} \in I_j$  from the interacting image. We describe three types of interactions: for interactions with subsequent images (*i.e.*,  $j > i$ ) we use

$$\Psi_{\text{fwd}}(r_{i,k}, r_{j,k'}) = \left( \frac{L_{\text{fwd}} r_{i,k}}{\|L_{\text{fwd}} r_{i,k}\|_2} \right)^{\top} \left( \frac{R_{\text{fwd}} r_{j,k'}}{\|R_{\text{fwd}} r_{j,k'}\|_2} \right). \quad (8)$$

For interactions with preceding images (*i.e.*,  $j < i$ ) we use

$$\Psi_{\text{bwd}}(r_{i,k}, r_{j,k'}) = \left( \frac{L_{\text{bwd}} r_{i,k}}{\|L_{\text{bwd}} r_{i,k}\|_2} \right)^{\top} \left( \frac{R_{\text{bwd}} r_{j,k'}}{\|R_{\text{bwd}} r_{j,k'}\|_2} \right). \quad (9)$$

For interactions within the image (*i.e.*,  $j = i$ ) we have

$$\Psi_{i,i}(r_{i,k}, r_{i,k'}) = \left( \frac{L_{i,i} r_{i,k}}{\|L_{i,i} r_{i,k}\|_2} \right)^{\top} \left( \frac{R_{i,i} r_{i,k'}}{\|R_{i,i} r_{i,k'}\|_2} \right). \quad (10)$$

Note,  $L_{\text{fwd}}, R_{\text{fwd}}, L_{\text{bwd}}, R_{\text{bwd}}, L_{i,i}, R_{i,i} \in \mathbb{R}^{d \times d}$  are trainable shared weights across the entire image sequence. Also, the object from the sentence-corresponding image will always be on the left side of the factor equation. Thus, the factor embeddings preserve the order.

**Local factor:** Differently from the previous interactions the following factor captures how important an object is based solely on the object representation. Given an object  $r_{i,k} \in I_i$ , we define the local factor as,

$$\Psi_i(r_{i,k}) = v^{\top} \text{ReLU}(V r_{i,k}), \quad (11)$$

where  $v \in \mathbb{R}^d$ ,  $V \in \mathbb{R}^{d \times d}$  are trainable weights.

### 3.2. Image-Sentence Attention (ISA)

In a next step we summarize the attended image representations  $a_i^s$  produced by OIA to compute the context embedding  $c_s$  for the sentence  $s$  that we wish to generate. For this we use the Image-Sentence Attention (ISA) unit. It picks the relevant image context for generating the specific sentence. Formally we obtain the context embedding via

$$c_s = \sum_{i=1}^N \hat{b}_{s,i} a_i^s, \quad (12)$$

where attention factors

$$\hat{b}_{s,i} \propto \exp \left( \hat{\alpha}_s \hat{\Psi}_i(a_i^s) + \hat{\alpha}_{s,s} \hat{\mu}_{s \rightarrow s}(a_i^s) \right), \quad (13)$$

and where  $\hat{\alpha}_s, \hat{\alpha}_{s,s} \in \mathbb{R}$  are scalars. To avoid spurious correlations between sentences, we consider only self interactions and a local factor. This is illustrated in Fig. 4. The self-message of the attended image representation  $a_i^s$  is

$$\hat{\mu}_{s \rightarrow s}(a_i^s) = \sum_{j=1}^N \hat{\Psi}(a_i^s, a_j^s). \quad (14)$$

Finally, the self and local factors are defined with a different set of weights following Eq. (10) and Eq. (11) respectively.

### 3.3. Story Decoding

The goal at each timestep of decoding is to compute the conditional probability  $p(y_{s,t}|y_{s,t-1}, c_s)$  where  $y_{s,t} \in \mathcal{Y}$  is the  $t$ -th word in sentence  $y_s$ ,  $\mathcal{Y}$  is the vocabulary and  $c_s$  is the context embedding detailed in Sec. 3.2. For this we use a GRU recurrent unit, tasked with generating probabilities over the vocabulary conditioned on the context embedding  $c_s$  and the previously generated token  $y_{s,t-1}$ :

$$p(y_{s,t} = w | y_{s,t-1}, c_s) \propto \exp(\beta_{s,t} \cdot g_w(y_{s,t-1}, h_{s,t-1}, c_s) + (1 - \beta_{s,t}) \cdot f_w(\phi_{s,t})), \quad (15)$$

where  $g_w$  is the output of a GRU unit for the word  $w$ . We set the GRU hidden dimension to  $d$ .  $h_{s,t-1} \in \mathbb{R}^d$  is the hidden state at timestep  $t-1$  for sentence  $s$ .  $f : \mathbb{R}^{|\mathcal{Y}|} \rightarrow \mathbb{R}^{|\mathcal{Y}|}$  is a learned prior over the vocabulary based on a bag-of-words prior histogram  $\phi_{s,t}$ , which we describe in the next paragraph. The purpose of  $f$  is to reduce text repetitions.  $f_w$  denotes the value of  $f$  for a word  $w$ . We also incorporate a calibration gate  $\beta_{s,t} : \mathbb{R}^d \rightarrow [0, 1]$  for functions  $f$  and  $g$  using

$$\beta_{s,t} = \sigma(v_{\beta}^{\top} \tanh(G_g h_{s,t} + G_f W_1(\phi_{s,t}))). \quad (16)$$

Here,  $G_g \in \mathbb{R}^{d \times d}$  and  $G_f \in \mathbb{R}^{\gamma \times d}$  are trained projections of the GRU hidden state and the bottleneck layer respectively,  $v_{\beta} \in \mathbb{R}^d$  are learned weights and  $\sigma$  is the sigmoid function.  $W_1$  is obtained from the prior as discussed next.

Method	M	B-1	B-2	B-3	B-4	R	C	Img Feat
seq2seq [14]	31.4	-	-	-	3.5	-	6.84	FC
h-attn-rank [41]	33.9	-	29.8	-	-	29.8	7.4	FC
Contextualize, Show & Tell [9]	34.4	60.1	36.5	21.1	12.7	29.2	7.1	FC
AREL [36]	35.0	63.8	39.1	23.2	14.1	29.5	9.4	FC
KnowledgeableStoryteller [39]	35.2	66.4	39.2	23.1	12.8	29.9	<b>12.1</b>	FC
HSRL [13]	35.2	-	-	-	12.3	29.5	8.4	Spatial
StoryAnchor [43]	35.5	65.1	40.0	23.4	14.0	30.0	9.9	FC
SGVST [35]	35.8	65.1	40.1	23.8	14.7	29.9	9.8	F-RCNN
Ours (ResNet)	36.3	66.3	41.5	23.7	14.5	30.0	9.8	Spatial
<b>Ours</b>	<b>36.8±0.1</b>	<b>68.4±0.7</b>	<b>42.7±0.3</b>	<b>25.2±0.2</b>	<b>15.3±0.2</b>	<b>30.2±0.1</b>	<b>10.1±0.2</b>	F-RCNN

Table 1: Quantitative results on the VIST dataset for METEOR, BLEU-1...4, ROUGE-L and CIDEr. The primary metric is METEOR. The ‘Img Feat’ column describes the pretrained image features. All models utilize a ResNet [10] backbone except CS&T which employs an Inception v3 model [30]. FC and Spatial refer to features extracted from the penultimate layer and the preceding one accordingly. F-RCNN are bottom up features [1].

**Bag-of-words (BOW) prior:** Remembering history during storytelling permits to stay on topic and advance the story in the desired direction. Although quite intuitive, mimicking this ability is not trivial. *E.g.*, most approaches for VST generate all the sentences in parallel. Converting the parallel sentence generation into a sequential one implies a major computational overhead during training.

To address this, we propose a simple yet effective learnable framework that does not require sequential training while still exploiting information found in prior sentences. The history is represented via a bag-of-words histogram  $\phi_{s,t}$ , which includes all words that have been used until timestep  $t$  for the  $s$ -th sentence. During training, we initialize  $\phi_{s,t=0}$  with the ground truth history counts found in the previous  $s - 1$  sentences. We update the statistics at each timestep with the predicted word  $y_{s',t}$  for  $s' < s$ , and produce the next state of the counter  $\phi_{s,t+1}$ . At inference we generate sentences sequentially and update  $\phi_{s,t}$  with the predicted words.  $\phi_{s,t}$  is fed through a shallow bottleneck network to obtain the prior  $f$ , composed of two layers  $W_1 \in \mathbb{R}^{|\mathcal{Y}| \times \gamma}$  and  $W_2 \in \mathbb{R}^{\gamma \times |\mathcal{Y}|}$  without activation, where  $\gamma$  is the bottleneck dimension:

$$f(\phi_{s,t}) = W_2(W_1(\phi_{s,t})). \quad (17)$$

Also note the use of  $W_1(\phi_{s,t})$  in the gate (Eq. (16)).

**Intra-repetition regularization:** To regularize intra-repetitions, we decay the probability of previously used words during sentence generation. A critical aspect of this approach is to exclude words that appear frequently in the language (*e.g.*, was, were, am). For this we preprocess the training set to calculate the average story frequency  $\rho(w)$  of a word  $w$  via  $\rho(w) = \frac{\# \text{ appearances of word } w}{\# \text{ stories } w \text{ was used}}$ . The final count for word  $w$  at timestep  $t$  is calculated as  $\phi'_{s,t}(w) = \max[0, (\phi_{s,t}(w) - \rho(w) + 1)]$ . Intuitively, a word will not be penalized before it is used more than the prior belief average  $\rho(w)$ . The final probability for word  $w$  being

used is given by

$$p'(y_{s,t} = w | y_{s,t-1}, c_s) = \frac{p(y_{s,t} = w | y_{s,t-1}, c_s)}{\pi \cdot \phi'_{s,t}(w) + 1}, \quad (18)$$

where  $\pi \geq 0$  is a constant hyper-parameter. A penalty of 2 proved to work best on the validation set.

## 4. Results

**Dataset:** To train and test the model we use the VIST dataset [14]. This dataset is composed of stories. Each story has 5 images and  $N = 5$  corresponding sentences. All images were collected from Flickr albums. Sequences of images belong to the same album. Each image sequence is annotated with 5 ground-truth reference stories. On average, around 2.5 stories are based on the images, and the rest are rewrites. The overall numbers are 40,098 training stories, 4,988 validation stories, and 5,050 test stories.

**Training Setup:** We extracted the image features using a pre-trained F-RCNN model with a ResNet152 backbone [10, 26, 1]. We set the number of extracted objects  $K = 36$ . Bounding box coordinates were normalized between 0 and 1. Words that appear less than 3 times in the training set are represented by an  $<\text{UNK}>$  token. The vocabulary size is 12,210 words. Word representations were initialized using GloVe embeddings [25]. We set the decay parameter  $\pi = 2$  and the image representation dimension  $d = 512$ . We set the dropout parameter to 0.3. We use cross-entropy loss to maximize likelihood of ground-truth stories. At decoding time we employ a beam search algorithm, with beam width set to 3. We use Adam [18] optimizer with a learning-rate of 0.0004, which is decayed by a factor of 0.8 if the validation score (METEOR) does not improve after 4 epochs. The total amount of trainable parameters is 13,092,194. Training converges after  $\sim 20$  epochs. Each epoch needs 20 minutes on an Nvidia V100 GPU.

### 4.1. Quantitative analysis

**Evaluation metrics:** As suggested by the creators of VIST [14], METEOR [2] correlates best with human judge-

Model	M	B-4	R	C	#Params
<b>attention</b>					
w/o OIA	36.0	14.1	30.0	8.4	11M
w/o ISA	35.9	14.2	29.9	9.3	11M
w/o attention	35.8	13.6	29.7	7.2	11M
no-direction	36.1	14.5	28.9	8.4	12M
<b>decoding</b>					
w/o rep. regularization	36.2	14.5	29.8	8.7	13M
w/o count norm	36.4	14.6	29.9	9.4	13M
w/o BOW prior	36.4	14.5	30.0	9.7	13M
Transformer	36.7	<b>15.7</b>	30.0	9.9	13M
<b>Ours</b>	<b>36.8</b>	15.3	<b>30.2</b>	<b>10.1</b>	13M

Table 2: Components ablation analysis.

Local	Self	Directional	M	B-1	B-2	B-3	B-4	R	C
✗	✓	✓	36.2	67.4	42.4	24.2	14.5	30.0	9.3
✓	✗	✓	36.0	67.8	42.3	24.2	14.4	29.8	9.2
✓	✓	✗	35.9	67.6	42.2	24.0	14.2	29.9	8.5
✓	✓	✓	<b>36.8</b>	<b>68.4</b>	<b>42.7</b>	<b>25.2</b>	<b>15.3</b>	<b>30.2</b>	<b>10.1</b>

Table 3: Factor ablation analysis.

ment. Following their example, we use METEOR as the primary metric. We also compute BLEU [23], ROUGE [20], and CIDEr [32] and compare to prior work where available. For evaluation we use the evaluation script of Yu *et al.* [41]<sup>1</sup>.

**Comparison to state-of-the-art:** In Tab. 1 we compare the method to recent baselines. Early methods did not take into account visual-spatial information, which harms the performance (*e.g.*, 35.5% *vs.* 36.8% on METEOR) [14, 36, 9]. Wang *et al.* [35] utilize image representations similar to our approach but do not consider relations between different images, resulting in a 1% drop on METEOR, showing that ordered structure encoding with OIA is beneficial. SGVST and StoryAnchor [41, 43] use different methods for mapping the image sequence to distinct topics. Differently, our approach is trained end-to-end. Further, our image representations are dependant on all the images in the sequence. Notably, unlike our method, SGVST uses scene graphs. This requires an additional model pre-trained on external scene graph data. Finally, Yang *et al.* [39] utilize an external commonsense dataset to enrich the input. Their CIDEr score is significantly higher, yet this improvement does not translate to all other metrics. The approach improves upon the current state-of-the-art by a margin (36.8% *vs.* 35.8% on METEOR). Note, the ROUGE-L metric is based on finding the longest subsequence matched to human generated stories. However, this score is almost identical for all prior works, indicating that this metric doesn't capture story generation improvements. We also report the performance with spatial ResNet152 features [10], which outperforms the state-of-the-art as well. This shows that the method is stable irrespective of image features.

**Ablation study:** In Tab. 2 we show the importance of different components via an ablation study. In 'w/o OIA,' we

<sup>1</sup>[http://github.com/lichengunc/vist\\_eval](http://github.com/lichengunc/vist_eval) - Codebase for commonly used evaluation scripts.

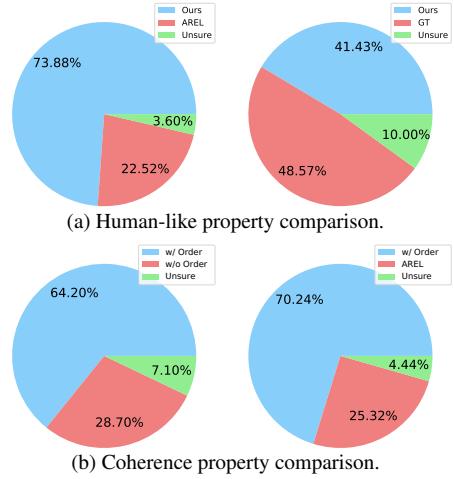


Figure 5: Human evaluation to compare properties.

Method	Focused	Coherent	Share	Human-like	Grounded	Detailed
AREL	3.49	3.18	3.18	3.26	3.32	3.15
Ours	3.67	3.52	3.20	3.56	3.54	3.32
GT	<b>3.72</b>	<b>3.57</b>	<b>3.34</b>	<b>3.64</b>	<b>3.56</b>	<b>3.53</b>

Table 4: Human evaluation results for rating survey (scores are between 1-5).

replace the OIA module (Sec. 3.1) with simple averaging of the  $K$  object representations of image  $I_i$ , resulting in a 0.8% drop on METEOR. Similarly, in 'w/o ISA,' we replace the ISA unit (Sec. 3.2) with averaging, leading to a 0.9% drop on METEOR. In 'w/o attention,' we removed both OIA and ISA, which dropped the METEOR score to 35.8%. For the method referred to as 'no-direction,' we use the same factor for preceding and proceeding interaction (*i.e.*,  $L_{\text{bwd}} = L_{\text{fwd}}$  and  $R_{\text{bwd}} = R_{\text{fwd}}$ ). Here, METEOR results drop by 0.7%. Hence, ordered interactions are beneficial. Next, we assess the decoding components (Sec. 3.3). We first remove the intra-repetition regularization (*i.e.*,  $\rho(w)$ ), which causes METEOR score to drop by 0.6%. Removing the popular words count ( $\phi'_{s,t}$ ), results in a 0.4% drop on METEOR. The METEOR score drops by 0.4% when we remove the BOW prior. Last, we replace the GRU decoding layer with a Transformer, which did not change results a lot.

In Tab. 3 we show an ablation analysis of the different factors used in OIA. We found that each factor contributes to the model's performance, and the directional factors (*i.e.*,  $\Psi_{\text{fwd}}$  and  $\Psi_{\text{bwd}}$ ) have the biggest impact.

In Tab. 5, we show the ability to reduce repetitions. As proposed by Bertoldi *et al.* [5], text repetitiveness is measured by the repetition rate of non-singleton n-grams within each story. In our experiment, we use up to 4-grams. The use of intra-repetition regularization reduces text repetition (0.14 to 0.04). Combined with the trainable bag-of-words prior module, we further improve this measure (0.008 *vs.* 0.14). We also report sentence repetitiveness, *i.e.*, the average number of repeated sentences in a story.

## 4.2. Human Evaluation

The subjective nature of the VST task calls for a human evaluation. We use a sample of 150 image sequences and test different story qualities by asking 3 MTurk annotators to rank or compare them to other methods. We compare our results to the AREL baseline since none of the more recent baselines are publicly available. Note that we also compare coherency against a model without ordered-factors, which already improves upon the prior state-of-the-art.

In Fig. 5a we provide the results when asking annotators to pick the most human-like story. We use the majority vote to decide the best model per story. The generated stories outperform the AREL baseline (73.87% vs. 22.53%). Surprisingly, in many cases, the annotators found the generated stories to be more human-like than the ground truth stories (41% vs. 48.57%). In Fig. 5b, we assess coherency. An important aspect of our work are the directional factors for coherency. To validate their effectiveness, we compared to a model that does not incorporate direction into the attention representation (*i.e.*, we use the same factor for preceding and proceeding interactions). The comparison shows a significant coherency improvement (64.2% vs. 28.7%). Also, a comparison against the AREL baseline demonstrates a more significant improvement (70.24% vs. 25.32%).

To further evaluate the quality of the stories, we follow the criteria set by the Visual Storytelling Challenge<sup>2</sup> and conduct a survey where judges are asked to rate six categories between 1-5: 1. *Focused*: the story contains information that is “naturally” relevant to the rest of the story; 2. *Coherence*: the sentences in the story are related and consistent; 3. *Share*: the inclination to share the story; 4. *Human-like*: the story was likely written by a human; 5. *Grounded*: the story directly reflects concrete entities in the image; and 6. *Detailed*: the story provides an appropriate level of detail. To obtain the final score, we average the annotators’ scores per sample, followed by averaging across the entire sample set. From Tab. 4 we observe: the model improved on all the criteria compared to the AREL model. Importantly, the generated stories are comparable to the ground-truth stories, indicating success in reducing the shortcomings found in prior methods. Nonetheless, the level of detail is still lacking, supporting the observation of Holtzman *et al.* [11] that current decoding strategies tend to generate well-formed yet somewhat generic text.

## 4.3. Qualitative evaluation

In Fig. 6, we show the ability of the method in reducing repetitions. We observe the AREL baseline to repeat the same sentences, for example, “...had a great time at...”. We also observe this repetitiveness when we remove the bag-of-words prior and the intra-sentence regularization (*i.e.*, No History column). Nevertheless, the method remains on topic, *i.e.*, family in the pool.

<sup>2</sup><http://visionandlanguage.net/workshop2018>



AREL	The kids <b>had a great time at</b> the pool. The little boy was excited to see the kids. We <b>had a great time at</b> the park. We <b>had a great time at</b> the pool. We <b>had a great time at</b> the park.
No History	The kids <b>had a great time at</b> the beach. The <b>baby</b> was happy to see the <b>baby</b> . We <b>had a great time at</b> the park. The <b>had a great time at</b> the pool. We <b>had a great time at</b> the park.
With History	The family went to the pool. The <b>baby</b> was very happy. The kids had a great time. The <b>kids</b> played in the pool. The little girl is having a good time.

Figure 6: An illustration of an image sequence along with three different stories generated by: (1) AREL baseline [36], (2) No History: a model without intra-repetition regularization and BOW prior (see Sec. 3.3); and (3) With History: the final model. Repeated sentences are highlighted with a yellow colored marker. Repeated words in a sentence are emphasized in red color.

Model		Text Rep.	Sent. Rep.
AREL [36]		0.16	0.4
BOG prior	Intra-repetition reg.		
No	No	0.14	0.33
Yes	No	0.10	0.18
No	Yes	0.04	0.04
Yes	Yes	<b>0.008</b>	<b>0.0</b>

Table 5: Story generation ablation analysis.

In Fig. 7 we sketch the attention maps along with the generated story. The first sentence, “We went to the mountains,” sets the theme for the story, which requires the processing of subsequent images. Notably, the ISA module picked the proceeding images. In contrast, for the second sentence, the attention focuses mostly on the second image resulting in a description of the lake observed exclusively in this image. The third sentence relates to the scenery. Hence the attention focuses on preceding and proceeding images.

## 5. Conclusion

We present a novel approach for VST, which encourages coherency of generated story. We incorporate structure between images with a new attention method that selects the important objects in an ordered image sequence. Human evaluation and quantitative analysis demonstrate that the approach outperforms existing methods. Further, we perform ablation and qualitative analysis to show effectiveness.

## References

- [1] 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 *CVPR*, 2017. **2, 6**
- [2] Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *ACL*, 2005. **6**
- [3] Kobus Barnard, Pinar Duygulu Sahin, David A. Forsyth, Nando de Freitas, David M. Blei, and Michael I. Jordan. Matching words and pictures. *JMLR*, 2003. **2**
- [4] Hedi Ben-Younes, Rémi Cadene, Matthieu Cord, and Nicolas Thome. Mutan: Multimodal tucker fusion for visual question answering. In *ICCV*, 2017. **3**

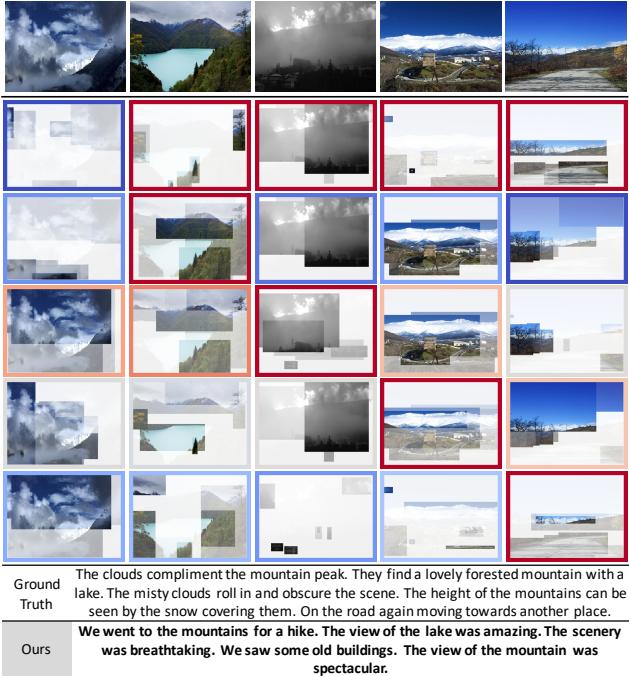


Figure 7: Illustration of OIA and ISA attention maps, the ground-truth story and the final generated story. Each row corresponds to a story sentence and shows objects OIA highlights. The attended images' border specifies the relevancy to sentence generation, from red (important) to blue (not important).

- [5] Nicola Bertoldi, Mauro Cettolo, and Marcello Federico. Cache-based online adaptation for machine translation enhanced computer assisted translation. In *MT Summit*, 2013. 7
- [6] Xinlei Chen and C. Lawrence Zitnick. Mind's eye: A recurrent visual representation for image caption generation. In *CVPR*, 2015. 1
- [7] Akira Fukui, Dong Huk Park, Daylen Yang, Anna Rohrbach, Trevor Darrell, and Marcus Rohrbach. Multimodal compact bilinear pooling for visual question answering and visual grounding. In *EMNLP*, 2016. 3
- [8] Peng Gao, Zhengkai Jiang, Haoxuan You, Pan Lu, Steven CH Hoi, Xiaogang Wang, and Hongsheng Li. Dynamic fusion with intra-and inter-modality attention flow for visual question answering. In *CVPR*, 2019. 5
- [9] Diana Gonzalez-Rico and Gibran Fuentes Pineda. Contextualize, show and tell: A neural visual storyteller. In *Storytelling Workshop, NAACL*, 2018. 1, 3, 6, 7
- [10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2015. 6, 7
- [11] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *ICLR*, 2020. 8
- [12] Chao-Chun Hsu, Zi-Yuan Chen, Chi-Yang Hsu, Chih-Chia Li, Tzu-Yuan Lin, Ting-Hao Huang, and Lun-Wei Ku. Knowledge-enriched visual storytelling. In *AAAI*, 2020. 3
- [13] Qiuyuan Huang, Zhe Gan, Asli Çelikyilmaz, Dapeng Wu, Jianfeng Wang, and Xiaodong He. Hierarchically structured

reinforcement learning for topically coherent visual story generation. In *AAAI*, 2018. 1, 3, 6

- [14] Ting-Hao Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Aishwarya Agrawal, Jacob Devlin, Ross B. Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, C. Lawrence Zitnick, Devi Parikh, Lucy Vanderwende, Michel Galley, and Margaret Mitchell. Visual storytelling. In *NAACL*, 2016. 1, 2, 3, 6, 7
- [15] Jin-Hwa Kim, Jaehyun Jun, and Byoung-Tak Zhang. Bilinear attention networks. In *NeurIPS*, 2018. 3
- [16] Jin-Hwa Kim, Kyoung-Woon On, Woosang Lim, Jeonghee Kim, Jung-Woo Ha, and Byoung-Tak Zhang. Hadamard product for low-rank bilinear pooling. In *ICLR*, 2017. 3
- [17] Taehyeong Kim, Min-Oh Heo, Seonil Son, Kyoung-Wha Park, and Byoung-Tak Zhang. Glac net: Glocal attention cascading networks for multi-image cued story generation. In *CoRR*, 2018. 3
- [18] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *CoRR*, 2014. 6
- [19] Jiacheng Li, Haizhou Shi, Siliang Tang, Fei Wu, and Yuet-ting Zhuang. Informative visual storytelling with cross-modal rules. In *MM*, 2019. 3
- [20] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *ACL*, 2004. 7
- [21] Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. Hierarchical question-image co-attention for visual question answering. In *NeurIPS*, 2016. 3
- [22] Yatri Modi and Natalie Parde. The steep road to happily ever after: an analysis of current visual storytelling models. In *Workshop on Shortcomings in Vision and Language, NAACL*, 2019. 1
- [23] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *ACL*, 2001. 7
- [24] Cesc C. Park and Gunhee Kim. Expressing an image stream with a sequence of natural sentences. In *NeurIPS*, 2015. 1
- [25] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *EMNLP*, pages 1532–1543, 2014. 6
- [26] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *IEEE PAMI*, 2015. 6
- [27] Idan Schwartz, Alexander G. Schwing, and Tamir Hazan. High-order attention models for visual question answering. In *NeurIPS*, 2017. 3, 5
- [28] Idan Schwartz, Seunghak Yu, Tamir Hazan, and Alexander G Schwing. Factor graph attention. In *CVPR*, 2019. 3, 5
- [29] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. Sequence to sequence learning with neural networks. In *NeurIPS*, 2014. 3
- [30] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2015. 6
- [31] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 5

- [32] Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *CVPR*, 2014. 7
- [33] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural image caption generator. In *CVPR*, 2014. 1, 3
- [34] Bairui Wang, Lin Ma, Wei Zhang, Wenhao Jiang, and Feng Zhang. Hierarchical photo-scene encoder for album storytelling. In *AAAI*, 2019. 1
- [35] Ruize Wang, Zhongyu Wei, Piji Li, Qi Zhang, and Xuanjing Huang. Storytelling from an image stream using scene graphs. In *AAAI*, 2019. 3, 6, 7
- [36] Xin Wang, Wenhua Chen, Yuanfang Wang, and William Yang Wang. No metrics are perfect: Adversarial reward learning for visual storytelling. In *ACL*, 2018. 1, 3, 6, 7, 8
- [37] Huijuan Xu and Kate Saenko. Ask, attend and answer: Exploring question-guided spatial attention for visual question answering. In *ECCV*, 2016. 2
- [38] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron C. Courville, Ruslan Salakhutdinov, Richard S. Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, 2015. 1, 2
- [39] Pengcheng Yang, Fuli Luo, Peng Chen, Lei Li, Zhiyi Yin, Xiaodong He, and Xu Sun. Knowledgeable storyteller: A commonsense-driven generative model for visual storytelling. In *IJCAI*, 2019. 3, 6, 7
- [40] Zichao Yang, Xiaodong He, Jianfeng Gao, Li Deng, and Alexander J Smola. Stacked attention networks for image question answering, 2015. 3
- [41] Licheng Yu, Mohit Bansal, and Tamara L. Berg. Hierarchically-attentive rnn for album summarization and storytelling. In *EMNLP*, 2017. 6, 7
- [42] Zhou Yu, Jun Yu, Chenchao Xiang, Jianping Fan, and Dacheng Tao. Beyond bilinear: Generalized multimodal factorized high-order pooling for visual question answering. In *NeurIPS*, 2018. 3
- [43] Bowen Zhang, Hexiang Hu, and Fei Sha. Visual storytelling via predicting anchor word embeddings in the stories. In *ICCV*. Workshop on Closing the Loop Between Vision and Language, 2020. 3, 6, 7