

# Transitional Adaptation of Pretrained Models for Visual Storytelling

Youngjae Yu<sup>1\*</sup>, Jiwan Chung<sup>2\*</sup>, Heeseung Yun<sup>2</sup>, Jongseok Kim<sup>3</sup>, Gunhee Kim<sup>2</sup>

<sup>1</sup>Allen Institute for AI, <sup>2</sup>Seoul National University, <sup>3</sup>Violet

{y.j.yu, jiwanchung, heeseung.yun, js.kim}@vision.snu.ac.kr, gunhee@snu.ac.kr

<https://vision.snu.ac.kr/projects/TAPM>

## Abstract

Previous models for vision-to-language generation tasks usually pretrain a visual encoder and a language generator in the respective domains and jointly finetune them with the target task. However, this direct transfer practice may suffer from the discord between visual specificity and language fluency since they are often separately trained from large corpora of visual and text data with no common ground. In this work, we claim that a transitional adaptation task is required between pretraining and finetuning to harmonize the visual encoder and the language model for challenging downstream target tasks like visual storytelling. We propose a novel approach named Transitional Adaptation of Pretrained Model (TAPM) that adapts the multi-modal modules to each other with a simpler alignment task between visual inputs only with no need for text labels. Through extensive experiments, we show that the adaptation step significantly improves the performance of multiple language models for sequential video and image captioning tasks. We achieve new state-of-the-art performance on both language metrics and human evaluation in the multi-sentence description task of LSMDC 2019 [50] and the image storytelling task of VIST [18]. Our experiments reveal that this improvement in caption quality does not depend on the specific choice of language models.

## 1. Introduction

Most models for vision-to-language generation tasks consist of a visual encoder to extract visual information from input images or videos, a language model to generate text sentences, and a mechanism to weld the two modules into one harmonized architecture. For example, recent models for visual captioning [7, 59] adopt a pretrained visual encoder and a pretrained language generator and then optimize the target cross-modal generation objective with the downstream datasets [63, 40, 49, 67, 69, 74]. In this pro-

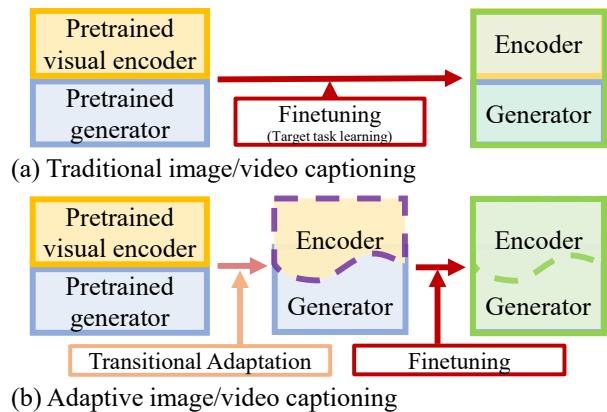


Figure 1. Comparison between existing captioning models and our Transitional Adaptation of Pretrained Model (TAPM). (a) Previous captioning models start from a pretrained visual encoder and a language generator and then directly finetune with the downstream datasets. (b) TAPM includes a simple pretext task as an adaptation process that harmonizes the generator with the visual encoder before optimizing the target objective.

cess, however, no transitional adaptation step has proposed to match the potentially substantial differences between the information stored in the visual encoder and the language generator, as they are separately trained from large sets of visual and text data with no common ground (*e.g.* images from ImageNet and text from Wikipedia).

This work is motivated by that this direct transfer of pretrained models to a downstream task may suffer from the dissonance between visual specificity and language fluency. For example, finetuning pretrained language models on another target task may result in catastrophic forgetting of the language generation capability [8, 66]. Moreover, existing captioning models have often been criticized for not sufficiently conditioning on the visual context and thus lack visual discriminability [34, 36].

Considering the potentially vast gap between the nature of the information stored in the visual encoder and the language decoder, it would be difficult for them to work in harmony at once for another challenging objective of vision-

\*Equal Contribution

to-language generation. In this light, we believe a simpler objective dedicated to improving coordination between the two separately pretrained models could help the model get prepared for the target objective eventually better and faster.

Therefore, we present *Transitional Adaptation of Pre-trained Model* (TAPM) for visual storytelling as the first approach that proposes an explicit visual adaptation step to harmonize the visual encoder with the pretrained language models as depicted in Fig. 1. Our adaptation step can be trained with only visual inputs, such as images or videos with no text label. We outline the contributions of this work as follows:

1. Our work is the first attempt to demonstrate an auxiliary adaptation loss’s effectiveness in welding a visual encoder with a pretrained language model. By extensive experiments, we show that this additional adaptation between pretraining and finetuning consistently improves the captioning quality of various language models such as GPT-2 [45], XLM [14], and QRNN [5].
2. We present the sequential coherence loss that can adapt the language generator using only sequential video/image inputs with no text label. We also introduce two recipes critical to TAPM’s success: (i) using the language model outputs for adaptation training and (ii) using the split-training process.
3. We evaluate TAPM in two storytelling tasks: sequential video captioning in the LSMDC 2019 [50] and sequential image captioning in VIST [18]. TAPM achieves new state-of-the-art performance in both tasks in terms of automatic language metrics and human evaluation.

## 2. Related Work

**Visual Storytelling.** Unlike direct and literal descriptions, visual storytelling aims to generate a more figurative and consistent narrative for consecutive images or videos [18]. Some earlier works [23, 24] explore the summarization of long videos into the storyline representation. Park *et al.* [41, 42] integrate an entity-based local coherence model to generate a coherent flow of multiple sentences for a photo album. Fan *et al.* [9] use a shorter prompt as the intermediate representation. Jain *et al.* [19] combine SMTs and RNNs to merge independent descriptions into a coherent story. Huang *et al.* [17] propose a two-level hierarchical RL-based decoder to plan a semantic topic first and then generate consistent sentences. Tang *et al.* [58] employ an attribute-based hierarchical decoder to create paragraphs using policy gradient with word-level rewards and adversarial training. Fan *et al.* [10] exploit a predicate-argument structure of the text to build coherent stories. Gella *et*

*al.* [11] introduce the VideoStory dataset for generating stories from social media videos. AdvInf [43] uses adversarial inference and MART [28] memory augmented transformer to generate paragraph-level captions.

Most previous works on visual storytelling require both visual encoder and language generator. Our work is orthogonally applicable to these approaches to better adapt the language decoder for visual context before training the models with the main vision-to-language objective, including Reinforcement and adversarial learning.

**Auxiliary Losses for Captioning.** Autoregressive language models trained with cross-entropy often suffer from exposure bias [3]. Several works on captioning have leveraged reinforcement learning by using rewards as auxiliary loss signals to ameliorate this bias. Zhang *et al.* [71] directly optimize language quality metrics with an actor-critic framework. Liu *et al.* [33] optimize a linear combination of language metrics using Monte Carlo rollouts. SCST [48] improves the REINFORCE algorithm to correctly normalize external rewards using the test-time inference algorithm’s output. Ren *et al.* [47] use the embedding similarity between generated sentences and image features as the reward. These reinforcement learning approaches have been extended to the video captioning problem [30, 65]. While reinforcement learning can help training non-differentiable objectives, it is known to be unstable [60]. Other types of auxiliary losses have also been adopted for captioning problems. Ma *et al.* [37] employ the cyclic reconstruction to enforce the localization of each word in an image. Zhou *et al.* [73] add visual grounding supervision to enhance the sentence generation quality. HINT [52] learns to match the attention map to human attention for grounded image captioning. VideoBERT [56] extends the text-based BERT to build bidirectional modeling between videos and captions. Compared to previous work, our work does not require additional visual caption data since it takes self-supervision losses with only sequential visual inputs.

**Pretrained Models for Vision-to-Language Tasks.** Recently, many works have demonstrated the power of self-supervision based representation learning in cross-modal settings. LXMERT [57] and ViLBERT [35] pretrain two-stream transformers on various tasks including masked cross-modal language model (LM) objectives. LXMERT is extended later with adaptive sparse attention [4]. VisualBERT [31] and VL-BERT [54] uses single-stream transformers. CMR [72] models the relevance between the textual entities and visual entities. UNITER [6] and Unicoder-VL [29] use object detection based objectives in addition to the masked LM loss. VideoBERT [56] trains a transformer for video-language tasks using vector quantization to categorize videos into discrete tokens. CBT [55] replaces the softmax loss of BERT with noise contrastive estimation.

These approaches aim to learn general representations,

and our method adapts the trained representations to the target cross-modal generation tasks. Thus, our model is orthogonal to the aforementioned self-supervised representations and consistently improves the final performance even with the self-supervised representation. Furthermore, they often use the masked cross-modal objectives that require both visual data and associated sentences (with blanks); contrarily, our method does not require text data at all for self-supervision.

### 3. Approach

We demonstrate our TAPM approach in visual storytelling tasks, which are a sequential extension of visual captioning. Its goal is to generate coherent  $C$  sentences for  $C$  visual inputs of video clips or images. We henceforth explain our model in the context of sequential video captioning because it subsumes sequential image captioning.

Fig. 2 illustrates the overall architecture, which consists of the visual encoder (section 3.1) and the language generator (section 3.2). We train the visual encoder and the language generator with the adaptation loss before finetuning them with the downstream captioning tasks (section 3.3). We employ the sequential coherence loss as the adaptation loss to encourage both distinctiveness and coherence in sequential captions. These losses are applied to the language model outputs in order to update the visual encoder in accordance with the language model (section 3.5). Finally, the encoder and the generator are trained with the target objective of visual storytelling.

For overall training, we use a split-training approach (section 3.5) that helps the decoder retain language generation capability. Since the adaptation loss is not a generation loss, it may degrade the language understanding of the pretrained language model. Hence, split-training fixes the language generator weights during the adaptation phase.

#### 3.1. The Visual Encoder

Given a video clip, we utilize pretrained feature extractors to extract vector feature  $v_{ij}$  for each frame  $j$ . The set of pretrained feature extractors varies depending on datasets and will be covered in section 4.1. We then reduce the vectors to  $M$  segments by mean-pooling them over temporal dimension.

With the extracted features of a video clip  $\mathbf{V}_i = \{\mathbf{v}_{i1}, \dots, \mathbf{v}_{iM}\}$  as inputs, the visual encoder builds task-specific representations  $\bar{\mathbf{V}}_i = \{\bar{\mathbf{v}}_{i1}, \dots, \bar{\mathbf{v}}_{iM}\}$ . Our visual encoder consists of two fully connected (FC) layers followed by Leaky ReLU [38], three layers of residual blocks, and a final self-attention layer [61]. A residual block consists of two FC layers and a ReLU activation [12]. After processing the visual inputs, we mean-pool the previous and next frame representation and concatenate them to the current representation to encode the context information.

#### 3.2. The Language Generator

For the language generator, one can use any language model. In our experiments, it is implemented by (but not confined to) GPT-2 [45], GPT [44], XLM [14], QRNN [5], and LSTM [16]. We use GPT-2-small [45] pretrained on a corpus dataset of 8 million web pages as the default generator due to its best performance among other language models. We will report the results of other language models in section 4.3.

#### 3.3. Adaptation training

We train the visual encoder with a simple auxiliary objective to harmonize it with the language generator in the adaptation phase. Here, we describe how to encode the visual and text representations for calculating the adaptation loss given the video inputs. The adaptation loss for visual storytelling will be discussed in the next section.

The language generator takes the task-specific representation  $\bar{\mathbf{V}}_i$  from the visual encoder as inputs and generates the contextualized representation for visual  $\tilde{\mathbf{V}}_i$  and text  $\hat{\mathbf{s}}_i$ . Specifically, the input  $\mathbf{X}_i$  to the generator is

$$\mathbf{X}_i = [\bar{\mathbf{V}}_i, [sep], [dummy]], \quad (1)$$

where  $[sep]$  and  $[dummy]$  are respectively separation and dummy tokens. Remind that  $\bar{\mathbf{V}}_i$  is a sequence of vectors with the number of segments  $M$ . Then the generator outputs

$$\tilde{\mathbf{X}}_i = [\tilde{\mathbf{V}}_i, [sep], \hat{\mathbf{s}}_i], \quad (2)$$

where  $\hat{\mathbf{s}}_i$  can be regarded as the text representation that the generator predicts for a sequential video input  $\bar{\mathbf{V}}_i$ .

Finally, the visual representation  $\hat{\mathbf{v}}_i$  is obtained by mean-pooling the sequence representation  $\tilde{\mathbf{V}}_i$  to a single vector. Note that the adaptation step does not use the caption label but inputs a dummy token into the generator to obtain text information. Thus, we can train the language generator with video-only datasets. While the dummy token can be arbitrarily selected from the pretrained vocabulary, we resort to the start-of-sentence token for all reported experiments. As will be shown in Table 3, TAPM with the dummy token performs comparably with the ground truth captions.

#### 3.4. The Sequential Coherence Loss

Visual storytelling is the problem of generating expressive, aligned, and coherent captions from a sequence of semantically connected visual inputs (e.g. videos or photo streams). Consecutive images or video clips tend to share common backgrounds, characters, and objects.

This closeness makes those visual features similar, and as a result, the captions generated from them overlap one another. To make consecutive captions not too overlapped

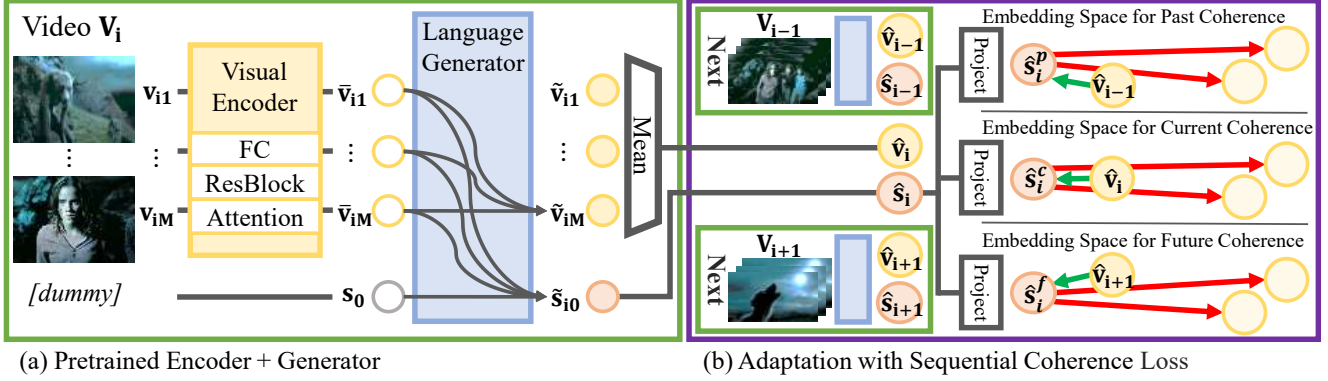


Figure 2. Illustration of the proposed TAPM framework. (a) TAPM harmonizes a pretrained visual encoder (section 3.1) with a pretrained language generator (section 3.2) to improve a target captioning task. In the adaptation phase, the model takes only videos (or images) as the input. Given a video, the language generator builds the corresponding video embedding ( $\hat{v}_i$ ) and text embedding ( $\hat{s}_i$ ) per each video. (b) We introduce sequential coherence loss to improve temporal coherence in visual storytelling tasks. We first use the respective FC layers ( $f^p$ ,  $f^c$  and  $f^f$ ) to project the text embedding ( $\hat{s}_i$ ) into the past, current, and future visual space. We then encourage the respective past, current, and future text embedding to be closer to their corresponding visual representations (**Pull (Green arrows)**) than the other visual representations (**Push (Red arrows)**).

but still coherent, we introduce the *sequential coherence loss* to build text representation of each visual input.

The sequential coherence loss enforces the text representation of a clip to predict the visual representations within its closed neighborhood well. We divide the sequential coherence loss into three parts of the past, current, and future matching loss for a better explanation. First, the past matching loss projects the text representation  $\hat{s}_i$  of video  $i$  by an FC layer  $f^p$  and makes it closer to the visual representation  $\hat{v}_{i-1}$  of the previous video  $i-1$  than the other videos, as in Figure 2. Second, the future matching loss is almost identical to the past matching loss except that it projects  $\hat{s}_i$  with a different FC layer  $f^f$  and matches with the next visual representation  $\hat{v}_{i+1}$ . Finally, the current matching loss matches the current visual representation  $\hat{v}_i$  with  $\hat{s}_i$  through an FC layer  $f^c$ . They are similar in that the text representation is projected in the past, future, current visual space by an FC layer and then drives the embeddings of correct matches closer (*pull*) and those of wrong matches farther away from each other (*push*).

To implement this notion, we employ margin ranking losses between correct matches and other wrong ones. The final loss is the sum of the past, current, and future matching losses as follows:

$$\begin{aligned}
 L_i = & \sum_{j \neq i-1} \max(0, 1 + \hat{v}_j * f^p(\hat{s}_i) - \hat{v}_{i-1} * f^p(\hat{s}_i)) \quad (3) \\
 & + \sum_{j \neq i} \max(0, 1 + \hat{v}_j * f^c(\hat{s}_i) - \hat{v}_i * f^c(\hat{s}_i)) \\
 & + \sum_{j \neq i+1} \max(0, 1 + \hat{v}_j * f^f(\hat{s}_i) - \hat{v}_{i+1} * f^f(\hat{s}_i)),
 \end{aligned}$$

where the operator  $*$  denotes the cosine similarity, and  $j$  indicates the index for wrong matches.

### 3.5. Training with the adaptation Loss

**Use of Language Model Outputs.** As described in the previous sections, our adaptation losses use the visual representation processed with the language model rather than the visual encoder outputs. Using the language model outputs enables the adaptation losses to update the visual encoder in accordance with the language model. On the other hand, using the encoder outputs would update the visual encoder in isolation. In Table 3, we will show that adaptation using the encoder outputs (TAPM+VisualA) does not improve upon the baseline (TAPM-A), while adaptation on the language model outputs (TAPM) does. Thus, this scheme is crucial to train the visual encoder in coordination with the language model to benefit the target task.

**Split-Training.** We split the training process into two phases: the adaptation loss step and the caption generation loss step. First, the visual encoder is updated for a given number of epochs by the adaptation loss, while the text encoder and the language generator are fixed. Then, we jointly update all the components with the generation loss. By splitting the training process, we give the model a chance to optimize the simpler adaptation task long enough before being presented with the harder generation objective. Fixing the language generator during the adaptation loss step prevents catastrophic forgetting of the language generation capability. Our ablation study in section 4.3 confirms that the split training leads to significant performance gains.

### 3.6. Finetuning and Inference

**Target-Task Training.** After adaptation training, we can finetune the language generator to the downstream captioning task with ground-truth data, where we input  $C$  pairs of video clips (or images) and text descriptions one by one:  $\{\mathbf{V}_1, \mathbf{S}_1, \dots, \mathbf{V}_C, \mathbf{S}_C\}$ . We use the teacher forcing as the training scheme with the cross-entropy loss:

$$L_i^G = - \sum_{l=1}^L \sum_{v=1}^V y_{il}^v \log p_{il}^v, \quad (4)$$

where  $v \in \{1, \dots, V\}$  is the vocabulary index,  $p_{il}$  is the prediction probability for the  $l$ -th token in  $\mathbf{S}_i$ , and  $y_{il}$  is the ground truth label. Finally, the language model head generates a caption output, consisting of a single FC layer that maps each vector of the language model outputs  $\tilde{\mathbf{S}}_i$  to a softmax layer to obtain the word probability  $p_i$  of each token over vocabulary.

**Cross-Modal Generation.** At inference, our goal is to generate a coherent sequence of  $C$  sentences for a visual test sample  $\{\mathbf{V}_1, \dots, \mathbf{V}_C\}$ . We first use the visual encoder to build the visual embedding  $\bar{\mathbf{V}}_i$  for  $i = 1, \dots, C$ . We then generate each sentence auto-regressively using the finetuned language generator. In the decoding step  $l$  for  $\bar{\mathbf{V}}_i$  (i.e., the  $i$ -th output sentence is generated up to  $l - 1$  words), the input to the language generator is  $[\bar{\mathbf{V}}_i, [sep], [dummy], s_{i1}, \dots, s_{i,l-1}]$ . We can obtain the word probability  $p_{il}$  with the output of the language generator  $\tilde{s}_{i,l-1}$ , and finally select the next word  $s_{il} = \arg \max_v p_{il}$ . We iterate this until the end-of-sentence token  $[eos]$  appears, or the output sentence reaches the predefined maximum length.

## 4. Experiments

We evaluate the TAPM approach in two visual storytelling tasks: sequential video captioning in LSMDC 2019 [50] and image captioning in VIST [18]. For both tasks, we achieve new state-of-the-art performance in both automatic evaluation (section 4.2) and human evaluation (section 4.4). We also perform various empirical analyses of our TAPM across various language models (section 4.3). Furthermore, we demonstrate that TAPM can benefit from additional visual-only datasets. TAPM is also extendable to other visual-linguistic tasks such as VQA and cross-modal retrieval, as shown in Appendix.

### 4.1. Experimental Setup

**Datasets.** The Multi-Sentence Description of LSMDC 2019 [50] is the task of generating consecutive captions for multiple short movie clips. For a given set of five clips, the model generates five sentences maintaining logical and contextual consistency. The dataset contains 128,085 clips

Table 1. Quantitative results on the LSMDC 2019 [50] public and blind test set. XE and AREL do not report the blind test score because they are not challenge participants. C stands for CIDEr and M for METEOR. All tests are done on the set level.

Models	Public Test		Blind Test	
	C	M	C	M
Official Baseline [43]	7.0	12.0	6.9	11.9
XE [64]	7.2	11.5	-	-
AREL [64]	7.3	11.4	-	-
TAPM (ours)	<b>10.0</b>	<b>12.3</b>	<b>8.8</b>	<b>12.4</b>

Table 2. Quantitative results on the VIST [18] test set. R stands for ROUGE-L.

Models	C	M	R
Huang et al.[18]	-	31.4	-
h-attn-rank[68]	7.5	34.1	29.5
GLACNet[25]	-	30.1	-
CST[13]	5.1	34.4	29.2
BLEU-RL[64]	8.9	34.6	29.0
CIDEr-RL[64]	8.1	34.9	29.7
GAN[64]	9.1	35.0	29.5
AREL[64]	9.4	35.0	29.5
StoryAnchor[70]	9.9	35.5	30.0
HSRL[17]	10.7	35.2	30.8
INet[21]	10.0	35.6	29.7
TAPM (ours)	<b>13.8</b>	<b>37.2</b>	<b>33.1</b>

Table 3. Ablation results of our TAPM model on the LSMDC 2019 public test set and the VIST test set. The evaluations for LSMDC are done on the sentence level.

Models	LSMDC			VIST		
	C	M	R	C	M	R
Baseline[43]	11.90	8.25	-	-	-	-
Baseline+GPT-2[45]	8.65	7.75	19.90	-	-	-
TAPM (ours)	<b>15.37</b>	8.41	<b>20.21</b>	<b>8.3</b>	<b>34.1</b>	<b>30.2</b>
-A	14.54	8.27	19.89	4.8	33.6	29.9
+Cap	15.29	<b>8.47</b>	20.19	6.7	33.8	29.8
+VisualA	14.59	8.37	20.00	4.9	33.0	29.9
-Split	14.28	8.34	19.71	4.5	32.8	29.8
-A+Split	14.01	8.28	19.60	6.5	33.8	30.0

from 200 movies and has four splits; 20,283 training, 1,486 validation, 2,018 public test, and 1,923 blind test samples. Following the challenge protocol, we combine the train and validation split as training data. The official performance is evaluated on the blind test split hidden from participants, while ablation studies are conducted on the public test split.

VIST [18] is a visual storytelling dataset, including 10,117 Flickr albums with 210,819 unique photos. Each story of VIST contains five sequential images with the corresponding captions. We use the SIS (Stories of Images in Sequence) tier that has more storytelling elements. Ignoring broken images, we use 40,071 training, 4,988 validation, and 5,050 testing story samples. In all experiments, we use the training/test split of [18, 68, 64]. As in [64], we evaluate

Table 4. Comparison between language models on LSMDC 2019 public test set. C, M, and R denote CIDEr, METEOR, and ROUGE-L, respectively. All evaluations are on the sentence level.

Models	No Adaptation			Adaptation (No split-training)			Adaptation (split-training)		
	C	M	R	C	M	R	C	M	R
Baseline [43]	11.90	8.25	-	-	-	-	-	-	-
LSTM-WT2	3.00	5.73	17.13	1.41	4.60	12.83	7.36	8.47	20.40
XLN [14]	10.05	7.09	19.01	7.50	6.95	17.66	13.11	8.00	20.01
GPT [44]	14.01	7.96	19.84	11.81	7.86	19.23	14.76	8.33	20.07
GPT-2	<b>14.54</b>	<b>8.27</b>	<b>19.89</b>	<b>14.28</b>	<b>8.34</b>	<b>19.71</b>	<b>15.37</b>	<b>8.41</b>	<b>20.21</b>

Table 5. (a) Official human evaluation results on the LSMDC 2019 blind test set. Lower is better. (b) Human evaluation results on VIST. Higher is better.

(a) Models		Scores	(b) TAPM vs XE			(b) TAPM vs AREL			
			Choice (%)	TAPM	XE	Tie	TAPM	AREL	Tie
Human		1.085	Relevance	<b>59.9</b>	34.1	6.0	<b>61.3</b>	32.8	5.9
Official Baseline [43]		4.015	Expressiveness	<b>57.3</b>	32.3	10.4	<b>57.3</b>	34.0	8.7
TAPM (ours)		<b>3.670</b>	Concreteness	<b>59.1</b>	30.3	10.7	<b>59.6</b>	30.4	10.0

Table 6. Results with additional visual-only data provided in the adaptation phase. The performance rises with the number of additional videos. C, M and R denotes CIDEr, METEOR and ROUGE-L, respectively.

Models	Videos	C	M	R
Baseline (Ours)	108,487	15.37	8.41	20.21
+ Additional LSMDC	10,053	15.49	8.51	20.26
+ Additional ActivityNet	480,860	16.48	8.67	20.35

at the album level by allowing only one story candidate per album regardless of photo sequences.

**Data Preprocessing.** For LSMDC 2019, we use the ResNet [15] pretrained on ImageNet [51] to extract frame features as in [48, 1]. For the challenge submission and human evaluation, we add the I3D feature [20] pretrained on Kinetics [22] as done in the official baseline [43]. We equally segment a video clip into three subshots and represent each by mean-pooling the features of frames. For the challenge results, we use set level evaluation by concatenating all captions within a set of 5 clips as dictated by the organizers. For ablation study, we use individual sentence level evaluation to compare with non-sequential generation models fairly. For VIST, we use the same ResNet extractor and additional features of object bounding boxes from Faster R-CNN [46] pretrained on Visual Genome [27]. We choose at most 20 objects with the highest likelihood per image from the R-CNN [46] detection results. After processing each feature through the visual encoder, we concatenate all features along the temporal dimension with a special separator token between them. We tokenize and numericalize the text using Byte Pair Encoding [53] for pretrained language models while using the whitespace tokenizer for the no pretrained models. In VIST, we use the default tokenizer to re-tokenize our generated samples for

evaluation. We generate each caption with beam search up to 30 tokens and cut every ground truth sentence to the maximum length of 50 tokens for all experiments.

**Metrics.** We use three n-gram based metrics to evaluate our approach: CIDEr [62], METEOR [2] and ROUGE-L [32]. CIDEr captures consensus by applying Term Frequency Inverse Document Frequency (TF-IDF) weighting for each n-gram. METEOR scores the sequence matches with explicit alignment at the sentence level. ROUGE-L is a recall-based metric computed with the length of the longest common subsequence. For computing METEOR in VIST, we use the official VIST challenge evaluation code <sup>1</sup>. All the other metric scores are computed with the pycocoevalcap library <sup>2</sup>.

**Baselines.** For LSMDC 2019, we compare our approach with the official baseline [50, 43]. We also adapt XE and AREL models [64] to LSMDC using the official codes. For VIST, we compare TAPM with eight state-of-the-art methods: GLACNet [25], h-attn-rank [68], Contextualize, Show and Tell (CST) [13], BLEU-RL [64], CIDEr-RL [64], GAN [64], AREL [64], StoryAnchor [70], HSRL [17], and INet [21]. The scores for BLEU-RL, CIDEr-RL, and AREL are referred from [64], while the results of GLACNet, CST, StoryAnchor, HSRL, and INet are referred from the respective papers. We use XE and AREL as baselines for human evaluation on the VIST dataset. XE shares the architecture of AREL except for the lack of adversarial rewards. We use the publicly available codes for both models.

**Hyperparameters.** Unless we mention it explicitly, we fix all random seeds to 0. For training, we use Adam optimizer [26] with linear learning rate decay. The learning rate is  $5e - 5$ , which is warmed up for the first 4000 steps.

<sup>1</sup><https://github.com/windx0303/VIST-Challenge-NAACL-2018>

<sup>2</sup><https://github.com/tylin/coco-caption>

We apply 0.5 dropout on the language generator outputs. In all experiments, we use the batch size of 8. For LSMDC dataset we train the adaptation loss for 5 epochs, whereas we train for 3 epochs in case of VIST dataset. We train all models up to 30 epochs.

## 4.2. Quantitative Results

We use OpenAI GPT-2 [45] as our default language generator due to its best performance among other language models. We use beam search with a size of 3 for the results in this section and section 4.4 and use a greedy search for the results in section 4.3 for faster computation.

Table 1 outlines the results of sequential video captioning on the LSMDC 2019 blind test set. Our TAPM method outperforms the strong adversarial inference official baseline [43] as well as the XE and AREL model in all metrics. Notably, our method shows significant gaps in the CIDEr metric, which is designed to score human-likeness [62].

Table 2 compares the results of sequential image captioning on the VIST test set. We report the scores computed using only one story per album following previous works. Even without explicitly optimizing the language metrics, our method is competent in the automatic evaluation. In CIDEr, our approach exhibits significant performance gains over the best-performing model AREL [64]. Our model also achieves the highest ROUGE accuracy and on-par METEOR performance with the baselines.

## 4.3. Further Analyses

We perform various empirical analyses of our TAPM model, including (i) ablation study to inspect the contributions of key ingredients and use of (ii) six other language models beyond GPT-2.

**Ablation Study.** We conduct an ablation study for the TAPM model in both LSMDC 2019 and VIST dataset. We test six variants: (i) (-A) removes the adaptation loss training, (ii) (+Cap) uses the ground truth captions instead of the dummy token, (iii) (+VisualA) applies the adaptation loss to the visual encoder output instead of the language generator output, (iv) (-Split) uses naive joint training of the adaptation and generation loss, (v) (-A+Split) is (-A) that uses split-training between the visual encoder and the generator,

Table 3 compares the results of the ablation variants. The performance of TAPM is comparable to that of TAPM+Cap, suggesting that adaptation with videos only is as successful as the supervision with the caption labels.

The slight performance drop from TAPM to TAPM-Split shows that naive joint training can be even worse than training without the adaptation loss. Significant degradation from TAPM-A to TAPM-A+Split proves the split training without the adaptation loss performs the worst. The results of TAPM+VisualA show that applying adaptation loss to visual encoder outputs does not improve the caption quality.

Hence, using language model outputs for adaptation is crucial. Our model, TAPM, performs the best when used as proposed.

Additionally, we replace the backbone of the baseline model [43] from the RNN encoder to GPT-2 pretrained language generator [45]. As shown in the table’s first two rows, the modified model performs even worse than the original baseline. This performance drop verifies our claim that employing a stronger language model does not automatically lead to a better storytelling capability. A stronger textual prior may weaken the visual conditioning when the visually conditioned target data size is insufficient. Without a proper adaptation step, the model would generate less visually relevant captions when using a strong language model such as GPT-2. Hence, the performance improvement of TAPM is attributable to the adaptation step rather than the strength of the language model.

**Other Language Models.** We test the generalization capability of TAPM using three pretrained language models, including LSTM-WT2 [16], XLM [14], and GPT [44]. LSTM [16] is an extension of RNN enlarging its memory capacity. We pretrain an LSTM-based two-layer encoder-decoder architecture on the WikiText-2 dataset [39]. XLM [14] is a multilingual language model designed to exploit both monolingual data and aligned bilingual data. GPT [44] is the predecessor of GPT-2. Table 4 compares the result of different language models. For all models, split-training with the adaptation loss contributes to consistent improvement in the language metrics, while naive joint training results in performance drops in terms of CIDEr and METEOR. These results prove that our TAPM method can improve the visual storytelling performance of a wide range of language models. Furthermore, both the adaptation loss and the split training are necessary to achieve the enhancement.

**Additional Visual-Only Data.** By not relying on ground-truth captions during the adaptation phase, we can exploit additional visual-only data. In Table 6, we perform experiments using additional video-only dataset to further improve TAPM in LSMDC. The generation performance increases along with the number of videos used, indicating that TAPM can use visual-only data to improve cross-modal generation capability.

## 4.4. Human Evaluation Results

We opt for human evaluation to robustly evaluate the captioning quality of our approach. As pointed out in [64], the automatic metrics often fail to capture expressiveness and coherence within a story. Please refer to [64] for details on the limitations of the language metrics for story evaluation.

Table 5 (a) shows human evaluation results conducted by the LSMDC 2019 challenge organizers. For 150 ran-



<b>GT</b>	[female] is scared, she's getting married.	this is her first wedding photo.	and here she is with everyone in the wedding.	this is [male], her husband with location, the best man.	here's [male] and his family.
<b>XE</b>	the <b>bride</b> was so happy to be married.	it was a beautiful day for the wedding.	the bride and groom were very happy to be married.	the bride and groom were very happy to be married.	<b>the bride and groom pose for a picture.</b>
<b>AREL</b>	it was a beautiful day for the wedding.	i had a great time there.	the bride and groom were so happy to be married.	then the bride and groom walked down the aisle.	<b>the bride and groom pose for a picture.</b>
<b>TAPM (Ours)</b>	it was a beautiful day for the wedding.	they <b>went down the stairs to get to the reception.</b>	the bride and groom posed for pictures.	after the ceremony, the <b>groom and the groom's father</b> pose for a picture.	<b>all of the guests</b> were happy to be at the wedding.

Figure 3. Qualitative comparison of sequential image captioning between our method and selected baselines on the VIST dataset. Blue and red fonts indicate correct and erroneous descriptions, respectively. Green shows the coherence between sentences. In the second sentence generated by TAPM, the model explains why the couple is going down the stairs.

dom sets of clips, human annotators rate generated multi-sentence descriptions from 5 (worst) to 1 (best) based on how helpful they are for a blind person to understand what is happening in the movie. To account for variability in human decisions, they aggregate three human judgments per caption and report the median score. We observe that TAPM is superior to the strong adversarial baseline [43].

For VIST, we follow previous research [64] to perform the pairwise comparison test, comparing a pair of generated samples by two methods. We ask human annotators to choose a better story between the two models' outputs for three aspects: relevance, expressiveness, and concreteness. The judges can conclude that the two samples are equally good. We randomly select 150 photo sequences and collect the medians of scores from five workers per test sample. For baselines of XE and AREL, we reproduce the results using the code and parameters provided by the original authors.

Table 5 (b) shows that our TAPM outperforms the baselines in all three aspects by large margins. The performance gain of our model is the most significant in terms of relevance. The gain suggests that the captions generated by TAPM reflect the pictorial narrative better than the baselines since the relevance measures how accurately the story describes what is happening in the image sequence.

#### 4.5. Qualitative Results

Fig. 3 presents a VIST example to compare the captions of TAPM against the baselines. Our generated output can avoid using some wrong words like *bride*, unlike the baselines. Furthermore, TAPM notably captures the causal relationship between the images well. In the second picture, TAPM states the purpose of going down the stairs is to get to the reception and deduces that the ceremony is over with

the third picture. The readers can find more examples in Appendix.

## 5. Conclusion

We proposed the *Transitional Adaptation of Pretrained Model* (TAPM) method for harmonizing the pretrained language model with the visual encoder for vision-to-language generation tasks. Extensive experiments showed that the adaptation phase using the adaptation loss consistently improves the caption quality across several language models and loss types. Our model achieved competitive performance in both automatic metrics and human evaluation for two visual storytelling tasks: the multi-sentence description of LSMDC 2019 and the image storytelling of VIST. There are several directions beyond this work. First, we can explore other adaptation loss types to improve the visual understanding capability of the pretrained language models that have proven their strengths in many language tasks. Second, one can apply our method to other cross-modal generation tasks utilizing the pretrained language models beyond visual storytelling.

**Acknowledgement.** We thank the anonymous reviewers for their thoughtful suggestions on this work. This work was supported by AIRS Company in Hyundai Motor Company & Kia Corporation through HKMC-SNU AI Consortium Fund, Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (No.2017-0-01772, Video Turing Test, No.2019-0-01082, SW StarLab), and the international cooperation program by the NRF of Korea (NRF-2018K2A9A2A11080927). Gunhee Kim is the corresponding author.



## 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*, 2018. 6
- [2] Satyanjeev Banerjee and Alon Lavie. METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments. In *ACL*, 2005. 6
- [3] Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. Scheduled Sampling for Sequence Prediction with Recurrent Neural Networks. In *NIPS*, 2015. 2
- [4] Prajwal Bhargava. Adaptive transformers for learning multimodal representations. In *ACL SRW*, 2020. 2
- [5] James Bradbury, Stephen Merity, Caiming Xiong, and Richard Socher. Quasi-Recurrent Neural Networks. In *ICLR*, 2017. 2, 3
- [6] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. UNITER: UNiversal Image-Text Representation Learning. *arXiv:1909.11740*, 2019. 2
- [7] Jacob Devlin, Hao Cheng, Hao Fang, Saurabh Gupta, Li Deng, Xiaodong He, Geoffrey Zweig, and Margaret Mitchell. Language models for image captioning: The quirks and what works. In *ACL*, 2015. 1
- [8] Sergey Edunov, Alexei Baevski, and Michael Auli. Pre-trained language model representations for language generation. In *NAACL*, 2019. 1
- [9] Angela Fan, Mike Lewis, and Yann Dauphin. Hierarchical Neural Story Generation. In *ACL*, 2018. 2
- [10] Angela Fan, Mike Lewis, and Yann Dauphin. Strategies for structuring story generation. In *ACL*, 2019. 2
- [11] Spandana Gella, Mike Lewis, and Marcus Rohrbach. A dataset for telling the stories of social media videos. In *EMNLP*, 2018. 2
- [12] Xavier Glorot and Yoshua Bengio. Understanding the Difficulty of Training Deep Feedforward Neural Networks. In *AISTATS*, 2010. 3
- [13] Diana Gonzalez-Rico and Gibran Fuentes-Pineda. Contextualize, Show and Tell: A Neural Visual Storyteller. *arXiv:1806.00738*, 2018. 5, 6
- [14] Alexis Conneau Guillaume Lample. Cross-lingual Language Model Pretraining. In *NIPS*, 2019. 2, 3, 6, 7
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *CVPR*, 2016. 6
- [16] Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. *Neural Computation*, 1997. 3, 7
- [17] Qiuyuan Huang, Zhe Gan, Asli Celikyilmaz, Dapeng Wu, Jianfeng Wang, and Xiaodong He. Hierarchically Structured Reinforcement Learning for Topically Coherent Visual Story Generation. In *AAAI*, 2019. 2, 5, 6
- [18] Ting-Hao (Kenneth) Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Aishwarya Agrawal, Jacob Devlin, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, C. Lawrence Zitnick, Devi Parikh, Lucy Vanderwende, Michel Galley, and Margaret Mitchell. Visual Storytelling. In *NAACL-HLT*, 2016. 1, 2, 5
- [19] Parag Jain, Priyanka Agrawal, Abhijit Mishra, Mohak Sukhwani, Anirban Laha, and Karthik Sankaranarayanan. Story Generation from Sequence of Independent Short Descriptions. In *SIGKDD Workshop on Machine Learning for Creativity (MLCreativity)*, 2017. 2
- [20] Andrew Zisserman Joao Carreira. Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset. In *CVPR*, 2017. 6
- [21] Yunjae Jung, Dahun Kim, Sanghyun Woo, Kyungsu Kim, Sungjin Kim, and In So Kweon. Hide-and-Tell: Learning to Bridge Photo Streams for Visual Storytelling. In *AAAI*, 2020. 5, 6
- [22] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew Zisserman. The Kinetics Human Action Video Dataset. *arXiv:1705.06950*, 2017. 6
- [23] Gunhee Kim, Leonid Sigal, and Eric P. Xing. Joint Summarization of Large-scale Collections of Web Images and Videos for Storyline Reconstruction. In *CVPR*, 2014. 2
- [24] Gunhee Kim and Eric P. Xing. Jointly Aligning and Segmenting Multiple Web Photo Streams for the Inference of Collective Photo Storylines. In *CVPR*, 2013. 2
- [25] 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. *CoRR*, 2018. 5, 6
- [26] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In *ICLR*, 2015. 6
- [27] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, Michael Bernstein, and Li Fei-Fei. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations. *IJCV*, 2017. 6
- [28] Jie Lei, Liwei Wang, Yelong Shen, Dong Yu, Tamara L. Berg, and Mohit Bansal. Mart: Memory-augmented recurrent transformer for coherent video paragraph captioning. In *ACL*, 2020. 2
- [29] Gen Li, Nan Duan, Yuejian Fang, Ming Gong, Daxin Jiang, and Ming Zhou. Unicoder-VL: A Universal Encoder for Vision and Language by Cross-modal Pre-training. In *AAAI*, 2020. 2
- [30] Lijun Li and Boqing Gong. End-to-End Video Captioning with Multitask Reinforcement Learning. In *WACV*, 2019. 2
- [31] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. What does bert with vision look at? In *ACL (short)*, 2020. 2
- [32] Chin-Yew Lin. ROUGE: A Package for Automatic Evaluation of Summaries. In *WAS*, 2004. 6
- [33] Siqi Liu, Zhenhai Zhu, Ning Ye, Sergio Guadarrama, and Kevin Murphy. Improved Image Captioning via Policy Gradient optimization of SPIDER. In *ICCV*, 2017. 2
- [34] Xihui Liu, Hongsheng Li, Jing Shao, Dapeng Chen, and Xiaogang Wang. Show, tell and discriminate: Image captioning by self-retrieval with partially labeled data. In *ECCV*, 2018. 1

- [35] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks. In *NIPS*, 2019. 2
- [36] Ruotian Luo, Brian Price, Scott Cohen, and Gregory Shakhnarovich. Discriminability objective for training descriptive captions. In *CVPR*, 2018. 1
- [37] Chih-Yao Ma, Yannis Kalantidis, Ghassan AlRegib, Peter Vajda, Marcus Rohrbach, and Zsolt Kira. Learning to Generate Grounded Image Captions without Localization Supervision. *arXiv:1906.00283*, 2019. 2
- [38] Andrew L. Maas, Awni Y. Hannun, and Andrew Y. Ng. Rectifier Nonlinearities Improve Neural Network Acoustic Models. In *ICML*, 2013. 3
- [39] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer Sentinel Mixture Models. In *ICLR*, 2017. 7
- [40] Pingbo Pan, Zhongwen Xu, Yi Yang, Fei Wu, and Yueting Zhuang. Hierarchical Recurrent Neural Encoder for Video Representation with Application to Captioning. In *CVPR*, 2016. 1
- [41] Cesc Chunseong Park and Gunhee Kim. Expressing an Image Stream with a Sequence of Natural Sentences. In *NIPS*, 2015. 2
- [42] Cesc Chunseong Park, Youngjin Kim, and Gunhee Kim. Retrieval of Sentence Sequences for an Image Stream via Coherence Recurrent Convolutional Networks. *IEEE TPAMI*, 2018. 2
- [43] Jae Sung Park, Marcus Rohrbach, Trevor Darrell, and Anna Rohrbach. Adversarial Inference for Multi-Sentence Video Description. In *CVPR*, 2019. 2, 5, 6, 7, 8
- [44] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving Language Understanding by Generative Pre-Training. 2018. 3, 6, 7
- [45] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language Models are Unsupervised Multitask Learners. 2019. 2, 3, 5, 7
- [46] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, 2015. 6
- [47] Zhou Ren, Xiaoyu Wang, Ning Zhang, Xutao Lv, and Li-Jia Li. Deep Reinforcement Learning-based Image Captioning with Embedding Reward. In *CVPR*, 2017. 2
- [48] Steven J. Rennie, Etienne Marcheret, Youssef Mroueh, Jarret Ross, and Vaibhava Goel. Self-Critical Sequence Training for Image Captioning. In *CVPR*, 2017. 2, 6
- [49] Anna Rohrbach, Marcus Rohrbach, and Bernt Schiele. The Long-Short Story of Movie Description. In *DAGM*, 2016. 1
- [50] Anna Rohrbach, Atousa Torabi, Marcus Rohrbach, Niket Tandon, Chris Pal, Hugo Larochelle, Aaron Courville, and Bernt Schiele. Movie description. *IJCV*, 2017. 1, 2, 5, 6
- [51] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *IJCV*, 2015. 6
- [52] Ramprasaath R. Selvaraju, Stefan Lee, Yilin Shen, Hongxia Jin, Dhruv Batra, and Devi Parikh. Taking a HINT: Leveraging Explanations to Make Vision and Language Models More Grounded. In *ICCV*, 2019. 2
- [53] Rico Sennrich, Barry Haddow, , and Alexandra Birch. Neural Machine Translation of Rare Words with Subword Units. In *ACL*, 2016. 6
- [54] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. Vi-bert: Pre-training of generic visiolinguistic representations. In *ICLR*, 2020. 2
- [55] Chen Sun, Fabien Baradel, Kevin Murphy, and Cordelia Schmid. Learning Video Representations using Contrastive Bidirectional Transformer. *arXiv:1906.05743*, 2019. 2
- [56] Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. VideoBERT: A Joint Model for Video and Language Representation Learning. In *ICCV*, 2019. 2
- [57] Hao Tan and Mohit Bansal. LXMERT: Learning Cross-Modality Encoder Representations from Transformers. In *EMNLP*, 2019. 2
- [58] Jinhui Tang, Jing Wang, Zechao Li, Jianlong Fu, and Tao Mei. Show, Reward, and Tell: Adversarial Visual Story Generation. In *AAAI*, 2018. 2
- [59] Marc Tanti, Albert Gatt, and Kenneth P. Camilleri. Transfer learning from language models to image caption generators: Better models may not transfer better. *arxiv:1901.01216*, 2019. 1
- [60] John N. Tsitsiklis and Benjamin Van Roy. Analysis of temporal-difference learning with function approximation. *NIPS*, 1996. 2
- [61] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention Is All You Need. In *NIPS*, 2017. 3
- [62] Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. CIDEr: Consensus-based Image Description Evaluation. In *CVPR*, 2015. 6, 7
- [63] Subhashini Venugopalan, Marcus Rohrbach, Jeff Donahue, Raymond Mooney, Trevor Darrell, and Kate Saenko. Sequence to Sequence – Video to Text. In *ICCV*, 2015. 1
- [64] Xin Wang, Wenhui Chen, Yuan-Fang Wang, and William Yang Wang. No Metrics Are Perfect: Adversarial Reward Learning for Visual Storytelling. In *ACL*, 2018. 5, 6, 7, 8
- [65] Xin Wang, Wenhui Chen, Jiawei Wu, Yuan-Fang Wang, and William Yang Wang. Video Captioning via Hierarchical Reinforcement Learning. In *CVPR*, 2018. 2
- [66] Jiacheng Yang, Mingxuan Wang, Hao Zhou, Chengqi Zhao, Weinan Zhang, Yong Yu, and Lei Li. Towards making the most of bert in neural machine translation. In *AAAI*, 2020. 1
- [67] Li Yao, Atousa Torabi, Kyunghyun Cho, Nicolas Ballas, Christopher Pal, Hugo Larochelle, and Aaron Courville. Describing Videos by Exploiting Temporal Structure. In *ICCV*, 2015. 1
- [68] Licheng Yu, Mohit Bansal, and Tamara L. Berg. Hierarchically-Attentive RNN for Album Summarization and Storytelling. In *EMNLP*, 2017. 5, 6
- [69] Youngjae Yu, Hyungjin Ko, Jongwook Choi, and Gunhee Kim. End-to-end Concept Word Detection for Video Captioning, Retrieval, and Question Answering. In *CVPR*, 2017. 1

- [70] Bowen Zhang, Hexiang Hu, and Fei Sha. The Steep Road to Happily Ever After: An Analysis of Current Visual Storytelling Models. In *ICCV19 CLVL workshop: 3rd Workshop on Closing the Loop Between Vision and Language*, 2019. 5, 6
- [71] Li Zhang, Flood Sung, Feng Liu, Tao Xiang, Shaogang Gong, Yongxin Yang, and Timothy M. Hospedales. Actor-Critic Sequence Training for Image Captioning. In *NIPS*, 2017. 2
- [72] Chen Zheng, Quan Guo, and Parisa Kordjamshidi. Cross-modality relevance for reasoning on language and vision. In *ACL*, 2020. 2
- [73] Luowei Zhou, Yannis Kalantidis, Xinlei Chen, Jason J. Corso, and Rohrbach Marcus. Grounded Video Description. In *CVPR*, 2019. 2
- [74] Luowei Zhou, Yingbo Zhou, Jason J. Corso, Richard Socher, and Caiming Xiong. End-to-End Dense Video Captioning with Masked Transformer. In *CVPR*, 2018. 1