
BYPASS NETWORK FOR SEMANTICS DRIVEN IMAGE PARAGRAPH CAPTIONING

Qi Zheng¹, Chaoyue Wang², and Dadong Wang³

¹University of Sydney, NSW 2008, Australia

²JD Explore Academy

³DATA61, CSIRO, NSW 2122, Australia

ABSTRACT

Image paragraph captioning aims to describe a given image with a sequence of coherent sentences. Most existing methods model the coherence through the topic transition that dynamically infers a topic vector from preceding sentences. However, these methods still suffer from immediate or delayed repetitions in generated paragraphs because (i) the entanglement of syntax and semantics distracts the topic vector from attending pertinent visual regions; (ii) there are few constraints or rewards for learning long-range transitions. In this paper, we propose a bypass network that separately models semantics and linguistic syntax of preceding sentences. Specifically, the proposed model consists of two main modules, i.e. a topic transition module and a sentence generation module. The former takes previous semantic vectors as queries and applies attention mechanism on regional features to acquire the next topic vector, which reduces immediate repetition by eliminating linguistics. The latter decodes the topic vector and the preceding syntax state to produce the following sentence. To further reduce delayed repetition in generated paragraphs, we devise a replacement-based reward for the REINFORCE training. Comprehensive experiments on the widely used benchmark demonstrate the superiority of the proposed model over the state of the art for coherence while maintaining high accuracy.

1 Introduction

Visual captioning explores machines' capability to comprehend the visual world by describing it with natural language. Traditional image captioning [18, 38, 23, 43, 36, 35, 41, 24, 31] depicts an image with a single sentence that covers the whole scene or the most salient object. Dense image captioning [11, 40, 44, 14] describes all the detected salient regions without regard to the order of descriptions. Later, [16] introduced image paragraph captioning that aims to give a coherent fine-grained paragraph.

The primary prerequisite of generating a coherent paragraph description for an image is to model coherence by its linguistic definition. Referring to linguistics [4, 7], the coherence exists when nearby sentences (a sentence-cluster) are "about" someone or something, and the whole paragraph does not jump back or forth among multiple entities. It indicates that the coherence flowing through a paragraph can be modeled by topic transition among sentences. Most existing methods [21, 3, 37, 16] attempt to guide the topic transition by dynamically inferring a topic vector at the beginning of each sentence. For instance, conditioning on preceding sentences, [21] transfers attended visual features into a topic vector. In the work of [37], the topic vector is predicted from the state of a hierarchical attention module, which integrates both regional features and preceding sentences.

Compared with directly transferring image captioning to paragraph captioning, these topic vector-based methods can generate more relative and diverse sentences owing to considering the information contained in the former sentence. However, directly using the former sentence as conditional input signals may only achieve a semblance of coherence. In other words, though maintaining linguistic consistency, the generated sentences fail to capture the variance in descriptive content. As a result, these methods still suffer from immediate or delayed repetition in synthesized paragraphs. In *immediate repetition*, the subsequent sentence is a duplicate of the previous one, either by copying words or expressing

the same stuff. As for *delayed repetition*, the subsequent sentence may go back to earlier descriptive content regardless of the coherence within the generated partial paragraph.

According to our observation, the topic transition that is directly controlled by preceding sentences fails to guarantee the desired coherence due to two reasons. On the one hand, such a transition does not distinguish the linguistic consistency from the content variance, which results in a compromise. Ideally, the topic transition should be driven by the development of descriptive content. However, in existing methods, the entanglement of syntax and semantics in preceding sentences distracts the topic vector from attending pertinent visual regions. On the other hand, most existing methods lack direct supervision to guide the long-range topic transition. The widely used maximum likelihood estimation promotes accuracy in word prediction but provides little feedback to sentence generation in a given context.

In this paper, we propose a bypass network to solve these issues. The proposed bypass structure separates semantics from the linguistic syntax of preceding sentences to produce a topic vector. The semantic stream drives topic transition, and the syntax is maintained by a bypass Part-Of-Speech (POS) stream. The decoder integrates these two elements to generate the following sentence. In this way, the topic vector can attend more accurately to pertinent visual regions and reduce immediate repetition. Given the simplicity of the proposed model, we provide a detailed deduction of the disentangling property in the appendix material. In addition, we devise a replacement-based reward for REINFORCE training. It enhances the quality of long-range transitions and alleviates the delayed repetition in generated paragraphs. Finally, we conduct comprehensive experiments to demonstrate the efficacy of the proposed model. The main contributions of our paper are listed below:

- We propose a Bypass network to model topic transition in image paragraph captioning. It drives topic transition by maintaining a semantic stream and reduces immediate repetition.
- We further devise a replacement-based reward to alleviate delayed repetition, which takes effect during REINFORCE training.
- We compare paragraph-level and sentence-level captioning performance to assess coarse- and fine-grain coherence and accuracy, respectively.

2 Related Work

2.1 Image Captioning and Dense Captioning

The image captioning task aims to depict a given image using one sentence that covers the semantic meaning of the whole scene or describes a salient object in the scene. Recent deep learning-based approaches [5, 13, 1, 10, 23, 33, 38] generally employ an encoder-decoder framework, which first extracts features from the image using CNNs [19] and then inputs them to a Recurrent Neural Network such as LSTMs [6]. Visual features range from CNN features, grid features to object proposals. There are also some attempts that introduce object attributes or visual relationships to boost the semantic quality of the generated caption [43, 42]. As for the decoder, attention mechanism [38, 1] has been introduced as a popular tool to attend and visualize where the model “sees” when it depicts.

Dense image captioning was introduced to enrich the descriptive details in an image. [11] propose a fully convolutional localization network (FCLN) that detects object regions and generates regional annotations in a single forward pass. To reduce redundancy in detected bounding boxes and increase semantic saliency of object regions, [40] incorporate joint inference and context fusion in the localization and captioning process, where the image-level feature served as the contextual cues. To capture more details during the captioning, [44] further utilize the neighboring feature and the attribute of each object region.

2.2 Image Paragraph Captioning

Paragraph descriptions can be generated either as a “big” sentence or as a sequence of sentences. The former works in an image-captioning way, where the most significant problem is the internal repetition. Upon the well-known Up-Down [1] model, [26] suggest suppressing the repetition by blocking the trigrams that have appeared before, denoted by TDC (Training for Diversity in image paragraph Captioning). This simple yet efficient technique gains a vast improvement in generated sentences with respect to automatic evaluation metrics. TEB [8] further generates a paragraph-vector using a Text Embedding Bank (TEB) to guide the captioning process.

We conclude that the success of TDC is twofold. First, since the decoder is trained to highlight words that co-occur frequently in the training corpus, suppression on the trigrams that have been generated gives way to the remaining highly correlated words. This way, TDC avoids producing repeated words while ensuring the quality of predicted words. Second, existing automatic evaluation metrics on image paragraph captioning take the generated description as one

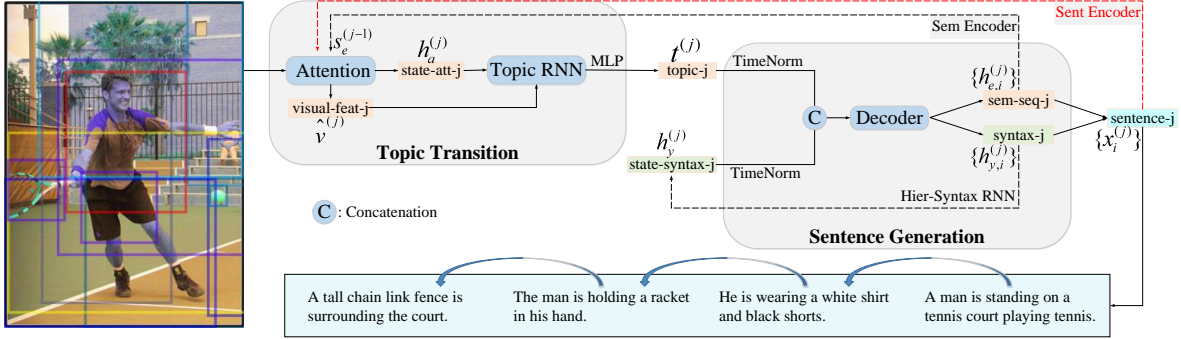


Figure 1: Illustration of generating the j -th sentence by our Bypass model. The model consists of i) the Topic-Transition module that predicts a topic vector based on preceding semantic stream, and ii) the Sentence Generation module that synthesizes a descriptive sentence given the topic vector and a syntactic state.

“big” sentence, where the order of sentences makes no difference. However, a primary concern is that sentence-level coherence is not guaranteed in the highly-scoring TDC.

Alternatively, many recent paragraph captioning methods learn to generate a sequence of sentences. For instance, [16] introduce a hierarchical framework (Regions-Hierarchical), where the sentence-level RNNs produce K topic vectors for K sentences and predict the probability that each vector indicates the ending of the paragraph. Unlike Regions-Hierarchical, [21] propose a generator where the two-level RNNs work as a whole and adopt a sentence discriminator and a topic-transition discriminator to assess sentence plausibility and topic coherence. [3] produce all the topic vectors in advance similar to Regions-Hierarchical, but then combine the topic vector with the semantics of the preceding sentence as the input to the word-level decoder. [25] introduce the LDA-Topic model to generate a topic distribution for each sentence, where the topic is represented by an average of several discrete abstract conceptual embeddings clustered from the training corpus. To increase the supervision of topic transition, [37] propose a densely supervised hierarchical policy-value (DHPV) network that dynamically assesses the contribution of the currently generated sentence or word. [39] use a hierarchical CNN architecture (i.e. ParaCNN) to extract contextual information between sentences and generate a visual paragraph. Similarly, Li et al [20] also use CNNs for paragraph captioning.

Though these approaches learn sentence-level topic transition, most of them ignore the gap between linguistic consistency and content transition and suffer from immediate or delayed repetition. Our method follows the branch of modeling topic transition among all generated sentences, but differently, we separate the semantics from the syntax of preceding sentences in predicting the topic vector, which reduces the immediate repetition. Furthermore, we devise a replacement-based reward to alleviate the delayed repetition.

3 Method

To enhance topic transition, we design a bypass network that separates semantics from syntax in producing the topic vector and combines them in sentence generation. We first provide an overview of the proposed method and then explain the two parts of our model in detail, namely topic transition and sentence generation. Finally, we introduce the optimization process and discuss the replacement-based reward.

3.1 Overview

Similar to [16, 21, 26], we first detect object regions in a given image I using offline object detectors. This produces a set of regional features $\{v_1, v_2, \dots, v_K\}$, as the input to the pipeline in Fig. 1. At the j -th time step, the Topic-Transition module first predicts an attention distribution $\{a_1, a_2, \dots, a_K\}^{(j)}$ over the regional features, conditioned on the semantics $s_e^{(j-1)}$ from the previous sentence. This gives the attended visual feature $\hat{v}^{(j)}$ and also the state $h_a^{(j)}$ of the attention network.

The Topic RNN takes in these two elements as input and produces a topic hidden state $h_t^{(j)}$. The topic vector t_j is generated by applying an MLP block on state $h_t^{(j)}$. The Sentence-Generation module applies the proposed Time-Normalization on the topic vector t_j and the syntax state $h_y^{(j)}$ from preceding sentences, respectively. Then the results are concatenated and input into the decoder to produce a sequence of words. At each step, the hidden state from the

decoder evolves into two streams, i.e. the semantic stream and the syntactic stream. The semantic sequence is encoded by Semantic Encoder and forms the condition to the Topic-Transition module, while the syntactic stream is encoded by hierarchical-syntax (Hier-Syntax) Encoder and outputs the syntactic state of the $(j+1)$ -th sentence.

3.2 Topic Transition

The Topic-Transition module determines the semantic content of the subsequent sentence, which can be seen as an extension of the Up-Down model [1]. In traditional image captioning, the Up-Down model dynamically predicts attention distribution in generating each word.

Differently, in image paragraph captioning, the dynamic is generally in sentence level, i.e. all the words in one sentence share the same visual attention. Specifically, given the regional features $\{v_1, v_2, \dots, v_K\}$, the Attention RNN predicts the attention distribution for the j -th sentence,

$$h_a^{(j)} = \text{RNN}_a([s_e^{(j-1)}; \bar{v}], h_a^{(j-1)}), \quad (1)$$

$$a_k^{(j)} = \text{softmax}\left(w_a^T \tanh(W_{va}v_k + W_{ha}h_a^{(j)})\right), \quad (2)$$

where $s_e^{(j-1)}$ is the semantic information of the $(j-1)$ -th sentence and $\bar{v} = 1/K \sum_k v_k$ is the mean-pooled regional feature. This gives the attended visual feature $\hat{v}^{(j)} = \sum a_k^{(j)} v_k$ for generating the j -th sentence.

Instead of directly being the topic vector as in the Up-Down model, $\hat{v}^{(j)} = \sum a_k^{(j)} v_k$ further passes a Topic RNN followed by an MLP block, and finally produces the topic vector $t^{(j)}$ and a probability $p^{(j)}$ of the j -th sentence being the last one,

$$h_t^{(j)} = \text{RNN}_t([\hat{v}; h_a^{(j)}], h_t^{(j-1)}), \quad (3)$$

$$t^{(j)} = \text{MLP}(h_t^{(j)}) \quad \text{and} \quad p^{(j)} = \text{FC}(h_t^{(j)}). \quad (4)$$

The necessity of the Topic RNN comes from the fact that some descriptions in a paragraph may depend more on an abstract concept than the specific objects in a scene, which is beyond the coverage of the attended visual features. More details and the verification are provided in our experimental section. It is worth noting that this module is very close to the hierarchical policy network in [37], except that the information of the preceding sentence is from a semantic stream. We consider this a simple yet efficient baseline module and made minimum modifications to it.

3.3 Sentence Generation

To reduce internal covariate shift of the vectors from different time steps, inspired by the well-known Batch Normalization [9], we propose Time Normalization for the topic vector $t^{(j)} \in \mathbb{R}^D$ and the syntax state $h_y^{(j)} \in \mathbb{R}^D$,

$$\hat{t}^{(j)} = \gamma_t \frac{t^{(j)} - \mu_t^{(j)}}{\sqrt{\sigma_t^{(j)} + \epsilon}}, \quad \hat{h}_y^{(j)} = \gamma_y \frac{h_y^{(j)} - \mu_y^{(j)}}{\sqrt{\sigma_y^{(j)} + \epsilon}}, \quad (5)$$

where $\gamma, \beta \in \mathbb{R}^D$ are parameters to be learned. During multiple time steps, the mean of the topic vector is $\mu_t^{(j)} = 1/j \sum_{m=0}^j t^{(m)}$ and the variance is $\sigma_t^{(j)} = 1/j \sum_{m=0}^j (t^{(m)} - \mu_t^{(j)})^2$, which is similar for the syntax state. For efficiency, the mean and the variance are updated in an iterative way,

$$\mu_t^{(j)} = \frac{1}{j} \left((j-1)\mu_t^{(j-1)} + t^{(j)} \right), \quad (6)$$

$$\sigma_t^{(j)} = \frac{j-1}{j} \left(\sigma_t^{(j-1)} + \frac{(\mu_t^{(j-1)} - t^{(j)})^2}{j} \right). \quad (7)$$

Given the normalized topic vector $\hat{t}^{(j)}$ and syntax state $\hat{h}_y^{(j)}$, the decoder produces a hidden state for each word. Then the hidden state $h_{d,i}^{(j)}$ splits into two branches, i.e. the semantic stream $h_{e,i}^{(j)} = \text{MLP}(h_{d,i}^{(j)})$ and the syntax stream $h_{y,i}^{(j)} = \text{MLP}(h_{d,i}^{(j)})$, and the word $x_i^{(j)}$ is predicted by

$$x_i^{(j)} = \text{softmax}\left(\text{FC}(h_{e,i}^{(j)} + h_{y,i}^{(j)})\right). \quad (8)$$

The sequence of semantic elements $h_{e,i}^{(j)}$ is encoded to form the semantic vector $s_e^{(j)} = \text{RNN}_e(h_{e,i}^{(j)})$, as the $(j+1)$ -th input to the Topic-Transition module. The sequence of syntax elements $h_{y,i}^{(j)}$ is decoded into POS tags $y_i^{(j)}$ of the j -th sentence, which then passes a sentence-level (Hier-Syntax) RNN to produce the unnormalized syntax state $h_y^{(j+1)}$.

3.4 Optimization

3.4.1 Training by cross-entropy loss.

Given the available text \mathbf{x}_* and POS tags \mathbf{y}_* , gram-level cross-entropy loss is applied to the generated word-sequence and tag-sequence, respectively,

$$\ell_x = -\log(P(\mathbf{x}^{(j)} = \mathbf{x}_*^{(j)})) = -\sum_i \log(P(x_i^{(j)} = x_{i,*}^{(j)})) \quad (9)$$

$$\ell_y = -\log(P(\mathbf{y}^{(j)} = \mathbf{y}_*^{(j)})) = -\sum_i \log(P(y_i^{(j)} = y_{i,*}^{(j)})) \quad (10)$$

A sentence-level cross-entropy loss is applied to the probability $p^{(j)}$ that indicates the ending sentence,

$$\ell_p = -\sum_j \log(P(p^{(j)} = p_*^{(j)})), \quad (11)$$

where $p_*^{(j)} = 1$ if the j -th sentence is the last one, otherwise $p_*^{(j)} = 0$. Finally, the total loss is a weighted sum of the three items,

$$\mathcal{L}_{\text{xe}} = \ell_x + \lambda_y \ell_y + \lambda_p \ell_p, \quad (12)$$

where λ_y and λ_p are hyper-parameters.

It is worth noting that we add no extra constraint on $h_{e,i}^{(j)}$ and $h_{y,i}^{(j)}$ to reduce their correlation. We proved that minimizing the loss \mathcal{L}_{xe} is equal to minimizing the off-diagonal entries of the function's Hessian matrix *w.r.t.* its inputs \mathbf{x} and \mathbf{y} , which encourages the disentanglement between the them. Therefore, no extra constraint is compulsory to reduce the correlation of $h_{e,i}^{(j)}$ and $h_{y,i}^{(j)}$.

3.4.2 Training by REINFORCE algorithm.

From the viewpoint of Reinforcement Learning, topic transition in image paragraph captioning can be modeled as a Markov reward process, because the transition (from *state* to *action*) probability is deterministic. In [37], an incremental sentence-level reward $r^{(j)}$ is proposed to assess the contribution of each sentence to the description paragraph,

$$r(\mathbf{y}^{(j)}) = \phi(\mathbf{y}^{(1:j)}, G) - \phi(\mathbf{y}^{(1:j-1)}, G), \quad (13)$$

where $\phi(\cdot)$ is a function that returns the value of the generated partial paragraph by comparing it with the annotated paragraph G .

For further training, an efficient solution is using the REINFORCE algorithm [34] to maximize the expected total reward,

$$\mathcal{J}(\theta) = \sum_j \sum_{\mathbf{y}^{(j)}} p_\theta(\mathbf{y}^{(j)} | v_{1:K}, \mathbf{y}^{(1:j-1)}) r(\mathbf{y}^{(j)}), \quad (14)$$

where θ is the set of parameters that have been defined. For simplicity, we denote $z = [v_{1:K}, \mathbf{y}^{(1:j-1)}]$ and then derive

$$\begin{aligned} \nabla_\theta \mathcal{J}(\theta) &= \sum_j \sum_{\mathbf{y}^{(j)}} r(\mathbf{y}^{(j)}) \nabla_\theta \log p_\theta(\mathbf{y}^{(j)} | z) \\ &\approx \sum_j \frac{1}{N} \sum_{n=1}^N r(\mathbf{y}^{(j,n)}) \nabla_\theta \log p_\theta(\mathbf{y}^{(j,n)} | z), \end{aligned} \quad (15)$$

where $\mathbf{y}^{(j,n)}$, $1 \leq n \leq N$ are the sentence sampled by Monte Carlo methods [12].

Replacement-based reward. One of the main drawbacks of the vanilla algorithm in Eq. (15) is its high variance. In the well-known SCST [30] algorithm for single-sentence captioning, a greedy-search baseline is introduced to reduce the variance. We could similarly set a baseline for the reward in Eq. (13) as $r(y^{(j)}) - r_b(y^{(j)})$,

$$r_b(y^{(j)}) = \phi([\mathbf{y}^{(1:j-1)}; \hat{\mathbf{y}}^{(j)}], G) - \phi(\mathbf{y}^{(1:j-1)}, G), \quad (16)$$

where $\hat{\mathbf{y}}^{(j)}$ is sampled by greedy search from the state policy $p_\theta(\mathbf{y}^{(j)}|z)$. However, we observed that such high variance comes much more from the reward itself. Specifically, an obvious difference of the length (i.e. numbers of words) between the generated partial paragraph $\mathbf{y}^{(1:j)}$ and the annotated paragraph G seriously undermines the estimation of the real contribution of $\mathbf{y}^{(j)}$.

On the other hand, Eq. (13) only counts the topic transition between adjacent sentences, i.e. $\mathbf{y}^{(j)}$ and $\mathbf{y}^{(j-1)}$, which is effective for immediate repetition but not the delayed repetition. Therefore, we devise a replacement-based reward to address both issues,

$$r(\mathbf{y}^{(j)}) = \phi(\mathbf{y}^{(1:j)}, \mathbf{y}_*^{(1:j)}) - \phi(\tilde{\mathbf{y}}^{(1:j)}, \mathbf{y}_*^{(1:j)}), \quad (17)$$

where $\tilde{\mathbf{y}}^{(1:j)} = [\mathbf{y}^{(1:j-1)}; \tilde{\mathbf{y}}^j]$ is the mixed partial description and $\phi(\tilde{\mathbf{y}}^{(1:j)}, \mathbf{y}_*^{(1:j)})$ serves as the baseline reward. $\tilde{\mathbf{y}}^j = \text{sample}(\mathbf{y}^{1:j-1})$ is the j -th sentence randomly sampled from previously generated $j-1$ sentences. This can also be optimized by Eq. (15).

4 Experiments

We conduct experiments on the Stanford image-paragraph dataset¹ [16]. This dataset contains 19,551 images from MS COCO dataset [22] and Visual Genome [17], each of which is labeled with one paragraph of about 67.50 words. Each paragraph consists of multiple sentences with 11.91 words/sentence on average. It has been officially divided into training (14,575), validation (2,487) and test (2,489) splits.

4.1 Implementation Details

Following [1], we use the Faster R-CNN [29] object detector to obtain $K = 36$ object proposals. The dimension of extracted regional features is 2,048. Unless otherwise specified, all the RNNs are one-layer GRU blocks, with 512 hidden units. Similar to [16, 21, 37], we set the number of sentences to 6 and the maximum words in a sentence to 30. All the MLP blocks consist of two linear layers, connected by one ReLU activation layer and one dropout layer. The dropout probability during training is 0.5. The dimension of word embeddings is 512, which are trained from scratch. We replace the words that appear less than four times in the dataset with ‘unk’ token and build a vocabulary of 4636 words. POS tags of the annotated paragraphs are extracted by NLTK tools² and then merged into 17 categories. The embedding size of POS tags is set as 64.

We train the proposed model to minimize Eq. (12) using the Adam [15] optimizer with an initial learning rate of 0.0005 for 70 epochs. We set the weights of ℓ_x , ℓ_y and ℓ_p as 5:5:1. The batch size throughout training is set to 10. The learning rate decays every three epochs, with a decay rate of 0.8. We further train the model to maximize Eq. (14) for 30 epochs with the same decay rate. One sample is used for the approximation in Eq. (15). We fix the beam-search size at 1 for our model throughout the following experiments.

4.2 Comparison Methods and Metrics

We compare our model with previous approaches reported in [16], including Sentence-Concat, Template, DenseCap-Concat and Image-Flat [13] and the latest state-of-the-art methods, including Regions-Hier [16], Up-Down [1], RTT-GAN [21], TDC [26], ParaCNN [39] and TEB [8]. Since the pre-trained models are not provided, we followed the officially released codes for training TDC³ [26], ParaCNN⁴ [39], and TEB⁵ [8]. Codes for the remaining works are unavailable, so we implemented the basic form of DHPV as our baseline model. The implementation follows the original paper [37] and training is under the same setting for a fair comparison.

¹<https://cs.stanford.edu/people/ranjaykrishna/im2p/index.html>

²<http://www.nltk.org/book/ch05.html>

³<https://github.com/lukemelas/image-paragraph-captioning>

⁴<https://github.com/Shiyang-Yan/ParaCNN>

⁵<https://github.com/arjung128/image-paragraph-captioning>

	METEOR	CIDEr*	Bleu@4	Bleu@3	Bleu@2	Bleu@1
TDC [26]	8.81 \pm 0.17	72.20 \pm 1.89	3.82 \pm 0.13	6.47 \pm 0.18	11.17 \pm 0.25	21.39 \pm 0.41
TEB [8]	8.10 \pm 0.17	66.10 \pm 1.75	3.63 \pm 0.13	6.07 \pm 0.17	10.50 \pm 0.24	19.84 \pm 0.31
ParaCNN [39]	8.52 \pm 0.16	61.06 \pm 1.54	3.89 \pm 0.11	6.87 \pm 0.16	12.11 \pm 0.22	23.11 \pm 0.29
Bypass	10.48 \pm 0.16	73.07 \pm 1.48	3.94 \pm 0.11	7.09 \pm 0.16	12.65 \pm 0.23	24.44 \pm 0.29

Table 1: Sentence-level scores of different methods. The margin of error (i.e. $\pm\epsilon$) is calculated at confidence interval 95%.

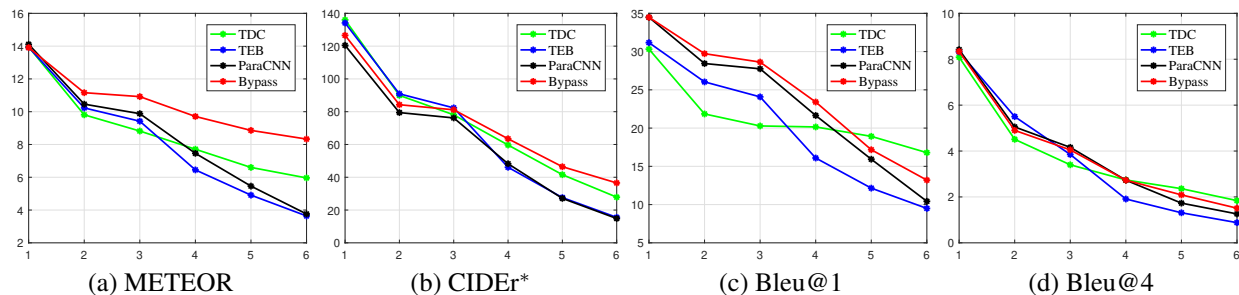


Figure 2: Sentence-level score curves of different methods. The index of the x-axis denotes the first to the sixth sentence.

We report six automatic evaluation metrics, i.e. Bleu@1-4 [27], METEOR [2] and CIDEr [32]. Among them, Bleu@1-4 scores measure the precision of n-gram matching. METEOR assesses the unigram matching based on the surface form, stemmed form, and meaning of the unigrams. CIDEr evaluates the consensus of generated captions based on TF-IDF weighting for each n-gram. First, following existing sentence/dense captioning methods, we treat the generated paragraph as a big sentence to compare these metrics. However, such measurement is insufficient to reflect the coherence within the generated paragraph, since it ignores the order of sentences within a paragraph. Therefore, we further measure sentence-level scores of these six metrics. Specifically, each sentence in a generated paragraph is compared with the corresponding sentence of the annotated paragraph, and the average over six sentences gives the final (Bleu@n, METEOR or CIDEr) score of a sample. Note that, since the distribution error of document frequency in the single-sentence ground truth, we replace CIDEr with CIDEr* that removes the IDF term.

4.3 Comparisons with Previous Methods

Paragraph-level evaluation As shown in Table 2, the methods Sentence-Concat, Template, DenseCap-Concat, Image-Flat [13] and Up-Down [1] perform poorly *w.r.t* both the n-gram metrics and the census-based metric. We examined the generated paragraphs from the Up-Down model and found that the description sentences in a paragraph are heavily repeated. Under cross-entropy training, Regions-Hier [16] achieves the highest Bleu@1 and METEOR scores, but a low CIDEr score. This reveals that the model learns overall accuracy but fails to predict those semantic words. On the contrary, TDC gains a high CIDEr score and a low Bleu@1 score. Our model integrates the advantages of the “big” sentence methods and the topic-based methods, and presents relatively stable paragraph-level coherence and accuracy. With adversarial training and extra data, RTT-GAN performs best *w.r.t* METEOR and Bleu scores, followed by our model. However, our model surpasses RTT-GAN by 27.90% on the CIDEr score, which validates the effectiveness of the replacement-based reward.

Sentence-level evaluation Results are listed in Table 1. Compared with the “big” sentence model TDC that suppresses immediate repetition, the TEB module contributes nothing to sentence accuracy and coherency. The ParaCNN achieves higher Bleu scores than TDC while lower METEOR and CIDEr* scores, which indicates increased accuracy and undermined semantic coherence. Our model surpasses TDC by 18.96% METEOR and 1.23% CIDEr* scores, and outperforms ParaCNN by 3.14% Bleu@4 and 5.76% Bleu@1 scores. It indicates that our model improves n-gram accuracy and enhances sentence coherence. This is further verified in Fig. 2. The proposed Bypass model consistently performs much better than the others *w.r.t* METEOR on all the six sentences, slightly better than TDC *w.r.t* CIDEr* and slightly better than ParaCNN *w.r.t* Bleu scores. The improvements contribute to simultaneously overcoming immediate repetition and delayed repetition.

	METER	CIDEr	Bleu@1	Bleu@4	RL
Sentence-Concat	12.05	6.82	31.11	3.98	✗
Template	14.31	12.15	37.47	7.38	✗
DenseCap-Concat	12.66	12.51	33.18	4.54	✗
Image-Flat	12.82	11.06	34.04	7.71	✗
Up-Down	13.66	12.89	32.78	6.89	✗
Regions-Hier*	15.95	13.52	41.90	8.69	✗
TDC	15.63	23.96	37.29	8.98	✗
Bypass	15.77	18.35	40.23	8.65	✗
Up-Down	13.63	13.77	29.67	5.88	✓
RTT-GAN*	18.39	20.36	42.06	9.21	✓
TEB	15.11	22.13	35.90	8.44	✓
ParaCNN	14.96	18.77	37.37	7.69	✓
Bypass	17.27	26.04	40.75	9.05	✓

Table 2: Paragraph-level scores of different methods on the test split of the Stanford image-paragraph dataset. ✗ and ✓ represent training by cross-entropy loss and finetuning by REINFORCE algorithm, respectively. Among the ✓, RTT-GAN and ParaCNN use adversarial training instead of RL. * indicates that extra data is used in training.

	METER	CIDEr	Bleu@1	Bleu@4
Baseline	15.48	16.40	38.83	8.34
Baseline + TimeNorm	15.39	17.01	38.73	8.30
Bypass	15.77	18.35	40.23	8.65
Baseline + R _{increment}	17.36	18.51	39.86	8.46
Bypass + R _{increment}	17.11	20.74	41.14	9.10
Bypass + R _{replacement}	17.27	26.04	40.75	9.05

Table 3: Paragraph-level scores *w.r.t* METEOR, CIDEr, Bleu@1 and Bleu@4.

4.4 Ablation Studies

Paragraph-level evaluation As shown in Table 3, equipped with Time Normalization, our model surpasses the original baseline model on all metrics. We observed faster convergence in experiments and improved diversity in generated descriptions. We attribute the accelerated convergence to the reduced internal shift, where Time Normalization reduces the difference among the distributions of semantic vector or syntax vector at different time steps. As for diversity, the normalized variance avoids mode-collapse in the input vectors for each sentence. Trained by R_{increment}, our model achieves higher scores than the baseline model *w.r.t* all the metrics. Comparing our model trained with the two rewards, we can see a boost on the CIDEr metric when using the replacement-based reward.

	METEOR	CIDEr*	Bleu@4	Bleu@3	Bleu@2	Bleu@1
Baseline	9.11 ± 0.17	71.39 ± 1.78	3.87 ± 0.13	6.79 ± 0.18	11.99 ± 0.25	22.79 ± 0.32
Baseline+TimeNorm	8.95 ± 0.16	70.53 ± 1.70	3.74 ± 0.12	6.62 ± 0.17	11.80 ± 0.24	22.71 ± 0.32
Bypass	9.32 ± 0.16	71.93 ± 1.68	3.84 ± 0.11	6.95 ± 0.17	12.30 ± 0.24	23.51 ± 0.32
Baseline+R _{increment}	9.91 ± 0.15	64.80 ± 1.35	3.36 ± 0.09	6.31 ± 0.15	11.68 ± 0.21	23.23 ± 0.27
Bypass+R _{increment}	10.15 ± 0.16	71.64 ± 1.57	3.99 ± 0.11	7.20 ± 0.17	12.76 ± 0.24	24.24 ± 0.30
Bypass+R _{replacement}	10.48 ± 0.16	73.07 ± 1.48	3.94 ± 0.11	7.09 ± 0.16	12.65 ± 0.23	24.44 ± 0.29

Table 4: Sentence-level scores of different methods. The margin of error (i.e. $\pm\epsilon$) is calculated at confidence interval 95%.

Sentence-level evaluation Table 4 lists the sentence-level scores under cross-entropy training and REINFORCE training. Interestingly, Time Normalization undermines the performance of the Baseline. It verifies our hypothesis that linguistic information harms topic transition. Time Normalization cannot function effectively on the integration of linguistics and semantics. However, the improvement is obvious after these two aspects are distinguished by our Bypass model. From the results of REINFORCE training, the replacement-based reward enhances the sentence-level coherence while the increment-based reward promotes the sentence-level accuracy, complementary to the Bypass architecture.

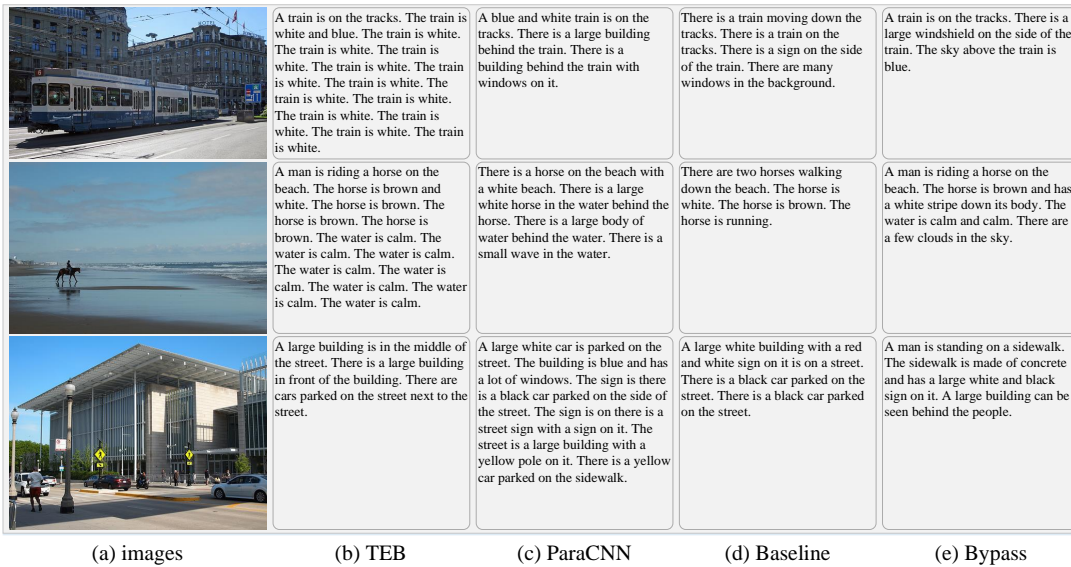


Figure 3: Descriptive captioning paragraphs generated by TEB, ParaCNN, the Baseline, and the Bypass model.

4.5 Qualitative Results

We list the descriptions generated by TEB, ParaCNN, the Baseline, and our model in Fig. 3. For the first image, the paragraph from TEB contains heavy immediate repetition, which keeps describing “the train is white”. ParaCNN improves the coherence of the paragraph by explicitly modeling topic transition, i.e. from *train*, *building* to *windows*. The Baseline model describes the visual scene fluently but has immediate repetition in the second sentence. Our model depicts around the *train*, from *track*, *windshield* to *sky*. As for the second image, similar repetition is observed in the captions from TEB, which is partially reduced by ParaCNN. The description from the Baseline is around the *horse*, but not accurate enough and somewhat redundant. Our model seamlessly transits from *man* and *horse* to *water* and *sky*, with little repeated words. The third image is one of the very few samples that TEB depicts with little repetition. However, ParaCNN suffers from delayed repetition in this case. The same topic goes back and forth in different sentences, breaking the semantic coherence in the paragraph, e.g. both the first and the last sentences are about *car*. The caption generated by the Baseline is accurate yet contains immediate repetition in the last sentence. Our model resembles the accuracy of the Baseline and avoids repetition.

Visualization of topic transition As shown in Fig. 4, top visual areas are masked among the bounding boxes for each sentence in the paragraph captions by the Baseline (right up) and by our model (right down).

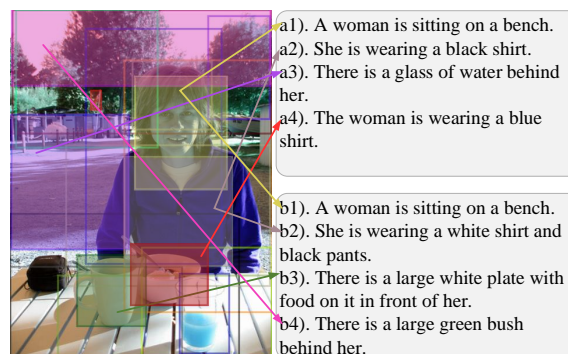


Figure 4: Visual attention during paragraph generation by the Baseline model (a1~a4) and the Bypass model (b1~b4).

We can see that both models capture accurate regions when producing the first two sentences, i.e., the *woman* and the *shirt*. However, for the third sentence, the Baseline model focuses on a *street* area and depicts the item on the table.

It continues describing *woman* and *shirt* in the fourth sentence while attending to the *bowl*. Instead, our model first focuses and describes the *cup* on the table in front of the woman in the third sentence, and then turns to the *bush* behind her. The visual attention in our model moves more smoothly and has better control of the generated descriptions.

Visualization of topic clustering Fig. 5 shows the clustering space of the learned topic vectors from the Baseline and the proposed model. We can observe that the topic vectors of the Bypass model present a much better clustering effect *w.r.t* the semantic. On the contrary, the topic vectors from the same semantic category spread more widely, and those from different categories mix with each other.

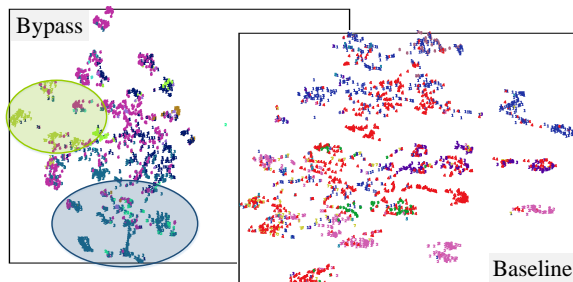


Figure 5: Clustering of topic vectors. Samples in the same color are from the same semantic category. The overlapped area covers very few points in the Bypass clustering space.

We circle two sets of clustered samples in space given by the Bypass model for illustration. The dark blue circle represents an outdoor scene, where the semantic words include {‘water’, ‘ocean’, ‘dirt’, ‘trees’, ‘tarmac’, ‘snow’, ‘rock’, ‘sky’, ‘grass’, ‘road’, ‘sand’, ‘leaves’, ‘gravel’, ‘forest’, ‘wooded’, ‘mountain’, ‘fence’, ‘tracks’, ‘ground’, ‘sun’, ‘enclosure’, ‘hill’, ‘slope’, ‘boulder’}. The light cyan represents an indoor scene, where the semantic words include {‘bottle’, ‘luggage’, ‘boxes’, ‘attached’, ‘window’, ‘meter’, ‘room’, ‘handles’, ‘bowls’, ‘glass’, ‘desk’, ‘box’, ‘door’, ‘toilet’, ‘plates’, ‘suitcases’, ‘shelf’, ‘suitcase’, ‘vase’, ‘hanging’, ‘bag’, ‘sink’, ‘bathroom’, ‘kitchen’}.

5 Conclusion

In this paper, we proposed a Bypass architecture that separates the semantic stream from sentences for a more coherent topic transition without extra disentangling loss. Moreover, a novel replacement-based reward was devised to reduce delayed repetition and the variance of the model when trained by the REINFORCE algorithm. Experiments showed improved sentence-level and paragraph-level coherence and accuracy of the proposed method without extra training data, and verified the validity of different modules. In the future, we would like to explore more on the syntactic stream to increase the linguistic diversity in generated paragraph descriptions.

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*, pages 6077–6086, 2018.
- [2] Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *ACL Workshop*, pages 65–72, 2005.
- [3] Moitrey Chatterjee and Alexander G Schwing. Diverse and coherent paragraph generation from images. In *ECCV*, pages 729–744, 2018.
- [4] Jurafsky Daniel and Martin James H. Speech and language processing. *Draft*, 2019.
- [5] Jeffrey Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, and Trevor Darrell. Long-term recurrent convolutional networks for visual recognition and description. In *CVPR*, pages 2625–2634, 2015.
- [6] Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. Speech recognition with deep recurrent neural networks. In *ICASSP*, pages 6645–6649. IEEE, 2013.

- [7] Barbara J Grosz, Scott Weinstein, and Aravind K Joshi. Centering: A framework for modeling the local coherence of discourse. *Computational linguistics*, 21(2):203–225, 1995.
- [8] Arjun Gupta, Zengming Shen, and Thomas Huang. Text embedding bank for detailed image paragraph captioning. In *AAAI*, volume 35, pages 15791–15792, 2021.
- [9] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015.
- [10] Wenhao Jiang, Lin Ma, Yu-Gang Jiang, Wei Liu, and Tong Zhang. Recurrent fusion network for image captioning. In *ECCV*, pages 499–515, 2018.
- [11] Justin Johnson, Andrej Karpathy, and Li Fei-Fei. Densecap: Fully convolutional localization networks for dense captioning. In *CVPR*, pages 4565–4574, 2016.
- [12] Malvin H Kalos and Paula A Whitlock. *Monte carlo methods*. John Wiley & Sons, 2009.
- [13] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *CVPR*, pages 3128–3137, 2015.
- [14] Dong-Jin Kim, Jinsoo Choi, Tae-Hyun Oh, and In So Kweon. Dense relational captioning: Triple-stream networks for relationship-based captioning. In *CVPR*, pages 6271–6280, 2019.
- [15] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [16] Jonathan Krause, Justin Johnson, Ranjay Krishna, and Li Fei-Fei. A hierarchical approach for generating descriptive image paragraphs. In *CVPR*, pages 317–325, 2017.
- [17] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *IJCV*, 123(1):32–73, 2017.
- [18] Girish Kulkarni, Visruth Premraj, Vicente Ordonez, Sagnik Dhar, Siming Li, Yejin Choi, Alexander C Berg, and Tamara L Berg. Babytalk: Understanding and generating simple image descriptions. *PAMI*, 35(12):2891–2903, 2013.
- [19] Yann LeCun, Yoshua Bengio, et al. Convolutional networks for images, speech, and time series. *The handbook of brain theory and neural networks*, 3361(10):1995, 1995.
- [20] Ruifan Li, Haoyu Liang, Yihui Shi, Fangxiang Feng, and Xiaojie Wang. Dual-cnn: A convolutional language decoder for paragraph image captioning. *Neurocomputing*, 396:92–101, 2020.
- [21] Xiaodan Liang, Zhiting Hu, Hao Zhang, Chuang Gan, and Eric P Xing. Recurrent topic-transition gan for visual paragraph generation. In *ICCV*, pages 3362–3371, 2017.
- [22] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, pages 740–755. Springer, 2014.
- [23] Jiasen Lu, Caiming Xiong, Devi Parikh, and Richard Socher. Knowing when to look: Adaptive attention via a visual sentinel for image captioning. In *CVPR*, pages 375–383, 2017.
- [24] Yue Lu, Chao Guo, Xingyuan Dai, and Fei-Yue Wang. Data-efficient image captioning of fine art paintings via virtual-real semantic alignment training. *Neurocomputing*, 2022.
- [25] Yuzhao Mao, Chang Zhou, Xiaojie Wang, and Ruifan Li. Show and tell more: Topic-oriented multi-sentence image captioning. In *IJCAI*, pages 4258–4264, 2018.
- [26] Luke Melas-Kyriazi, Alexander M Rush, and George Han. Training for diversity in image paragraph captioning. In *EMNLP*, pages 757–761, 2018.
- [27] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *ACL*, pages 311–318. Association for Computational Linguistics, 2002.
- [28] William Peebles, John Peebles, Jun-Yan Zhu, Alexei Efros, and Antonio Torralba. The hessian penalty: A weak prior for unsupervised disentanglement. In *ECCV*, pages 581–597. Springer, 2020.
- [29] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015.
- [30] Steven J Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. Self-critical sequence training for image captioning. In *CVPR*, pages 7008–7024, 2017.
- [31] Jia Huei Tan, Ying Hua Tan, Chee Seng Chan, and Joon Huang Chuah. Acort: A compact object relation transformer for parameter efficient image captioning. *Neurocomputing*, 2022.

- [32] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *CVPR*, pages 4566–4575, 2015.
- [33] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural image caption generator. In *CVPR*, pages 3156–3164, 2015.
- [34] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- [35] Jie Wu, Tianshui Chen, Hefeng Wu, Zhi Yang, Guangchun Luo, and Liang Lin. Fine-grained image captioning with global-local discriminative objective. *IEEE Transactions on Multimedia*, 23:2413–2427, 2021.
- [36] Lingxiang Wu, Min Xu, Jinqiao Wang, and Stuart Perry. Recall what you see continually using gridstm in image captioning. *IEEE Transactions on Multimedia*, 22(3):808–818, 2020.
- [37] Siying Wu, Zheng-Jun Zha, Zilei Wang, Houqiang Li, and Feng Wu. Densely supervised hierarchical policy-value network for image paragraph generation. In *IJCAI*, pages 975–981. AAAI Press, 2019.
- [38] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, pages 2048–2057, 2015.
- [39] Shiyang Yan, Yang Hua, and Neil Robertson. Paracnn: Visual paragraph generation via adversarial twin contextual cnns. *arXiv preprint arXiv:2004.10258*, 2020.
- [40] Linjie Yang, Kevin Tang, Jianchao Yang, and Li-Jia Li. Dense captioning with joint inference and visual context. In *CVPR*, pages 2193–2202, 2017.
- [41] Min Yang, Wei Zhao, Wei Xu, Yabing Feng, Zhou Zhao, Xiaojun Chen, and Kai Lei. Multitask learning for cross-domain image captioning. *IEEE Transactions on Multimedia*, 21(4):1047–1061, 2019.
- [42] Ting Yao, Yingwei Pan, Yehao Li, and Tao Mei. Exploring visual relationship for image captioning. In *ECCV*, pages 684–699, 2018.
- [43] Ting Yao, Yingwei Pan, Yehao Li, Zhaofan Qiu, and Tao Mei. Boosting image captioning with attributes. In *ICCV*, pages 4894–4902, 2017.
- [44] Guojun Yin, Lu Sheng, Bin Liu, Nenghai Yu, Xiaogang Wang, and Jing Shao. Context and attribute grounded dense captioning. In *CVPR*, pages 6241–6250, 2019.

A Deduction of the disentanglement

We add no extra constraint on $h_{e,i}^{(j)}$ and $h_{y,i}^{(j)}$ to reduce their correlation since the employed single addition layer under word-level cross-entropy loss facilitates the disentanglement property.

For simplicity, we reuse the symbol $\mathbf{x} \in \mathbb{R}^D$ and $\mathbf{y} \in \mathbb{R}^D$ to denote the semantic element $h_{e,i}^{(j)}$ and the syntactic element $h_{y,i}^{(j)}$, respectively. The input vector to the *softmax* layer is \mathbf{z} and the output is \mathbf{f} . Thus we have,

$$\mathbf{z} = \mathbf{W}(\mathbf{x} + \mathbf{y}) \text{ and } \mathbf{f} = \sigma(\mathbf{z}), \tag{18}$$

where \mathbf{W} is the parameter matrix and $\sigma(\cdot)$ is *softmax* function.

Given that the i -th ground-truth word in the j -th sentence is an one-hot word vector \mathbf{f}^* with $f_k = 1$, the cross-entropy loss is

$$\ell = -\log f_k, \tag{19}$$

and we have that

$$\frac{\partial f_k}{\partial \mathbf{x}} = \frac{\partial f_k}{\partial \mathbf{z}} \cdot \frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \frac{\partial f_k}{\partial \mathbf{z}} \cdot \mathbf{W}, \tag{20}$$

$$\frac{\partial f_k}{\partial \mathbf{y}} = \frac{\partial f_k}{\partial \mathbf{z}} \cdot \frac{\partial \mathbf{z}}{\partial \mathbf{y}} = \frac{\partial f_k}{\partial \mathbf{z}} \cdot \mathbf{W}. \tag{21}$$

Then according to gradient chain rule, we can derive

$$\frac{\partial \ell}{\partial \mathbf{x}} = -\frac{\partial \log f_k}{\partial \mathbf{x}} = -\frac{1}{f_k} \frac{\partial f_k}{\partial \mathbf{x}}. \tag{22}$$

To analyze the correlation between \mathbf{x} and \mathbf{y} , we calculate the second-order partial derivatives as [28] does

$$\frac{\partial^2 \ell}{\partial \mathbf{x} \partial \mathbf{y}} = \frac{\partial}{\partial \mathbf{y}} \frac{\partial \ell}{\partial \mathbf{x}} = -\frac{\partial}{\partial \mathbf{y}} \left(\frac{1}{f_k} \frac{\partial f_k}{\partial \mathbf{x}} \right) \quad (23)$$

$$= \frac{1}{f_k^2} \frac{\partial f_k}{\partial \mathbf{y}} \frac{\partial f_k}{\partial \mathbf{x}} - \frac{1}{f_k} \frac{\partial^2 f_k}{\partial \mathbf{x} \partial \mathbf{y}}. \quad (24)$$

The two terms in Eq. (24) regardless of coefficient are

$$\begin{aligned} \frac{\partial f_k}{\partial \mathbf{y}} \frac{\partial f_k}{\partial \mathbf{x}} &= \left(\frac{\partial f_k}{\partial x_t} \frac{\partial f_k}{\partial y_s} \right) \\ &= \left(\left(\sum_i \frac{\partial f_k}{\partial z_i} w_{it} \right) \left(\sum_j \frac{\partial f_k}{\partial z_j} w_{js} \right) \right) \\ &= \left(\sum_{i,j} \frac{\partial f_k}{\partial z_i} \frac{\partial f_k}{\partial z_j} w_{it} w_{js} \right), \end{aligned} \quad (25)$$

and

$$\begin{aligned} \frac{\partial^2 f_k}{\partial \mathbf{x} \partial \mathbf{y}} &= \frac{\partial}{\partial \mathbf{y}} \left(\frac{\partial f_k}{\partial \mathbf{x}} \right) = \frac{\partial}{\partial \mathbf{y}} \left(\frac{\partial f_k}{\partial \mathbf{z}} \cdot \mathbf{W} \right) \\ &= \left(\frac{\partial}{\partial y_s} \frac{\partial f_k}{\partial x_t} \right) = \left(\frac{\partial}{\partial y_s} \left(\sum_i \frac{\partial f_k}{\partial z_i} w_{it} \right) \right) \\ &= \left(\sum_{i,j} \frac{\partial^2 f_k}{\partial z_i \partial z_j} \frac{\partial z_j}{\partial y_s} w_{it} \right) \\ &= \left(\sum_{i,j} \frac{\partial^2 f_k}{\partial z_i \partial z_j} w_{js} w_{it} \right) = \mathbf{W}^T \left(\frac{\partial^2 f_k}{\partial \mathbf{z}^2} \right) \mathbf{W}. \end{aligned} \quad (26)$$

Replacing Eq. (23) with Eq. (25) and Eq. (26), we have that

$$\begin{aligned} \frac{\partial^2 \ell}{\partial \mathbf{x} \partial \mathbf{y}} &= \frac{1}{f_k^2} \left(\sum_{i,j} w_{it} w_{js} \left(\frac{\partial f_k}{\partial z_i} \frac{\partial f_k}{\partial z_j} - f_k \frac{\partial^2 f_k}{\partial z_i \partial z_j} \right) \right) \\ &\doteq \frac{1}{f_k^2} \mathbf{W}^T \mathbf{A} \mathbf{W}, \end{aligned} \quad (27)$$

where $\mathbf{A} = (a_{ij}) = \left(\frac{\partial f_k}{\partial z_i} \frac{\partial f_k}{\partial z_j} - f_k \frac{\partial^2 f_k}{\partial z_i \partial z_j} \right)$.

Since $\mathbf{f} = \sigma(\mathbf{z})$, we have $f_k = \frac{e^{z_k}}{\sum_m e^{z_m}}$. The partial derivative of the *softmax* function gives

$$\frac{\partial f_k}{\partial z_i} = \begin{cases} -f_i f_k & \text{if } i \neq k \\ f_k(1 - f_k) & \text{if } i = k. \end{cases} \quad (28)$$

Based on Eq. (28), the second-order partial derivative is

$$\frac{\partial^2 f_k}{\partial z_i \partial z_j} = \begin{cases} f_k(2f_k - 1)(f_k - 1) & \text{if } i = j = k \\ f_j f_k(2f_k - 1) & \text{if } i = k \neq j \\ f_i f_k(2f_k - 1) & \text{if } i \neq j = k \\ f_k f_i(2f_i - 1) & \text{if } i = j \neq k \\ 2f_i f_j f_k & \text{if } i \neq j \neq k. \end{cases} \quad (29)$$

Correspondingly, the elements of \mathbf{A} can be calculated as

$$a_{ij} = \begin{cases} f_k^3(1 - f_k) & \text{if } i = j = k \text{ (①)} \\ -f_k^3 f_j & \text{if } i = k \neq j \text{ (②)} \\ -f_i f_k^3 & \text{if } i \neq j = k \text{ (③)} \\ f_k^2 f_i(1 - f_i) & \text{if } i = j \neq k \text{ (④)} \\ -f_i f_j f_k^2 & \text{if } i \neq j \neq k \text{ (⑤)}. \end{cases} \quad (30)$$

Supposing $k = 1$, a visualization of \mathbf{A} is

$$\mathbf{A} = \left[\begin{array}{c|ccc} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{array} \right] \doteq \left[\begin{array}{c|cc} \textcircled{1} & & \textcircled{2} \\ \textcircled{4} & & \textcircled{5} \\ \textcircled{3} & & \ddots \\ & \textcircled{5} & & \textcircled{4} \end{array} \right].$$

Minimizing the loss ℓ leads to $f_k \rightarrow 1$ and $f_m \rightarrow 0$ for $m \neq k$. Thus, $a_{ij} \rightarrow 0$ (by Eq. (30)) and $\mathbf{A} \rightarrow \mathbf{0}$. We assume that the parameter matrix \mathbf{W} is bounded (i.e., w_{ij} is finite), then it derives $\frac{\partial^2 \ell}{\partial \mathbf{x} \partial \mathbf{y}} \rightarrow \mathbf{0}$ and $\frac{\partial^2 \ell}{\partial \mathbf{y} \partial \mathbf{x}} \rightarrow \mathbf{0}$ by Eq. (27).

Thus, minimizing the loss ℓ also minimizes the off-diagonal entries of the function’s Hessian matrix *w.r.t.* its inputs \mathbf{x} and \mathbf{y} , which encourages the disentanglement between the inputs. Therefore, no extra constraint is compulsory to reduce the correlation of $h_{e,i}^{(j)}$ and $h_{y,i}^{(j)}$.

B More qualitative examples

Fig. 6 is a more complete version of the quality comparison. We can see that ‘big’ sentence methods Up-Down and TEB suffer from both immediate repetition and delayed repetition. TDC reduces immediate repetition by blocking trigrams, while ParaCNN and the Baseline model address it by explicitly modeling topic transition. However, they still suffer from delayed repetition. Our Bypass reduces both types of repetition and enhances the coherence of generated descriptive paragraphs.




<p>(a) images</p>				
	<p>(b) Up-Down</p>	<p>A train is on the tracks. There is a red and white train on the train. There is a white building behind the train.</p>	<p>A man is riding a horse. The horse is brown and white. The horse is brown. The horse is brown. The horse is brown. The horse is brown. The horse is brown. The horse is brown. The sky is blue. The sky is blue. The sky is blue.</p>	<p>A man is walking on the sidewalk. There is a man walking on the sidewalk. There is a man walking on the sidewalk. There are cars parked on the street. There are cars parked on the street. There are many cars on the street. There are many cars on the street.</p>
		<p>(c) TEB</p>	<p>A train is on the tracks. The train is white and blue. The train is white. The train is white. The train is white. The train is white. The train is white. The train is white. The train is white.</p>	<p>A man is riding a horse on the beach. The horse is brown and white. The horse is brown. The horse is brown. The water is calm. The water is calm. The water is calm. The water is calm. The water is calm. The water is calm.</p>
	<p>(d) TDC</p>	<p>A train is on the tracks. There is a large building behind the train. There are many windows on the side of the train.</p>	<p>A man is riding a horse on the beach. The horse is brown and white. The man is wearing a black helmet. The horse is white. The water is calm. The sky is blue. The clouds are white. There are many clouds in the sky. The sky is blue and white. There is a lot of clouds in the sky. There are clouds in the sky.</p>	<p>A street is lined with cars. There are cars parked on the street. There is a man in a white shirt walking on the sidewalk. There is a white car parked on the side of the street.</p>
	<p>(e) ParaCNN</p>	<p>A blue and white train is on the tracks. There is a large building behind the train. There is a building behind the train with windows on it.</p>	<p>There is a horse on the beach with a white beach. There is a large white horse in the water behind the horse. There is a large body of water behind the water. There is a small wave in the water.</p>	<p>A large white car is parked on the street. The building is blue and has a lot of windows. The sign is there is a black car parked on the side of the street. The sign is on there is a street sign with a sign on it. The street is a large building with a yellow pole on it. There is a yellow car parked on the sidewalk.</p>
	<p>(f) Baseline</p>	<p>There is a train moving down the tracks. There is a train on the tracks. There is a sign on the side of the train. There are many windows in the background.</p>	<p>There are two horses walking down the beach. The horse is white. The horse is brown. The horse is running.</p>	<p>A large white building with a red and white sign on it is on a street. There is a black car parked on the street. There is a black car parked on the street.</p>
	<p>(g) Bypass</p>	<p>A train is on the tracks. There is a large windshield on the side of the train. The sky above the train is blue.</p>	<p>A man is riding a horse on the beach. The horse is brown and has a white stripe down its body. The water is calm and calm. There are a few clouds in the sky.</p>	<p>A man is standing on a sidewalk. The sidewalk is made of concrete and has a large white and black sign on it. A large building can be seen behind the people.</p>

Figure 6: Descriptive captioning paragraphs generated by Up-Down, TDC, TEB, ParaCNN, the Baseline, and the Bypass model.