

1 We thank all reviewers for their positive comments on idea novelty, technical quality, paper writing, and promising
2 directions. We respond to the concerns point-by-point as below.

3 **R#1.Q1. Apart from Fig. 6, would be useful to see if syntactic/lexical (S./L.) effects are disentangled.**

4 Following this suggestion, we provide quantitative comparisons to analyze the disentanglement of the L. and S. effects
5 in Tab. I. After modifying L./S., we have found: (1) the word and POS sequences are both changed ($ED \neq 0$ and $B-1,3$
6 $\neq 100\%$), which reveals the inherent correlation between L. and S., (2) the change on word/POS sequences is bigger
7 than that of POS/word, *i.e.*, higher ED and lower B-1,3, which indicates that the L./S. variables have more effect on the
8 L./S., and (3) the change on POS sequences is smaller than that on word sequences, which is probably due to smaller
9 S. (POS) vocabulary. We will add the above results and discussions in our paper to further enrich the insights.

10 **R#1.Q2. On technical details.**

11 For the question whether beam search is used in [28][3], com-
12 mon image captioning methods [28] and [3] use beam search
13 for sampling. For the question which captions are used for
14 computing metrics in Tab. 1, we use likelihood to sample top
15 5 captions of each testing image, compute their metrics, and
16 choose the top caption for the evaluation in Tab. 1 to ensure a
17 fair comparison. We will clarify the above details in our paper.

18 **R#2.Q1. On the originality compared to [13] & VAE.**

19 Thanks, we agree that our work does share certain intersection with [13] and VAEs. However, we rebut that our novelty
20 are fundamentally sufficient comparing to [13] and VAEs, which are detailed in two aspects:

21 *New problem formulation:* Conventional encoder-decoder depends solely on sampling to import randomness [28],
22 which limits the diversity among the outputs (see Tab. 2). VAE based encoder-decoder introduces the latent vari-
23 able and makes a two-stage inference for latent variables and words, which enhances the diversity (also see Tab. 2).
24 However, the latent variables in VAEs have a very general prior (standard Gaussians), which does not consider any
25 domain-specific knowledge. We argue that it may waste model capacity, and one should consider the unique problem
26 structures of image captioning instead of using the VAE as is. Therefore, we introduce the domain knowledge from
27 NLP and decompose the latent variables into lexicon and syntax variables. This fundamental change in the problem
28 representation is the core novelty of our approach. Under this guidance, it’s a straightforward thinking to model a
29 structured variational inferrer with the assistance of VP-Tree [13] and the adaptation of VAE. However, we kindly ar-
30 gue that such assistance/adaptation is not our core contribution. One can replace VP-Tree with other visual structured
31 representations, or use generative models other than VAE. However, in order to directly demonstrate the effectiveness
32 of the core idea, we intentionally chose these straightforward adaptations, which help the readers directly catch our
33 main innovation and not get distracted by the complicated adaptation.

34 *New technical design:* VP-Tree [13] can provide the lexicon/syntax probabilities, which, however, does not involve
35 the construction and the prior/posterior inference of the latent variables (Sec. 3.2). For VAEs, though commonly used,
36 they never consider modeling the structured latent variables with the domain-specific knowledge, as well as jointly
37 optimizing the reconstruction and the prior-posterior distance with the structured latent variables (Sec. 3.3).

38 **R#2.Q2. Notations involving “ \cdot ” and “ $\cdot\cdot$ ” can be replaced with sub/super-scripts ℓ and s .**

39 Thank you for this suggestion. We will modify accordingly in our paper.

40 **R#3.Q1. Confusion on model generalization for longer captions or paragraphs.**

41 Thanks for this inspiring question. VarMI-tree can be easily expanded to the case with more tree nodes for long
42 captions. Our model itself has no such limitation on caption length, and we just set the node number as 7 to cover
43 the captions in the COCO dataset. As for image paragraph description, it is quite different from the sentence-level
44 captioning due to different topics of sentences [A][B]. However, as long as the topic feature of each sentence can
45 be extracted from RNN [A][B] to construct VarMI-tree, it does not restrict our model to be generalized to image
46 paragraph description.

47 [A] M. Chatterjee *et al.* Diverse and Coherent Paragraph Generation from Images. ECCV 2018.

48 [B] J. Krause *et al.* A Hierarchical Approach for Generating Descriptive Image Paragraphs. CVPR 2017.

49 **R#3.Q2. Reorganizing Section 3 for easier follow.**

50 Thanks for the suggestion. We will strengthen the method overview with a table of notations and an algorithm flow.

51 **R#3.Q3. Providing more abstract analysis and empirical discussion.**

52 Thanks for this suggestion. Please kindly refer to our response to Q1 of R#2, which will be added in our paper.

Table I: Consistency between the word/POS output sequences of the modified and original versions. Lexicon/syntax (L./S.) is modified by assigning random probabilities to words/POSS in VP-tree. Edit Distance (ED), Bleu-1 (B-1, %) and Bleu-3 (B-3, %) [29] are taken to measure the consistency (B-1,3 are for the word-level evaluation).

	Evaluation on word			Evaluation on POS		
	ED	B-1	B-3	ED	B-1	B-3
Modify L.	4.08	51.24	32.20	5.48	85.36	73.27
Modify S.	3.82	55.14	36.44	5.83	84.22	71.38