

1 We thank all the reviewers for their helpful comments and for recognizing the novelty of our approach (R2–4) and its
2 potential for positive impact (R2, R4). We are glad that the reviewers found our experimental setup exhaustive (R1–4).

3 **Technical contribution (R1, R3).** The novelty of our approach lies in the manner in which we formulate our COS latent
4 space in Eq. (3) to encode the object and context information, allowing to better capture the underlying distribution.
5 This enables the exploitation of contextual similarities between images and captions beyond that from paired training
6 data in our pseudo-supervised setup. This is not feasible with prior work, *e.g.* Seq-CVAE [6]. This is the first work that
7 addresses diversity *and* novel object captioning in a single framework (Sec. 3.2, Eq. (7), ll. 198–212).

8 **Oracle evaluation and consensus re-ranking (R2).** The oracle evaluation has been used extensively in prior work
9 on multi-prediction (Guzmán-Rivera et al., NIPS 2012). The oracle evaluation does score the top sentence, but is
10 evaluated for a limited number of sentences per image and averaged across the test set. Thus, for a high oracle score the
11 model needs to put high probability mass on the likely captions for an image without relying on “chance”. The oracle
12 score provides an upper bound to any re-ranking method. While consensus re-ranking can indeed be improved using
13 stronger networks to close the gap between oracle performance and consensus re-ranking, we follow the same consensus
14 re-ranking procedure as [6,14] for fairness and achieve much better accuracy (Tab. 9, suppl.). Furthermore, the accuracy
15 after re-ranking (Fig. 4) is competitive to methods that generate a single caption, *e.g.* [4], using beam search. We will
16 move Tab. 9 (suppl.) to the main paper. We compute $p(x|z, I)$ and, therefore, ranking captions x for an image I by
17 decoder probabilities $p(x|I)$ as in beam search would require marginalization over z , which is computationally difficult.

18 **Accuracy scores (R2, R3, R4).** Consensus re-ranking provides a realistic evaluation
19 against the competing methods (Tab. 9). To further validate our approach, we compare the
20 best-10 sentences after consensus re-ranking in Fig. 4 against Seq-CVAE [6] and beam
21 search on a standard captioning method of [4], which has the same base network. The
22 accuracy results show that COS-CVAE is competitive to beam search and thus diverse
23 *yet* accurate. The average(worst) CIDEr for all the 20 samples with COS-CVAE is
24 0.875(0.457) compared to 0.792(0.351) for Seq-CVAE. We will add this.

25 **F1-score, Fig. 3b (R2, R4).** Ranking captions by decoder probabilities is not feasible
26 as discussed above. Consensus re-ranking is not possible as the training data does not
27 contain novel objects. The test data statistics for held-out COCO [17] show that $\sim 57\%$
28 of the images with novel objects have at least one ground-truth caption *without* the mention of novel object. For
29 ground-truth annotations, mention of novel object in any of the sentences is counted towards the F1-score [3]. Therefore,
30 it would be limiting to expect every caption from a diverse captioning approach, which learns the distribution of the
31 likely captions, to mention the novel object (Tab. 11). We will clarify and highlight these challenges in the final version.
32 We thus show the F1-score and the highest CIDEr for sample sizes of z in $\{5, \dots, 100\}$ in Fig. 3b and demonstrate
33 competitive results even for a small sample size of 5. Our COS-CVAE with 5 samples has the highest CIDEr of 1.111
34 and an average CIDEr of 0.775, which is competitive to NBT [30] with a similar base network. We will add this.

35 **Effect of context-object split (R2).** Verbs, adjectives, and spatial relationships are crucial components of the context
36 (ll. 105). Random samples in Tab. 2, 11, and 12 show that the captions from COS-CVAE are coherent (ll. 287–288).

37 **Self-CIDEr (R3).** COS-CVAE has a score of 0.742 while Seq-CVAE(attn) has 0.714. We will add a human evaluation.

38 **Diversity metrics, Tab. 4 (R4).** Consistent with prior work, *e.g.* [6], Unique and Novel are computed with 100 samples
39 and Div-1, Div-2, and m-Bleu are computed with top-5 generated captions after consensus re-ranking (ll. 251–253).

40 **Benefit of pseudo supervision (R3).** Pseudo supervision allows to leverage additional contextual descriptions from
41 human annotations that can describe the image, thus capturing the true diversity of the underlying distribution (ll. 94–96).
42 We use pseudo supervision along with ground-truth captions for best performance because ground-truth captions help
43 accurately guide the attention of the decoder (ll. 128–130). The application of pseudo supervision to a retrieval task is
44 an interesting idea and can be considered for future work since the extension is non-trivial.

45 **CBS constraints (R4).** Seq-CVAE (attn)+CBS yields an F1-score of 55.4 with 5 samples compared to an F1-score of
46 70.1 with COS-CVAE+CBS. However, CBS decoding leads to less diverse captions (one of the main goals here).

47 **Clarifications (R2).** The results of Div-BS are on the same test split as [6, 14] (*cf.* Tab. 1 in [6]). The base vision
48 network for different methods can be different as various models impose different constraints in the latent space [14,
49 41]. We have included the Seq-CVAE(attn) baseline with the same base vision network as our COS-CVAE.

50 **Nocaps dataset (R1, R2, R4).** Our variational COS-CVAE framework models the distribution of likely captions for a
51 given image. The evaluation server for the nocaps dataset (Agrawal et al.), however, accepts only one caption per image
52 and does not support methods modeling one-to-many relationships for images and captions. We thus provide the results
53 on the real world held-out COCO dataset for novel object captioning [17] instead.

54 **Additional references (R2, R3).** Thank you. We will add this.

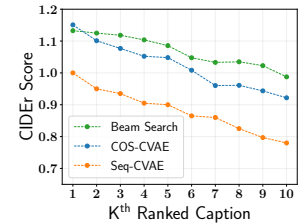


Figure 4: CIDEr score of consensus re-ranked best-10 captions from 20 samples.