Response to "Visual Dialogue without Vision or Dialogue" [7]

Abhishek Das Devi Parikh Dhruv Batra

Georgia Institute of Technology {abhshkdz, parikh, dbatra}@gatech.edu

Abstract

In a recent workshop paper, Massiceti *et al.* [7] presented a baseline model and subsequent critique of Visual Dialog [3] that raises what we believe to be unfounded concerns about the dataset and evaluation. This article intends to rebut the critique and clarify potential confusions for practitioners and future participants in the Visual Dialog challenge.

1 Introduction

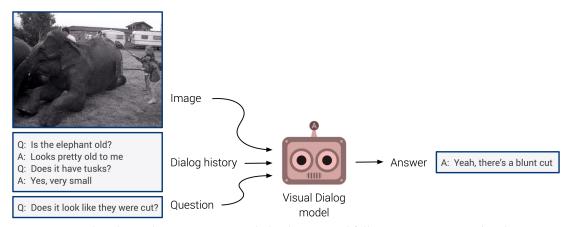


Figure 1: Visual Dialog task: given an image, dialog history, and follow-up question, predict the answer.

Task. The goal of Visual Dialog is to develop conversation agents that can talk about images. Towards this end, in previous work [3], we proposed a task – given an image, dialog history, and follow-up question, predict a free-form natural language answer to the question (Fig. 1) – and a large-scale dataset¹, evaluation metrics and server², and baseline models³ for this task.

Key challenge. A fundamental challenge in dialog systems is automatic evaluation of long free-form answers since existing metrics such as BLEU, METEOR, and ROUGE are known to correlate poorly with human judgement [6]. Thus, as proposed in our initial paper [3], to evaluate Visual Dialog, models are provided a list of 100 candidate answers for each question – consisting of the ground-truth answer from the dataset mixed with nearest neighbors, popular, and random answers – and evaluated on how well they rank the ground-truth answer on retrieval metrics such as mean reciprocal rank (MRR), recall (R@1, 5, 10), and mean rank.

¹visualdialog.org/data

²evalai.cloudcv.org/web/challenges/challenge-page/103/overview

 $^{^3}$ github.com/batra-mlp-lab/visdial

As we describe in our paper [3], these candidate answers for each question are programmatically curated from other answers in the dataset and not human-generated, and so, some candidate answers may be semantically identical (e.g. 'yeah' and 'yes'). Thus, more recently, we conducted new human studies – asking four human subjects to annotate whether each of the 100 candidate answers is correct or not for all questions in the VisDial test split. For evaluation, we report the normalized discounted cumulative gain (NDCG) over the top K ranked options, where K is the number of answers marked as correct by at least one annotator. For this computation, we consider the relevance of an answer to be the fraction of annotators that marked it as correct. This was the primary evaluation criterion for the 1st Visual Dialog Challenge⁴.

As described in [3], there are two broad families of dialog models (unfortunately with names that are overloaded in machine learning) – 'generative' models (that produce a response word-by-word given some context and are evaluated on the ranking of the likelihood scores they assign to candidate answers), and 'discriminative' models (that simply learn to rank a list of candidate answers and cannot produce a new response). This retrieval-based evaluation holds for both families. Compatibility of the evaluation metric with generative models is crucial, since they are more useful for real-world applications where answer options are not available.

2 Concern 1: Suitability of NDCG evaluation

Massiceti et al. [7] note that 'the VisDial dataset was recently updated to version 1.0, where the curators try to ameliorate some of the issues with the single-"ground-truth" answer approach. They incorporate a human-agreement scores for candidate answers, and introduce a modified evaluation which weighs the predicted rankings by these scores. However, in making this change, the primary evaluation for this data has now become an explicit classification task on the candidate answers – requiring access, at train time, to all 100 candidates for every question-image pair. For the stated goals of Visual Dialog, this change can be construed as unsuitable as it falls into the category of redefining the problem to match a potentially unsuitable evaluation measure – how can one get better ranks in the candidate-answer-ranking task.'

The claim that "the primary evaluation for this data has now become an explicit classification task on the candidate answers" is incorrect and thus the conclusion drawn from it is inaccurate and confusing. First, the task has not changed, only the evaluation metric (from MRR to NDCG). The task did not and does not "require access, at train time, to all 100 candidates". Discriminative models use 100 candidate answers at train time; generative models do not. This was discussed in our initial paper [3] and continues to be true.

Perhaps what the authors [7] are trying to say and express concern for is – this metric (NDCG) will favor one kind of model family over another. This is possible and something we have given a lot of thought to. Empirical findings from the 1st Visual Dialog Challenge⁵ indicate that these generative models perform comparably (or even better sometimes) than discriminative models on the NDCG metric – for example, 53.67 vs. 49.58 on VisDial v1.0 test-std for Memory Network + Attention with generative vs. discriminative decoding respectively. Code and models available here: https://github.com/batra-mlp-lab/visdial#pretrained-models-1. While this is still a potentially weak surrogate for human-in-the-loop evaluation of Visual Dialog models, it is encouraging that there now seems to be an automatic evaluation criterion on which generative models, which do not have access to candidate answers during training, outperform discriminative models. As we describe on visualdialog.org, the reason why we chose a single track for the challenge was that in practice, the distinction between the two model families can get blurry (e.g., non-parametric models that internally maintain a large list of answer options), and the separation would be difficult to enforce. Note that our choice of ranking for evaluation isn't an endorsement of either approach (generative or discriminative).

⁴visualdialog.org/challenge/2018#evaluation

⁵visualdialog.org/challenge/2018#winners

3 Concern 2: Comparison to proposed CCA baseline [7]

Massiceti *et al.* [7] proposed a simple CCA baseline with two variants – 1) question-only (ignoring image and dialog history), 2) question + image (ignoring dialog history), which they show outperforms state-of-the-art models on the mean rank metric. They further note that *'an important takeaway from our analyses is that it is highly effective to begin exploration with the simplest possible tools one has at one's disposal. This is particularly apposite in the era of deep neural networks, where the prevailing attitude appears to be that it is preferable to start exploration with complicated methods that aren't well understood, as opposed to older, perhaps even less fashionable methods that have the benefit of being rigorously understood.'*

We agree that simple and strong baselines are important, and are pleasantly surprised to see that a CCA baseline performs so well on mean rank. However, there are a few problems with this analysis. First, the baseline proposed by Massiceti *et al.* [7] is not close to state-of-the-art – the authors cherry-pick the mean rank metric and ignore trends on *all other metrics* (see Tab. 1). Second, it ignores that a similar finding has already been presented in the original Visual Dialog paper [3], that question-only and question + image models perform close to but slightly worse than full Q+I+H models. We recreate Tab. 1 from [3]. Third, the authors [7] ignore that the CCA baselines perform worse than not just state-of-the-art models, but also these Q and Q+I ablations [3], and comparable to answer prior and nearest neighbor (NN) baselines [3] on MRR and R@k. Finally, the results presented in [7] are not directly comparable. The proposed CCA baselines use Resnet-34 [5] features and FastText [1] embeddings, while the baselines in [3] use VGG-16 [8] and learn word embeddings from scratch respectively.

	Model	NDCG	MRR	R@1	R@5	R@10	Mean Rank
v0.9 val	' Answer prior	-	0.3735	23.55	48.52	53.23	26.50
	NN-Q	-	0.4570	35.93	54.07	60.26	18.93
	NN-QI	-	0.4274	33.13	50.83	58.69	19.62
	LF-Q-G	-	0.5048	39.78	60.58	66.33	17.89
	LF-QI-G	-	0.5204	42.04	61.65	67.66	16.84
	LF-QIH-G	-	0.5199	41.83	61.78	67.59	17.07
	HRE-QIH-G	-	0.5237	42.29	62.18	67.92	17.07
	HREA-QIH-G	-	0.5242	42.28	62.33	68.17	16.79
	MN-QIH-G	-	0.5259	42.29	62.85	68.88	17.06
	A-Q (Massiceti et al. [7])		0.3031	16.77	44.86	58.06	16.21
	A-QI (Massiceti <i>et al</i> . [7])	-	0.2427	12.17	35.38	50.57	18.29
v1.0 test-std	LF-QIH-G	0.5121	0.4568	35.08	55.92	64.02	18.81
	HRE-QIH-G	0.5245	0.4561	34.78	56.18	63.72	18.78
	MN-QIH-G	0.5280	0.4580	35.05	56.35	63.92	19.31
	A-Q (Massiceti et al. [7])	<u>-</u>	0.2832	15.95	40.10	55.10	17.08
	A-QI (Massiceti et al. [7])	-	0.2393	12.73	33.05	48.68	19.24

Table 1: Performance of methods on VisDial v0.9 and v1.0, measured by normalized discounted cumulative gain (NDCG), mean reciprocal rank (MRR), recall@k and mean rank. Higher is better for NDCG, MRR, and recall@k, while lower is better for mean rank.

4 Conclusion

To summarize:

- In an attempt to make evaluation for Visual Dialog more reliable, we have recently had multiple human subjects indicate whether each of the candidate answers for a question is correct, which is then used as the reference score while computing the NDCG metric. Massiceti *et al.* [7] claim that this changes the Visual Dialog task to an "explicit classification task on the candidate answers" which is incorrect. The task remains the same as before [3], only the evaluation has changed.
- Further, NDCG evaluation using dense annotations does not favor a particular family of Visual Dialog models (between discriminative and generative), as evidenced by findings from the 1st Visual Dialog challenge noted on visualdialog.org.

- While we welcome simple and strong baselines, the CCA baseline for Visual Dialog proposed by Massiceti *et al.* [7] is not close to state-of-the-art. The authors solely focus on one metric (mean rank) while ignoring all other metrics (MRR, R@k, NDCG) on which their approach is significantly worse, not just against state-of-the-art models, but also against ablations from [3] (see Tab. 1).
- Finally, the VisDial dataset [3] and evaluation are not perfect unbiased testbeds. VisDial likely has many biases and trivial correlations models can pick up on, as has been previously observed in other unstructured (or loosely structured) real-world datasets [4]. Further, automatic evaluation of dialog is an open research problem, and our NDCG evaluation protocol for Visual Dialog is an attempt at making it more robust. Alternatively, evaluation with humans paired with dialog models, conversing for the human to be able to achieve a downstream goal (e.g. understand their visual surroundings, book a flight ticket, etc.) would perhaps be the truest form of evaluation, as has been explored in [2], although this is expensive. There is scope for improvement across all axes task/dataset, evaluation, as well as methods.

References

- [1] P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov. Enriching word vectors with subword information. *TACL*, 2017. 3
- [2] P. Chattopadhyay, D. Yadav, V. Prabhu, A. Chandrasekaran, A. Das, S. Lee, D. Batra, and D. Parikh. Evaluating Visual Dialog Agents via Cooperative Human-AI games. In *HCOMP*, 2017. 4
- [3] A. Das, S. Kottur, K. Gupta, A. Singh, D. Yadav, J. M. F. Moura, D. Parikh, and D. Batra. Visual Dialog. In *CVPR*, 2017. 1, 2, 3, 4
- [4] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *CVPR*, 2017. 4
- [5] K. He, X. Zhang, S. Ren, and J. Sun. Deep Residual Learning for Image Recognition. In CVPR, 2016.
- [6] C.-W. Liu, R. Lowe, I. V. Serban, M. Noseworthy, L. Charlin, and J. Pineau. How NOT To Evaluate Your Dialogue System: An Empirical Study of Unsupervised Evaluation Metrics for Dialogue Response Generation. In EMNLP, 2016. 1
- [7] D. Massiceti, P. K. Dokania, N. Siddharth, and P. H. Torr. Visual Dialogue without Vision or Dialogue. In NeurIPS Workshop on Critiquing and Correcting Trends in Machine Learning, 2018. 1, 2, 3, 4
- [8] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In ICLR, 2015. 3