

University of Stuttgart
Germany



Visual Dialog System - Interactive mode

Mohamed Adnen ABDESSAIED

Study program: Simulation Technology

Matriculation number: 3158373

supervised by:

Prof. Dr. Ngoc Thang VU

&

Dirk VÄTH

SimTech research project 2

*University of Stuttgart
Institute for Natural Language Processing (IMS)
Winter Semester 2019-2020*

CONTENTS

| | | |
|------------|---------------------------------------|---|
| I | Abstract | 1 |
| II | Introduction | 1 |
| III | Related work | 2 |
| IV | Approach | 2 |
| IV-A | Pre-training | 2 |
| IV-A1 | Hyperparameters | 2 |
| IV-A2 | Results | 2 |
| IV-B | Interactive mode | 3 |
| IV-B1 | Usage description | 3 |
| IV-B2 | Our data | 3 |
| IV-C | Fine tuning | 4 |
| V | Results | 4 |
| V-A | On our data | 4 |
| V-B | On a subset of VisDial v1.0 | 5 |
| VI | Conclusion and future work | 5 |
| VII | Acknowledgment | 6 |
| | References | 7 |

I. ABSTRACT

One of the ultimate goals of artificial intelligence is the field of NLP is to train agents capable of maintaining a conversation with a human for some predefined period of time. Although this goal is still so far from achieving, many attempts were made to push the boundaries of what such agents are capable of doing. Das et al. proposed a visual dialog task in [1] where the AI-agent is challenged to answer questions about a picture provided first with only its caption and then, as the conversation evolves, with the dialog history. The images used for this task are collected from the COCO dataset [2] and the conversations are generated using the Amazon Mechanical Turk (AMT).

In [1], a variety of encoder-decoder network architectures were suggested: 3 encoders - Late Fusion, Hierarchical Recurrent Encoder and Memory Network- and 2 decoders - discriminative and generative. All of the combinations of the aforementioned encoders and decoders outperformed many established and sophisticated baselines. However, until now the agents were trained offline on a huge dataset of 120K images and 1,2M question-answer pairs.

We tried in this work to introduce an interactive training mode that allows us to integrate the human-in-the-loop concept into the visual dialog task. This interactive mode lets humans evaluate the agent and give it feedback that can be used for online training purposes. In addition, we investigate the possibility of improving the quality of learning using only limited amounts of data.

Further information about the dataset and the overall settings of the task can be found at <https://visdialdialog.org/>.

Keywords: Visual dialog, encoder-decoder networks, interactive training, human-in-the-loop.

II. INTRODUCTION

In the last decade, machine learning and AI in general dominated many fields in science and engineering as they outperformed the most sophisticated non learning-based methods. Whether in computer vision (CV), natural language processing (NLP), classification or even high-level AI tasks like playing chess or GO, machine learning established itself as the number one option when trying to deal with and solve such tasks. In the last years, new problems started to emerge dealing with multi-modal data as a consequence of many fields becoming more and more intertwined with one another. In fact, some of the most abundant types of data are images and text which happen to be the main constituents of the task we will address in this work: Visual Dialog. It is one of the most daunting problems at the intersection between computer vision and natural language processing. The main objective here is to train intelligent agents capable of conversing with humans based on some image by accurately answering questions with respect to it. This has tremendous application potential and can be used to improve and facilitate people's lives. Some of the applications may include:

- Helping the visually impaired people to get through their day-to-day life.
Agent: Your friend Sarah has posted a new photo on Instagram.
Human: Is she smiling?
Agent: No, but she seems happy.
- Improving our experience with AI assistants.
Human: OK Google! can you see the TV in the living-room?
Agent: Yes.
Human: Is it on?
Agent: Yes.
Human: Turn it off.

Visual dialog is not the only task in the intersection between computer vision and natural language processing. Image captioning [3] and visual question answering (VQA) [4] are classic examples that fall under the same category. However, visual dialog has an utterly different definition than these two.

Visual dialog VS image captioning: The main difference between image captioning and visual dialog is that the former does not portray any form of human conversation whatsoever. It simply provides a human with an over-all description of some image.

Visual dialog VS visual question answering: Despite the fact that visual question answering comes a little bit closer to achieving human-machine interaction than image captioning, it falls short at maintaining a somewhat satisfactory conversation as long as its duration/length is concerned. VQA is in fact a special case of visual dialog consisting of only one dialog round, i.e. one question-answer pair. Although the visual dialog data and the VQA data seem to appear similar, they have two fundamental differences:

- **Coreference in dialog:** As each dialog has 10 rounds, i.e. 10 question-answer pairs, the presence of many pronouns becomes inevitable. Thus, the agent has to overcome some ambiguities and figure out

what/whom these pronouns refer to. Some statistics made in [1] show that the amount of such pronouns gradually increase as the conversation evolves which is in accordance with what actually happens in real world scenarios.

- **Temporal continuity in topics:** Das et al. showed that the conversations also have some continuity if the topics being discussed. In fact, they conducted a human study to annotate the questions of 40 random images, i.e. 400 questions in total. The annotation was based on predefined topics, e.g. *asking about an object, a scene, the weather or even the image itself*. Across the 10 rounds, approximately 4.55 topics were discussed suggesting that they are not independent of each other. This implies that the data has some sort of temporal topic continuity.

As you may have already guessed, acquiring data for the visual dialog task is a tedious and long process, especially if one needs huge amounts of data. This may put major limitations to the potential of this task as not many people have the appropriate hardware and the resources to collect such data. Another major issue of the current learning strategies is the dominance of offline training. The latter requires also special hardware such as powerful GPUs or the access to some cloud.

In this work we try to address these two issues at once by introducing an interactive mode that enables us to take action and actively contribute to the training of such agents whilst generating new data.

III. RELATED WORK

Many attempts were made to improve on the work of Das et al. [1]. In [5], Niu et al. introduced a novel recursive visual attention strategy to tackle the most daunting problem of visual dialog: co-reference between questions and dialog histories. [6] introduced a goal-driven training for visual question answering and dialog agents using deep reinforcement learning. Two agents Q-Bot and A-Bot communicate in natural language dialog so that Q-Bot can select an unseen image from a lineup of images. Although this improvement led to better performance at the downstream dialog-conditioned image-guessing task, this improvement saturates and starts degrading after a few rounds of interactions, and thus, cannot lead to a better visual dialog model. [7] tried to improve the latter and explained its failure by the fact that the repeated interactions between Q-Bot and A-Bot during self-talk are not informative with respect to the image. Their solution to this problem was devising a simple auxiliary objective that incentivizes Q-Bot to ask diverse questions. This reduced question repetitions and enabled A-Bot to explore a larger state space during reinforced learning.

IV. APPROACH

Our approach addresses one key improvement that was not considered by all of the aforementioned methods as the human interaction with the system was always left out of the equation. With other words, there was no human supervision during or after training. Our approach aims at solving this issue by introducing an interactive mode based

on the human-in-the-loop concept and online training/fine tuning.

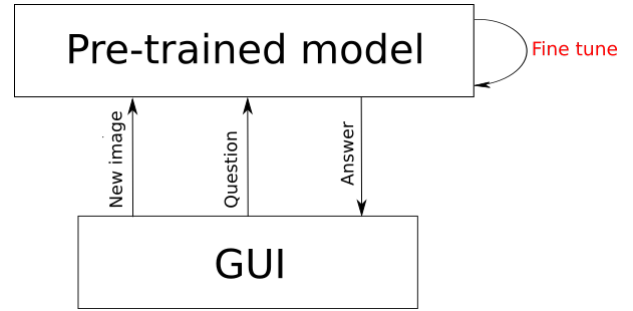


Figure 2: Interactive mode

Our approach consists of 3 main steps as Fig. 2 shows: pre-training, interactive mode, fine tuning (FT).

A. Pre-training

Since many efforts were made to collect valuable data for the visual dialog task - Factoid Question-Answer corpus [8], the Q&A dataset of DeepMind [9], 100K SimpleQuestions dataset [10], the bAbI dataset [11], the SQuAD dataset for reading comprehension [12] and the VisDial dataset [1] - it would not be such a wise choice not to take advantage of them. In this work, we used only v-0.9 of the VisDial dataset to pre-train some of the baseline networks introduced by Das et al. whose architectures are very well explained in [1]. More details about the different versions of the dataset can be found at <https://visualdialog.org/data>. In the following we adapt the same network naming convention as in [1], i.e. `<encoder>-<input>-<decoder>`. E.g. MN-QIH-G means that we use a Memory Network encoder with question-image-history input and generative decoder. Throughout this work, we only consider the generative decoders as they are easier to train on newly generated data, i.e. we do not need to provide a list of 100 candidate answers at each dialog round like in the discriminative setting. Furthermore, we think that generative decoders resemble more the way human interact with each other: one person asks a question and the other one generates an answer rather than picking up the best response from a set of possible answers.

1) *Hyperparameters:* All the networks have LSTMs with 2 layers and 512-dim hidden states. The word embedding is trained from scratch on the entire word corpus of the training data. The latter is shared across both the encoder and the decoder. We use the l2-normalized activations from the penultimate layer of VGG16 [13] to represent the image. As optimizer, we use Adam [14] and clamp the gradients clamped to $[-5, 5]$ at each step to avoid gradient explosion. Furthermore, we set the dropout value to 0.2 and we train the network for 20 epochs with batch size 8.

2) *Results:* Table I shows the performance of the pre-trained LF-QIH-G and MN-QIH-G networks measured on the validation set of the VisDial v0.9 containing 40K samples.

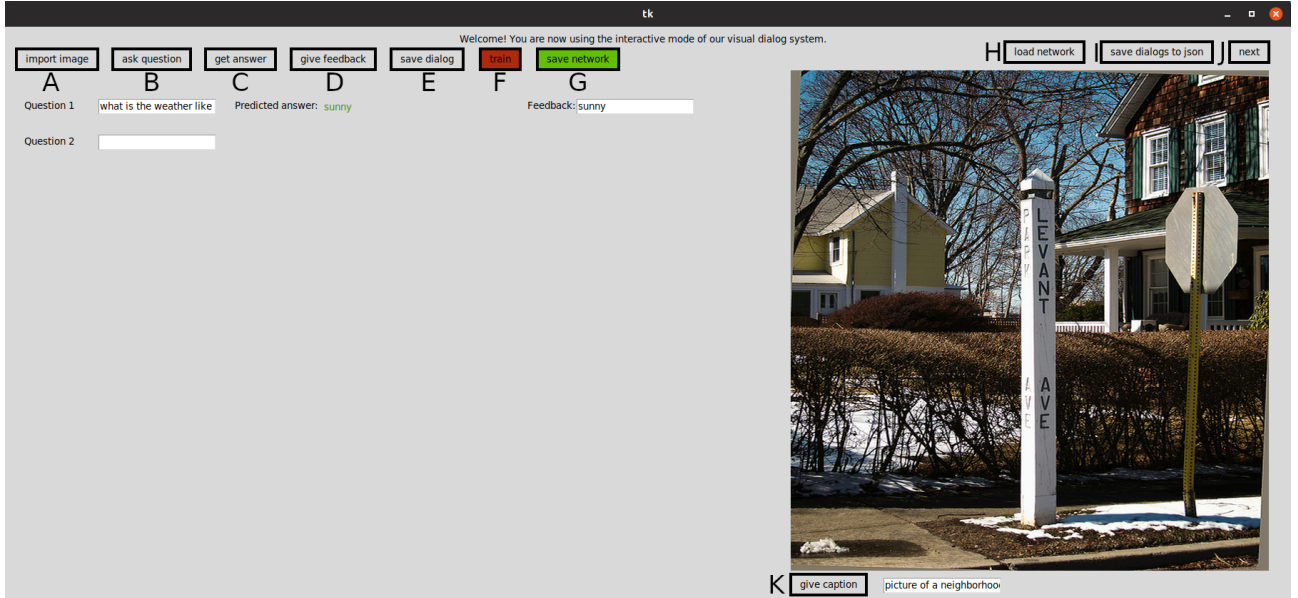


Figure 3: Annotated GUI of the interactive mode

| Model | MRR | R@1 | R@5 | R@10 |
|----------|--------|-------|-------|-------|
| LF-QIH-G | 0.4791 | 35.48 | 60.96 | 66.93 |
| MN-QIH-G | 0.4903 | 37.80 | 60.33 | 66.10 |

Table I: Performance on VisDial v0.9 measured by mean reciprocal rank (MRR) and recall @k.

As we can see, our networks achieved comparable results to those of [1]. In addition, based on the R@1 metric, which happens to be the most significant one, MN-QIH-G outperformed LF-QIH-G as [1] suggested. Thus, we can safely assume that our networks are ready to be deployed within our interactive mode. Furthermore, we expect them to produce reliable answers in real world scenarios. In fact, this will be solidified afterwards in section IV-B2 by examining their performance on our dataset. It is worth mentioning, that since MN-QIH-G performed better than LF-QIH-G, it was used in the rest of this work for generating our dataset as well as for FT.

B. Interactive mode

1) *Usage description:* After pre-training, the network is used for inference on unseen data by our interactive mode. The latter takes the form of a graphical user interface (GUI) as illustrated in Fig. 3. The interactive mode resembles somehow the process of collecting data described in [1]. The main difference consists in generating the data whilst assessing the network performance. Our goal here is to collect new data and use them to incrementally fine tune the last linear layer of the decoder net and therefore improve the overall performance of the system.

The user starts off by loading a pre-trained network and then importing an image using the button H and A of Fig. 3 respectively. If the image is successfully imported, it will be displayed on the right-hand-side of the GUI alongside with the "give caption" button underneath it, i.e. button K of Fig. 3. The user has to annotate the image with a suitable caption and then hit the "Enter" key to save it and be able to start asking questions. This can be done with the

help of button B of the GUI. After the user has asked the question, he/she has to save it by hitting "Enter". At this stage, the user is able to get the generated response of the network using the button C of the interface. Since we want to compute some kind of accuracy afterwards, the user has to give a feedback to the network using the button D. This feedback will play the role of the ground truth of the corresponding question. The dialog can be saved using the button E and only after saving it can we train, i.e. fine tune the network (button F). This option makes online learning possible and allows us to assess the gradual improvements of our system as we continue interacting with it. The fine tuned version of the network can be saved at any time using the button G of our interface. After the user finishes asking questions about one particular images, he/she can proceed to the next one using the button J. It is worth mentioning that one can ask 10 questions at maximum on one particular picture. At the end, the generated dialogs can be logged into a json file using the button I of the interface.

2) *Our data:* In this section, we will describe the data we collected using our interactive mode. We used the first 100 images from the visual dialog testing set that can be found at <https://visualdialog.org/data>. The dialogs generated have an average length of 7.08 rounds. The longest dialog consists of 9 rounds and the shortest of only 3. The following figures illustrate some examples.



Figure 4: Example 1

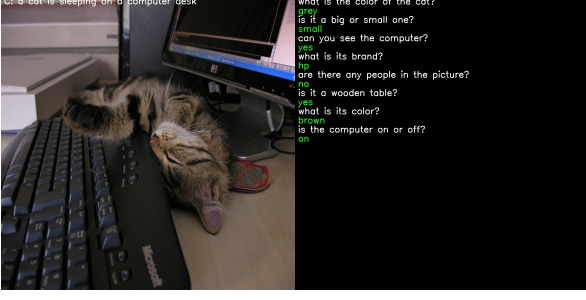


Figure 5: Example 2

C. Fine tuning

The fine tuning step consists of updating the weights of the last fully connected layer of the decoder network. This layer is responsible for generating words from the LSTM outputs. For further details concerning the architecture of the network, please refer to [1]. Thus, fine tuning the last layer will update the distribution of the words given all the previous ones as well as the encoder output, i.e.

$$P(w_i|w_{i-1}, \dots, w_1, out_{encoder}),$$

where w_i and $out_{encoder}$ denote the word generated at step i and the encoder output respectively. The word at step i is the one that maximizes the previous probability, i.e.

$$w_i = \arg \max_{w \in V} P(w|w_{i-1}, \dots, w_1, out_{encoder}),$$

where V is our vocabulary introduced in section IV-A1. Fig. 6 illustrates how the answers are generated using the generative decoder. In this work, we fix $w_1 = \langle \text{SOS} \rangle$, and generate words until $w_i = \langle \text{EOS} \rangle$ or $n = n_{max}$. We used a value of 10 for n_{max} .

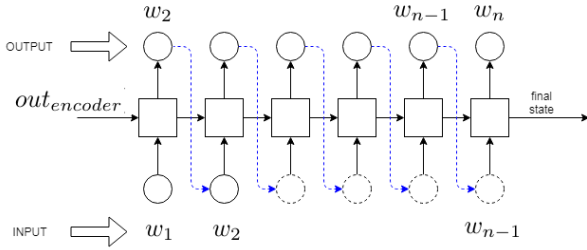


Figure 6: Answer generation diagram

The hyperparameters used for fine tuning are the same as described in IV-A1. The only difference is that we train only for a total of 4 epochs with batch size 1.

V. RESULTS

A. On our data

We first of all compare the performances of the network on the data we generated via the interactive mode, i.e. the data used for FT, before and after FT. This is a good way to see whether the network has learned something from the FT step or not. To do so, we investigated the generated answers in both cases and computed the corresponding network accuracies. This was done on a per dialog basis, i.e. a human assessed the goodness of the answers of the system and gave it 1 point if the answer was acceptable

and 0 points otherwise. Note here that there is not a unique right/acceptable answer for each question. Thus, basing the accuracy solely on the ground truths, i.e. the feedback of the user as mentioned in section IV-B1, does not make much sense. E.g. if the question and its ground truth were "How big is the room?" and "It is big" respectively than it would be unfair to reject an answer from the system like "The room is big" or just simply "big". Some datasets like VisDial [1] v.1.0 addressed this issue differently by introducing the so called "dense annotations". However, we opted for the simpler manual evaluation as our dataset only consists of 100 samples. The accuracies were then used to fit a PDF as illustrated in Fig. 7.

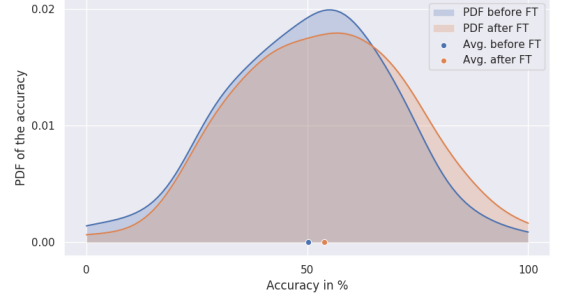


Figure 7: PDFs of the accuracy

As we can see, the PDF of the accuracy has shifted with a considerable amount towards 100% which indicates that the number of dialogs with high accuracy rates has increased after FT at the expense of those with low accuracy rates. This suggests that we can indeed improve the performance of the system by simply updating the last layer only. In fact, the network after FT scored an average accuracy of approx. 53.80%. That is a 3.55% increase in comparison with the 50.25% average accuracy before FT. This difference as small as it may appear is not to be underestimated especially due to the fact that the FT data is negligible in comparison with the original data used to pre-train the network.

We also want to emphasize on the fact that the fine tuned system did not simply reproduce the same answers given as ground truth by the user. In fact, it produced utterly different but acceptable answers in many cases. The following example illustrates this very clearly. We use Q_i, GT_i, A_i, A_i^{ft} to denote the question, the ground truth, the answer before FT and the answer after FT of dialog round i respectively.



Figure 8: A picture of a mother zebra and its son

Caption: A picture of a mother zebra and its son.

Q₁: Where are the zebras?

GT₁: In the wild.

A₁: They are standing in the water.

A₁^{ft}: In a field.

Q₂: Do you see green plants?

GT₂: Yes.

A₂: no

A₂^{ft}: No.

Q₃: Are they running or standing?

GT₃: Standing.

A₃: Standing.

A₃^{ft}: Standing.

Q₄: Are they eating grass?

GT₄: No.

A₄: No.

A₄^{ft}: Yes.

Q₅: Is the sun shining?

GT₅: Yes.

A₅: Yes.

A₅^{ft}: Yes.

Q₆: Do they look friendly?

GT₆: Yes.

A₆: Yes.

A₆^{ft}: Yes.

We are mostly interested in the first dialog round where the system initially failed to answer the question correctly. However, after FT it succeeded in generating a totally acceptable response which is different than the ground truth: *"In a field"* vs *"in the wild"*.

B. On a subset of VisDial v1.0

In the previous section, we showed that our FT method did indeed improve the performance of the system on our data, i.e. the data used for FT, and the network succeeded not to reproduce the same answers given as ground truths. However, we still wanted to see how the network performs on unseen data, i.e. data not seen during both pre-training and FT.

Although we do not expect the fine tuned network to outperform the pre-trained one since the amount of FT data is very limited and only the last layer was updated, we must not witness any huge drop in the generalizability of the system as this defies the whole purpose of our FT

approach. In order to assess the generalizability of our fine tuned network, we used a subset of the VisDial v1.0 dataset. This subset was randomly chosen and contains 100 samples and was not used for pre-training the network as we used VisDial v0.9 for that purpose. Another reason that made us opt for this dataset is the fact that it comes with the dense annotations used to compute the normalized discounted cumulative gain (NDCG). It is a very fair metric to measure the performance of such AI agents since it very elegantly solves the problem of the non-uniqueness of answers addressed in section V-A. Detailed information about the NDCG metric can be found at <https://visualdialog.org/challenge/2019#evaluation>. In the following, we took the same approach to assess the networks as we did before, i.e. we investigate the PDF of each metric separately before and after FT.

| Model | NDCG ×100 | MRR | R@1 | R@5 | R@10 |
|-----------|--------------|-------|-------|-------|-------|
| Before FT | 57.89 | 46.01 | 34.90 | 56.30 | 62.60 |
| After FT | 57.67 | 45.56 | 34.10 | 56.10 | 62.40 |

Table II: Performance on a subset of VisDial v1.0 measured by normalized discounted cumulative gain (NDCG), mean reciprocal rank (MRR), recall @k and mean rank.

Table II and Fig. 9 - 13 show the influence of FT when the network is tested on unseen data. Although the density functions of the different metrics are slightly modified after FT, their average value are practically exactly the same as table II shows. This suggests that FT does not deteriorate the generalizability of the system. However, by using this interactive mode more intensively and thus generating richer data we can outperform the baseline system and achieve higher results after FT.

VI. CONCLUSION AND FUTURE WORK

In this work, we have addressed the issues of efficiently acquiring valuable data for the visual dialog problem while incorporating human supervision. This was done by introducing an interactive mode that helps generate dialogs while interacting with a pre-trained system. Moreover, we have shown that incrementally fine tuning a pre-trained model can push its performance and helps it improve its affinity and creativity in answering questions as we have seen in the example of section IV-B2. However, FT only the last layer of the decoder and with limited amounts of data seems not to have a huge impact when the network is tested on unseen data, i.e. data not used in both pre-training and FT. In fact, this offers a significant margin to improve our method in the future. Moreover, the choice of the optimal weights to fine tune remains a very daunting yet crucial task that we aim to solve in future endeavors in order to improve the generalizability of our agent. Our generated data using the interactive mode and our code can be found at <https://github.com/adnenabdessaied/VisDiag>.

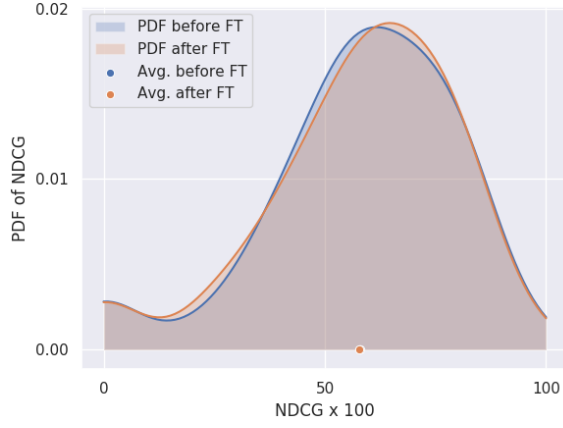


Figure 9: PDFs of the NDCG

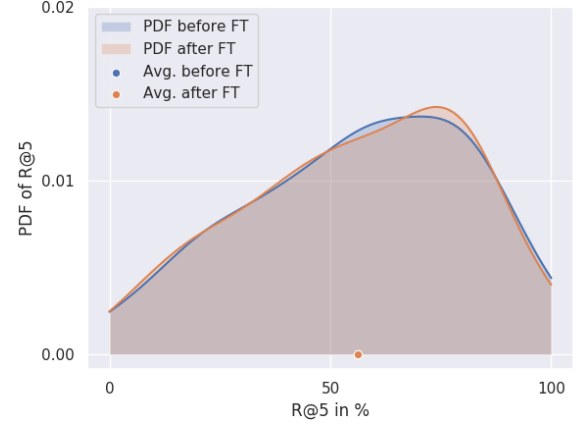


Figure 12: PDFs of the R@5

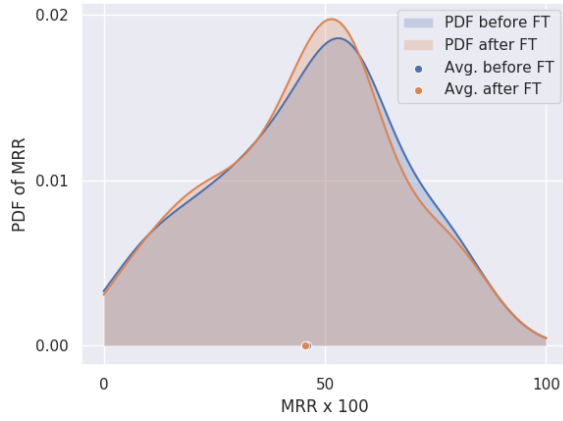


Figure 10: PDFs of the MRR

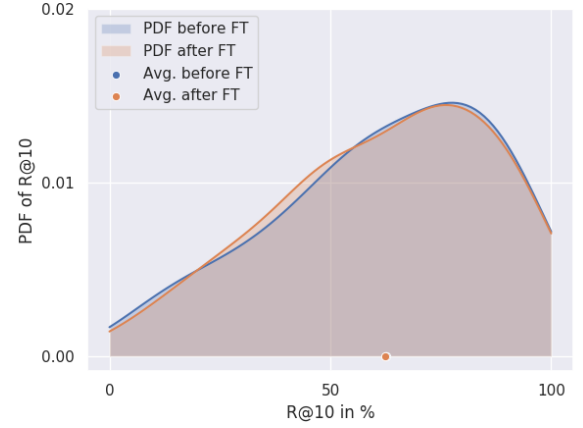


Figure 13: PDFs of the R@10

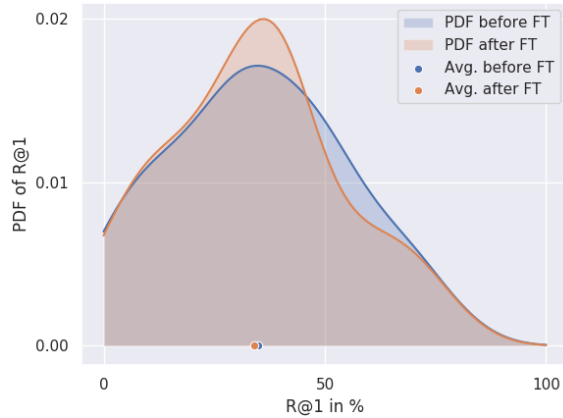


Figure 11: PDFs of the R@1

VII. ACKNOWLEDGMENT

At the end, we want to thank the Visual Dialog challenge creators who made huge efforts in collecting the data we used to pre-train our networks as well as providing good deep learning architectures capable of solving such a challenging task [1].

REFERENCES

- [1] Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, Jose M. F. Moura, Devi Parikh, and Dhruv Batra. Visual dialog. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [2] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014.
- [3] R. Subash, R. Jebakumar, Yash Kamdar, and Nishit Bhatt. Automatic image captioning using convolution neural networks and LSTM. *Journal of Physics: Conference Series*, 1362:012096, nov 2019.
- [4] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *The IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- [5] Yulei Niu, Hanwang Zhang, Manli Zhang, Jianhong Zhang, Zhiwu Lu, and Ji-Rong Wen. Recursive visual attention in visual dialog. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6679–6688, 2019.
- [6] Abhishek Das, Satwik Kottur, José M.F. Moura, Stefan Lee, and Dhruv Batra. Learning cooperative visual dialog agents with deep reinforcement learning. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [7] Dhruv Batra Devi Parikh Abhishek Das Vishvak Murahari, Prithvijit Chattopadhyay. Improving generative visual dialog by answering diverse questions. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*, 2019.
- [8] Kelvin Jiang, Dekun Wu, and Hui Jiang. FreebaseQA: A new factoid QA data set matching trivia-style question-answer pairs with Freebase. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 318–323, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
- [9] Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. Teaching machines to read and comprehend. In *Advances in neural information processing systems*, pages 1693–1701, 2015.
- [10] Antoine Bordes, Nicolas Usunier, Sumit Chopra, and Jason Weston. Large-scale simple question answering with memory networks. *arXiv preprint arXiv:1506.02075*, 2015.
- [11] Jason Weston, Antoine Bordes, Sumit Chopra, Alexander M Rush, Bart van Merriënboer, Armand Joulin, and Tomas Mikolov. Towards ai-complete question answering: A set of prerequisite toy tasks. *arXiv preprint arXiv:1502.05698*, 2015.
- [12] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.
- [13] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [14] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.