

Look, Read and Ask: Learning to Ask Questions by Reading Text in Images

Soumya Jahagirdar¹[0000–0002–3460–9151], Shankar Gangisetty¹[0000–0003–4448–5794], and Anand Mishra²[0000–0002–7806–2557]

¹ KLE Technological University, Hubballi, India
{01fe17bcs212,shankar}@kletech.ac.in

² Vision, Language, and Learning Group (VL2G)
IIT Jodhpur, India
mishra@iitj.ac.in

Abstract. We present a novel problem of text-based visual question generation or TextVQG in short. Given the recent growing interest of the document image analysis community in combining text understanding with conversational artificial intelligence, e.g., text-based visual question answering, TextVQG becomes an important task. TextVQG aims to generate a natural language question for a given input image and an automatically extracted text also known as OCR token from it such that the OCR token is an answer to the generated question. TextVQG is an essential ability for a conversational agent. However, it is challenging as it requires an in-depth understanding of the scene and the ability to semantically bridge the visual content with the text present in the image. To address TextVQG, we present an QCR-consistent visual question generation model that Looks into the visual content, Reads the scene text, and Ask a relevant and meaningful natural language question. We refer to our proposed model as OLRA. We perform an extensive evaluation of OLRA on two public benchmarks and compare them against baselines. Our model – OLRA automatically generates questions similar to the public text-based visual question answering datasets that were curated manually. Moreover, we significantly outperform baseline approaches on the performance measures popularly used in text generation literature.

Keywords: Visual Question Generation (VQG) · Conversational AI · Visual Question Answering (VQA).

1 Introduction

“To seek truth requires one to ask the right questions.”

Suzy Kassem

Developing agents that can communicate with a human has been an active area of research in artificial intelligence (AI). The document image analysis community has also started showing interest in this problem lately as evident from

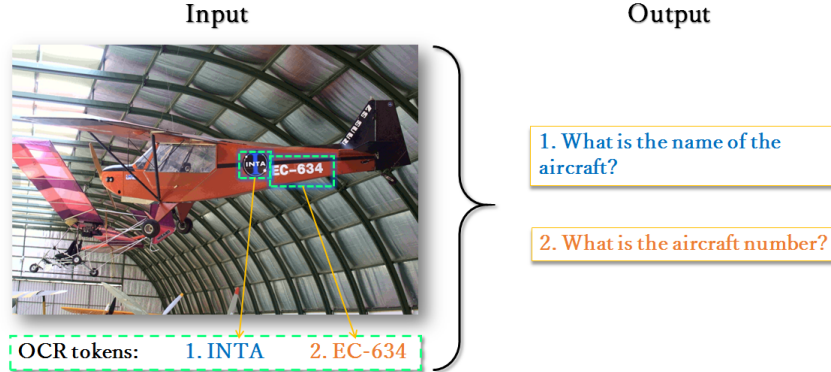


Fig. 1. TextVQG. We introduce a novel problem of visual question generation by leveraging text in the image. Given an image and an OCR token automatically extracted from it, our goal is to generate a meaningful natural language question whose answer is the OCR token.

the efforts of the community on text-based visual question answering [4, 25, 28, 39], and ICDAR 2019 robust reading challenge on scene text visual question answering [1]. Visual question answering (VQA) is only one desired characteristic of a conversational agent where the agent answers a natural language question about the image. An intelligent agent should also have the ability to ask meaningful and relevant questions with respect to its current visual perception. Given an image, generating meaningful questions, also known as visual question generation (VQG) is an essential component of a conversational agent [30]. VQG is a precursor of visual dialogue systems and might help in building large-scale VQA datasets automatically. To the best of our knowledge, the document image analysis community has not yet looked into the important problem of VQG leveraging text. In this work, we fill up this gap prevailing in the literature by introducing a novel problem of text-based visual question generation or TextVQG in short.

The TextVQG has the following goal: given a natural image containing text, and an OCR token extracted from it, generate a natural language question whose answer is the OCR token. This problem is challenging because it requires in-depth semantic interpretation of the text as well as the visual content to generate meaningful and relevant questions with respect to the image. VQG has been a well-explored area in vision and language community [13, 17, 30, 43]. However, these works often ignore the text appearing in the image, and only restrict themselves to the visual content while generating questions. It should be noted that text in the image helps in asking semantically meaningful questions connecting visual and textual content. For example, consider the image shown in Fig. 1. Given an image of an aircraft and two words – *INTA* and *EC-634* written on it, we aim to automatically generate questions such as “What is the name of the

aircraft?” and “What is the aircraft number?”.

Baselines: Motivated by the baselines presented in the VQG literature [30], a few of the plausible approaches to address TextVQG are as follows: (i) **maximum-entropy language model (MELM)**: which uses the extracted OCR tokens based on their confidence scores along with detected objects to generate question, (ii) **seq2seq model**: which first generates a caption for the input image, and then, this caption is fed into a seq2seq model to generate a question, and (iii) **GRNN model**: where the CNN feature of the input image is passed to a gated recurrent unit (GRU) to generate questions. We empirically observe that these methods often fall short in performance primarily due to their inability (a) to semantically interpret the text and visual content jointly, and (b) to establish consistency between generated question and the OCR token which is supposed to be the answer. To overcome these shortcomings, we propose a novel TextVQG model as described below.

Our Approach: To encode visual content, scene text, its position and bringing consistency between generated question and OCR token; we propose an OCR-consistent visual question generation model that Looks into the visual content, Reads the scene text, and Ask a relevant and meaningful natural language question. We refer to this architecture for TextVQG as OLRA. OLRA begins by representing visual features using pretrained CNN, extracted OCR token using FastText [5] followed by an LSTM, and positions of extracted OCR using positional representations. Further, these representations are fused using a multimodal fusion scheme. The joint representation is further passed to a one-layered LSTM-based module and a maximum likelihood estimation-based loss is computed between generated and reference question. Moreover, to ensure that the generated question and corresponding OCR tokens are consistent with each other, we add an OCR token reconstruction loss that is computed by taking l_2 -loss between the original OCR token feature and the representation obtained after passing joint feature to a multi-layer perception. The proposed model OLRA is trained in a multi-task learning paradigm where a weighted combination of the OCR-token reconstruction and question generation loss is minimized to generate a meaningful question. We evaluate the performance of OLRA on TextVQG, and compare against the baselines. OLRA significantly outperforms baseline methods on two public benchmarks, namely, ST-VQA [4] and TextVQA [39].

Contributions of this paper: The major contributions of this paper are two folds:

1. We draw the attention of the document image analysis community to the problem of text-based visual question generation by introducing a novel task referred to as TextVQG. TextVQG is an important and unexplored problem in the literature with potential downstream applications in building visual dialogue systems and augmenting training sets of text-based visual question answering models. We firmly believe that our work will boost ongoing re-

search efforts [4, 25, 39, 37] in the broader area of conversational AI and text understanding.

2. We propose OLRA – an OCR-consistent visual question generation model that looks into the visual content, reads the text and asks a meaningful and relevant question for addressing TextVQG. OLRA automatically generates questions similar to the datasets that are manually curated. Our model viz. OLRA significantly outperforms the baselines and achieves a BLEU score of 0.47 and 0.40 on ST-VQA [4] and TextVQA [39] datasets respectively.

2 Related Work

The performance of scene text recognition and understanding has significantly improved over the last decade [12, 21, 27, 31, 36, 42]. It has also started influencing other computer vision areas such as scene understanding [15], cross-modal image retrieval [24], image captioning [37], and visual question answering (VQA) [4, 25, 39]. Among these, text-based VQA works [4, 25, 39] in the literature can be considered one of the major steps by document image analysis community towards conversational AI. In the follow-up sections, we shall first review the text-based VQA followed by VQG which is the primary focus of this work.

2.1 Text-based Visual Question Answering

Traditionally, VQA works in the literature focus only on the visual content [2, 14, 32], and ironically fall short in answering the questions that require reading the text in the image. Keeping the importance of text in the images for answering visual questions, researchers have started focusing on text-based VQA [4, 20, 28, 39]. Among these works, authors in [39] use top-down and bottom-up attention on text and visual objects to select an answer from the OCR token or a vocabulary. The scene text VQA model [4] aims to answer the questions by performing reasoning over the text present in the natural scene. In OCR-VQA [28], OCR is used to read the text in the book cover images, and a baseline VQA model is proposed for answering questions enquiring about author name, book title, book genre, publisher name, etc. Typically, these methods are based on a convolutional neural network (CNN) to perform the visual content analysis, a state-of-the-art text recognition engine for detecting and reading text, and a long short-term memory (LSTM) network to encode the questions. More recently, T-VQA [20] presents a progressive attention module and a multimodal reasoning graph for reading and reasoning. Transformer and graph neural network-based models have also started gaining popularity for addressing text-based VQA [9, 11]. Another direction of work in text-based visual question answering is text-KVQA [38] where authors propose a VQA model that reads textual content in the image, connects it with a knowledge graph, and performs reasoning using a gated graph neural network to arrive at an accurate answer.

2.2 Visual Question Generation

VQG is a dual task of VQA and is essential for building visual dialogue systems [8, 13, 17, 19, 29, 30, 43, 44]. Moreover, the ability to generate relevant question to the image is a core to in-depth visual understanding. Question generation from images as well as from raw text have been a well-studied problem in the literature [10, 13, 17, 22, 35, 43]. Automatic question generation techniques in NLP have also enabled chat-bots [10, 22, 35].

Among visual question generation works, authors in [30] focused on generating natural and engaging questions from the image, and provided three distinct datasets, each covering object to event-centric images. They proposed three generative models for tackling the task of VQG which we believe are the baselines for any novel VQG tasks, and we adopt these models for TextVQG and compare it with the proposed model. In [8], authors generated diverse questions of different types [8], such as, when, what, where, which, and how questions. In [17], a goal-driven variational auto-encoder model is used to generate questions by maximizing the mutual information between visual content as well as the expected answer category. In [19], authors posed VQG as a dual-task to VQA and jointly addressed both the task.

These works in the literature restrict their scope to asking questions only about visual content and ignore any scene text present in the image. We argue that conversational agents must have the ability to ask questions by semantically bridging text in the image with visual content. Despite its fundamental and applied importance, the problem of TextVQG - text-based visual question generation has not been looked into in the literature. We fill this gap in the literature through our work.

3 Proposed Model: OLRA

Given a natural image containing text and an OCR token extracted from the image, TextVQG aims to generate a natural language question such that the answer is the OCR token. To solve this problem, the proposed method should be able to successfully recognize text and visual content in the image, semantically bridge the textual and visual content and generate meaningful questions. For text detection and recognition, we rely on one of the successful modern scene text detection and recognition methods [3, 36].³ Further, the pre-trained CNN is used for computing the visual features and an LSTM is used to generate a question. The overall architecture of the proposed model viz. OLRA is illustrated in Fig. 2. OLRA has the following three modules:

(i) Look and Read Module: In this module, we extract convolutional features from the given input image I . We use pre-trained CNN ResNet-50 which gives us 512-dimensional features ϕ_I . Further, since our objective is to generate questions by relating visual content and text appearing in the image, we use [3, 36] for

³ For one of the datasets namely TextVQA, OCR-tokens extracted from Rosetta [6] are provided with the dataset.

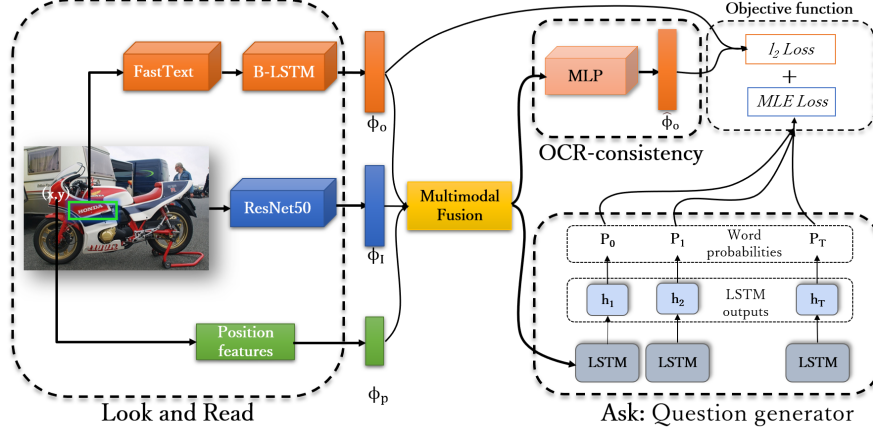


Fig. 2. OLRA. Our proposed visual question generator architecture viz. OLRA has three modules, namely, (i) Look and Read, (ii) OCR-consistency, and (iii) Ask module. Note: since OCR tokens can be more than a single word, we pass their FastText embeddings to a B-LSTM. Please refer to Section 3 for more details. **[Best viewed in color]**.

detecting and recognizing text. We prefer to use [3, 36] for text detection and recognition due to its empirical performance. However, any other scene text recognition module can be plugged in here. Once we detect and recognize the text, we obtain its FastText [5] embedding. Since OCR tokens can be proper nouns or noisy, and therefore, can be out of vocabulary; we prefer FastText over Word2Vec [26] or Glove [33] word embeddings. The FastText embedding of the OCR token is fed to a bi-directional long short-term memory (B-LSTM) to obtain a 512-dimensional OCR-token feature or ϕ_o .

Further, positions of OCR can play a vital role in visual question generation. For example as illustrated in Fig. 3, in a sports scene, the number appearing in a sportsman’s jersey and an advertisement board will require us to generate questions such as “What is the jersey number of the player who is pitching the ball?” and “What is the number written on advertisement board?” respectively. In order to use OCR token positions in our framework, we use 8 features i.e., topleft-x, topleft-y, width, height, rotation, yaw, roll, and pitch of the bounding box as ϕ_p .⁴

Once these three features, i.e., ϕ_o : OCR features, ϕ_I : image features, and ϕ_p : positions features are obtained, our next task is to learn joint representation by judiciously combining them. To this end, we use a multi-layer perceptron \mathcal{F} to obtain joint representation Ψ as follows:

⁴ When we use CRAFT [3] for text detection, we use only first four positional features i.e., topleft-x, topleft-y, width and height of the bounding box.

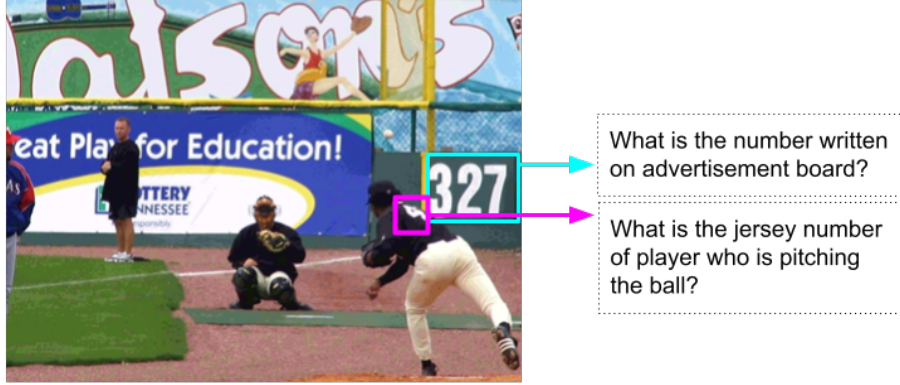


Fig. 3. Importance of positions of OCR tokens for question generation. Numbers detected in different positions in the image i.e., number 4 and 327 in this example demands generating different questions.

$$\Psi = \mathcal{F}(W, \phi). \quad (1)$$

Here, W is a learnable weight matrix and $\phi = [\phi_o; \phi_I; \phi_p]$ is concatenation of features. The joint feature obtained at this module Ψ is a 512-dimensional vector that is fed to OCR token-reconstruction and question generator as described next.

(ii) OCR-consistency Module: After learning the joint representation of both image and OCR features, the model should be able to generate questions given specific tokens. To ensure this consistency, it is important that joint representation learned after fusion (Ψ) preserves the OCR representation (ϕ_o). We reconstruct $\hat{\phi}_o$ by passing joint feature Ψ to a multi-layer perceptron to obtain the reconstructed answer representation. We minimize the l_2 -loss between the original answer and the reconstructed answer representations, i.e.,

$$\mathcal{L}_a = \|\phi_o - \hat{\phi}_o\|_2^2. \quad (2)$$

(iii) Ask Module: Our final module generates the relevant questions to the image. We use a decoder LSTM to generate the question q' from combined space as explained in the previous section. The minimization of maximum-likelihood error (MLE) between the generated question and the true question in the training set makes the model learn to generate appropriate questions. The 512-dimensional joint feature vector serves as the initial state to an LSTM. We produce one word of the output question at a time using the LSTM. This happens until we encounter an end-of-sentence token. At each step during decoding, we feed in an embedding of the previously predicted word and predict the next word. The loss between generated and true questions guides the generator to generate better

questions. The obtained joint feature Ψ from the previous module acts as an input to the question generator. This question generator captures the time-varying characteristics and outputs questions related to the image and text present in the image. As mentioned above, we minimize the MLE objective \mathcal{L}_q between generated question q' and ground truth question q .

Training: OLRA is trained in a multitask learning paradigm for two tasks, i.e., OCR-token reconstruction and question generation. We use training set of text-based visual question answering datasets to train the model by combining both MLE loss \mathcal{L}_q and OCR-token reconstruction loss \mathcal{L}_a , as follows:

$$\mathcal{L} = \mathcal{L}_q + \lambda \mathcal{L}_a, \quad (3)$$

where λ is a hyperparameter which controls relative importance in optimizing these two losses, and ensures that the generated questions have better semantic structure with respect to the textual content along with the scene of the image.

4 Experiments and Results

In this section, we experimentally validate our proposed model for TextVQG. We first discuss the datasets and performance measures in Section 4.1 and Section 4.2 respectively. We then describe the baseline models for TextVQG in Section 4.3 followed by implementation details in Section 4.4. The quantitative and qualitative experimental results of the proposed model and the comparative analysis is provided in Section 4.5.

4.1 Datasets

As this is the first work towards a visual question generation model that can read the text in the image, there is no dedicated dataset available for this task. We, therefore, make use of the two popular text-based VQA datasets, namely, TextVQA [39] and ST-VQA [4]. Note that, unlike these datasets where originally the task is to answer a question about the image, our aim is to generate questions. In other words, given an image and an OCR token (often answers in these datasets), our proposed method learns to automatically generate questions similar to these manually curated datasets. A brief statistical description related to these datasets are as follows:

1. **ST-VQA** [4]: consists of 23,038 images with 31,791 question-answer pairs obtained from different public datasets. A total of 16,063 images with 22,162 questions and 2,834 images with 3,910 questions are used for training and testing respectively, considering only the question-answer pairs having OCR tokens as their answers.
2. **TextVQA** [39]: comprises of 28,418 images obtained from openImages [16] with 45,336 questions. It also provides OCR tokens extracted from Rosetta [6]. We use 21,953 images with 25,786 questions and 3,166 images with 3,702 questions in all as training and testing set respectively, considering only the question-answer pairs having OCR tokens as their answers.

4.2 Performance metrics

Following the text-generation literature [7], we use popular evaluation measures such as bilingual evaluation understudy (BLEU), recall-oriented understudy for gisting evaluation - longest common sub-sequence (ROUGE-L), and metrics for evaluation of translation with explicit ordering (METEOR). The BLEU score compares n -grams of the generated question with the n -grams of the reference question and counts the number of matches. The ROUGE-L metric indicates similarity between two sequences based on the length of the longest common sub-sequence even though the sequences are not contiguous. The METEOR is based on the harmonic mean of uni-gram precision and recall and is considered a better performance measure in the text generation literature [18]. Higher values of all these performance measures imply better matching of generated questions with the reference questions.

4.3 Baseline models

Inspired by the VQG models in literature [30], we present three baseline models by adopting them for TextVQG task.

Maximum Entropy Language Model (MELM): Here, a set of OCR tokens extracted from the image along with a set of detected objects using Faster-RCNN [34] are fed to a maximum entropy language model to generate a question.

Seq2seq model: In this baseline, we first obtain caption of the image using a method proposed in [23]. Then, the generated caption and the extracted OCR tokens from the images are passed to a seq2seq model [40]. In seq2seq model, the encoder contains an embedding layer followed by an LSTM layer and for decoding, we use an LSTM layer followed by a dense layer. The seq2seq model is trained for question generation tasks using the training set of datasets described earlier.

Gated Recurrent Neural Network (GRNN): In this model, we obtain visual features using InceptionV3 [41]. This yields a feature vector of $1 \times 1 \times 4096$ dimensions that are then passed to a GRU to generate a question. We train GRU for question generation by keeping all the layers of InceptionV3 unchanged.

4.4 Implementation details

We train our proposed network using the Adam optimizer with a learning rate of $1e-4$, batch size of 32, and a decay rate of 0.05. The value of λ in Equation 3 is set to 0.001. The maximum length of the generated questions is set to 20. We train the model for 10 epochs. In multi-modal fusion, a two-layer attention network is used with feature sizes of 1032 and 512 respectively. The model is trained on a single NVIDIA Quadro P5000.

4.5 Results and discussions

Quantitative Analysis. We evaluate the performance of our proposed model i.e., OLRA, and compare it against three baseline approaches, namely, MELM,

Table 1. Comparison of OLRA with baselines. We observe that the proposed model viz. OLRA clearly outperforms baselines on both the datasets for TextVQG.

Method	<i>ST-VQA</i> [4] dataset			<i>TextVQA</i> [39] dataset		
	BLEU	METEOR	ROUGE-L	BLEU	METEOR	ROUGE-L
MELM	0.34	0.12	0.31	0.30	0.11	0.30
Seq2seq	0.29	0.12	0.30	0.27	0.11	0.29
GRNN	0.36	0.12	0.32	0.33	0.12	0.30
Ours (OLRA)	0.47	0.17	0.46	0.40	0.14	0.40

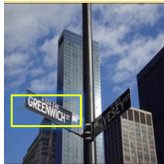


	Ground Truth	MELM	seq2seq	GRNN	Ours (OLRA)
	What does the street sign say?	What is the street name?	What is another name for a street?	What is the sign in the bottom?	What is the name of the city on street sign?
	What name of the author of the book?	Which book is it?	Who reads book?	What is name to the left?	What is the name of the book?
	What is the brand of the milk used?	What is the brand?	Bananas, bananas, and what else are on a plate of food?	What is the name of the company?	What is brand of the milk?

Fig. 4. Qualitative comparison. Visual question generation by MELM, seq2seq, and GRNN baselines, and ours on a set of images from ST-VQA dataset.

seq2seq, and GRNN on ST-VQA and TextVQA datasets. The comparative result is shown in Table 1 using performance measures discussed in Section 4.2. Note that higher values for all these popularly used performance measures are considered superior.

Among the three baseline approaches, GRNN generates comparatively better questions. Our proposed model i.e., OLRA significantly outperforms all the baseline models. For example, on ST-VQA dataset, OLRA improves BLEU score by 0.11, METEOR score by 0.05, and ROUGE-L score by 0.14 as compared to the most competitive baseline i.e., GRNN. It should be noted that under these performance measures these gains are considered significant [7]. We observe similar performance improvement on TextVQA dataset as well.

Qualitative Analysis. We perform a detailed qualitative analysis of the baselines as well as our proposed model. We first show a comparison of generated questions using all the three baselines versus OLRA in Fig. 4. We observe that

Table 2. Ablation study. BLEU scores analysis (i) with inclusion of positional information and (ii) with n-word answers on both the datasets.

OLRA	<i>ST-VQA</i> [4] dataset	<i>TextVQA</i> [39] dataset
w/o position	0.44	0.39
w/ position	0.45	0.39
w/ position and OCR-consistency	0.47	0.40
w/ 1-word answers	0.48	0.41
w/ 2-word answers	0.47	0.39
w/ 3-word answers	0.46	0.40

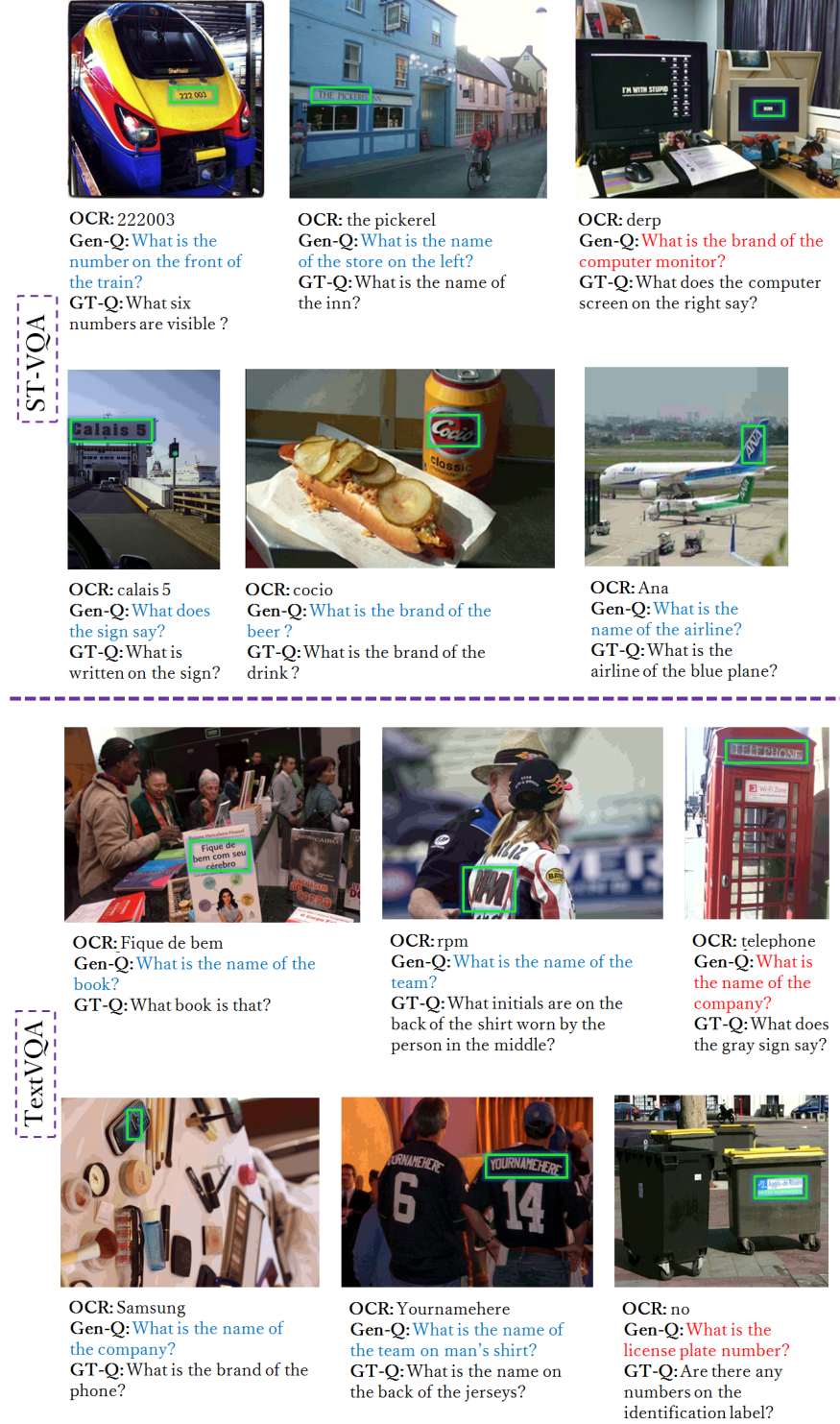
the baselines are capable of generating linguistically meaningful questions. However, they do not fulfill the sole purpose of **Look**, **Read** and **Ask**-based question generation. For the first example in Fig. 4 the expected question is “What does the street sign say?”, the baseline approaches and the proposed model generate nearly the same question. But, as the complexity of the scene increases, the baseline models fail to generate appropriate questions, and the proposed model due to its well-designed look, read, ask, and OCR-consistency modules, generates better questions. For example, in the last row of Fig. 4, MELM generates “What is the brand?” but it fails to specify the kind of the product in the scene. Whereas the proposed model generates “What is the brand of the milk?” which is very close to the target question.

Further, more results of our model are shown in Fig. 5. Here, we represent OCR token using green bounding boxes, correctly generated questions in the blue color text, and incorrectly generated questions in the red color text. Consider the example of a train image in Fig. 5. Here, our method successfully generates the question “What is the number on the front of the train?”. Similar such example question generations can be seen in Fig. 5.

The failure of our model pronounced when either OCR-token is misinterpreted (for example, the word “derp” on the computer screen is misinterpreted as computer monitor brand, and the word “TELEPHONE” on a telephone booth is misinterpreted as company name) or there is a need of generating questions whose answer is not an OCR token, for example: “Are there any numbers on the identification label?”.

Ablation Study: We perform two ablation studies to demonstrate: (i) utility of the positional information and (ii) model’s capability to generate those questions which have multi-word answers.

Model without positional information considers image features ϕ_I + token features ϕ_o , and with positional information considers image features ϕ_I + token features ϕ_o + positional information ϕ_p to generate questions. We observe that model with positional information and OCR-consistency i.e., our full model (OLRA) as shown in Table 2 enhances the BLEU score and quality of generated question over other models on both ST-VQA and TextVQA datasets.



Further, in Table 2, we also show OLRA’s performance on generating those questions that have 1-word, 2-word, and 3-word answers respectively. Based on statistical analysis with respect to ST-VQA test set, there are, 69% 1-word, 19.6% 2-word, 7% 3-word, and 4.4% above 3-word question-answer pairs. While in TextVQA test set, there are, 63.7% 1-word, 21.3% 2-word, 8% 3-word, and 7% above 3-word question-answer pairs. We observe that BLEU scores are nearly the same and the model performs equally well in all the three cases (see Table 2). This ablation study indicates that our model is capable of generating even those questions which have two or three-word length as the answer.

5 Conclusions

We introduced the novel task of ‘Text-based Visual Question Generation’, where given an image containing text, the system is tasked with asking an appropriate question with respect to the OCR token. We proposed OLRA – an OCR-consistent visual question generation model to ask meaningful and relevant visual questions. OLRA outperformed three baseline approaches on two public benchmark datasets. As the first work towards developing a visual question generation model that can read, we restrict our scope to generating simple questions whose answer is the OCR token itself. Generating complex questions that require deeper semantic and commonsense reasoning, and improving text-based VQA by augmenting its training data using automatically generated questions are few tasks that we leave as future works.

We firmly believe that the captivating novel task and the benchmarks presented in this work will encourage researchers to develop better TextVQG models, and thereby gravitate ongoing research efforts of the document image analysis community towards conversational AI.

References

1. ICDAR 2019 Robust Reading Challenge on Scene Text Visual Question Answering. <https://rrc.cvc.uab.es/?ch=11>, accessed: 2021-02-01
2. Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C.L., Parikh, D.: VQA: visual question answering. In: ICCV (2015)
3. Baek, Y., Lee, B., Han, D., Yun, S., Lee, H.: Character region awareness for text detection. In: CVPR (2019)
4. Biten, A.F., Tito, R., Mafla, A., i Bigorda, L.G., Rusiñol, M., Jawahar, C.V., Valveny, E., Karatzas, D.: Scene text visual question answering. In: ICCV (2019)
5. Bojanowski, P., Grave, E., Joulin, A., Mikolov, T.: Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics* **5**, 135–146 (2017)
6. Borisjuk, F., Gordo, A., Sivakumar, V.: Rosetta: Large scale system for text detection and recognition in images. In: KDD (2018)
7. Celikyilmaz, A., Clark, E., Gao, J.: Evaluation of text generation: A survey. *CoRR abs/2006.14799* (2020)

8. Fan, Z., Wei, Z., Li, P., Lan, Y., Huang, X.: A question type driven framework to diversify visual question generation. In: IJCAI (2018)
9. Gao, D., Li, K., Wang, R., Shan, S., Chen, X.: Multi-modal graph neural network for joint reasoning on vision and scene text. In: CVPR (2020)
10. Gülçehre, Ç., Dutil, F., Trischler, A., Bengio, Y.: Plan, attend, generate: Planning for sequence-to-sequence models. In: NIPS (2017)
11. Hu, R., Singh, A., Darrell, T., Rohrbach, M.: Iterative answer prediction with pointer-augmented multimodal transformers for textvqa. In: CVPR (2020)
12. Jaderberg, M., Simonyan, K., Vedaldi, A., Zisserman, A.: Reading text in the wild with convolutional neural networks. *International Journal of Computer Vision* **116**(1), 1–20 (2016)
13. Jain, U., Zhang, Z., Schwing, A.G.: Creativity: Generating diverse questions using variational autoencoders. In: CVPR (2017)
14. Johnson, J., Hariharan, B., van der Maaten, L., Fei-Fei, L., Zitnick, C.L., Girshick, R.B.: CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning. In: CVPR (2017)
15. Karaoglu, S., Tao, R., Gevers, T., Smeulders, A.W.M.: Words matter: Scene text for image classification and retrieval. *IEEE Transaction on Multimedia* **19**(5), 1063–1076 (2017)
16. Krasin, I., Duerig, T., Alldrin, N., Ferrari, V., Abu-El-Haija, S., Kuznetsova, A., Rom, H., Uijlings, J., Popov, S., Veit, A., Belongie, S., Gomes, V., Gupta, A., Sun, C., Chechik, G., Cai, D., Feng, Z., Narayanan, D., Murphy, K.: Openimages: A public dataset for large-scale multi-label and multi-class image classification. Dataset available from <https://github.com/openimages> (2017)
17. Krishna, R., Bernstein, M., Fei-Fei, L.: Information maximizing visual question generation. In: CVPR (2019)
18. Lavie, A., Agarwal, A.: METEOR: an automatic metric for MT evaluation with high levels of correlation with human judgments. In: WMT@ACL (2007)
19. Li, Y., Duan, N., Zhou, B., Chu, X., Ouyang, W., Wang, X., Zhou, M.: Visual question generation as dual task of visual question answering. In: CVPR (2018)
20. Liu, F., Xu, G., Wu, Q., Du, Q., Jia, W., Tan, M.: Cascade reasoning network for text-based visual question answering. In: ACM Multimedia (2020)
21. Long, S., He, X., Yao, C.: Scene text detection and recognition: The deep learning era. *International Journal of Computer Vision* **129**(1), 161–184 (2021)
22. Lopez, L.E., Cruz, D.K., Cruz, J.C.B., Cheng, C.: Transformer-based end-to-end question generation. *CoRR abs/2005.01107* (2020)
23. Luo, R., Price, B.L., Cohen, S., Shakhnarovich, G.: Discriminability objective for training descriptive captions. In: CVPR (2018)
24. Mafla, A., de Rezende, R.S., Gómez, L., Larlus, D., Karatzas, D.: Stacmr: Scene-text aware cross-modal retrieval. *CoRR abs/2012.04329* (2020)
25. Mathew, M., Karatzas, D., Manmatha, R., Jawahar, C.V.: DocVQA: A dataset for VQA on document images. In: WACV (2021)
26. Mikolov, T., Yih, W., Zweig, G.: Linguistic regularities in continuous space word representations. In: HLT-NAACL (2013)
27. Mishra, A., Alahari, K., Jawahar, C.V.: Scene text recognition using higher order language priors. In: BMVC (2012)
28. Mishra, A., Shekhar, S., Singh, A.K., Chakraborty, A.: OCR-VQA: visual question answering by reading text in images. In: ICDAR (2019)
29. Misra, I., Girshick, R.B., Fergus, R., Hebert, M., Gupta, A., van der Maaten, L.: Learning by asking questions. In: CVPR (2018)

30. Mostafazadeh, N., Misra, I., Devlin, J., Mitchell, M., He, X., Vanderwende, L.: Generating natural questions about an image. In: ACL (2016)
31. Neumann, L., Matas, J.: Real-time lexicon-free scene text localization and recognition. *IEEE Transaction on Pattern Analysis and Machine Intelligence* **38**(9), 1872–1885 (2016)
32. Patro, B.N., Kurmi, V.K., Kumar, S., Namboodiri, V.P.: Deep bayesian network for visual question generation. In: WACV (2020)
33. Pennington, J., Socher, R., Manning, C.D.: Glove: Global vectors for word representation. In: EMNLP (2014)
34. Ren, S., He, K., Girshick, R.B., Sun, J.: Faster R-CNN: towards real-time object detection with region proposal networks. *IEEE Transaction on Pattern Analysis and Machine Intelligence* **39**(6), 1137–1149 (2017)
35. Serban, I.V., García-Durán, A., Gülçehre, Ç., Ahn, S., Chandar, S., Courville, A.C., Bengio, Y.: Generating factoid questions with recurrent neural networks: The 30m factoid question-answer corpus. In: ACL (2016)
36. Shi, B., Bai, X., Yao, C.: An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition. *IEEE Transaction on Pattern Analysis and Machine Intelligence* **39**(11), 2298–2304 (2017)
37. Sidorov, O., Hu, R., Rohrbach, M., Singh, A.: Textcaps: A dataset for image captioning with reading comprehension. In: ECCV (2020)
38. Singh, A.K., Mishra, A., Shekhar, S., Chakraborty, A.: From strings to things: Knowledge-enabled VQA model that can read and reason. In: ICCV (2019)
39. Singh, A., Natarajan, V., Shah, M., Jiang, Y., Chen, X., Batra, D., Parikh, D., Rohrbach, M.: Towards VQA models that can read. In: CVPR (2019)
40. Sutskever, I., Vinyals, O., Le, Q.V.: Sequence to sequence learning with neural networks. In: NIPS (2014)
41. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: CVPR (2016)
42. Wang, K., Babenko, B., Belongie, S.J.: End-to-end scene text recognition. In: ICCV (2011)
43. Yang, J., Lu, J., Lee, S., Batra, D., Parikh, D.: Visual curiosity: Learning to ask questions to learn visual recognition. In: CoRL. Proceedings of Machine Learning Research (2018)
44. Zhang, S., Qu, L., You, S., Yang, Z., Zhang, J.: Automatic generation of grounded visual questions. In: IJCAI (2017)