

ICDAR 2021 Competition on Document Visual Question Answering

Rubèn Tito^{*,1}, Minesh Mathew^{*,2}, C.V. Jawahar², Ernest Valveny¹, and Dimosthenis Karatzas¹

¹ Computer Vision Center, UAB, Spain
{rperez, ernest, dimos}@cvc.uab.cat

² CVIT, IIIT Hyderabad, India

Abstract. In this report we present results of the ICDAR 2021 edition of the Document Visual Question Challenges. This edition complements the previous tasks on Single Document VQA and Document Collection VQA with a newly introduced on Infographics VQA. Infographics VQA is based on a new dataset of more than 5,000 infographics images and 30,000 question-answer pairs. The winner methods have scored 0.6120 ANLS in Infographics VQA task, 0.7743 ANLSL in Document Collection VQA task and 0.8705 ANLS in Single Document VQA. We present a summary of the datasets used for each task, description of each of the submitted methods and the results and analysis of their performance. A summary of the progress made on Single Document VQA since the first edition of the DocVQA 2020 challenge is also presented.

Keywords: Infographics · Document Understanding · Visual Question Answering.

1 Introduction

Visual Question Answering (VQA) seeks to answer natural language questions asked on images. The initial works on VQA focused primarily on images in the wild or natural images [1,10]. Most models developed to perform VQA on natural images, make use of (i) deep features from whole images or objects (regions of interest within the image), (ii) a Question Embedding module which make use of word embeddings and Recurrent Neural Networks (RNN) or Transformers [29], and (iii) a fusion module which fuses the image and text modalities using attention [2,27].

Images in the wild often contain text and reading this text is critical to semantic understanding of natural images. For example, Veit et al. observe that nearly 50% of images in MS-COCO dataset have text present in it [30]. TextVQA [26] and ST-VQA [4] datasets evolved out of this actuality and both the datasets feature questions where reading text present on the images is important to arrive at the answer. Another track of VQA which require reading text on the images

* Equal contribution.

is VQA on charts. DVQA [13], FigureQA [14] and LEAF-QA [5] comprise a few types of standard charts such as bar charts or line plots. Images in these datasets are rendered using chart plotting libraries using either dummy data or real data and the datasets contain millions of questions created using question templates.

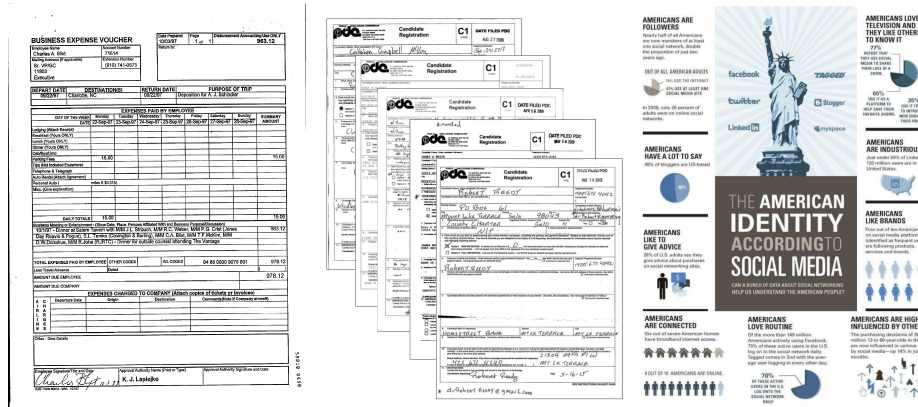
On the other hand, Reading Comprehension [25,12,7] and Open-domain Question Answering [22,16] on machine readable text is a popular research topic in Natural Language Processing (NLP) and Information Retrieval (IR). Unlike VQA where the context on which the question is asked is an image, here questions are asked on a given text passage or a text corpus, say the entire Wikipedia. Recent developments in large-scale pretraining of language models using huge corpora of unlabeled text and later transferring pretrained representations to downstream tasks enabled considerable progress in this space [6,35].

Similar to QA involving either images or text, there are multimodal QA tasks where questions are asked on a context of both images and text. Textbook QA [15] is a QA task where questions are based on a context of text and diagrams and images. RecipeQA [34] questions require models to read text associated with a recipe and understand related culinary pictures. Both datasets have multiple choice style questions and the text modality is presented as machine readable text passages, unlike unstructured text in VQA involving scene text or chart VQA, which need to be first recognized from the images using an OCR.

Document analysis and recognition aims to automatically extract information from document images. DAR research tends to be bottom up and focus on generic information extraction tasks (character recognition, layout analysis, table extraction, etc), disconnected from the final purpose the extracted information is used for. The DocVQA series of challenges aims to bridge this gap, by introducing Document Visual Question Answering (DocVQA) as a high-level semantic task dynamically driving DAR algorithms to conditionally interpret document images. Automatic document understanding is a complicated endeavour that implies much more than just reading the text. In designing a document the authors order information in specific layouts so that certain aspects “stand out”, organise numerical information into tables or diagrams, request information by designing forms, and validate information by adding signatures. All these aspects and more would need to be properly understood in order to devise a generic document VQA system.

The DocVQA challenge comprises three tasks. The Single Document VQA and Document Collection VQA tasks were introduced in the 2020 edition of the challenge as part of the ‘Text and Documents in the Deep Learning Era’ workshop in CVPR 2020 [21], and have received continuous interest since then. We discuss here important changes and advancements in these tasks since the 2020 edition. The Infographics VQA task is a newly introduced task for the 2021 edition.

The Single Document VQA task requires answering questions asked on a single-page document image. The dataset for Single Document VQA comprises business documents, which include letters, fax communications reports with tables and charts and others. These documents mostly contain text in the form of



Question: What is the extension number as per the voucher?
Answers: (910) 741-0673

(a) Single Document VQA

Question: In which years did Anna M. Rivers run for the State senator office?
Answers: 2016, 2020
Doc. Evidences: 454, 10901

(b) Document Collection VQA

Question: What percentage of Americans are online?
Answers: 90%, 90

(c) Infographics VQA

Fig. 1: Examples for the three different tasks in DocVQA. Single Document VQA task (a) is a standard VQA on business/industry documents. Instead, in Document Collection VQA task (b), questions are asked on a document collection comprising documents of same template. Finally, Infographics VQA task (c) is similar to Single Document VQA task, but the images in the dataset are infographics.

sentences and passages [20]. We discuss here new submissions for this task since the 2020 edition.

In the Document Collection VQA task, questions are posed over a whole collection of documents which share the same template but different content. More than responding with the right answer, in this scenario it is important to also provide the user with the right evidence in collection to support the response. For the 2021 edition, we have revisited the evaluation protocols to better deal with unordered list answers and to provide more insight into the submitted methods.

In the newly introduced Infographics VQA task, questions are asked on a single image, but unlike Single Document VQA where textual content is dominating the document images, here we focus on infographics, where the structure, graphical and numerical information are also important to convey the message. For this task we introduced a new dataset of infographics and associated questions and answers annotations [19].

More details about the datasets, the future challenges and updates can be found in <https://docvqa.org>.

2 Competition Protocol

The Challenge ran from November 2020 to April 2021. The setup of the Single Document VQA and Document Collection VQA tasks was not modified with respect to the 2020 edition, while for Infographics VQA, which is a completely new task, we released the training and validation sets between November 2020 and January 2021, and the test set on February 11, 2021, giving participants two months to submit results until April 10th. We relied at all times on the scientific integrity of the authors to follow the established rules of the challenge that they had to agree with upon registering at the Robust Reading Competition portal.

The Challenge is hosted at the Robust Reading Competition (RRC) portal³. All submitted results are evaluated automatically, and per-task ranking tables and visualization options to explore results are offered through the portal. The results presented in this report reflect the state of submissions at the closure of the 2021 edition of the challenge, but the challenge will remain open for new, out-of-competition, submissions. The RRC portal should be considered as the archival version of results, where any new results, submitted after this report was compiled will also appear.

3 Infographics VQA

The design of the new task on Infographics VQA was informed by our analysis of results from the 2020 edition of the Single Document VQA task. In particular, we noticed that most state of the art methods for Single Document VQA followed a simple pipeline consisting on applying OCR (that was also provided) followed by purely text-based QA approaches. Although such approaches give good overall results as running text dominates the kind of documents of the Single Document VQA task, a closer look at the different question categories reveals that those methods performed quite worse on questions that rely on interpreting graphical information, handwritten text or layout information. In addition, these methods work because the Single Document VQA questions were designed so that they can be answered in an extractive manner, which means that the answer text appears in the document image and usually can be inferred from the surrounding context.

For the new Infographics VQA task, we focused on a domain where running text is not dominating the document, while we downplayed the amount of questions based on purely textual evidence and defined more questions which require the models to interpret the layout and other visual aspects. We also made sure that a good number of questions require logical reasoning or elementary arithmetical operations to arrive at the final answer.

We have received and evaluated a total of 337 different submissions on this task from 18 different users. Those submissions include different versions from the same methods. In section 3.4 can be found the ranking of the 6 methods that finally participated in the competition.

³ <https://rrc.cvc.uab.es/?ch=17>

3.1 Evaluation Metric

To evaluate and rank the submitted methods on this task we use the standard ANLS [3] metric for reading-based VQA tasks. There exist a particular case in Infographics VQA task where the full answer is actually a set of items for which the order is not important, like a list of ingredients in a recipe. In the dataset description (section 3.2) we name this as 'Multi-Span' answer. In this case we create all possible permutations of the different items and accept it as correct answers.

3.2 The InfographicsVQA Dataset

Infographic images were downloaded from various online sources using the Google and Bing image search engines. We removed duplicate images using a Perceptual Hashing technique implemented in Imagededup library [11]. To gather the question-answer pairs, instead of using online crowdsourcing platforms we used an internal web based annotation tool, hired annotators and worked closely with them through continuous interaction to ensure the annotations quality and balance the amount of questions for each type defined. In cases where there are multiple correct answers due to language variability, we collected more than one answer per question. For example in the case of the example shown in Figure 1c the question has two valid ground-truth answers.

In total 30,035 question-answer pairs were annotated on 5,485 infographics. We split the data into train, validation and test splits in an 80-10-10 ratio, with no overlap of images between different splits. We also collected additional information regarding questions and answers in the validation and test splits which help us better analyze VQA performance. In particular, we collected the type of answer, which indicates whether the answer can be found as (i) a single text-span in the image ('image-span'), (ii) a concatenation of different text-spans in the image ('multi-span'), (iii) a span of the question text ('question span') or (iv) a 'non-span' answer. The evidence type for each answer is also annotated, and there could be multiple evidences for some answers. The different evidence types are Text, Table/list, Figure, Map (a geographical map) and Visual/layout. Additionally, for questions where a counting, arithmetic or sorting operation is required we collect the operation type as well. More details of the annotation process, statistics of the dataset, analysis of images, questions and answers and examples for each type of evidence, answer, and operation are presented in [19].

Along with the InfographicsVQA dataset, we also provided OCR transcriptions of each of the images in the dataset obtained using Amazon Textract OCR.

3.3 Baselines

In order to establish a baseline performance, we used two models built on state of the art methods. The first model is based on a layout aware language modelling approach — LayoutLM [33]. We take a pretrained LayoutLM model and continue pretraining it on the train split of the InfographicsVQA dataset using a Masked

Language Modelling task. Later we finetune this model for an extractive question answering using a span prediction head at the output. We name this baseline as “(Baseline) LayoutLM”.

The second baseline is the state of the art Scene Text VQA method M4C [9]. This method embeds all different modalities – question, image text and image, into a common space, where a stack of transformer layers is applied allowing each entity to attend to inter- and intra- modality features providing them with context. Then, the answer is produced by an iterative decoder with a dynamic pointer network that can provide an answer either from a fixed vocabulary or from the OCR tokens spotted on the image. This baseline is named as “(Baseline) M4C”. In [19] we give more details of these models, various ablations we try out and detailed results and analysis.

We also provide a human performance evaluation carried out by two volunteers “(Baseline) Human”.

3.4 Submitted methods

We received 6 submissions in total. All of them are based on pretrained language representations. However, we can appreciate that the top ranked methods use multi-modal pretrained architectures combining visual and textual information extracted from the image, while the lower ranked methods use representations based only on natural language.

- #1 - Applica.ai TILT [24]: The winning method learns simultaneously layout information, visual features and textual semantics with a multi-modal transformer based method, and rely on an encoder-decoder architecture that can generate values that are not included in the input text.
- #2 - IG-BERT (single model): A visual and language pretrained model on infographic text pairs. The model was initialized from BERT-large and trained on InfographicsVQA training and validation data. The visual features are extracted using a Faster-RCNN trained on Visually29K [18]. Also, they used OCR tokens extracted by the Google Vision API instead of the ones provided in the competition.
- #3 - NAVER CLOVA: This method uses an extractive QA method based on the HyperDQA [21] approach. First they pre-train BROS [8] model on the IIT-CDIP [17] dataset but sharing the parameters between projection matrices during self-attention. Then, they perform additional pretraining on SQuAD [25] and WikitableQA [23] datasets. Finally, they also fine-tune on the DocVQA [21] dataset.
- #4 - Ensemble LM and VLM: An ensemble between two different methods from extractive NLP QA method (Method 1) and scene-text VQA method (Method 2). The first one is based on BERT-large pretrained on SQuAD + DocVQA and fine-tuned on InfographicsVQA. The final output is vote-based by three

trained models with different hyper-parameters. The second model is the SS-Baseline [36] trained on TextVQA [26], ST-VQA [4] and InfographicsVQA. At the end a rule-based post-processing is performed for three types of questions, i.e., selection, inverse percentage, and summation to get the final result by filling the empty results of Model 1 with answers predicted by Model 2.

#5 - bert baseline: BERT-large model pretrained on SQuAD [25]. It also uses a fuzzy search algorithm to better find the start and end index of the answer span.

#6 - BERT (CPDP): BERT-large model pretrained on SQuAD [25].

Method	ANLS
(Baseline) Human	0.9800
Applica.ai TILT	0.6120
IG-BERT (single model)	0.3854
NAVER CLOVA	0.3219
Ensemble LM and VLM	0.2853
(Baseline) LayoutLM	0.2720
bert baseline	0.2078
BERT (CPDP)	0.1678
(Baseline) M4C	0.1470

Table 1: **Infographics VQA task results table.** Top section show methods withing the ICDAR 2021 competition. Bottom section shows baseline methods proposed by the competition organizers.

3.5 Performance analysis

In table 1 the final competition ranking of the Infographics VQA task is shown. In addition, in Figure 2 we show the methods' Average Normalized Levenshtein Similarity (ANLS) score breakdown by type of answer (top), type of evidence (middle) and type of operations and reasoning required to answer the posed questions (bottom). As it can be seen, the winning method Applica.ai TILT outperforms the rest in all the categories. In addition, the plots show a drop in performance of all methods when the text of the answer is a set of different text spans in the image (Multi-Span) or it does not appear in the image (Non-Span). The easiest questions are the ones that can be directly answered with running text from the documents (Figure 2-middle: Text), this should be expected given the extractive QA nature of most of the methods. Finally, the hardest ones are all the questions that require to perform an operation to come up with the answer as the last plot illustrates comparing the Counting, Sorting and Arithmetic against the performance on None.

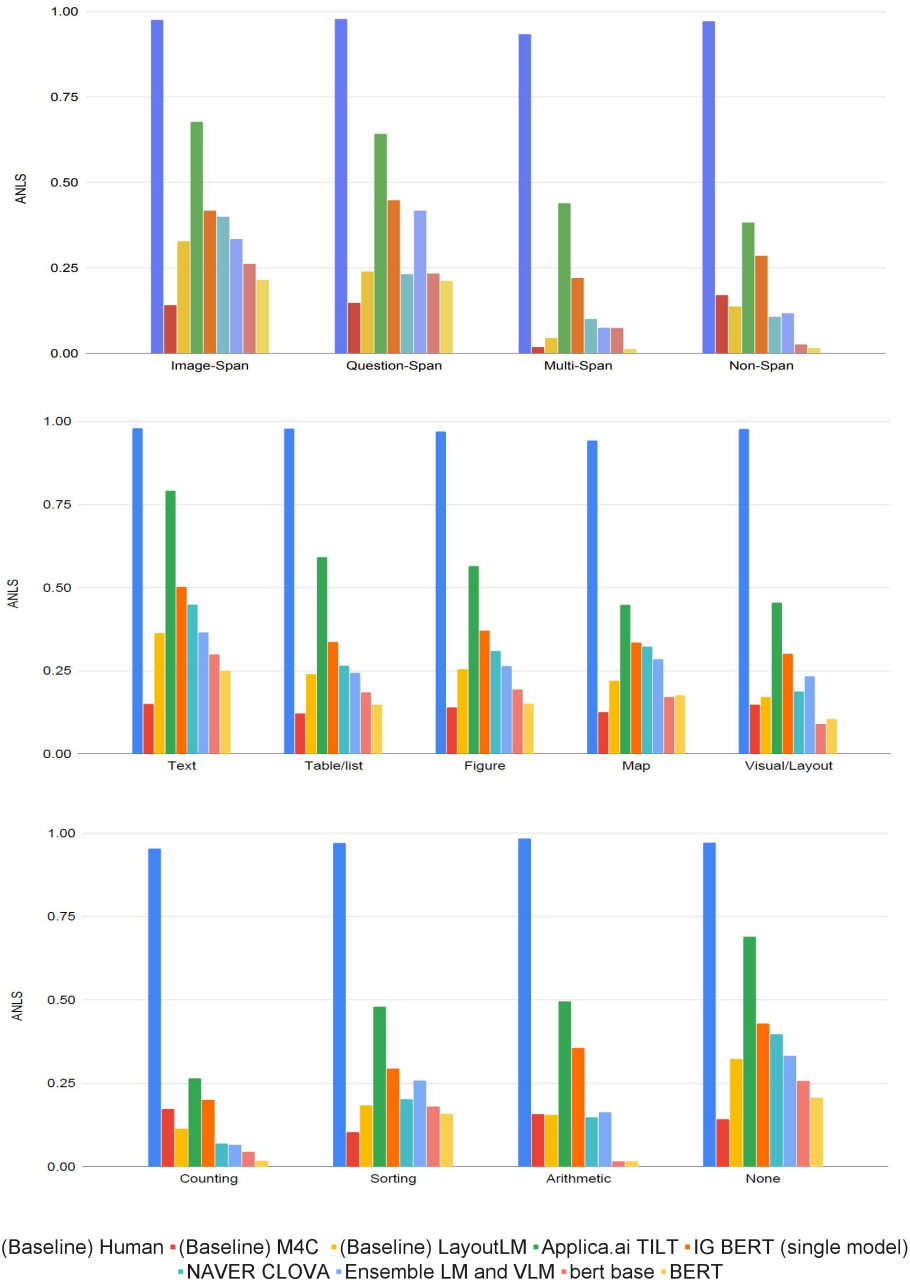


Fig. 2: ANLS score of Infographics VQA task methods breakdown by answer types (top), evidence type (middle) and operations required to answer the posed questions (bottom).

4 Document Collection VQA

In Document Collection VQA task, the questions are posed over a whole collection of 14K document images. The objective in this task is to come up with the answers to the given questions, but also with the positive evidences. We consider positive evidence the document IDs from which the answer can be inferred from.

This task was first introduced in the scope of the CVPR 2020 DocVQA Challenge edition, and we have received and evaluated 60 submissions from 7 different users, out of which 5 have been made public by their authors and feature in the ranking table. 3 of these submissions are new, tackling the problem of answering the questions and not only providing the evidences.

4.1 Evaluation Metric

In this task, the whole answer to a given question is usually a set of sub-answers extracted from different documents, and consequently the order in which those itemized sub-answers are provided is not relevant. Thus, it presents a problem similar to the Multi-Span answers in Infographics VQA. However, the amount of items in this task varies from 1 to > 200 depending on the question, and therefore is not smart to just create all possible permutations. For this reason the metric used to assess the answering performance in this task is the Average Normalized Levenshtein Similarity for Lists (ANLSL) presented in [28], an adaptation of the ANLS [3], to work with a set of answers for which the order is not important. Thus, captures smoothly OCR recognition errors and evaluates reasoning capability at the same time while handles unordered set of answers by performing the Hungarian Matching algorithm between the method’s provided answer and the ground truth.

In addition, even though the retrieval of the positive evidences is not used to rank the methods, we evaluate them according to the Mean Average Precision (MAP) which allows to better analyze the method’s performance.

4.2 Baselines

We defined two baseline methods to establish the base performance for this task. Both baselines work in two steps. First they rank the documents in the collection to find which ones are relevant to answer the given question and then, they extract the answer from those documents marked as relevant. The first method “(Baseline) TS - BERT” ranks the documents by performing text spotting of specific words in the question over the documents in the collection. Then using an extractive QA pretrained BERT model extracts the answers from the top ranked documents. The second method named “(Baseline) Database”, makes use of the commercial OCR Amazon Textract to extract the key-value pairs from the form’s fields for all the documents in the collection. With this information a database-like structure is built. Then, the questions are manually parsed into Structured Query Language (SQL) from which the relevant documents and answers are retrieved. Note that while the first baseline is generic in nature and

could potentially be applied on different collections, the second one relies on the nature of the documents (forms). Detailed information for both methods can be found in [28].

4.3 Submitted methods

In this task, only one new method has submitted. The winning method Infrd-RADAR (Retrieval of Answers by Document Analysis and Re-ranking) first applies OCR to extract textual information from all document images and the combines results with image information to extract key information. The question is parsed into SQL queries by using the spaCy library to split and categorize the question chunks into predefined categories. The SQL queries are then used to retrieve a set of relevant documents and BERT-Large is used to re-rank them again. Afterwards, the re-ranked document IDs are used to filter the extracted information and finally, based on the parsed questions a particular field is collected and posted as an answer.

4.4 Performance analysis

Table 2 shows the competition result ranking comparing the submitted method in this challenge with the baselines and the methods from the CVPR 2020 challenge version. The Infrd-RADAR method outperforms all baselines in the ranking metric (ANLSL). However, methods from the 2020 edition show better performance on the retrieval of positive evidence. This indicates a potential for improvement since given the collection nature of this dataset, a better performance in finding the relevant documents is expected to lead to better answering performance. In figure 3 we show a breakdown of the evidence and answer scores by query. On one hand there is a set of questions for which most of the methods performs well (Q10, Q15, Q16) while Infrd-RADAR is the only one to come up with the answers for the questions Q11 and Q13 and providing their evidences. However, it performs worse than Database in questions Q8, Q9 and Q18, which is probably a consequence of the performance drop when finding the relevant

Method	Answering ANLSL	Retrieval MAP
[†] Infrd-RADAR	0.7743	74.66%
(Baseline) Database [28]	0.7068	71.06%
(Baseline) TS-BERT [28]	0.4513	72.84%
DQA [21]	0.0000	80.90%
DOCR [21]	0.0000	79.15%

Table 2: **Document Collection VQA task results table.** New submissions after CVPR challenge edition [21] are indicated by [†] icon.

documents as it can be seen in the top plot. Notice also that some questions refer to only one single document (Q13, Q15, Q16, Q19), which facilitates methods to score either 0 or 1.

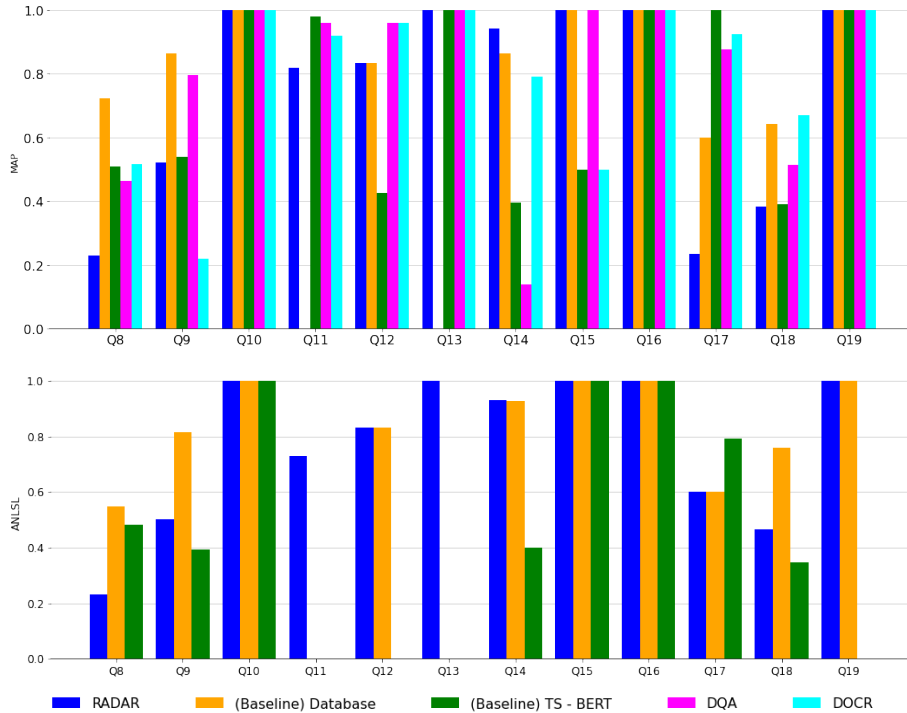


Fig. 3: **Document Collection VQA methods performance by questions.** Top figure shows the retrieval MAP score. Bottom figure shows ANLSL answering score. In bottom figure methods which didn't provide answers are omitted for clarity.

5 Single Document VQA

The Single Document VQA task was first introduced in the scope of the CVPR 2020 DocVQA Challenge edition. Since then, we have received and evaluated > 750 submissions from 44 different users, out of which 28 have been made public by their authors and feature in the ranking table. 21 of these submissions are new, and have pushed overall performance up by $\sim 2\%$. In this task questions were posed over single page document images and the answer to the questions was in most of the cases text that appears within the image. The documents of this dataset were sourced from the industry documents library⁴ and the dataset finally comprised 30,035 question-answer pairs and 5,485 images.

⁴ <https://www.industrydocuments.ucsf.edu/>

Although the documents feature complex layouts and include tables and figures, running text is dominant. Combined with way questions were defined, extractive approaches inspired from NLP perform quite well. To evaluate the submitted methods the standard ANLS [3] metric for reading-based VQA tasks was used.

5.1 Baselines

In order to establish a baseline performance, we used two models built on state of the art methods presented in [20]. The first model is the extractive question answering model from NLP BERT [6] pretrained on SQuAD and finetuned on Single Document VQA dataset. BERT is a language representation pre-trained from unlabeled text passages using transformers. To get the context necessary to extract the answer span we concatenate the detected OCR tokens following left-right, top-down direction. We name this baseline as “(Baseline) BERT Large”. The second baseline is the same state of the art Scene Text VQA method M4C [9] as in Infographics VQA task. However, since the notion of visual objects in real word images is not directly applicable in case of document images, the features of detected objects are omitted. This baseline is named as “(Baseline) M4C”. We also provide a human performance evaluation carried out by few volunteers “(Baseline) Human”.

5.2 Submitted methods

Here we briefly describe the current top 3 ranked methods (new submissions) while the previous edition’s methods description can be found in [21]. The complete ranking can be found in Table 3.

- #1 - Applica.ai TILT [24]: The same team and method as the winners of the Infographics VQA task.
- #2 - LayoutLM 2.0 [32]: Pretrained representation that also learns the text, layout and image in a multi-modal framework using new pretraining tasks where cross-modality interaction are better learned than in previous LayoutLM [33]. It also integrates a spatial-aware mechanism that allows the model to better exploit the relative positional relationship among different text blocks.
- #3 - Alibaba DAMO NLP: An ensemble of 30 models and a pretrained architecture based on StructBERT [31].

6 Conclusions and Future Work

The ICDAR 2021 edition of the DocVQA challenge is the continuation of a long-term effort towards Document Visual Question Answering. One year after

Method	ANLS
(Baseline) Human	0.9811
[†] Applica.ai TILT	0.8705
[†] LayoutLM 2.0 (single model)	0.8672
[†] Alibaba DAMO NLP	0.8506
PingAn-OneConnect-Gammalab-DQA	0.8484
Structural LM-v2	0.7674
[†] QA_Base_MRC_2	0.7415
QA_Base_MRC_1	0.7407
HyperDQA_V4	0.6893
(Baseline) Bert Large	0.6650
Bert fulldata fintuned	0.5900
[†] UGLIFT v0.1 (Clova OCR)	0.4417
(Baseline) M4C	0.3910
Plain BERT QA	0.3524
HDNet	0.3401
CLOVA OCR	0.3296
DocVQAQV_V0.1	0.3016

Table 3: **Single Document VQA task results table.** New submissions after CVPR challenge edition [21] are indicated by [†] icon.

we first introduced it, the DocVQA Challenge has received significant interest by the community. The challenge has evolved, by improving evaluation and analysis methods and by introducing a new complex task on Infographics VQA. In this report, we have summarised these changes, and presented new submissions to the different tasks. Importantly, we note that while the standard approach one year ago was a pipeline combining commercial OCR with NLP-inspired QA methods, state of the art approaches in 2021 are multi-modal. This confirms that the visual aspect of documents is indeed important, even in more mundane layouts such as the documents of Single Document VQA. This is especially true when analysing specific types of questions. The DocVQA challenge will remain open for new submissions in the future. We plan to extend this further with new tasks and insights in the results. We also hope that the VQA paradigm will give rise to more top-down, semantically-driven approaches to document image analysis.

Acknowledgments

This work was supported by an AWS Machine Learning Research Award, the CERCA Programme / Generalitat de Catalunya, and UAB PhD scholarship No B18P0070. We thank especially Dr. R. Manmatha for many useful inputs and discussions.

References

1. Agrawal, A., Lu, J., Antol, S., Mitchell, M., Zitnick, C.L., Batra, D., Parikh, D.: Vqa: Visual question answering (2016)
2. Anderson, P., He, X., Buehler, C., Teney, D., Johnson, M., Gould, S., Zhang, L.: Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering (2017)
3. Biten, A.F., Tito, R., Mafla, A., Gomez, L., Rusinol, M., Mathew, M., Jawahar, C., Valveny, E., Karatzas, D.: Icdar 2019 competition on scene text visual question answering. In: 2019 International Conference on Document Analysis and Recognition (ICDAR). pp. 1563–1570. IEEE (2019)
4. Biten, A.F., Tito, R., Mafla, A., Gomez, L., Rusinol, M., Valveny, E., Jawahar, C., Karatzas, D.: Scene text visual question answering. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 4291–4301 (2019)
5. Chaudhry, R., Shekhar, S., Gupta, U., Maneriker, P., Bansal, P., Joshi, A.: Leaf-qa: Locate, encode attend for figure question answering. In: WACV (2020)
6. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In: ACL (2019)
7. Dua, D., Wang, Y., Dasigi, P., Stanovsky, G., Singh, S., Gardner, M.: DROP: A Reading Comprehension Benchmark Requiring Discrete Reasoning Over Paragraphs. In: NAACL-HLT (2019)
8. Hong, T., Kim, D., Ji, M., Hwang, W., Nam, D., Park, S.: Bros: A pre-trained language model for understanding texts in document (2021) (2021)
9. Hu, R., Singh, A., Darrell, T., Rohrbach, M.: Iterative answer prediction with pointer-augmented multimodal transformers for textvqa. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2020)
10. Hudson, D.A., Manning, C.D.: GQA: a new dataset for compositional question answering over real-world images. CoRR **abs/1902.09506** (2019), <http://arxiv.org/abs/1902.09506>
11. Jain, T., Lennan, C., John, Z., Tran, D.: Imagededup. <https://github.com/idealo/imagededup> (2019)
12. Joshi, M., Choi, E., Weld, D., Zettlemoyer, L.: TriviaQA: A Large Scale Distantly Supervised Challenge Dataset for Reading Comprehension. In: ACL (2017)
13. Kafle, K., Price, B., Cohen, S., Kanan, C.: DVQA: Understanding Data Visualizations via Question Answering. In: CVPR (2018)
14. Kahou, S.E., Michalski, V., Atkinson, A., Kádár, Á., Trischler, A., Bengio, Y.: Figureqa: An annotated figure dataset for visual reasoning. arXiv preprint arXiv:1710.07300 (2017)
15. Kembhavi, A., Seo, M., Schwenk, D., Choi, J., Farhadi, A., Hajjishirzi, H.: Are You Smarter Than A Sixth Grader? Textbook Question Answering for Multimodal Machine Comprehension. In: CVPR (2017)
16. Kwiatkowski, T., Palomaki, J., Redfield, O., Collins, M., Parikh, A., Alberti, C., Epstein, D., Polosukhin, I., Kelcey, M., Devlin, J., Lee, K., Toutanova, K.N., Jones, L., Chang, M.W., Dai, A., Uszkoreit, J., Le, Q., Petrov, S.: Natural questions: a benchmark for question answering research. Transactions of the Association of Computational Linguistics (2019)
17. Lewis, D., Agam, G., Argamon, S., Frieder, O., Grossman, D., Heard, J.: Building a test collection for complex document information processing. In: Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval. pp. 665–666 (2006)

18. Madan, S., Bylinskii, Z., Tancik, M., Recasens, A., Zhong, K., Alsheikh, S., Pfister, H., Oliva, A., Durand, F.: Synthetically trained icon proposals for parsing and summarizing infographics. arXiv preprint arXiv:1807.10441 (2018)
19. Mathew, M., Bagal, V., Tito, R.P., Karatzas, D., Valveny, E., Jawahar, C.: Infographicvqa. arXiv preprint arXiv:2104.12756 (2021)
20. Mathew, M., Karatzas, D., Jawahar, C.V.: DocVQA: A Dataset for VQA on Document Images. In: WACV (2020)
21. Mathew, M., Tito, R., Karatzas, D., Manmatha, R., Jawahar, C.: Document visual question answering challenge 2020. arXiv preprint arXiv:2008.08899 (2020)
22. Nguyen, T., et al.: Ms marco: A human generated machine reading comprehension dataset. CoRR **abs/1611.09268** (2016)
23. Pasupat, P., Liang, P.: Compositional semantic parsing on semi-structured tables. In: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers). pp. 1470–1480 (2015)
24. Powalski, R., Borchmann, L., Jurkiewicz, D., Dwojak, T., Pietruszka, M., Palka, G.: Going full-tilt boogie on document understanding with text-image-layout transformer. arXiv preprint arXiv:2102.09550 (2021)
25. Rajpurkar, P., Zhang, J., Lopyrev, K., Liang, P.: Squad: 100,000+ questions for machine comprehension of text. In: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing. pp. 2383–2392 (2016)
26. Singh, A., Natarajan, V., Shah, M., Jiang, Y., Chen, X., Batra, D., Parikh, D., Rohrbach, M.: Towards vqa models that can read. In: Proceedings of the IEEE/CVF CVPR. pp. 8317–8326 (2019)
27. Teney, D., Anderson, P., He, X., van den Hengel, A.: Tips and tricks for visual question answering: Learnings from the 2017 challenge (2017)
28. Tito, R., Karatzas, D., Valveny, E.: Document collection visual question answering. arXiv preprint arXiv:2104.14336 (2021)
29. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: Proceedings of the 31st International Conference on NeurIPS Information Processing Systems. pp. 6000–6010 (2017)
30. Veit, A., Matera, T., Neumann, L., Matas, J., Belongie, S.: Coco-text: Dataset and benchmark for text detection and recognition in natural images (2016)
31. Wang, W., Bi, B., Yan, M., Wu, C., Bao, Z., Xia, J., Peng, L., Si, L.: Structbert: Incorporating language structures into pre-training for deep language understanding. arXiv preprint arXiv:1908.04577 (2019)
32. Xu, Y., Xu, Y., Lv, T., Cui, L., Wei, F., Wang, G., Lu, Y., Florencio, D., Zhang, C., Che, W., et al.: Layoutlmv2: Multi-modal pre-training for visually-rich document understanding. arXiv preprint arXiv:2012.14740 (2020)
33. Xu, Y., Li, M., Cui, L., Huang, S., Wei, F., Zhou, M.: Layoutlm: Pre-training of text and layout for document image understanding. In: Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. pp. 1192–1200 (2020)
34. Yagcioglu, S., Erdem, A., Erdem, E., Ikizler-Cinbis, N.: RecipeQA: A Challenge Dataset for Multimodal Comprehension of Cooking Recipes. In: EMNLP (2018)
35. Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R.R., Le, Q.V.: XLNet: Generalized Autoregressive Pretraining for Language Understanding. In: NeurIPS (2019)
36. Zhu, Q., Gao, C., Wang, P., Wu, Q.: Simple is not easy: A simple strong baseline for textvqa and textcaps. arXiv preprint arXiv:2012.05153 (2020)