# OCR-IDL: OCR Annotations for Industry Document Library Dataset

Ali Furkan Biten[1], Rubèn Tito[1], Lluis Gomez[1], Ernest Valveny[1], and Dimosthenis Karatzas[1]

[1] Computer Vision Center, UAB
[2] {abiten,rperez,lgomez,ernest,dimos}@cvc.uab.cat

**Abstract.** Pretraining has proven successful in Document Intelligence tasks where deluge of documents are used to pretrain the models only later to be finetuned on downstream tasks. One of the problems of the pretraining approaches is the inconsistent usage of pretraining data with different OCR engines leading to incomparable results between models. In other words, it is not obvious whether the performance gain is coming from diverse usage of amount of data and distinct OCR engines or from the proposed models. To remedy the problem, we make public the OCR annotations for IDL documents using commercial OCR engine given their superior performance over open source OCR models. The contributed dataset (OCR-IDL) has an estimated monetary value over 20K US$. It is our hope that OCR-IDL can be a starting point for future works on Document Intelligence. All of our data and its collection process with the annotations can be found in `https://github.com/furkanbiten/idl_data`.

## 1 Introduction

Analysis of masses of scanned documents is essential in intelligence, law, knowledge management, historical scholarship, and other areas [25]. The documents are often complex and varied in nature that can be digital or scanned born, containing elements such as forms, figures, tables, graphics and photos, while being produced by various printing and handwriting technologies. Some common examples of documents comprise of purchase orders, financial reports, business emails, sales agreements, vendor contracts, letters, invoices, receipts, resumes, and many others [54]. Processing various document types to user's intent is done with manual labor that is time-consuming and expensive, meanwhile requiring manual customization or configuration. In other words, each type of document demands hard-coded changes when there is a slight change in the rules or workflows of documents or even when dealing with multiple formats.

To address these problems, Document Intelligence models and algorithms are created to automatically structure, classify and extract information from documents, improving automated document processing. Particularly, Document Intelligence as a research field aims at creating models for automatically analyzing and understanding documents, reducing the time and the cost associated

with it. From a research perspective, what makes Document Intelligence especially challenging is the requirement of combining various disciplines such as optical character recognition (OCR), document structure analysis, named entity recognition, information retrieval, authorship attribution and many more.

Recent methods on Documents Intelligence utilize deep neural networks combining Computer Vision and Natural Language Processing. Hao *et al.* [17] proposed an end-to-end training using Convolutional Neural Networks to detect tables in documents. Several published works [46,42,56] exploit the advances in object detection [40,19] to further improve the accuracy in document layout analysis. Even though these works have advanced the Document Intelligence field, there are two main limitations to be recognized: (i) they rely on a small human annotated dataset and (ii) they use pre-trained networks that have never seen any documents, hence the interaction between text and layout. Inspired by BERT [11],Xu *et al.* [54] identified these problems and propose a pre-training strategy to unlock the potential of large-scale unlabeled documents. More specifically, they obtain OCR annotations from an open source OCR engine Tesseract [45] for 5 Million documents from IIT-CDIP [25] dataset. With the introduction of pre-training strategy and advances in modern OCR engine [1,12,20,28,34], many contemporary approaches [7,2,53] have utilized even more data to advance the Document Intelligence field.

In this work, we make public the OCR annotations for 26 Millions pages using a commercial OCR engine that has the monetary value over 20K US$. Our motivation for releasing a massive scale documents dataset annotated with a commercial OCR engine is two-fold. First of all, the usage of different amount of documents and different OCR engines across the papers makes it impossible to fairly compare their results and hence their architecture. By creating this dataset, we hope that the works in Document Intelligence will become more comparable and have better intuition on what the proposed architecture can actually accomplish.

Secondly, we decide to use a commercial OCR engine, specifically Amazon Textract[3], over Tesseract. It is because the performance of the OCR engines can significantly affect the model's performance which can be seen in fields that use OCR annotations, such as in fine-grained classification [29,30,31], in scene-text visual question answering [9,44,8,13], in document visual question answering (DocVQA) [50,33]. Apart from improving the annotation quality significantly, we want to level the differences between research groups and companies.

We provide the annotations for publicly available documents from Industry Documents Library (IDL). IDL is a digital archive of documents created by industries which influence public health, hosted by the University of California, San Francisco Library [4]. IDL has already been used in the literature for building datasets: IIT-CDIP [25], RVL-CDIP [18], DocVQA [50,33]. Hence, our OCR annotations can be used to further advance in these tasks.

---

[3] `https://aws.amazon.com/textract/`
[4] `https://www.industrydocuments.ucsf.edu`
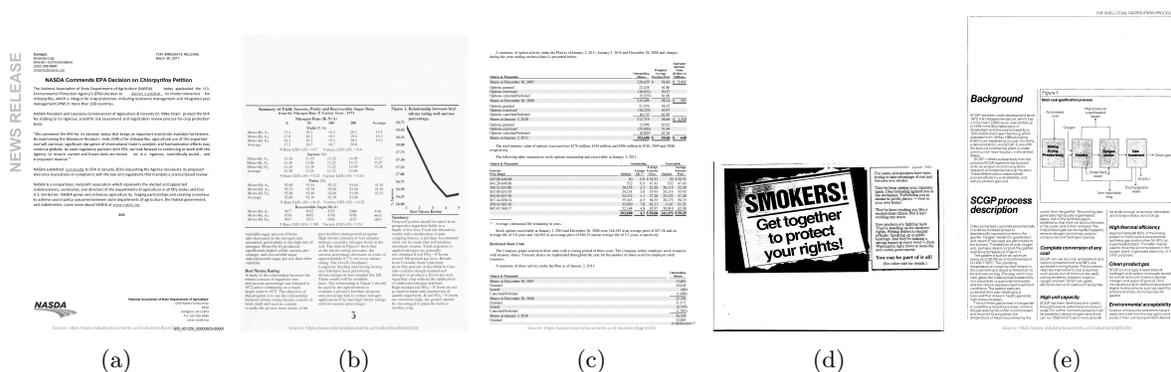
|  (a)  |  (b)  |  (c)  |  (d)  |  (e)  |

Fig. 1: **Document images of OCR-IDL.** The dataset includes a wide variety of documents with dense text (a), tables (b), figures (c), and complex layouts that combines different elements (d, e).

The rest of the paper is structured as follows. First, we briefly explain all the related works. Next, we will elaborate on our data collection and comparison to other datasets. Finally, we will provide various statistics of the annotations and conclude our paper.

## 2   Related Work

Document Intelligence can be considered as an umbrella term covering problems of Key Information Extraction [10,54], Table Detection [41,38] and Structure Recognition [39,55], Document Layout Segmentation [5,4] Document Layout Generation [6,36,3,48], Document Visual Question Answering [51,50,32], Document Image Enhancement [49,22,47] which involves the understanding of visually rich semantic information and structure of different layout entities of a whole page.

Early days of Document Intelligence has relied on rule-based handcrafted approaches mainly divided into bottom-up and top-down methods. Bottom-up methods [15,24,35,43] first detect connected components at the pixel level, later to be fused into higher level of structure through various heuristics and name depending on distinct structural features. While top-down methods [16] dissects a page into smaller units such as titles, text blocks, lines, and words.

Lately, the success of large-scale pre-training [11] in Natural Language Processing has been integrated into Document Intelligence, resulting in impressive performance gains. These methods follow a two step procedure where first they pretrain the models on unlabeled documents (OCR annotations are obtained by an off-the-shelf OCR engine), then they finetune it on specific downstream tasks. LayoutLM [54] is one of the first works that pretrain BERT based language model with document layout information, using masked language/vision modeling and multi label classification. BROS [21] is built on top of Span-BERT [23] with

| Dataset | # of Docs | # of Pages | Docs source | Docs. description | OCR-Text | OCR-BB | Layout | Doc. type |
|---|---|---|---|---|---|---|---|---|
| IIT-CDIP [25] | $6.5M$* | $35.5M$* | UCSF-LTD | Industry documents | Unknown | ✗ | ✗ | ✓ |
| RVL-CDIP [18] | -[†] | $400K$ | UCSF-LTD | Industry documents | ✗ | ✗ | ✗ | ✓ |
| PublayNet [56] | -[†] | $364K$ | PubMedCentral | Journals and articles | ✗ | ✗ | ✓ | ✗ |
| DocBank [26] | -[†] | $500K$ | arXiv | Journals and articles | ✗ | ✗ | ✓ | ✗ |
| DocVQA [51] | $6K$ | $12K$ | UCSF-IDL | Industry documents | Microsoft OCR | ✓ | ✗ | ✓ |
| OCR-IDL | $4.6M$ | $26M$ | UCSF-IDL | Industry documents | Amazon Textract | ✓ | ✗ | ✓ |

Table 1: Summary of other Document Intelligence Datasets. *We skipped $145K$ documents that gave xml parsing errors, didn't contain document ID or number of pages. [†]No traceability between different pages of the same document.

spatially aware graph decoder. For the pretraining loss, they use area-masked language model. Self-Doc [27] utilizes two separate transformer [52] encoders for visual and textual features and later to be fed to multi-modal transformers encoder. TILT [37] tries to encode the layout information by integrating pairwise 1D and 2D information into their models. Uni-Doc [14] is designed to do most document understanding tasks that takes words and visual features from a semantic region of a document image by combining three self-supervised losses. More recent methods, Doc-Former [2] and LayoutLMv2 [53] combine multiple pretraining losses such as image-text alignment, learning to construct image features and multi-modal masked language modeling together to achieve state-of-the-art results.

Yet, comparing all of these works are cumbersome since each work that performs pretraining uses different amounts of data with diverse OCR engines. Hence, this makes it especially hard to understand where the gain is coming from. In other words, we can not draw clear conclusions to questions such as: "What is the effect of the amount of pretraining data on the performance?", "Is the performance gain coming from a better/stronger OCR engine or from the proposed architecture?", "What is the effect of the pretraining loss on the downstream tasks keeping OCR and the amount of data identical?" To help answer these questions, we collect and annotate the largest public OCR annotated documents dataset (OCR-IDL).

## 3   OCR-IDL Dataset

In this section, we elaborate on various details regarding OCR-IDL. Firstly, we explain the process we follow on how we get the IDL data and use Amazon-Textract to obtain OCR annotations. Next, we compare OCR-IDL to other datasets that have proven useful for the document intelligence tasks. And finally, we provide in-depth statistics on the documents we use.
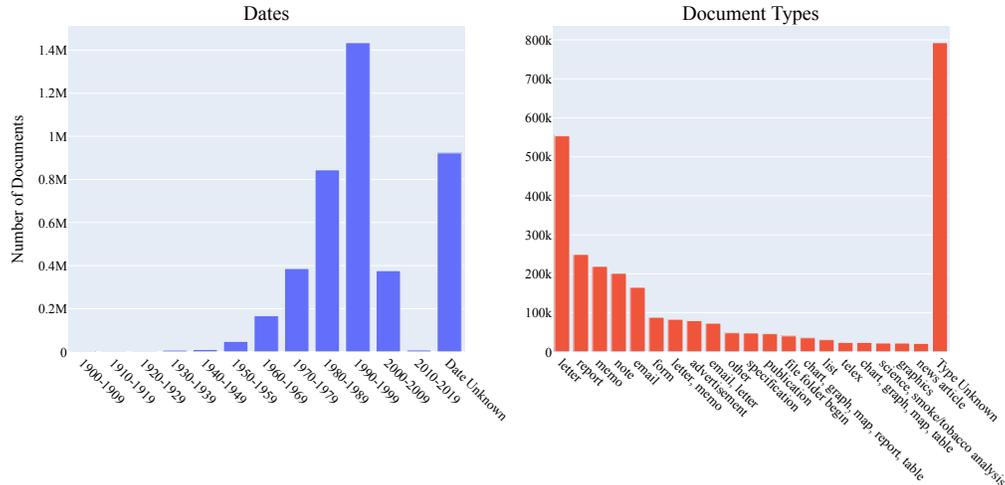
Fig. 2: Distribution of the annotated documents in terms of document types and dates.

### 3.1   Data Collection

As already mentioned, IDL is an industry documents library hosted by UCSF, its main purpose is to "identify, collect, curate, preserve, and make freely accessible internal documents created by industries and their partners which have an impact on public health, for the benefit and use of researchers, clinicians, educators, students, policymakers, media, and the general public at UCSF and internationally"[5]. IDL in total contains over 70 millions documents. We use the publicly available link[6] to downlad the 4.6 Million Documents from IDL dataset which comes in the format of PDFs. Some examples can be viewed in Figure 1. We appreciate that the documents are quite varied in terms of layout where there are tables, figures, ads and combination of them in a single page. Also, we see that the documents are quite text-rich containing many OCR words. Finally, the pages contain smudges, taints and other discoloration what is found in the in a real use-case scenario.

We choose to annotate only 4.6M documents since annotating 13M documents not only would be much costlier (42K dollars instead of 18K) but also in the literature it is shown that using more data have diminishing returns on the downstream task [7]. After obtaining the data, we preprocess the documents to remove empty, faulty and broken pdfs. This process resulted in the elimination of 6548 documents. Moreover, we remove also documents that have more than

---

[5] https://en.wikipedia.org/wiki/Industry_Documents_Library

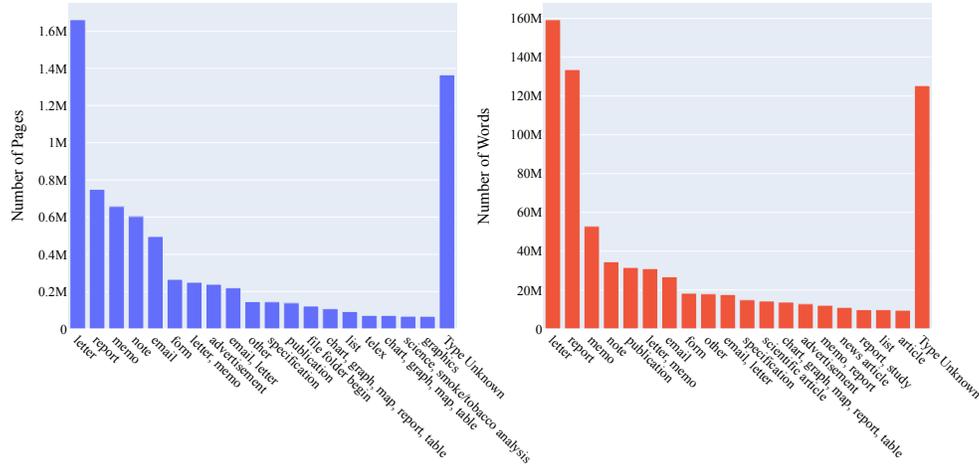[6] https://s3-us-west-2.amazonaws.com/edu.ucsf.industrydocuments.artifacts/

Fig. 3: Distribution of number of pages and number of words per document type.

2900 pages which are 71 in total. The necessity of removing such huge documents was because Amazon-Textract OCR only accepts up to 3000 pages. After pre-processing the documents, we feed all the documents to the OCR engine to obtain the annotations. The annotations provided by the Textract engine include transcription of words and lines and their corresponding bounding boxes and polygons with text type that can be printed or handwritten. Processing all 4.6M documents was done by a single machine with 16 parallelized cores and took about 1 month.

### 3.2   Comparison to existing datasets

In this section, we compare the statistics of the amount of documents and pages to other datasets that are used in Document Intelligence. To name a few, the Illinois Institute of Technology dataset for Complex Document Information Processing (IIT-CDIP) [25] is the biggest document dataset and it is designed for the task of information retrieval. The Ryerson Vision Lab Complex Document Information Processing (RVL-CDIP) [18] dataset used the IIT-CDIP metadata to create a new dataset for document classification. PublayNet [56] and DocBank [26] are datasets designed for layout analysis tasks and DocVQA [33,51] instead, is designed for Visual Question Answering task over document images.

We summarize all the key information for comparison in Table 1. First, we stress that OCR-IDL is the second biggest dataset in amount for pre-training and biggest dataset with annotations obtained from commercial OCR. This provides unique opportunity for the Document Intelligence research for utilizing the unlabeled documents in their research. Furthermore, even though OCR-IDL uses
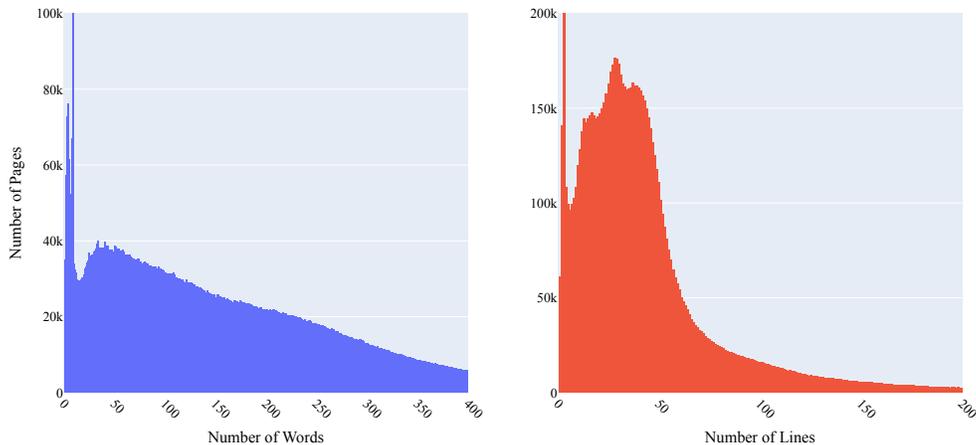
Fig. 4: Distribution of number of words (left) and lines (right) by pages.

the documents from the same source as IIT-CDIP and RVL-CDIP, it also contains other type of industrial documents. OCR-IDL contains documents from chemical, medical and drug industries, hence having more variety in terms of content as well as layout information.

### 3.3  Dataset Statistics

IDL documents come with metadata that is curated by human annotators. They include information about the date, industry, drugs, chemicals, document types and many more. We restrict our analysis on exploring what type of documents we have and the distribution of the dates they are created which can be found in Figure 2. In IDL metadata, there are 35k various document types from which we show only the most common 20 which include letters, report, email, memo, note, etc. As can be seen in the left side of Figure 2, most of the documents' type is unknown. Moreover, even though we have a skewed distribution on letter, report and memo, we also have very distinct distribution such as chart, graphics, news articles. This is especially encouraging because it provides diverse layouts within various contexts. Moreover, the documents are created spanning 100 years where most of the documents are in the range of 1980-2010 as can be seen in the right side of Figure 2. The date the documents are created contributes to the variability of the documents not only in terms of semantics but more importantly having visual artifacts (smudges, different resolution) with different printing technologies.

On top of the amount of documents, we give more details on the number of pages and words for each document type. The number of pages follows the same distribution as the amount of documents per type, as can be appreciated
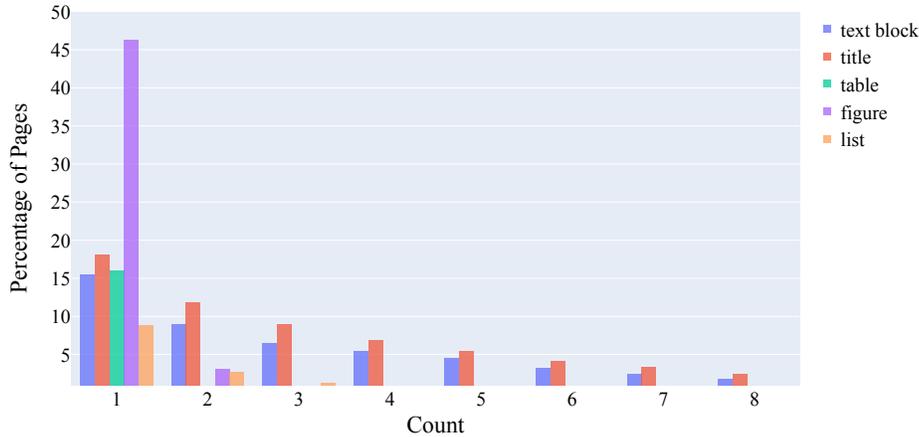
Fig. 5: Distribution of various layout blocks in OCR-IDL. Our documents contain a lot of diversity in terms of title, figure, text block, list and table.

in Figure 3. Also, we can see from Figure 3 that the report type is much richer in terms of text while the rest follows more or less the same distribution.

We turn our attention to OCR annotation statistics. In total, we obtain OCR annotations for 4614232 (4.6M) documents where we have 26621635 (26M) pages, averaging 6 pages per document. Moreover, since documents are known to be a text-rich environment, we provide extra details regarding the amount of words and lines per page and documents. We have 166M words with 46M lines, on average there are 62.5 words and 17.5 lines per page while 360.8 words and 101.25 lines per document. To have a better understanding of the distribution of words and lines per page, we present Figure 4. As shown in the figure, mean distribution is between 20 to 100 words per page while there is a significant amount of pages that contain more than 200 words per page. The distribution for lines follows a different distribution in which it can be observed that most of the pages contain from 10 to 50 lines. In either case, it is clearly observed that documents at hand are ideal for performing pretraining with their diverse layouts and text-rich settings.

Finally, to quantify the diversity of the documents in terms of layout, we run publicly available Faster-RCNN [40] trained on PubLayNet [56] to segment a document into text block, title, figure, list, table. To obtain the segmentation results, we randomly selected 40K pages and some segmentation examples can be found in Figure 6. It can be appreciated from Figure 5 that 40% of the documents have at least 1 figure. We also observe that 10-20% of the pages have at least 1 table and list, showing that documents at hand contain very diverse layout information. Moreover, more than 45% of the pages contain more than 1 text
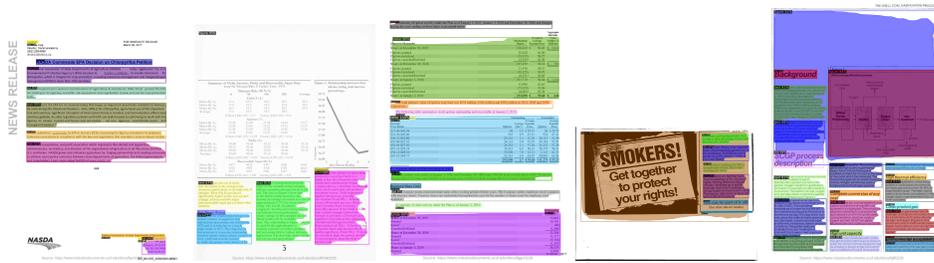
Fig. 6: Qualitative results for segmentation of layout information in OCR-IDL.

block and 1 title, making the documents a text rich environment for pre-training on Document Intelligence tasks.

## 4    Conclusion

In this paper, we have presented our effort to provide OCR annotations for the large-scale IDL document dataset called OCR-IDL. These annotations have a monetary value over $20,000 and are made publicly available with the aim of advancing the Document Intelligence research field. Our motivation is two-fold, first we make use of a commercial OCR engine to obtain high quality annotations, leading to reduce the noise provided by OCR on pretraining and downstream tasks. Secondly, it is our hope that OCR-IDL can be a starting point for future works on Document Intelligence to be more comparable. Throughout this article we have detailed the process that we have followed to obtain the annotations, we have presented a statistical analysis, and compared them with other datasets in the state of the art. The provided analysis shows that our contribution has a high potential to be used successfully in pre-training strategies for document intelligence models. All the code for data collection process and annotations can be accessed in `https://github.com/furkanbiten/idl_data`.

## References

1. Aberdam, A., Litman, R., Tsiper, S., Anschel, O., Slossberg, R., Mazor, S., Manmatha, R., Perona, P.: Sequence-to-sequence contrastive learning for text recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 15302–15312 (2021)
2. Appalaraju, S., Jasani, B., Kota, B.U., Xie, Y., Manmatha, R.: Docformer: End-to-end transformer for document understanding. arXiv preprint arXiv:2106.11539 (2021)
3. Arroyo, D.M., Postels, J., Tombari, F.: Variational transformer networks for layout generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 13642–13652 (2021)
4. Biswas, S., Banerjee, A., Lladós, J., Pal, U.: Docentr: An end-to-end document image enhancement transformer. arXiv preprint arXiv:2201.11438 (2022)

5. Biswas, S., Riba, P., Lladós, J., Pal, U.: Beyond document object detection: instance-level segmentation of complex layouts. International Journal on Document Analysis and Recognition (IJDAR) **24**(3), 269–281 (2021)
6. Biswas, S., Riba, P., Lladós, J., Pal, U.: Docsynth: a layout guided approach for controllable document image synthesis. In: International Conference on Document Analysis and Recognition. pp. 555–568. Springer (2021)
7. Biten, A.F., Litman, R., Xie, Y., Appalaraju, S., Manmatha, R.: Latr: Layout-aware transformer for scene-text vqa. arXiv preprint arXiv:2112.12494 (2021)
8. 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)
9. 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)
10. Carbonell, M., Riba, P., Villegas, M., Fornés, A., Lladós, J.: Named entity recognition and relation extraction with graph neural networks in semi structured documents. In: 2020 25th International Conference on Pattern Recognition (ICPR). pp. 9622–9627. IEEE (2021)
11. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018)
12. Fang, S., Xie, H., Wang, Y., Mao, Z., Zhang, Y.: Read like humans: Autonomous, bidirectional and iterative language modeling for scene text recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7098–7107 (2021)
13. Gómez, L., Biten, A.F., Tito, R., Mafla, A., Rusiñol, M., Valveny, E., Karatzas, D.: Multimodal grid features and cell pointers for scene text visual question answering. Pattern Recognition Letters **150**, 242–249 (2021)
14. Gu, J., Kuen, J., Morariu, V., Zhao, H., Jain, R., Barmpalios, N., Nenkova, A., Sun, T.: Unidoc: Unified pretraining framework for document understanding. Advances in Neural Information Processing Systems **34** (2021)
15. Ha, J., Haralick, R.M., Phillips, I.T.: Document page decomposition by the bounding-box project. In: Proceedings of 3rd International Conference on Document Analysis and Recognition. vol. 2, pp. 1119–1122. IEEE (1995)
16. Ha, J., Haralick, R.M., Phillips, I.T.: Recursive xy cut using bounding boxes of connected components. In: Proceedings of 3rd International Conference on Document Analysis and Recognition. vol. 2, pp. 952–955. IEEE (1995)
17. Hao, L., Gao, L., Yi, X., Tang, Z.: A table detection method for pdf documents based on convolutional neural networks. In: 2016 12th IAPR Workshop on Document Analysis Systems (DAS). pp. 287–292. IEEE (2016)
18. Harley, A.W., Ufkes, A., Derpanis, K.G.: Evaluation of deep convolutional nets for document image classification and retrieval. In: 2015 13th International Conference on Document Analysis and Recognition (ICDAR). pp. 991–995. IEEE (2015)
19. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: Proceedings of the IEEE international conference on computer vision. pp. 2961–2969 (2017)
20. He, Y., Chen, C., Zhang, J., Liu, J., He, F., Wang, C., Du, B.: Visual semantics allow for textual reasoning better in scene text recognition. arXiv preprint arXiv:2112.12916 (2021)
21. Hong, T., Kim, D., Ji, M., Hwang, W., Nam, D., Park, S.: Bros: A pre-trained language model for understanding texts in document (2020)

22. Jemni, S.K., Souibgui, M.A., Kessentini, Y., Fornés, A.: Enhance to read better: A multi-task adversarial network for handwritten document image enhancement. Pattern Recognition **123**, 108370 (2022)
23. Joshi, M., Chen, D., Liu, Y., Weld, D.S., Zettlemoyer, L., Levy, O.: Spanbert: Improving pre-training by representing and predicting spans. Transactions of the Association for Computational Linguistics **8**, 64–77 (2020)
24. Lebourgeois, F., Bublinski, Z., Emptoz, H.: A fast and efficient method for extracting text paragraphs and graphics from unconstrained documents. In: 11th IAPR International Conference on Pattern Recognition. Vol. II. Conference B: Pattern Recognition Methodology and Systems. vol. 1, pp. 272–273. IEEE Computer Society (1992)
25. 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)
26. Li, M., Xu, Y., Cui, L., Huang, S., Wei, F., Li, Z., Zhou, M.: Docbank: A benchmark dataset for document layout analysis. In: Proceedings of the 28th International Conference on Computational Linguistics. pp. 949–960 (2020)
27. Li, P., Gu, J., Kuen, J., Morariu, V.I., Zhao, H., Jain, R., Manjunatha, V., Liu, H.: Selfdoc: Self-supervised document representation learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5652–5660 (2021)
28. Litman, R., Anschel, O., Tsiper, S., Litman, R., Mazor, S., Manmatha, R.: Scatter: selective context attentional scene text recognizer. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 11962–11972 (2020)
29. Mafla, A., Dey, S., Biten, A.F., Gomez, L., Karatzas, D.: Fine-grained image classification and retrieval by combining visual and locally pooled textual features. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 2950–2959 (2020)
30. Mafla, A., Dey, S., Biten, A.F., Gomez, L., Karatzas, D.: Multi-modal reasoning graph for scene-text based fine-grained image classification and retrieval. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 4023–4033 (2021)
31. Mafla, A., Rezende, R.S., Gómez, L., Larlus, D., Karatzas, D.: Stacmr: Scene-text aware cross-modal retrieval. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 2220–2230 (2021)
32. Mathew, M., Bagal, V., Tito, R., Karatzas, D., Valveny, E., Jawahar, C.: Infographicvqa. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 1697–1706 (2022)
33. Mathew, M., Karatzas, D., Jawahar, C.: Docvqa: A dataset for vqa on document images. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 2200–2209 (2021)
34. Na, B., Kim, Y., Park, S.: Multi-modal text recognition networks: Interactive enhancements between visual and semantic features. arXiv preprint arXiv:2111.15263 (2021)
35. O'Gorman, L.: The document spectrum for page layout analysis. IEEE Transactions on pattern analysis and machine intelligence **15**(11), 1162–1173 (1993)
36. Patil, A.G., Ben-Eliezer, O., Perel, O., Averbuch-Elor, H.: Read: Recursive autoencoders for document layout generation. In: Proceedings of the IEEE/CVF

Conference on Computer Vision and Pattern Recognition Workshops. pp. 544–545 (2020)

37. Powalski, R., Borchmann, Ł., Jurkiewicz, D., Dwojak, T., Pietruszka, M., Pałka, G.: Going full-tilt boogie on document understanding with text-image-layout transformer. In: International Conference on Document Analysis and Recognition. pp. 732–747. Springer (2021)

38. Prasad, D., Gadpal, A., Kapadni, K., Visave, M., Sultanpure, K.: Cascadetabnet: An approach for end to end table detection and structure recognition from image-based documents. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops. pp. 572–573 (2020)

39. Raja, S., Mondal, A., Jawahar, C.: Table structure recognition using top-down and bottom-up cues. In: European Conference on Computer Vision. pp. 70–86. Springer (2020)

40. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. Advances in neural information processing systems **28**, 91–99 (2015)

41. Riba, P., Dutta, A., Goldmann, L., Fornés, A., Ramos, O., Lladós, J.: Table detection in invoice documents by graph neural networks. In: 2019 International Conference on Document Analysis and Recognition (ICDAR). pp. 122–127. IEEE (2019)

42. Schreiber, S., Agne, S., Wolf, I., Dengel, A., Ahmed, S.: Deepdesrt: Deep learning for detection and structure recognition of tables in document images. In: 2017 14th IAPR international conference on document analysis and recognition (ICDAR). vol. 1, pp. 1162–1167. IEEE (2017)

43. Simon, A., Pret, J.C., Johnson, A.P.: A fast algorithm for bottom-up document layout analysis. IEEE Transactions on Pattern Analysis and Machine Intelligence **19**(3), 273–277 (1997)

44. 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 Conference on Computer Vision and Pattern Recognition. pp. 8317–8326 (2019)

45. Smith, R.: An overview of the tesseract ocr engine. In: Ninth international conference on document analysis and recognition (ICDAR 2007). vol. 2, pp. 629–633. IEEE (2007)

46. Soto, C., Yoo, S.: Visual detection with context for document layout analysis. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). pp. 3464–3470 (2019)

47. Souibgui, M.A., Biswas, S., Jemni, S.K., Kessentini, Y., Fornés, A., Lladós, J., Pal, U.: Docentr: An end-to-end document image enhancement transformer. arXiv preprint arXiv:2201.10252 (2022)

48. Souibgui, M.A., Biten, A.F., Dey, S., Fornés, A., Kessentini, Y., Gomez, L., Karatzas, D., Lladós, J.: One-shot compositional data generation for low resource handwritten text recognition. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 935–943 (2022)

49. Souibgui, M.A., Kessentini, Y.: De-gan: A conditional generative adversarial network for document enhancement. IEEE Transactions on Pattern Analysis and Machine Intelligence (2020)

50. Tito, R., Karatzas, D., Valveny, E.: Document collection visual question answering. arXiv preprint arXiv:2104.14336 (2021)

51. Tito, R., Mathew, M., Jawahar, C., Valveny, E., Karatzas, D.: Icdar 2021 competition on document visual question answering. In: International Conference on Document Analysis and Recognition. pp. 635–649. Springer (2021)
52. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. Advances in neural information processing systems **30** (2017)
53. 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)
54. 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)
55. Zhong, X., ShafieiBavani, E., Jimeno Yepes, A.: Image-based table recognition: data, model, and evaluation. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXI 16. pp. 564–580. Springer (2020)
56. Zhong, X., Tang, J., Yepes, A.J.: Publaynet: largest dataset ever for document layout analysis. In: 2019 International Conference on Document Analysis and Recognition (ICDAR). pp. 1015–1022. IEEE (2019)