
Measuring the State of Document Understanding

Łukasz Borchmann*

Michał Pietruszka*

Tomasz Stanisławek*

Dawid Jurkiewicz

Michał Turski

Karolina Szyndler

Filip Graliński

Applica.ai
firstname.surname@applica.ai

Abstract

1 Understanding documents with rich-layouts plays a vital role in digitization and
2 hyper-automation but remains a challenging topic in the NLP research community.
3 Additionally, the lack of a commonly accepted benchmark made it difficult to
4 quantify progress in the domain. To empower research in Document Understanding,
5 we present a suite of tasks that fulfill the highest quality, difficulty, and licensing
6 criteria. The benchmark includes Visual Question Answering, Key Information
7 Extraction, and Machine Reading Comprehension tasks over various document
8 domains, and layouts featuring tables, graphs, lists, and infographics. The current
9 study reports systematic baselines making use of recent advances in layout-aware
10 language modeling. To support adoption by other researchers, both the benchmarks
11 and reference implementations will be shortly released.

12 1 Introduction

13 While mainstream Natural Language Processing focuses on plain text documents, content one
14 encounters when reading, e.g., scientific articles, company announcements, or even personal notes, is
15 rarely plain and purely sequential. In particular, the document’s visual and layout aspects that guide
16 our reading process and carry non-textual information appear to be an essential aspect that requires
17 comprehension. **These layout aspects, as we understand them, are prevalent in tasks that can be much
18 better solved when given not only sequence text on the input but pieces of multimodal information
19 covering aspects such as text-positioning (i.e., location of words on the 2D plane), text-formatting (e.g.,
20 different font sizes, colors), and graphical elements (e.g., lines, bars, presence of figure) among others.**

21 Over the decades, systems dealing with document understanding developed an inherent aspect of
22 multi-modality that nowadays revolves around the problems of integrating visual information with
23 spatial relationships and text [34, 1, 49, 11]. Within this frame of reference, Document Understanding
24 may involve the ability to comprehend documents by integrating information from different modalities,
25 e.g., to analyze the figure in the context of accompanying text [2].

26 In general, when document processing systems are considered, the term *understanding* is thought
27 of specifically as the capacity to convert a document into meaningful information [9, 56, 14]. The
28 exact nature of this information depends on the task under consideration and can range from the
29 location of document components to the answer valid for some content-related questions formulated
30 in natural language.

31 Despite its importance for digital transformation, the problem of measuring how well available
32 models obtain information from a wide range of document types and how suitable they are for

*Equal contribution

33 freeing workers from paperwork through process automation is not yet addressed. We intend to
 34 bridge this major gap by introducing the first Document Understanding benchmark (see Section 5
 35 for a review). It includes tasks that either originally had a vital layout understanding component
 36 or were reformulated in such a way that after our modification requires layout understanding. In
 37 particular, there is no structured representation of the underlying text, such as a database-like table
 38 given in advance, and it has to be determined as a part of the end-to-end process from the raw input
 39 file. Every time, there is only a PDF file provided as an input with accompanying textual tokens
 40 and their locations (bounding boxes). It is not enough to process the text in a sequential manner
 41 (token by token), and there is no ground truth reading order given in advance. There are also some
 42 common document understanding problems involved (Section 1.2).

43 **Contribution.** We review and evaluate available data to asses its quality, provide manually annotated
 44 diagnostic sets, measure the human performance, improve data splits, and correct the existing manual
 45 annotations. Importantly, part of the existing datasets is reformulated in a document understanding
 46 paradigm, such as a more competitive problem fitting real-world situations is derived. Additionally,
 47 we propose a novel format for storing data, and provide datasets in an unified form, making their joint
 48 processing and evaluation more accessible. Finally, we provide and open source baselines solving the
 49 task, to facilitate further research on the problem.

50 1.1 Importance and applications

51 As a means of end-to-end process automation, Document Understanding fits into the rapidly growing
 52 market of hyperautomation-enabling technologies, estimated to reach nearly \$600 billion in 2022, up
 53 24% from 2020 [40]. One of the core bottlenecks for this growth is a demand for structured data that
 54 serves technologies such as, e.g., big data analytics and workflow automation. Considering that un-
 55 structured data is orders of magnitude more abundant than structured data, the lack of necessary tools
 56 to extract and analyze it can limit the performance of these intelligent services. The process of struc-
 57 turing data and content must be robust to various document domains. It should also not assume a static
 58 or template layout since the diversity of documents and formats is increasing. A good document un-
 59 derstanding *benchmark* should measure to what extent these technologies are supported in their tasks.

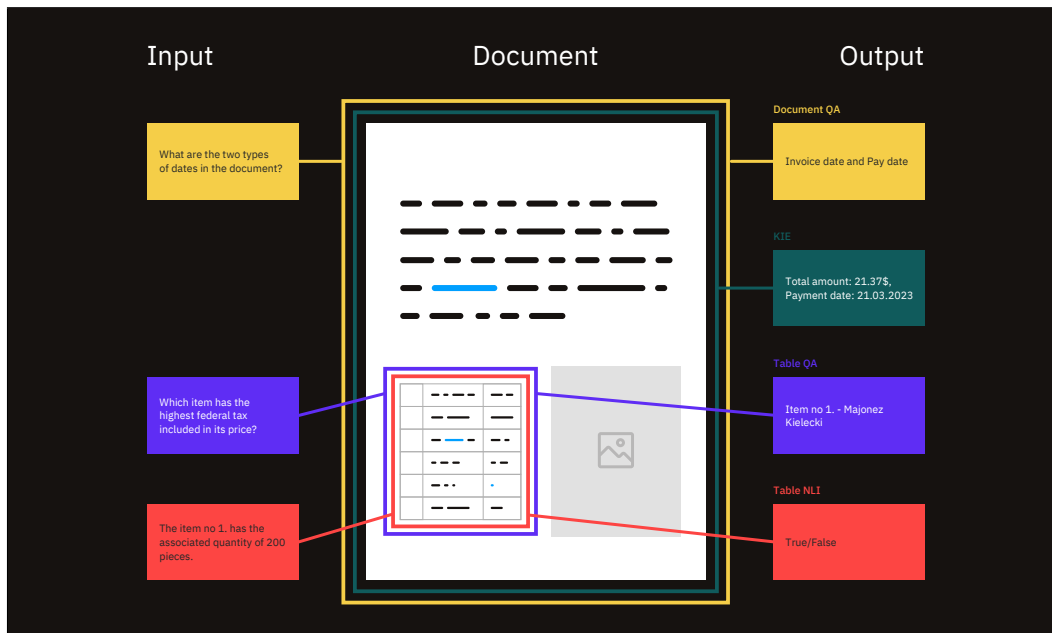


Figure 1: Document Understanding covers problems ranging from the extraction of key information, through verification statements related to rich content, to answering open questions regarding an entire file. It may involve the comprehension of multi-modal information conveyed by a document.

60 1.2 Challenges

61 Owing to its end-to-end nature and heterogeneity, Document Understanding is the touchstone of
62 Machine Learning. The challenges begin to pile up due to the mere form a document is available in,
63 as there is a widespread presence of analog materials such as scanned paper records.

64 **C1.** Consequently, architectures evaluated with different OCR engines are incomparable, e.g., it has
65 been shown that the choice of an OCR engine may impact results more than the choice of model
66 architecture [41]. The overall performance is affected by the noise resulting from OCR errors and
67 incorrectly detected reading order, which impacts commonly used solutions based on sequence
68 labeling. The latter problem is currently being investigated with models independent of sequential
69 order [19, 38].

70 **C2.** Yet another layout-related problem can be exemplified by trying to understand a value from
71 a table cell. In Document Understanding, contrary to the problems of QA over tables, there is no
72 parsed table given in advance as both born-digital and analog documents lack such information. As a
73 result, the application of a particular architecture might require layout analysis and inference of table
74 structure as a necessary step. Nevertheless, it is common to rely on end-to-end models comprehending
75 spatial relationships between the bounding boxes of words instead [63, 55, 38].

76 **C3.** In addition to layout and textual semantics, part of the covered problems demand a Computer
77 Vision component, e.g., to detect a logo, analyze a figure, recognize text style, determine whether the
78 document was signed or the checkbox nearby was selected. Thus, Document Understanding naturally
79 incorporates challenges of both multi-modality and each modality individually.

80 **C4.** Moreover, it is common that token-level annotation is not available, and one receives merely key-
81 value or question-answer pairs assigned to the document. Even in problems of extractive nature, token
82 spans cannot be easily obtained, and consequently, the application of state-of-the-art architectures
83 from other tasks is not straightforward. In particular, authors attempting Document Understanding
84 problems in sequence labeling paradigms were forced to rely on faulty handcrafted heuristics [38].

85 While part of the mentioned challenges are either task- or dataset-specific, they are widespread across
86 Document Understanding problems. A good Document Understanding model should achieve high
87 accuracy and work robustly for documents with the challenges mentioned above.

88 1.3 Desiderata

89 **Gather.** We define our desiderata as follows. We wish to gather both the sparse Document Under-
90 standing datasets published over the years and datasets from related fields that can be reformulated
91 for Document Understanding. Even though a plethora of commercial solutions deal with the problem
92 and it is an object of increasing interest, the availability of public datasets in this field is limited [42].
93 It results from the common practice of publishing works with an evaluation performed on a private,
94 presumably confidential dataset.

95 **Examine.** Then, we intend to examine the value of gathered resources, select the most promising,
96 eliminate their identified disadvantages and provide missing information wherever applicable. Our
97 changes may include improvements of annotation quality and dataset splits, elimination of biases, or
98 preparation of human baselines.

99 **Unify.** To eliminate some of the barriers in future experiments, we wish to propose a format to unify
100 varied Document Understanding tasks and convert all of the datasets included in the benchmark.
101 Additionally, to address challenge C1 (Section 1.2), we provide versioned OCR layers for scanned
102 documents to make models evaluated in the future directly comparable.

103 **Evaluate.** To show there is much space for improvement, we intend to evaluate state-of-the-art
104 models and comment on their result comparisons and human baselines. By precisely diagnosing
105 aspects where these models underperform, we wish to aid the community in identifying where to
106 focus their efforts to conduct valuable research and development.

107 **Open.** Our stance on the future of the Document Understanding Benchmark is to be open and
108 evolving. With further deep learning advances, some tasks may be considered solved, and benchmarks
109 need to hold to that pace. Given the scarcity of datasets available that conform to our Design Process
110 criteria of quality, difficulty, and licensing (see Section 3.1 for detailed analysis), we intend to mimic

111 changes in the publicly available datasets. Specifically, we view an extension of our suite as a
112 continuous process when there is a new dataset complying with defined standards.

113 2 Landscape of Document Understanding tasks

114 For the purposes of the present work, we treat Document Understanding as an umbrella term covering
115 problems of Key Information Extraction, Classification, Question Answering, Layout Analysis, and
116 Machine Reading Comprehension whenever they involve rich documents in contrast to plain-text or
117 image-text pairs (Figure 1).

118 In addition to the problems strictly classified as Document Understanding, several related tasks can
119 be reformulated as such. These provide either text-figure pairs instead of real-world documents or
120 parsed tables given in their structured form. Since both can be rendered as synthetic documents with
121 some loss of information involved, they are worth considering bearing in mind the low availability
122 of proper Document Understanding tasks. **Importantly, such reformulated tasks share an important
123 aspect of C2 challenge we outlined in Section 1.2, as content interpretation is no longer available.**

124 **KIE.** Key Information Extraction, also referred to as Property Extraction, is a task where (properties,
125 document) tuple values are to be provided. Contrary to QA problems, there is no question in natural
126 language but rather a phrase or keyword, such as *total amount*, or *place of birth*. Public datasets in
127 the field include extraction performed on receipts [17, 35], invoices, reports [41], and forms [21].
128 Documents within each of the mentioned tasks are homogeneous, whereas the set of properties to
129 extract is limited and known in advance – in particular, the same type-specific property names appear
130 in both test and train sets. In contrast to Name Entity Recognition, KIE typically does not assume
131 token-level annotations are available, and may require to normalize values found within the document.
132 Moreover, accurate prediction of property values requires some form of layout comprehension.

133 **QA and MRC.** At first glance, Question Answering and Machine Reading Comprehension over
134 Documents is simply the KIE scenario where a question in natural language replaced a property
135 name. More differences become evident when one notices that QA and MRC involve an open set of
136 questions and various document types. Consequently, there is pressure to interpret the question and
137 to possess better generalization abilities. Furthermore, a specific content to analyze demands a much
138 stronger comprehension of visual aspects, as the questions commonly relate to figures and graphics
139 accompanying the formatted text [29, 28, 46].

140 **Classification.** Though document image classification was initially approached using solely the
141 methods of Computer Vision, it has recently become evident that multi-modal models can achieve
142 significantly higher accuracy [54, 55, 38]. Similar conclusions were recently reached in other
143 tasks, e.g., assigning labels to excerpts from biomedical papers depending on the used experiment
144 method [53]. Classification in our context involves rich content, where comprehension of both visual
145 and textual aspects is required since unimodal models underperform.

146 **Layout analysis.** Document Layout Analysis, performed to determine a document’s components,
147 is the oldest Document Understanding problem, initially motivated by the need to optimize storage
148 and the transmission of large information volumes [34]. Even though the motivation behind it has
149 changed over the years, it is rarely an end itself but rather a means to achieve a different goal, such as
150 improving OCR systems. A typical dataset in the field assumes detection and classification of page
151 regions or tokens, depending on the area they belong to [62, 26].

152 **QA over figures.** Question Answering over Figures is, to some extent, comparable with QA and
153 MRC over documents described above. The difference is that a ‘document’ here consists of a single
154 born-digital plot, reflecting information from chosen, desirably real-world data. Because questions
155 are typically templated and figures generated by authors of the task, regular datasets in this category
156 contain millions of examples [31, 4]. Interestingly, questions here can be demanding, e.g., require the
157 estimation of a line chart value at some point.

158 **QA and NLI over tables.** Question Answering and Natural Language Inference over Tables are
159 similar, though in the case of NLI, there is a statement to verify instead of a question to answer.
160 There is never a need to analyze the actual layout, as both assume comprehension of a provided data
161 structure in a way that is equivalent to a database table. Consequently, the methods proposed here
162 are distinct from those used in Document Understanding. There are, however, similarities to exploit,
163 i.e., every task of this type can be reformulated as Document Understanding by simply rendering

164 the table and treating it as an actual document. Apart from the NLI specific, the resulting scenario is
165 similar to QA and MRC over documents described earlier, with the difference that a 'document' now
166 consists of a single born-digital table.

167 3 Benchmark overview

168 Many of the datasets existing in the previously analyzed landscape cannot, on their own, provide
169 enough information that would allow scientists to generalize results to other tasks within the
170 document understanding. Here we describe the proposed suite of tasks and the designed process
171 that led to its creation. **Extensive documentation of the process, including the datasheet, is available**
172 **in Appendices A-H and supplementary materials.**

173 3.1 Desired characteristics

174 As the value and importance of Document Understanding result from its application to process
175 automation, a good benchmark should measure to which degree workers can be supported in their
176 tasks. Though Layout Analysis is oldest of the Document Understanding problems, its output is often
177 not an end in itself but rather a half-measure disconnected from the final information the system is
178 used for. Consequently, we excluded all datasets of this kind by design and restricted ourselves to
179 English tasks of classification, KIE, QA, MRC, and NLI over complex documents, figures, and tables.

180 Candidate tasks resulted from an extensive review of both literature and data science challenges
181 without accompanying publication. Gathered proposals were filtered according to the criteria of
182 quality, difficulty, and licensing.

183 **Quality.** Availability of high-quality annotation was a condition *sine qua non* for a task to qualify.
184 To ensure the highest annotation quality, we excluded resources prepared using a distant annotation
185 procedure, e.g., classification tasks where entire sources were labeled instead of individual instances,
186 or templated question-answer pairs.

187 **Difficulty.** As it makes no sense to measure progress on solved problems, only tasks with a
188 substantial gap between human performance and state-of-the-art models were considered. In the case
189 of promising tasks lacking a human baseline, we provided our estimation.

190 **Licensing.** In publishing our benchmark, we are making efforts to ensure the highest standards for
191 the future of the machine learning community. Only tasks with a permissive license to use annotations
192 and data for further research can be considered.

193 At the same time, we recognized it is essential to approach the benchmark construction holistically, i.e.,
194 to carefully select tasks from diverse domains and types in the rare cases where datasets are abundant.

195 3.2 Selected tasks

196 Table 1 summarizes the selected tasks described in detail below, whereas Appendix G covers the
197 complete list of considered datasets and reasons we omitted them. Lack of the classification and
198 figure QA tasks in this selection results from the fact that none of the available fulfills the assumed
199 selection criteria.

200 The ★ symbol denotes that the dataset was reformulated or modified to improve its quality or
201 align with the Document Understanding paradigm (See Table 2 and Appendix B). We do not
202 distinguish with this mark minor changes, such as data deduplication introduced in multiple datasets
203 (Appendix A).

204 **DocVQA.** Dataset for Question Answering over single-page excerpts from various real-world
205 industry documents. Typical questions present here might require comprehension of images, free
206 text, tables, lists, forms, or their combination [29]. The best-performing solutions so far make
207 use of layout-aware multi-modal models employing either encoder-decoder or sequence labeling
208 architectures [38, 55]. We take the dataset as is without introducing any modification.

209 **InfographicsVQA.** The task of answering questions about visualized data from a diverse collection
210 of infographics, where the information needed to answer a question may be conveyed by text, plots,
211 graphical or layout elements. Currently, the best result is obtained by an encoder-decoder model.

212 **Kleister Charity**. A task for extracting information about charity organizations from their published
 213 reports is considered, as it is characterized by careful manual annotation by linguists and a significant
 214 gap to human performance. It addresses important areas, namely high layout variability (lack of
 215 templates), need for performing an OCR, the appearance of long documents, and multiple spatial
 216 features (e.g., tables, lists, and titles). We take the dataset as is without introducing any modification.

217 **PWC★**. Papers with Code Leaderboards dataset was designed to extract result tuples from machine
 218 learning papers, including information on task, dataset, metric name, score. The best performing ap-
 219 proach involves a multi-step pipeline, with modules trained separately on identified subproblems [24].
 220 In contrast to the original formulation, we provide a complete paper as input instead of the table.
 221 This approach allows us to treat the problem as an end-to-end Key Information Extraction task with
 222 grouped variables (Appendix B).

223 **DeepForm★**. KIE dataset consisting of socially important documents related to election spending.
 224 The task is to extract contract number, advertiser name, amount paid, and air dates from advertising
 225 disclosure forms submitted to the Federal Communications Commission [44]. We use a subset of
 226 distributed datasets and improve annotations errors (Appendix B).

227 **WikiTableQuestions★**. Dataset for QA over semi-structured HTML tables sourced from Wikipedia.
 228 The authors intended to provide complex questions, demanding multi-step reasoning on a series
 229 of entries in the given table, including comparison and arithmetic operations [36]. The problem is
 230 commonly approached assuming a semantic parsing paradigm, with an intermediate state of formal
 231 meaning representation, e.g., inferred query or predicted operand to apply on selected cells [59, 16].
 232 We reformulate the task as document QA by rendering the original HTML and restrict available
 233 information to layout given by visible lines and token positions (Appendix B). It is forbidden for
 234 participating systems to use the HTML source.

235 **TabFact★**. To study fact verification with semi-structured evidence over relatively clean and simple
 236 tables collected from Wikipedia, entailed and refuted statements corresponding to a single row
 237 or cell were prepared by the authors of TabFact [6]. Without being affected by the simplicity of
 238 binary classification, this task poses challenges due to the complex linguistic and symbolic reasoning
 239 required to perform with high accuracy. Analogously to WikiTableQuestion, we render tables and
 240 reformulate the task as document NLI (Appendix B).

241 **Challenges we outlined in Section 1.2 are prevalent in this selection. In particular, scanned documents**
 242 **and infographics have a crucial component of OCR-related problems (C1), no task provides easily**
 243 **accessible information on content interpretation (C2) or token-level annotations (C4), and they may**
 244 **demand comprehension of visual clues to perform well (C3).**

245 3.3 Diagnostic subsets

246 We propose several auxiliary validation subsets, spanning across all the tasks, to improve result
 247 analysis and aid the community in identifying where to focus its efforts. A detailed description of
 248 these categories and related annotation procedures is provided in Appendix E.

249 **Answer characteristic**. We consider four features regarding the answer shallow characteristic. First,
 250 we indicate whether the answer is provided in the text explicitly in exact form (*extractive* data point)

Task	Size (thousands)			Type	Metric	Features		Domain
	Train	Dev	Test			Input	Scanned	
DocVQA	10.2	1.3	1.3	Visual QA	ANLS	}Doc.	+	Business
InfographicsVQA	4.4	.5	.6	Visual QA	ANLS		-	Open
Kleister Charity	1.7	.4	.6	KIE	F1		+/-	Legal
PWC★	.2	.06	.12	KIE*	F1		-	Scientific
DeepForm★	.7	.1	.3	KIE	F1		+/-	Finances
WikiTableQuestions★	1.4	.3	.4	Table QA	Acc.	}Exc.	-	Open
TabFact★	13.2	1.7	1.7	Table NLI	Acc.		-	Open

Table 1: Comparison of selected tasks with their base characteristic, including information regarding whether an input is an entire document (Doc.) or document excerpt (Exc.)

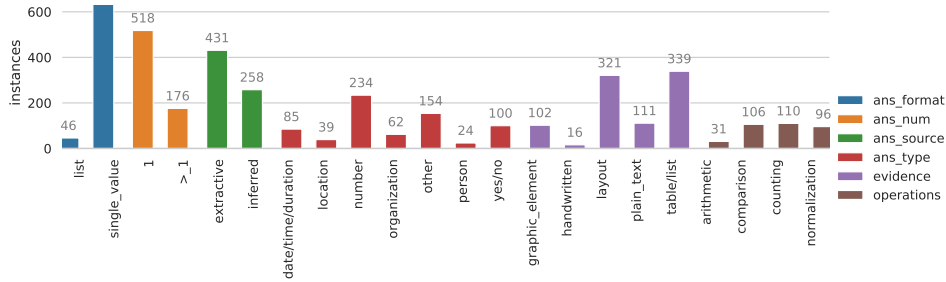


Figure 2: Number of annotated instances in each diagnostic subset category.

251 or has to be inferred from the document content (*abstractive* one). The second category includes,
 252 e.g., all the cases where value requires normalization before being returned (e.g., changing the date
 253 format). Next, we distinguish expected answers depending on whether they contain a *single value* or
 254 *list* of values. Finally, we decided to recognize several popular data types depending on shapes or
 255 class of expected named entity, i.e., to distinguish *date*, *number*, *yes/no*, *organization*, *location*, and
 256 *person* classes.

257 **Evidence form.** As we intend to analyze systems dealing with rich data, it is natural to study
 258 the performance w.r.t. the form that evidence is presented within the analyzed document. We
 259 distinguished *table/list*, *plain text*, *graphic element*, *layout*, and *handwritten* categories.

260 **Required operation.** Finally, we distinguish whether i.e., *arithmetic operation*, *counting*,
 261 *normalization* or some form of *comparison* has to be performed to answer correctly.

262 3.4 Unified format

263 We propose a unified format for storing information in the Document Understanding domain and
 264 deliver converted datasets as part of the released benchmark. It assumes three interconnected dataset,
 265 document annotation and document content levels. Please refer to the repository for examples and
 266 formal specifications of the schemes.

267 **Dataset.** The dataset level is intended for storing the general metadata, e.g., name, version, license,
 268 and source. Here, the JSON-LD format based on the well-known schema.org web standard is used.²

269 **Document.** The documents annotation level is intended to store annotations available for individual
 270 documents within datasets and related metadata (e.g., external identifiers). Our format, valid for all
 271 of the Document Understanding tasks, is specified using the JSON-Schema standard. This ensures
 272 that every record is well-documented and makes automatic validation possible. Additionally, to make
 273 the processing of large datasets efficient, we provide JSON Lines file for each split, thus it is possible
 274 to read one record at a time.

²See <https://json-ld.org/> for information on the JSON-LD standard, and <https://developers.google.com/search/docs/data-types/dataset> for the description of adapted schema.

Dataset	Diagnostic sets	Unified format	Human performance	Manual annotation	Reformulation as DU	Improved split
DocVQA	+	+	-	-	-	-
InfographicsVQA	+	+	-	-	-	-
Kleister Charity	+	+	-	-	-	-
PWC	+	+	+	+	+	-
DeepForm	+	+	+	+	-	+
WikiTableQuestions	+	+	+	-	+	+
TabFact	+	+	-	-	+	-

Table 2: Brief characteristics of our contribution, major fixes and modifications introduced to particular datasets. See Appendix B for a full description.

275 **Content.** As part of the original annotation or additional data we provide is related to document
276 content (e.g., the output of a particular OCR engine **that is of critical importance due to C2**), we
277 introduce the document’s content level. Similarly to the document level, we propose an adequate
278 JSON Schema and provide the JSON Lines files in addition. PDF files with the source document
279 accompany dataset -, document-, and content-level annotations. If the source PDF was not available,
280 a lossless conversion was performed.

281 3.5 Human performance

282 Estimation of human performance for PWC, WikiTableQuestions, DeepForm was performed in-
283 house by professional annotators who are full-time employees of our company after completing
284 the task-specific training (See Appendix D). Each dataset was approached with two annotators; the
285 average of their scores, when validated against the gold standard, is treated as the human performance
286 (See Table 3). **Interestingly, human scores on PWC are relatively low in terms of F1 value – we**
287 **explained this and justified keeping the task in Appendix B.**

288 3.6 Evaluation protocol

289 All the benchmark submissions are expected to conform to the following rules to guarantee fair
290 comparison, reproducibility, and transparency.

- 291 1. All results should be automatically obtainable starting from either raw PDF documents or the
292 JSON files we provide. In particular, it is not permitted to rely on the potentially available source
293 file that our PDFs were generated from or in-house manual annotation.
- 294 2. Despite the fact that we provide an output of various OCR mechanisms wherever applicable, it is
295 allowed to use software from outside the list. In such cases, participants are highly encouraged to
296 donate OCR results to the community, and we declare to host them along with other variants. It is
297 expected to provide detailed information on used software and its version.
- 298 3. Any dataset can be used for unsupervised pretraining. The use of supervised pretraining is limited
299 to datasets where there is no risk of information leakage, e.g., one cannot train models on datasets
300 constructed from Wikipedia tables unless it is guaranteed that the same data does not appear in
301 WikiTableQuestions and TabFact.
- 302 4. It is encouraged to either use datasets already publicly available or to release private data used
303 for pretraining. The minimum requirement for a private dataset is the description of its size and
304 creation sufficient to reproduce the results.
- 305 5. Training performed on a development set is not allowed. We assume participants select the model
306 to submit using training loss or validation score. We do not release test sets and keep them secret
307 by introducing a daily limit of evaluations performed on the benchmark’s website.
- 308 6. Although we allow submissions limited to one category, e.g., QA or KIE, complete evaluations of
309 models that are able to comprehend all of the tasks with one architecture are highly encouraged.
- 310 7. Since different random initialization or data order can result in considerably higher scores, we
311 require the bulk submission of at least three results with different random seeds.
- 312 8. Every submission is required to have an accompanying description. It is recommended to include
313 the link to the source code.

314 **Scoring.** To provide an objective means of comparison with the previously published results, we
315 decided to retain the initially formulated metrics.

316 Regarding the overall score, we consider resorting to an arithmetic mean of different metrics desirable
317 due to its simplicity and straightforward calculation and analysis. Moreover, it is partially justified as
318 all ANLS, F1, and accuracy are interrelated variants of the same measure.

319 To discount for a different number of tasks within each type, we perform a two-step averaging, i.e.,
320 average scores within each category and then average aggregated scores of Document QA, KIE and
321 Table QA groups.³

³Scores on the DocVQA and InfographicsVQA test sets are calculated using the official website.

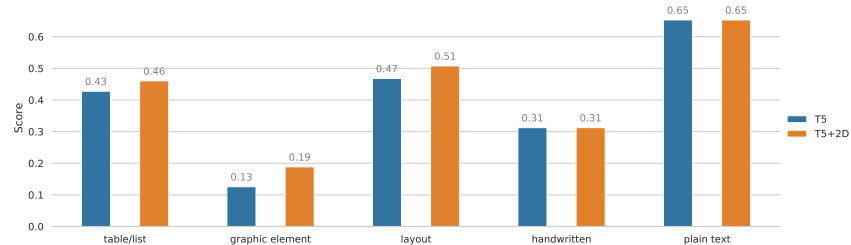


Figure 3: Demonstration of one of the diagnostic subsets we introduced. Here, the impact of 2D bias on cases with particular evidence types can be easily verified.

322 4 Experiments

323 We intended to facilitate future research by providing extensible model with a straightforward training
 324 procedure which can be applied to all of the proposed task in an end-to-end manner. Consequently,
 325 we decided to rely on the extended T5 model to ensure these and to identify the current level of
 326 performance on the chosen tasks [39]. The introduced modifications is the addition of 2D positional
 327 bias [55, 38, 63] that has been shown to perform well on tasks that demand layout understanding.
 328 Both models are released with the benchmark on the MIT license. Following the evaluation protocol,
 329 the training is run three times for each configuration of model size, architecture, and OCR engine.

330 Comparison of the best-performing baselines in relation to human performance and top results
 331 reported in the literature is presented in Table 3. In several cases there is a significant gap between
 332 performance of our baselines and the external best. It can be attributed to several factors. First of all,
 333 WTQ and TabFact were reformulated in a document understanding paradigm. External bests for these
 334 are no longer applicable to the benchmark in its present, more demanding form. Moreover, they were
 335 task-specific, i.e., were explicitly designed for particular task and do not support other datasets within
 336 the benchmark. Secondly, there are differences between the evaluation protocol that we assume and
 337 what the previous authors assumed (e.g., we do not allow training models on the development sets,
 338 we require reporting an average of multiple runs, we disallow pretraining on datasets that might lead
 339 to information leak). Thirdly, to simplify the process and ensure easier reproducibility, we did not
 340 conduct any unsupervised pretraining (contrary to the most state-of-the-art models). Fourthly, there
 341 is no aspect of vision comprehension in our baseline that could possibly address the C3 challenge.
 342 Finally, there is the case of Kleister Charity. An encoder-decoder model we relied on as a one-to-fit-all
 343 baseline cannot process an entire document due to memory limitations. As a result, the score was
 344 lower as we consumed only a part of the document.

345 Irrespective of the task and whether our competitive baselines or external results are considered,
 346 there is still a large gap to humans, which is desired for novel baselines. Moreover, one can notice
 347 that the addition of 2D positional bias to the T5 architecture leads to better scores, which is yet
 348 another result we anticipated as it suggests that considered tasks have an essential component of
 349 layout comprehension. Availability of the diagnostic subsets we introduced allows one to verify this
 350 assumption. Figure 3 compares baselines w.r.t. the evidence form and shows that the difference
 351 can be attributed to a better comprehension of tables, lists, layout, and even graphic elements where
 352 spatial relationships play a pivotal role. We hope the research community will use these and other
 353 diagnostic subsets to investigate particular approaches’ trade-offs, locate current bottlenecks and
 354 answer the question of where should we look for improvement?

355 5 Relation to existing benchmarks and evaluation campaigns

356 The benchmarks in existence consider either well-established NLP tasks in separation (e.g., Question
 357 Answering, Language Modeling, or Natural Language Inference) or a particular aspect of models
 358 applied in the field, such as performance w.r.t. the input sequence length. Consequently, recurring
 359 evaluation campaigns in Document Understanding can be considered to be related works despite
 360 being distinct from a benchmark *per se*.

361 **NLP benchmarks.** The most recognizable NLP benchmarks of GLUE and SuperGLUE cover a
 362 wide range of problems related to language understanding, such as semantic similarity and Natural
 363 Language Inference [52, 51]. In contrast to the present work, they are related to short text excerpts
 364 lacking layout and accompanying graphical materials. Longer documents and documents collections
 365 are covered by decaNLP casting a variety tasks as question answering over a context [30], KILT
 366 assuming comprehension grounded in real-world knowledge [37], or Long Range Arena focused on
 367 the computational efficiency of the models [48]. All of the mentioned consider plain-text documents
 368 without the rich structure that we are aiming at.

369 Recently, there is a growing interest in dynamic benchmarks enabling customizable model comparison
 370 or with tasks changing through time [33, 58, 27]. Our approach is to some extent related as we focus
 371 on customization, e.g., multiple leaderboards are available, and it is up to the participant to decide
 372 whether to evaluate the model on an entire benchmark or particular category. Secondly, we place
 373 attention on the explanation by providing means to analyze the performance concerning document
 374 or problem types (e.g., using the diagnostic sets we provide). Finally, we intend to gather datasets
 375 not included in the present version of the benchmark to facilitate evaluations in an entire field of
 376 Document Understanding, regardless of if they are included in the current version of the leaderboard.

377 **Evaluation campaigns.** So far, efforts of the Document Understanding community were focused on
 378 recurring shared tasks collocated with the major conferences in the field. The most prominent of them
 379 is the Robust Reading Competition collocated with the ICDAR. Though not all of the RRC fit into the
 380 Document Understanding as we define it, this is where tasks involving information extraction from
 381 historical handwritten records and receipts, text extraction from biomedical figures, and document
 382 visual question answering have been proposed [29, 18, 57, 12]. An essential difference between the
 383 recurring events and a benchmark is that tasks from historical competition require the re-assessment of
 384 difficulty in spite of the current state-of-the-art. Additionally, they are often considered in separation,
 385 while the benchmark intends to measure system abilities on various tasks at once.

386 6 Conclusions

387 To efficiently pass information to the reader, writers often assume that structured forms such as tables,
 388 graphs, or infographics are more accessible than sequential text due to human visual perception and
 389 our ability to understand a text’s spatial surroundings. We investigate the problem of correctly mea-
 390 suring the progress of models able to comprehend such complex documents and propose a benchmark
 391 – a suite of tasks that balance factors such as quality of a document, importance of layout information,
 392 type and source of documents, task goal, and the potential usability in modern applications.

393 We aim to track the future progress on them with the website prepared for transparent verification
 394 and analysis of the results. The former is facilitated by the diagnostics subsets we derived to measure
 395 vital features of the Document Understanding systems. Finally, we provide a set of solid baselines,
 396 datasets in the unified format, and released source code to bootstrap the research on the topic.

Dataset / Task type	Score (task-specific metric)			
	T5	T5+2D	External best	Human
DocVQA	72.5	74.1	87.1 [38]	98.1
InfographicsVQA	37.8	43.1	61.2 [38]	98.0
Kleister Charity	57.9	57.7	83.6 [63]	97.5
PWC★	24.2	25.2	—	51.1
DeepForm★	73.4	74.8	—	98.5
WikiTableQuestions★	32.5	33.4	51.8 [60]	76.7
TabFact★	52.2	53.7	83.9 [10]	92.1
Visual QA	55.2	58.6	—	98.1
KIE	51.8	52.6	—	82.4
Table QA/NLI	42.4	43.6	—	84.4
Overall	49.8	51.6	—	88.3

Table 3: Best results of the T5+2D model in relation to human performance and external best.

397 **Checklist**

- 398 1. For all authors...
- 399 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
400 contributions and scope? [Yes] Since this is a benchmark paper, the main claims are
401 described in Section 1.3 as a Desiderata.
- 402 (b) Did you describe the limitations of your work? [Yes] See 2 where we discuss the
403 broader landscape of available tasks and why we consider part of them.
- 404 (c) Did you discuss any potential negative societal impacts of your work? [N/A] Since this
405 is a benchmark paper, we do not see any negative societal impacts.
- 406 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
407 them? [Yes]
- 408 2. If you are including theoretical results...
- 409 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 410 (b) Did you include complete proofs of all theoretical results? [N/A]
- 411 3. If you ran experiments...
- 412 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
413 imental results (either in the supplemental material or as a URL)? [Yes] Please see
414 Supplementary Materials.
- 415 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
416 were chosen)? [Yes] Please see Appendix F and Supplementary Materials
- 417 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
418 ments multiple times)? [N/A] Since we are providing baselines to roughly estimate
419 whether the task is solved or not, or what type of information the documents contain,
420 multiple runs are not necessary.
- 421 (d) Did you include the total amount of compute and the type of resources used (e.g., type
422 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix F.
- 423 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 424 (a) If your work uses existing assets, did you cite the creators? [Yes] See, e.g., Section 2.
- 425 (b) Did you mention the license of the assets? [Yes] Only the datasets with permissive
426 licenses were chosen, see Section 3.1.
- 427 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
428 Yes, we provide models, data and code.
- 429 (d) Did you discuss whether and how consent was obtained from people whose data you’re
430 using/curating? [N/A]
- 431 (e) Did you discuss whether the data you are using/curating contains personally identifiable
432 information or offensive content? [N/A]
- 433 5. If you used crowdsourcing or conducted research with human subjects...
- 434 (a) Did you include the full text of instructions given to participants and screenshots, if
435 applicable? [Yes] We annotated data by ourselves, based on instructions given in the
436 Appendix E.
- 437 (b) Did you describe any potential participant risks, with links to Institutional Review
438 Board (IRB) approvals, if applicable? [N/A]
- 439 (c) Did you include the estimated hourly wage paid to participants and the total amount
440 spent on participant compensation? [N/A]

441 **References**

- 442 [1] T. Bayer, J. Franke, U. Kressel, E. Mandler, M. Oberländer, and J. Schürmann. *Towards*
443 *the Understanding of Printed Documents*, pages 3–35. Springer Berlin Heidelberg, Berlin,
444 Heidelberg, 1992.
- 445 [2] Z. Bylinskii. *Computational perception for multi-modal document understanding*. PhD thesis,
446 Massachusetts Institute of Technology, 2018.
- 447 [3] I. Chalkidis, E. Fergadiotis, P. Malakasiotis, and I. Androutsopoulos. Large-scale multi-label text
448 classification on EU legislation. In *Proceedings of the 57th Annual Meeting of the Association*
449 *for Computational Linguistics*, pages 6314–6322, Florence, Italy, July 2019. Association for
450 Computational Linguistics.
- 451 [4] R. Chaudhry, S. Shekhar, U. Gupta, P. Maneriker, P. Bansal, and A. Joshi. Leaf-qa: Locate,
452 encode attend for figure question answering. In *2020 IEEE Winter Conference on Applications*
453 *of Computer Vision (WACV)*, pages 3501–3510, 2020.
- 454 [5] L. Chen, X. Chen, Z. Zhao, D. Zhang, J. Ji, A. Luo, Y. Xiong, and K. Yu. Websrc: A dataset for
455 web-based structural reading comprehension, 2021.
- 456 [6] W. Chen, H. Wang, J. Chen, Y. Zhang, H. Wang, S. Li, X. Zhou, and W. Y. Wang. Tabfact : A
457 large-scale dataset for table-based fact verification. In *International Conference on Learning*
458 *Representations (ICLR)*, Addis Ababa, Ethiopia, April 2020.
- 459 [7] W. Chen, H. Zha, Z. Chen, W. Xiong, H. Wang, and W. Wang. Hybridqa: A dataset of multi-hop
460 question answering over tabular and textual data, 2021.
- 461 [8] M. Cho, R. K. Amplayo, S. won Hwang, and J. Park. Adversarial tableqa: Attention supervision
462 for question answering on tables, 2018.
- 463 [9] M. Dehghani. Toward document understanding for information retrieval. *SIGIR Forum*,
464 51(3):27–31, Feb. 2018.
- 465 [10] J. Eisenschlos, S. Krichene, and T. Müller. Understanding tables with intermediate pre-training.
466 In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 281–296,
467 Online, Nov. 2020. Association for Computational Linguistics.
- 468 [11] F. Esposito, D. Malerba, G. Semeraro, and S. Ferilli. Knowledge revision for document
469 understanding. In *ISMIS*, 1997.
- 470 [12] A. Fornés, V. Romero, A. Baró, J. I. Toledo, J. A. Sánchez, E. Vidal, and J. Lladós. Icdar2017
471 competition on information extraction in historical handwritten records. In *2017 14th IAPR*
472 *International Conference on Document Analysis and Recognition (ICDAR)*, volume 01, pages
473 1389–1394, 2017.
- 474 [13] T. Gebru, J. Morgenstern, B. Vecchione, J. W. Vaughan, H. M. Wallach, H. D. III, and K. Craw-
475 ford. Datasheets for datasets. *CoRR*, abs/1803.09010, 2018.
- 476 [14] R. M. Haralick. Document image understanding: Geometric and logical layout. In *CVPR*,
477 volume 94, pages 385–390, 1994.
- 478 [15] A. W. Harley, A. Ufkes, and K. G. Derpanis. Evaluation of deep convolutional nets for
479 document image classification and retrieval. In *International Conference on Document Analysis*
480 *and Recognition (ICDAR)*, 2015.
- 481 [16] J. Herzig, P. K. Nowak, T. Müller, F. Piccinno, and J. Eisenschlos. TaPas: Weakly supervised
482 table parsing via pre-training. In *Proceedings of the 58th Annual Meeting of the Association for*
483 *Computational Linguistics*, pages 4320–4333, Online, July 2020. Association for Computational
484 Linguistics.
- 485 [17] Z. Huang, K. Chen, J. He, X. Bai, D. Karatzas, S. Lu, and C. Jawahar. ICDAR2019 competition
486 on scanned receipt OCR and information extraction. In *ICDAR*, 2019.

- 487 [18] Z. Huang, K. Chen, J. He, X. Bai, D. Karatzas, S. Lu, and C. V. Jawahar. Icdar2019 competition
488 on scanned receipt ocr and information extraction. In *2019 International Conference on*
489 *Document Analysis and Recognition (ICDAR)*, pages 1516–1520, 2019.
- 490 [19] W. Hwang, J. Yim, S. Park, S. Yang, and M. Seo. Spatial dependency parsing for semi-structured
491 document information extraction, 2020.
- 492 [20] S. K. Jauhar, P. D. Turney, and E. H. Hovy. Tabmcq: A dataset of general knowledge tables and
493 multiple-choice questions. *CoRR*, abs/1602.03960, 2016.
- 494 [21] G. Jaume, H. K. Ekenel, and J.-P. Thiran. FUNSD: A dataset for form understanding in noisy
495 scanned documents. In *ICDAR-OST*, 2019.
- 496 [22] G. Jaume, H. K. Ekenel, and J.-P. Thiran. Funsd: A dataset for form understanding in noisy
497 scanned documents, 2019.
- 498 [23] K. V. Jobin, A. Mondal, and C. V. Jawahar. In *2019 International Conference on Document*
499 *Analysis and Recognition Workshops (ICDARW)*, volume 1, pages 74–79, 2019.
- 500 [24] M. Kardas, P. Czapla, P. Stenetorp, S. Ruder, S. Riedel, R. Taylor, and R. Stojnic. Axcell:
501 Automatic extraction of results from machine learning papers, 2020.
- 502 [25] A. Kembhavi, M. Seo, D. Schwenk, J. Choi, A. Farhadi, and H. Hajishirzi. Are you smarter than
503 a sixth grader? textbook question answering for multimodal machine comprehension. pages
504 5376–5384, 07 2017.
- 505 [26] M. Li, Y. Xu, L. Cui, S. Huang, F. Wei, Z. Li, and M. Zhou. Docbank: A benchmark dataset for
506 document layout analysis, 2020.
- 507 [27] P. Liu, J. Fu, Y. Xiao, W. Yuan, S. Chang, J. Dai, Y. Liu, Z. Ye, and G. Neubig. Explainboard:
508 An explainable leaderboard for nlp. *arXiv preprint arXiv:2104.06387*, 2021.
- 509 [28] M. Mathew, V. Bagal, R. P. Tito, D. Karatzas, E. Valveny, and C. V. Jawahar. Infographicvqa,
510 2021.
- 511 [29] M. Mathew, D. Karatzas, and C. Jawahar. Docvqa: A dataset for vqa on document images. In
512 *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*,
513 pages 2200–2209, January 2021.
- 514 [30] B. McCann, N. S. Keskar, C. Xiong, and R. Socher. The natural language decathlon: Multitask
515 learning as question answering. *CoRR*, abs/1806.08730, 2018.
- 516 [31] N. Methani, P. Ganguly, M. M. Khapra, and P. Kumar. Plotqa: Reasoning over scientific plots.
517 In *The IEEE Winter Conference on Applications of Computer Vision (WACV)*, March 2020.
- 518 [32] L. Nan, C. Hsieh, Z. Mao, X. V. Lin, N. Verma, R. Zhang, W. Kryściński, N. Schoelkopf,
519 R. Kong, X. Tang, M. Mutuma, B. Rosand, I. Trindade, R. Bandaru, J. Cunningham, C. Xiong,
520 and D. Radev. Fetaqa: Free-form table question answering, 2021.
- 521 [33] Y. Nie, A. Williams, E. Dinan, M. Bansal, J. Weston, and D. Kiela. Adversarial NLI: A new
522 benchmark for natural language understanding. In *Proceedings of the 58th Annual Meeting of the*
523 *Association for Computational Linguistics*, pages 4885–4901, Online, July 2020. Association
524 for Computational Linguistics.
- 525 [34] D. Niyogi and S. N. Srihari. A rule-based system for document understanding. In *Proceedings*
526 *of the Fifth AAAI National Conference on Artificial Intelligence*, pages 789–793, 1986.
- 527 [35] S. Park, S. Shin, B. Lee, J. Lee, J. Surh, M. Seo, and H. Lee. CORD: A consolidated receipt
528 dataset for post-ocr parsing. In *Document Intelligence Workshop at NeurIPS*, 2019.
- 529 [36] P. Pasupat and P. Liang. Compositional semantic parsing on semi-structured tables. *CoRR*,
530 abs/1508.00305, 2015.

- 531 [37] F. Petroni, A. Piktus, A. Fan, P. Lewis, M. Yazdani, N. D. Cao, J. Thorne, Y. Jernite, V. Pla-
532 chouras, T. Rocktäschel, and S. Riedel. Kilt: a benchmark for knowledge intensive language
533 tasks. In *arXiv:2009.02252*, 2020.
- 534 [38] R. Powalski, L. Borchmann, D. Jurkiewicz, T. Dwojak, M. Pietruszka, and G. Palka. Go-
535 ing full-tilt boogie on document understanding with text-image-layout transformer. *CoRR*,
536 abs/2102.09550, 2021.
- 537 [39] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J.
538 Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of*
539 *Machine Learning Research*, 21(140):1–67, 2020.
- 540 [40] M. Rimol. Gartner Forecasts Worldwide Hyperautomation-Enabling Software Market to Reach
541 Nearly \$600 Billion by 2022. [https://www.gartner.com/en/newsroom/press-releas](https://www.gartner.com/en/newsroom/press-releases/2021-04-28-gartner-forecasts-worldwide-hyperautomation-enabling-software-market-to-reach-nearly-600-billion-by-2022)
542 [es/2021-04-28-gartner-forecasts-worldwide-hyperautomation-enabling-sof](https://www.gartner.com/en/newsroom/press-releases/2021-04-28-gartner-forecasts-worldwide-hyperautomation-enabling-software-market-to-reach-nearly-600-billion-by-2022)
543 [tware-market-to-reach-nearly-600-billion-by-2022](https://www.gartner.com/en/newsroom/press-releases/2021-04-28-gartner-forecasts-worldwide-hyperautomation-enabling-software-market-to-reach-nearly-600-billion-by-2022), 2021.
- 544 [41] T. Stanisławek, F. Graliński, A. Wróblewska, D. Lipiński, A. Kaliska, P. Rosalska, B. Topolski,
545 and P. Biecek. Kleister: Key information extraction datasets involving long documents with
546 complex layouts, 2021.
- 547 [42] N. Subramani, A. Matton, M. Greaves, and A. Lam. A survey of deep learning approaches for
548 ocr and document understanding, 2021.
- 549 [43] H. Sun, Z. Kuang, X. Yue, C. Lin, and W. Zhang. Spatial dual-modality graph reasoning for key
550 information extraction, 2021.
- 551 [44] S. Svetlichnaya. DeepForm: Understand structured documents at scale. [https://wandb.ai](https://wandb.ai/stacey/deepform_v1/reports/DeepForm-Understand-Structured-Documents-at-Scale--VmlldzoyODQ3Njg)
552 [/stacey/deepform_v1/reports/DeepForm-Understand-Structured-Documents-a](https://wandb.ai/stacey/deepform_v1/reports/DeepForm-Understand-Structured-Documents-at-Scale--VmlldzoyODQ3Njg)
553 [t-Scale--VmlldzoyODQ3Njg](https://wandb.ai/stacey/deepform_v1/reports/DeepForm-Understand-Structured-Documents-at-Scale--VmlldzoyODQ3Njg), 2020.
- 554 [45] A. Talmor, O. Yoran, A. Catav, D. Lahav, Y. Wang, A. Asai, G. Ilharco, H. Hajishirzi, and
555 J. Berant. Multimodalqa: Complex question answering over text, tables and images. *CoRR*,
556 abs/2104.06039, 2021.
- 557 [46] R. Tanaka, K. Nishida, and S. Yoshida. Visualmrc: Machine reading comprehension on
558 document images. In *AAAI*, 2021.
- 559 [47] R. Tanaka, K. Nishida, and S. Yoshida. Visualmrc: Machine reading comprehension on
560 document images. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(15):13878–
561 13888, May 2021.
- 562 [48] Y. Tay, M. Dehghani, S. Abnar, Y. Shen, D. Bahri, P. Pham, J. Rao, L. Yang, S. Ruder, and
563 D. Metzler. Long range arena: A benchmark for efficient transformers. *ArXiv*, abs/2011.04006,
564 2020.
- 565 [49] S. L. Taylor, D. Dahl, M. Lipshutz, C. Weir, L. M. Norton, R. Nilson, and M. Linebarger.
566 Integrated text and image understanding for document understanding. In *HLT*, 1994.
- 567 [50] H. M. Vu and D. T. Nguyen. Revising FUNSD dataset for key-value detection in document
568 images. *CoRR*, abs/2010.05322, 2020.
- 569 [51] A. Wang, Y. Pruksachatkun, N. Nangia, A. Singh, J. Michael, F. Hill, O. Levy, and S. Bowman.
570 SuperGlue: A stickier benchmark for general-purpose language understanding systems. In
571 H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors,
572 *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- 573 [52] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. Bowman. GLUE: A multi-task
574 benchmark and analysis platform for natural language understanding. In *Proceedings of the*
575 *2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*,
576 pages 353–355, Brussels, Belgium, Nov. 2018. Association for Computational Linguistics.
- 577 [53] T.-L. Wu, S. Singh, S. Paul, G. Burns, and N. Peng. Melinda: A multimodal dataset for
578 biomedical experiment method classification. *ArXiv*, abs/2012.09216, 2020.

- 579 [54] Y. Xu, M. Li, L. Cui, S. Huang, F. Wei, and M. Zhou. Layoutlm: Pre-training of text and layout
580 for document image understanding, 2019.
- 581 [55] Y. Xu, Y. Xu, T. Lv, L. Cui, F. Wei, G. Wang, Y. Lu, D. Florencio, C. Zhang, W. Che, M. Zhang,
582 and L. Zhou. Layoutlmv2: Multi-modal pre-training for visually-rich document understanding,
583 2020.
- 584 [56] S. Yacoub. Automated quality assurance for document understanding systems. *IEEE Software*,
585 20(3):76–82, 2003.
- 586 [57] C. Yang, X.-C. Yin, H. Yu, D. Karatzas, and Y. Cao. Icdar2017 robust reading challenge on
587 text extraction from biomedical literature figures (detext). In *2017 14th IAPR International
588 Conference on Document Analysis and Recognition (ICDAR)*, volume 01, pages 1444–1447,
589 2017.
- 590 [58] Z. Yang, S. Zhang, J. Urbanek, W. Feng, A. H. Miller, A. Szlam, D. Kiela, and J. Weston.
591 Mastering the dungeon: Grounded language learning by mechanical turker descent. In *6th
592 International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada,
593 April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.
- 594 [59] P. Yin, G. Neubig, W. tau Yih, and S. Riedel. TaBERT: Pretraining for joint understanding of
595 textual and tabular data. In *Annual Conference of the Association for Computational Linguistics
596 (ACL)*, July 2020.
- 597 [60] P. Yin, G. Neubig, W.-t. Yih, and S. Riedel. TaBERT: Pretraining for joint understanding of
598 textual and tabular data. In *Proceedings of the 58th Annual Meeting of the Association for
599 Computational Linguistics*, pages 8413–8426, Online, July 2020. Association for Computational
600 Linguistics.
- 601 [61] K. Zaporozhets, J. Deleu, C. Develder, and T. Demeester. Dwie: an entity-centric dataset for
602 multi-task document-level information extraction, 2021.
- 603 [62] X. Zhong, J. Tang, and A. J. Yepes. Publaynet: largest dataset ever for document layout analysis.
604 In *2019 International Conference on Document Analysis and Recognition (ICDAR)*, pages
605 1015–1022. IEEE, Sep. 2019.
- 606 [63] Łukasz Garncarek, R. Powalski, T. Stanisławek, B. Topolski, P. Halama, and F. Graliński.
607 LAMBERT: Layout-aware (language) modeling using bert for information extraction, 2020.