

Form-NLU: Dataset for the Form Language Understanding

Yihao Ding
The University of Sydney
Sydney, NSW, Australia

Siqu Long
The University of Sydney
Sydney, NSW, Australia

Jiabin Huang
The University of Sydney
Sydney, NSW, Australia

Kaixuan Ren
The University of Sydney
Sydney, NSW, Australia

Xingxiang Luo
The University of Sydney
Sydney, NSW, Australia

Hyunsuk Chung
FortifyEdge
Sydney, NSW, Australia

Soyeon Caren Han
The University of Sydney
Sydney, NSW, Australia

ABSTRACT

Compared to general document analysis tasks, form document structure understanding and retrieval are challenging. Form documents are typically made by two types of authors; A form designer, who develops the form structure and keys, and a form user, who fills out form values based on the provided keys. Hence, the form values may not be aligned with the form designer's intention (structure and keys) if a form user gets confused. In this paper, we introduce Form-NLU, the first novel dataset for form structure understanding and its key and value information extraction, interpreting the form designer's intent and the alignment of user-written value on it. It consists of 857 form images, 6k form keys and values, and 4k table keys and values. Our dataset also includes three form types: digital, printed, and handwritten, which cover diverse form appearances and layouts. We propose a robust positional and logical relation-based form key-value information extraction framework. Using this dataset, Form-NLU, we first examine strong object detection models for the form layout understanding, then evaluate the key information extraction task on the dataset, providing fine-grained results for different types of forms and keys. Furthermore, we examine it with the off-the-shelf pdf layout extraction tool and prove its feasibility in real-world cases.

CCS CONCEPTS

• Information systems → Information retrieval.

KEYWORDS

Datasets, Form understanding, Natural language understanding

ACM Reference Format:

Yihao Ding, Siqu Long, Jiabin Huang, Kaixuan Ren, Xingxiang Luo, Hyunsuk Chung, and Soyeon Caren Han. 2023. Form-NLU: Dataset for the Form Language Understanding. In *Proceedings of The 46th International ACM*

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
SIGIR '23, July 23–27, 2023, Taipei, Taiwan

© 2023 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00
<https://doi.org/XXX>

SIGIR Conference on Research and Development in Information Retrieval (SIGIR '23). ACM, New York, NY, USA, 10 pages. <https://doi.org/XXX>

1 INTRODUCTION

The structural information and its value extraction in a document can be a valuable source for Natural Language Processing (NLP) tasks, especially information extraction and retrieval. Recently, NLP communities and industries, like IBM and Microsoft, have proposed a range of techniques to *understand* positional and/or logical structure of documents[14, 17, 29, 32], and *extract* the essential information[4]. Those researches have mainly focused on Visually Rich Document(VRD)-based tasks, such as academic papers[17, 32], receipts[9, 16], and forms [10, 27], and many benchmark problems has been solved, including layout analysing [13, 32], table structure recognition [2, 30, 31], document question answering [15, 21].

Most VRD-based problems have been successfully solved, however, the form-understanding task is relatively challenging. This is mainly because of two reasons: two types of authors in a form and the combination of diverse visual cues. First, unlike general VRDs, the main aim of forms is very clear, *collecting data from the form users*, and this aim is applied from the medical domain to administrative data collection. According to this aim, we can assume that there are two main authors; A form designer and a form user. A form designer focuses on developing a form structure to collect the required information by defining the clear key point. The developed form would be used as a user interface so a form user can supply the form value based on their understanding. Unfortunately, not every form is clear and easy to understand. To collect the diverse required information, several form designers tend to make forms with diverse layouts/structures, which have complex logical and positional relationships between semantic entities. The form user can easily get confused with the designer's intention, and this derives a wrong alignment of key-value pairs. The confusion about the form designer's intention and the uncertainty of the form user would raise the difficulty of form document understanding and information extraction. Secondly, due to the involvement of *form developers-to-users* relationship, the form has a high possibility of having a combination of different natures, such as digital, printed, or handwritten. For example, a designer may provide users with electrical paper forms, while users may submit filled forms via various carriers, such as digital, printed or handwritten versions. It also commonly happens that users provide various types of noise

(low resolution or uneven scanning or bad handwriting) in the submitted forms. This causes huge difficulties in understanding form document structure and extracting the essential key-value pairs. Regarding form understanding, several datasets have been released in recent years, collected from scanned receipts [9, 16], contracts [19], and cross-domain forms [10, 27] (shown in Table 1). However, those datasets produce the form developer’s intention as relatively simple and general, which does not deal with the confusion about the form designer’s intention and the form user’s uncertainty. Moreover, most datasets do not cover the various carriers of document versions and their noises. This would worsen understanding of the form structure and extracting the key information.

In this paper, we introduce a new dataset for form structure understanding and key information extraction. The dataset would enable the interpretation of the form designer’s specific intention and the alignment of user-written value on it. Our dataset also includes three form types: digital, printed, and handwritten, which cover diverse form appearances/layouts and deal with their noises. In addition to this, we propose a new baseline for form structure understanding and key-value information extraction, which applies robust positional and logical relations. To do this, we cover user-caused form diversities to precisely extract key information from forms, which is an emerging industrial demand currently. Our model covers the hierarchical structure of documents, from words to sentences, sentences to a semantic entities like a paragraph, and entities to a document page. Note that exiting transformer-based models mainly focus on token [5, 8, 25, 26, 28] or entity level [12, 20, 29] independently, ignoring the contextual dependency between different level elements. In the evaluation, we first examine strong object detection models for the form layout understanding, then test the proposed model for the key information extraction task. Moreover, we also examine our key information extraction dataset and proposed model with the off-the-shelf pdf layout extraction tool and prove the feasibility in real-world cases.

The main contribution of this research can be summarised as follows: 1) We introduce Form-NLU, a new form structure understanding and key information extraction dataset that covers specific form designers’ intentions and makes an alignment with the user’s values. 2) We propose a new baseline model that handles positional and logical relations and hierarchical structure in form documents. The proposed model has outperformed other SOTA form key information extraction models on Form-NLU. 3) We apply the proposed dataset and the model with the off-the-shelf pdf layout extraction tool and prove the feasibility in real-world form document cases.

2 RELATED WORK

There are general documents understanding benchmarks introduced in [17, 21, 32], but we briefly introduce the form document understanding and information extraction, which involves multi-party interaction resulting in more complicated positional and logical relationships, as shown in Table 1. There are three major points that we would like to discuss and compare with previous benchmarks. First, most benchmarks [10, 24, 27] do not cover the various carriers of document versions and their noises. For example, a form designer may provide the form with digital forms, while the form user may submit the filled form with various carriers, such as digital

Name	Source	Type			Features		Designer Intention
		\mathcal{D}	\mathcal{P}	\mathcal{H}	B.Box	Text	
FUNSD [10]	Noise Form	×	○	○	○	○	General
XFUND [27]	Synthetic Form	○	×	○	○	○	General
EPHOIE [24]	Exam Paper	×	×	○	○	○	General
Charity [19]	Annual Report	○	○	○	×	○	Specific
NDA [19]	Agreements	○	×	×	×	○	Specific
Form-NLU (Ours)	Financial Form	○	○	○	○	○	Specific

Table 1: Summary of Form Understanding Datasets.

(\mathcal{D}), printed (\mathcal{P}) or handwritten (\mathcal{H}) version. This trend commonly affects the quality of form structure understanding and information extraction due to the various types of noise, including resolution or scanning issues. Hence, the successful benchmark should cover real-world cases with various carriers. Secondly, in order to conduct the form structure understanding and information extraction, it is crucial to interpret the positional and logical relationships between form components. Most benchmarks enable understanding positional and logical/semantic relations by using bounding boxes (B.Box) and Textual information. However, Stanisławek et al. [19] releases two form information extraction datasets, NDA and Charity, without providing bounding box coordinates of semantic entities. Finally, the form document has two main authors, form designers and form users. The form designers design the form structure to collect the required information (designer’s intent), and the form users try to understand the designer’s intention but easily get confused. Hence, handling the designer’s intention is crucial to deal with understanding the form structure and extracting key information. Several popular benchmarks [10, 27] produce the form developer’s intention as relatively simple and general, which does not deal with the confusion about the form designer’s intention and the form user’s uncertainty. For example, those datasets cover simple form component types, keys and values. Scanned exam paper [24] datasets have a relatively simple layout with horizontal key-value pair structures, which do not include any dynamic layout components, such as tables, paragraphs, or complex key-value pairs. **Form-NLU** is the first visual-linguistics form language understanding dataset for supporting researchers in interpreting specific designer intentions under noises from user’s input with various types of form carriers.

3 DATASET

3.1 Data Collection

Our **Form-NLU** is a subset of the publicly available financial form data source for Form 604 (notice of change of interests of the substantial holder)¹ collected by SIRCA. It provides the text records of each substantial shareholder notice form submitted to the Australian Stock Exchange (ASX)² from 2003 to 2015. To better comprehend the form designer’s intentions, we included the twelve most essential form fields [1]. Each form field expects a **Value** from the user, which should contain the specific information being asked by the corresponding **Key** that expresses the form designer’s intention. For example, the key of “*company name*” anticipates a string value that gives the name of a company, whereas “*voting percentage*” asks

¹<https://asic.gov.au/regulatory-resources/forms/forms-folder/604-notice-of-change-of-interests-of-substantial-holder/>

²<https://www2.asx.com.au/>

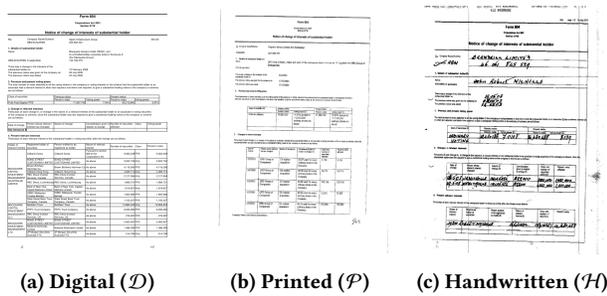


Figure 1: Digital, Printed and Handwritten Form Samples.

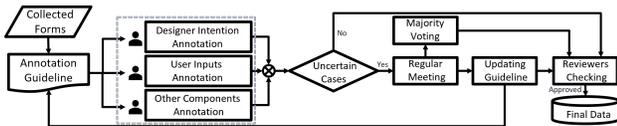


Figure 2: Overall Annotation Workflow.

for a number that tells the portion of the voting. In order to address the aforementioned developer-to-user relationships in real-world scenarios, we also covered three variants of format depending on how the forms were filled and submitted by different users: (1) Digital (\mathcal{D}), (2) Printed (\mathcal{P}) and (3) Handwritten (\mathcal{H}). The Digital forms are filled in digitally and directly submitted as PDF files. The Printed and Handwritten forms are filled in digitally and by hand respectively and then scanned first before being saved as PDF files for submission. A sample of each type of the form is provided in Figure 1, which demonstrates the potential diversity in terms of the form format, e.g., font styles and rotations.

3.2 Annotation Guideline

Based on our task goals (see Section 5), we first developed the annotation schema and guidelines. For Layout Analysis (as Task A³), we created seven distinct form layout component categories (*Title*, *Section*, *Form Key*, *Form Value*, *Table Key*, *Table Value* and *Others*) that encompass all possible components from the collected forms. Each form layout component represents a semantic segment of the text area. The *Title* and *Section* correspond to the title of the form document and the sections, respectively. The *Form/Table Key* and *Form/Table Value* refer to the intended questions of the form designer (e.g. Company Name) and the corresponding user input fields. The difference between the *Key/Value* of form and table is that, by our definition, the *Form Key-Value* pairs are horizontally aligned, whereas the *Table Key-Value* pairs are vertically aligned. Any components that do not fall into these four types are defined as *Others*, such as a customized paragraph and cells added and modified by some companies on their own purposes.

While the Layout Analysis (as Task A) focuses on understanding the overall form structure via those general semantic components, the Key Information Extraction (as Task B⁴) aims to comprehend different intentions of the form designers. Thus, in Task B, the *Key* and *Value* components from Task A are further distinguished based

on the twelve Key intentions, which include: (1) Company Name (*com_nm*), (2) Company ID (*com_id*), (3) Holder Name (*hold_nm*), (4) Holder ID (*hold_id*), (5) Change Date (*chg_date*), (6) Previous Notice Given Date (*gvn_date*), (7) Previous Notice Date (*ntc_date*), (8) Class of Securities (*class*), (9) Previous Share (*pre_shr*), (10) Previous Voting Percentage (*pre_pct*), (11) Share (*new_shr*), (12) New Voting Percentage (*new_pct*). Key (1)-(2) and Key (3)-(4) ask for the identity of the listed company and the substantial holder, respectively, whereas Key (5)-(12) query the value of the shares being changed. We simply appended the identifier *Key* and *Value* separately to the twelve(12) *Key* names, and created twenty-four(24) labels for the key-value pairs accordingly, e.g., Each key-value pair would produce 2 labels; for example, "Company Name Key" and "Company Name Value". Based on these different key intentions, we inspected the value patterns to explore the potential noise caused by the diversified user understanding and uncertainty. The statistical analysis can be found in Section 4 (Figure 3).

3.3 Annotation Procedure

We recruited three human annotators and two human reviewers⁵ and performed an iterative annotation for each form component using the DataTorch⁶. Specifically, as shown in Figure 2, we split the annotation into three specialized sub-tasks: (1) *Key annotation* for the paired *Values*, and (2) *Value annotation* for the rest of the non-Key and non-Value components, including *Title*, *Section* and *Others*. Each annotator was assigned one sub-tasks respectively. They identified all possible targeted components and annotated the corresponding rectangular bounding boxes and labels based on the annotation guideline. During the annotation, any uncertain cases would be marked by the annotator and further decided by the discussion and majority voting in the regular meeting. The annotation guidelines were updated whenever needed to handle similar cases later. For instance, one typical type of uncertain case was the annotation of *Form Key* due to its various formats made by different form users, such as "The previous notice was dated 12 Feb 2003" as an example for the Key (7) Previous Notice Date, where the final agreed Key "The previous notice was dated" is expressed in a complete sentence with the Value "12 Feb 2003". This iterative annotation process produced 757 annotated Digital form documents. In addition, we also randomly selected 50 Printed and 50 Handwritten forms and conducted the annotation process the same as the Digital forms. We included these Printed and Handwritten forms as our additional test set for exploring the possibility of handling different form natures in real-world scenarios.

To ensure the final quality of the human annotation, we visualized the annotated bounding boxes and labels of all forms and assigned them to the two human reviewers for parallel manual-check. They checked each form one by one with reference to the up-to-date annotation guideline and annotated the correctness. They reviewed any form annotated with the incorrect label by both reviewers or received disagreement between them (i.e. one of them annotated as correct while the other annotated as incorrect). Then,

³See Section 5.1 Task A - Form Layout Analysing

⁴See Section 5.2 Task B - Key Information Extraction

⁵The three annotators are with a background in Computer Science or Financial at the University of Sydney while the two human reviewers are financial domain experts.

⁶<https://datatorch.io/features/annotator>

those are modified based on the final agreed decision of the two reviewers. To measure the quality of the annotations, we calculated the annotation agreement rate using both Cohen’s Kappa [3] and the Hamming Loss, which derived the overall scores of **0.998** and **0.003**, respectively, indicating a high annotation quality.

3.4 Annotation Format

The final annotation of our proposed **Form-NLU** is provided in *.json* format, separately for Task A and B⁷. The annotations for each form segment are stored as a *dictionary object* containing multiple key-value pairs with indicative key names in the *.json* file. The main attributes for each form segment shared by the two tasks include *bbox* for bounding box coordinates, *text* for textual tokens of this specific segment (e.g. from OCR or pdfminer⁸), and the *label* for this segment based on the specific task, e.g., “*Section*” for a segment in Task A or “*Company ID (com_id)*” for a segment in Task B. Besides, we also include some auxiliary attributes derived from our experiments for potential development in the future, such as the *segmentation* coordinates for Task A as well as the *visual_feature* and *bert_cls* for Task B (See Section 6.4, 6.5).

4 DATASET ANALYSIS

4.1 Component Distribution

The final version of our annotated **Form-NLU** consists of 857 forms, including 757 Digital, 50 Printed and 50 Handwritten forms. Table 2 shows the overall statistics of the data splits and the breakdown distribution of form components. We randomly split the Digital forms using the ratio of 70/10/20, resulting in 535, 76 and 146 forms for training, validation and testing. Besides, the 50 Printed and 50 Handwritten forms are included as our additional test splits. The breakdown distribution shows that the *Key* and *Value* persistently dominate since they are the main contents of the form. *Values* can be less than *Keys* due to the cases of empty values. The *Title* and *Section* are the least and demonstrate similar occurrences because each form document normally contains one main title with optional one or two subtitles while always having the two essential sections. The *Others* are slightly more frequent than the *Title* or *Section*, and the occurrence depends on the different form fillers. This consistent overall distribution across forms is attributed to the use of the Form 604 template. It makes our **Form-NLU** a promising benchmark dataset for exploring solutions to the real-world financial form layout understanding problem when the standard form template is available (our Task A). In addition, we further differentiate the keys and values for comprehending the specific form intentions (our Task B). These key-value pairs are contained in either *Table* (mostly vertical key-value alignment) or *Form* (mostly horizontal key-value alignment), which indicates the variety of spatial relationships that requires the model to learn in order to accurately identify the values and align with each key. In Table 3, we provide the average bounding box size and textual tokens of each form component, which reflects both their visual and linguistic features that can be

⁷We provide several real examples of the *.json* files for both tasks in https://github.com/adlnp/form_nlu#dataset-loading-and-samples. The complete dataset will be released upon acceptance of this paper.

⁸We use PDFMiner to extract the text of digital-born forms \mathcal{D} and Google Cloud Vision to extract the text of printed and handwritten sets (\mathcal{P} , \mathcal{H})

Type	Usage	Form Images	Title		Form		Table		Others
			Section	Key	Value	Key	Value		
Train	Digital	535	1068	1070	3708	3568	2669	2669	1691
Val	Digital	76	152	152	524	510	380	379	246
Test	Digital	146	292	292	1009	978	730	730	458
	Printed	50	98	100	346	332	250	249	152
	Handwritten	50	100	100	348	315	249	226	149
Total Number		857	1710	1714	5935	5703	4278	4253	2696

Table 2: Number of Form Components for Each Data Split.

Type	Title	Section	Form		Table	
			Key	Value	Key	Value
Bounding Box Average Width	186.15	136.82	95.46	85.98	57.29	41.92
Bounding Box Average Height	30.44	12.47	12.47	12.47	10.50	10.45
Bounding Box Average PX	4275.94	1720.62	1204.96	1183.10	609.75	456.94
Average Number of Tokens	7.35	7.10	5.16	4.14	4.29	1.74

Table 3: Average Bounding Box Width, Height, Number of Pixel(PX), and Number of Tokens for Each Form Component

potentially helpful for specific key and value identification. *Title* and *Section* tend to have longer textual content and bigger component size with larger font sizes. In comparison, *Key* and *Value* are much shorter and are contained in smaller bounding boxes. Especially the *Value* in the *Table* has the shortest length as they are mostly numeric values.

4.2 Value Pattern Proportion

As mentioned before, the filled-in content varies and involves potential noises due to the diversified user understanding and uncertainty. To illustrate the variety of the filled-in content, we summarise the main value patterns for each key intention and provide the proportions in Figure 3. Throughout all keys, the largest proportion of value follows the common value pattern that the form users tend to provide minimal information to fulfil the intention of the keys. For example, they put simply the company name only to the key *com_nm* in Figure 3a or ID only for 3b. However, there are several other value patterns available, including additional information shown in Figure 3a, 3b and 3e, or having different data format or writing style, such as the various date formats in Figure 3c. These human-caused variants well reflect the real-world form understanding scenarios and imply potential challenges for achieving precise information retrieval in our Task B. Overall, the Share Class (Figure 3d) and the date type values (Figure 3c) show comparatively more patterns while the numerical type values in Figure 3e tend to be more consistent with the common pattern. This pattern similarity among the keys with the same data type values indicates that solely relying on the linguistic cues may not be enough as understanding the relative spatial location of the values in the form is also required in order to distinguish these similar values of different keys.

4.3 Key-value Pair Comparison Analysis

To provide an in-depth analysis of the nature of key-value pairs, we further summarize the average character number and the ratio of spatial relation for the twelve key-value pairs in Figure 4 and 5, respectively. It can be observed that typically there are seven fixed key-value pairs in Form (i.e., Figure 4a/5a) and another fixed five in Table (i.e., Figure 4b/5b) formatted by the standard Form 604 template. Two groups of key-value pairs demonstrate their own features. As shown in Figure 4a for the form-based key-value pairs,

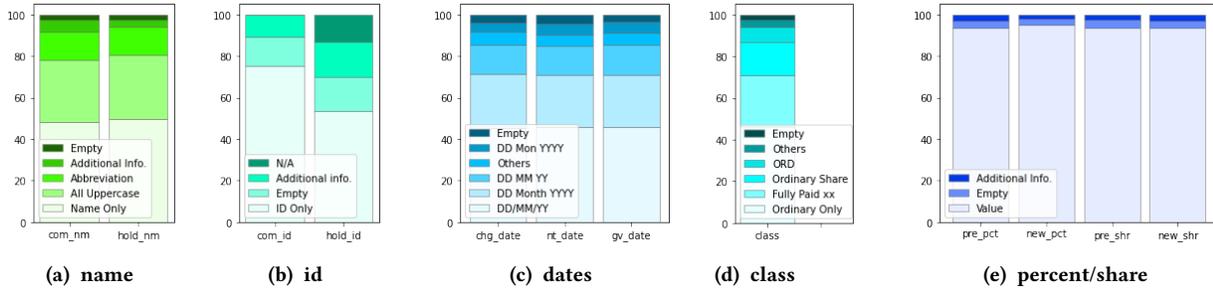


Figure 3: Value Patterns Distributions for Each Key Group.

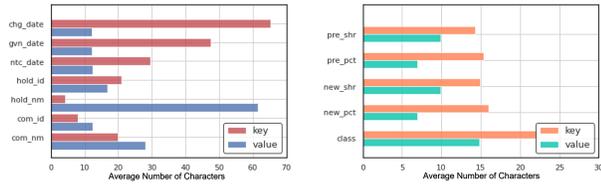


Figure 4: Average # of Characters in Each Key Value Pair

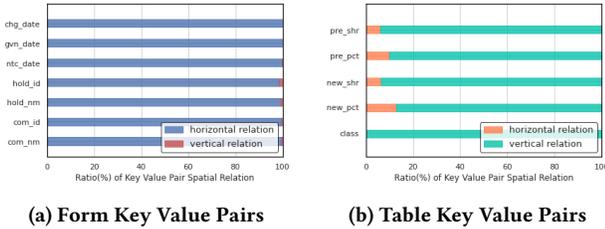


Figure 5: Ratio (%) of Horizontal and Vertical Relation for Each Key Value Pair

the three date type pairs tend to have much longer value than the key where as the other four pairs may have either longer or shorter values than keys. The *hold_nm* key-value shows the largest character length gap among all. In comparison, the five Table-based pairs in Figure 4b all have shorter values than keys, and the length gaps are similar. Overall, most keys have short values (e.g., < 20 characters) from which the semantic context could be scarce. On the other hand, as seen from Figure 5, the seven Form-based pairs are uniformly aligned horizontally, whereas the other five pairs in Table are mostly vertical. However, some special cases with the opposite alignment are also observed in both groups, such as the vertical key-value pairs for the *hold_id/hold_nm* and *com_id/com_nm* in the Form-based group (Figure 5a), as well as the horizontal pairs for the four percentage type keys and values in the Table-based group (Figure 5b). Thus, simply memorizing the spatial alignment for each key cannot lead to optimized value retrieval in Task B.

5 TASKS OVERVIEW

5.1 Task A - Form Layout Analysing

The purpose of Task A is to detect the semantic entities (*Title*, *Section*, *Form_key*, *Form_value*, *Table_key*, *Table_value*, *Others*) of

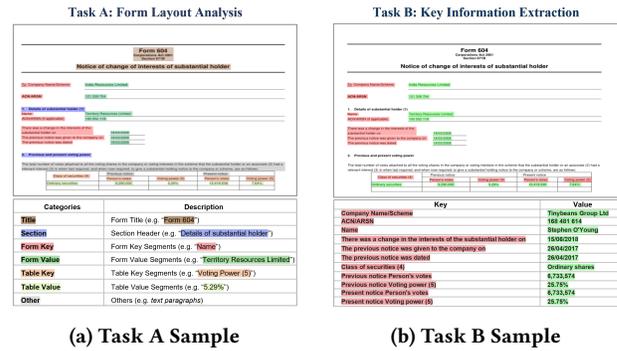


Figure 6: Examples for Task A and B. For Task A, the users need to detect each layout component’s bounding box and recognize the associated category. Task B asks the user to feed the fixed key text (red area) into the model to predict the corresponding value of RoI’s index (green area)

forms. This is a prior task to understand the form layout by detecting the position of each semantic entity, including key, value and other objects. Given a form image I , an object detection model is used to detect a set of RoIs (r_1, r_2, \dots, r_n) . Each r_i contains bounding box b_i coordinates (x_i, y_i, w_i, h_i) with semantic category c_i . As shown in Figure 6a, the bounding box of each semantic entity in an input form image is detected and coloured based on the recognized categories. The model should effectively detect layout components such as form *Title* (like "Form 604") and *Section* headers (such as "Details of substantial holder"). Additionally, we also expect the models can differentiate keys/values located in *Form* or *Table*; for example, "Name" should be detected as a *Form_key* instance, while "Voting Power (5)" is a *Table_key* instance.

5.2 Task B - Key Information Extraction

Task B aims to evaluate whether the proposed models could comprehensively understand designers’ intentions and user uncertainties to extract valuable information from input forms. It allows using ground truth RoIs’ set $\mathcal{R}_{gt} = (R_1, R_2, \dots, R_n)$ during the training and inference stage. Given key text information t , a document image I and a set of ground truth RoIs (r_1, r_2, \dots, r_n) , a model H can output the RoI’s index number aligned with input t . As Figure 6b shown, each highlighted entity with red background colour is the key t we

need to feed into the proposed model, and the green paired RoI's index is the desired output from the model based on the current input key. For example, inputting the required key text content "Company Name/Scheme", the desired output from the model is the RoI's index of paired values where we can easily get the text content ("Tinybeans Group Ltd") from it.

6 EVALUATION SETUP

6.1 Implementation Detail

For Task A, we fine-tune Faster-RCNN and Mask-RCNN models with two backbones⁹ respectively on our dataset based on Detectron2 platform. We set 5000, 128, and 0.02 as the maximum iteration times, batch size and base learning rate, and other setups are the same as Detectron2 official tutorial¹⁰. Regarding Task B, we employ various approaches to encode the vision and language features. Firstly, all Task B adopted baselines use pretrained BERT to encode key textual content. Moreover, for the visual aspect, VisualBERT, LXMERT, and M4C models utilize 2048-d features extracted from the Res5 layer of ResNet101. The maximum number of input key text tokens and the number of segments on each page are all defined as 50 and 41, respectively. Task A and B experiments are conducted on 51 GB Tesla V100-SXM2 with CUDA11.2.

6.2 Task A Baselines and Metric

We use two popular object detection models, **FasterRCNN** [18] and **MaskRCNN** [6] with different ImageNet pretrained backbones testing on our dataset. To evaluate the performance of the object detection model, we apply a *mAP* (mean Average Precision), which is commonly used in object detection tasks. The *mAP* score is calculated by averaging *APs* over all categories' overall pre-defined IoU thresholds¹¹.

6.3 Task B Baselines and Metric

Task B mainly focus on predicting the corresponding value of the RoI index based on input text content (as shown in Figure 6b). Thus, several multi-modal transformer frameworks are adopted as baselines on **Form-NLU**, of which inputs are multi-aspect RoI features, including large pre-trained **VisualBERT** [12], **LXMERT** [20] and non-pre-trained **M4C** [7] models¹². Regarding evaluation, we use weighted F1-score as the primary evaluation metrics for representing overall and breakdown performance. Note that some visual language pre-trained models, such as ViLT [11], LayoutLM [25], were excluded since those are mainly based on the image patches or pieces and do not fit into the task demands.

6.4 Task A Model

For Task A - Form Layout Analysing task (Section 7.1), we use Faster-RCNN and Mask-RCNN with various depth backbones (ResNet-50 and ResNet-101) as baselines and contact experiments to check the effects of diverse model architecture and model size. Additionally,

⁹Faster-RCNN backbones are faster_rcnn_R_50_FPN_3x, faster_rcnn_R_101_FPN_3x and Mask-RCNNs are mask_rcnn_R_50_FPN_3x, mask_rcnn_R_101_FPN_3x

¹⁰https://colab.research.google.com/drive/16jcaJoc6bCFAQ96jDe2HwtXj7BMD_-m5

¹¹We refer [32] to adopt *mAP* as metrics and follow their thresholds for Task A.

¹²Detailed baseline setup can be found in https://github.com/adlnlp/form_nlu#baseline-model-description

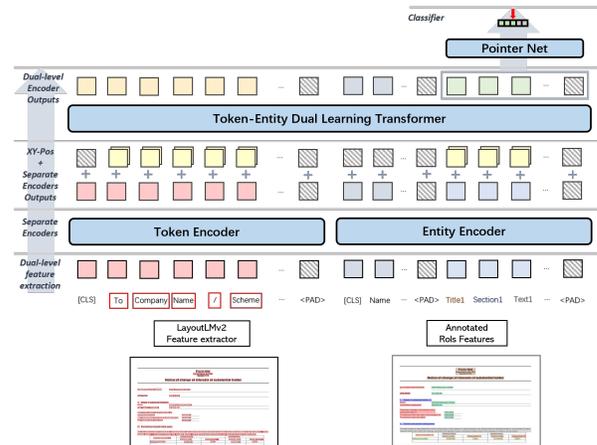


Figure 7: Our Proposed Key Information Extraction Model Architecture for Task B

for the real-world case analysis (Section 7.3), we adopted some methods using open-source PDF parsers (such as PDFMiner) to analyse the form layout without training the deep learning models. The PDF parser outputs can also be used as inputs for Task B.

6.5 Task B Model

For Task B, we propose a new document key information extraction model (as Figure 7 shown) to predict the corresponding values from input key content. To achieve this, we investigate positional and logical relations and their hierarchical structure with two components: multiple aspect features extraction/integration and entity-token dual-level Model.

6.5.1 Multi-Aspect Features. We investigate five aspect features, including Visual (*V*), Textual (*T*), Positional (*P*), Density (*D*), and Gap Distance (*G*) features, with the corresponding encoding approaches to comprehensively explore which features benefit understanding designer intentions and argument form-like document segment representations. Many traditional methods for document understanding tasks have shown the effectiveness and significance of visual, textual, and positional features only [25, 28, 29]. Except for those three aspect features, we also aim to explore the effectiveness of text density for form understanding, which has been demonstrated by [14] for document layout analysis. Moreover, based on the dataset analysis results, the gap distance between entities is crucial for understanding the form structure. Hence, we introduce two layout-related features, including normalized positional features (bounding box coordinates) and the gap distance between an entity to neighbours¹³.

6.5.2 Entity-Token Dual Level Model. The proposed new entity-token dual-level form information extraction model contains four components, including Entity Encoder $Encoder_{entity}$, Token Encoder $Encoder_{token}$, Dual Encoder $Encoder_{dual}$ and Pointer Net-based classifier ϕ .

¹³Our Form-NLU Github (https://github.com/adlnlp/form_nlu#multi-aspect-features) provides more detailed feature representation approaches.

1) **Entity Encoder** aims to learn the semantic relations between entities for enhancing the vanilla entity representations. Based on the preliminary results, we select the pretrained LXMERT to citetan2019lxmert as the entity encoder. The inputs of the LXMERT-based entity encoder include key text T_{key} , Entities' visual features V_{entity} , and the bounding box coordinates B_{entity} . After feeding those features into pretrained LXMERT, the enhanced Entity Rols' representations $E \in \mathbb{R}^{768}$ can be extracted.

2) **Token Encoder** aims to acquire cross-modal token representations. The SoTA layout-aware pre-trained visual language transformer, LayoutLMv2 [28], is adopted as the token encoder. We firstly employ *LayoutLMv2FeatureExtractor* to extract and encode the token level features, including token text T_{token} , token Rols' bounding box B_{token} with pixel features I_{pixel} . Then, we follow the original LayoutLMv2 setup to process the extracted token-level features and feed them into our token-level encoder *Encoder_{token}*, pretrained LayoutLMv2 to get the token-level representation $T \in \mathbb{R}^{768}$.

3) **XY-Pos Dual Encoder** is designed for learning the geometrically sensitive multi-grained and multi-modality feature representations. After getting entity and token level feature representations from pretrained models, we treat them as the sequence inputs of a dual-level mutual learning encoder *Encoder_{dual}*, where a 6-layer transformer with 8-heads self-attention is used as the basic framework. Unlike the original transformer that adopts sinusoidal positional encoding [22], we proposed a new geometric-sensitive positional encoding *XY-Pos*. It flattens the stacked sets of normalized bounding box coordinates along with the X or Y axis to get X_i^{pos} and Y_i^{pos} . Supposing $b_i = (x_i, y_i, w_i, h_i)$ is the bounding box of r_i of *Encoder_{dual}* in document page D_j , the size of D_j is (W_j, H_j) .

$$pos_i^x = \left[\frac{x_i + d_{i_1}}{W_j}, \frac{x_i + d_{i_2}}{W_j}, \dots, \frac{x_i + d_{i_m}}{W_j} \right] \quad (1)$$

where the l -th step $d_{i_l} = \frac{w_i \times l}{m}$

$$X_i^{pos} = flatten([pos_i^x] \times n) \quad (2)$$

We set m and n as 32 and 24 to get a 768-d vector. Finally X_i^{pos} is flattened into a $m \times n$ dimensional vector, we have $X_i^{pos} \in \mathbb{R}^{768}$. Similar procedures are used to generate positional encoding Y_i^{pos} along *Y-axis* of r_i . The final input representations before feeding into *Encoder_{dual}* can be represented as:

$$E_{pos_{xy}} = E + X_{entity}^{pos} + Y_{entity}^{pos} \quad (3)$$

$$T_{pos_{xy}} = T + X_{token}^{pos} + Y_{token}^{pos} \quad (4)$$

Unlike most VLPMS [20, 28] adopted positional encoding methods through linear projecting 4-d bounding box coordinates into high dimensional vectors. *XY-pos* could capture more geometric features between entities and tokens. Then we will feed them into *Encoder_{dual}* and get the updated token and entity representations T_{dual}, E_{dual} .

4) **Pointer Net based Classifier** The multi-aspect features are concatenated with E_{dual} to get E_{multi} . E_{multi} is fed into the pointer net [23] based classifier ϕ to get a score vector which following a *softmax* to retrieve the final prediction results y_{pre} .

Test set	Model	mAP	Breakdowns (Precision)						
			Title	Section	Form_key	Form_value	Table_key	Table_value	Others
\mathcal{D}	F-50	68.98	73.54	69.36	67.93	70.37	69.18	66.76	65.73
	F-101	69.99	71.32	68.30	72.37	71.70	68.74	72.84	64.65
	M-50	71.74	77.29	69.06	73.01	70.30	68.40	73.25	70.90
	M-101	71.93	79.54	69.64	71.25	72.38	67.02	74.38	69.33
\mathcal{P}	F-50	56.46	57.44	53.98	54.49	51.46	59.31	62.86	55.67
	F-101	59.54	55.72	52.81	63.13	61.32	59.42	68.11	56.26
	M-50	63.06	59.83	60.98	63.74	61.91	70.75	65.93	58.29
	M-101	65.47	63.05	62.49	66.84	61.80	71.82	70.19	62.10
\mathcal{H}	F-50	49.87	58.73	43.85	63.14	25.15	63.48	40.14	54.56
	F-101	50.39	53.92	33.36	65.28	36.68	66.97	41.06	55.43
	M-50	57.99	61.45	48.46	68.16	43.14	71.76	56.95	56.05
	M-101	60.22	64.81	54.46	69.97	46.82	73.34	52.81	59.32

Table 4: Overall Performance and Breakdown Results for Layout Analysing Task (Task A). F-50 and F-100 represent the Faster-RCNN with ResNet50 and ResNet101 as backbones. The same patterns occur for Mask-RCNN (M-50 or M-101).

7 RESULTS

7.1 Task A: Form Layout Analysing

7.1.1 *Overall and Breakdown Performance.* This section shows the test performance of layout analysing (Task A) models on digital (\mathcal{D}), printed (\mathcal{P}), and handwritten (\mathcal{H}) sets. From Table 4, we observe that the *mAP* of Mask-RCNNs can achieve around 2%, 6% and 10% higher than Faster-RCNNs with identical backbones on \mathcal{D} , \mathcal{P} and \mathcal{H} , respectively. This may result from the finer spatial localisation of auxiliary instance segmentation tasks adopted by Mask-RCNNs. Then, we explore the effects of various backbones, like Faster-RCNN with ResNet-101 (F-101) or Mask-RCNN with ResNet-50 (M-50). From Table 4, we can find F-101 and M-101 consistently achieve around 1% to 2% higher than F-50 and M-50. It demonstrates that the deeper or large-scale pretrained backbones may generate more comprehensive visual representations. Notably, the overall performance of \mathcal{D} is around 5% higher than \mathcal{P} and even about 11% more than \mathcal{H} . The main reason may result from the apparent difference between scanned (both printed \mathcal{P} and handwritten \mathcal{H}) and digital (\mathcal{D}) forms, especially the handwritten forms involving more user-uncertainties such as writing mistakes or scanning rotation.

As our Form-NLU provides fine-grained layout component types such as *Value* subdivided into *Form_value* and *Table_value*, it enables us to observe and analyse the performance of subdivided components affected by distinct layout position distribution. For example, certain layout components show stable performance among the three test sets, such as *Form_key* and *Table_key*. It may result from those components designed by form designers with more shared layout and visual patterns like font and layout arrangement. However, *Form_value* and *Table_value* components may involve more uncertainties and noise from users, electric devices or propagation process, leading to *mAP* on \mathcal{H} (46.82% and 52.81%) being lower than \mathcal{P} (61.80% and 70.19%) and much lower than \mathcal{D} (72.38% and 74.38%). In addition, for *Title* and *Others* components, the performance of M-101 on \mathcal{D} (79.54% and 69.33%) is apparently higher than \mathcal{P} (63.05% and 62.10%) and \mathcal{H} (64.81% and 59.32%), while \mathcal{P} and \mathcal{H} almost have similar performance. The reason may come from the difference between scanned (\mathcal{P} and \mathcal{H}) and digital forms, such as the rotation or lower resolution of scanned forms.

7.1.2 *Stepped Training Set Ratios.* We set stepped ratios of training set size (\mathcal{T}) (10%, 50% and 100%) to train M-101 and evaluate it on \mathcal{D} , \mathcal{P} , and \mathcal{H} sets to explore the effects of training size on Task A. It can be seen from Figure 8 as \mathcal{T} size increases, overall and

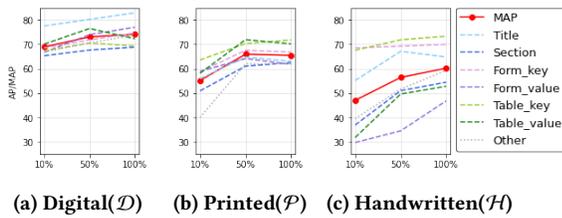


Figure 8: Different Training Ratio mAP on Digital (\mathcal{D}), Printed (\mathcal{P}), and Handwritten (\mathcal{H}) Test Sets for Task A

breakdown performance improve, whereas this trend becomes inapparent when \mathcal{T} reaches 50% on \mathcal{D} and \mathcal{P} . Regarding performance on \mathcal{H} , increasing the training size could boost the performance of all types of layout components, significantly for *Table_value*, *Section* and *Others*. Especially, as one significant layout component type for downstream tasks, *Table_value* hold around 30% on \mathcal{H} . In general, the trends shown in Figure 8 demonstrate that our training set (\mathcal{T}) contains various form types to enable handling distinct application scenarios, even with the small ratio of \mathcal{T} . Moreover, it also represents increasing the size of digital training sets (\mathcal{T}) that could effectively improve printed and handwritten performance. More training samples may enable catching more shared feature patterns, such as spatial distributions, typographical similarity, font type and size of specific components.

7.2 Task B: Key Information Extraction

7.2.1 Overall and Multi-aspect Feature Performance. This section compares and analyses the performance between our model and three widely used baselines under different configurations on digital (\mathcal{D}), printed (\mathcal{P}) and handwritten (\mathcal{H}) test sets. Firstly, we focus on comparing the performance among vanilla models (the first row of each model group in Table 5) in which input features and model architectures are the same as their paper described. Our model achieves better performance than all vanilla baselines on \mathcal{D} , \mathcal{P} and \mathcal{H} . It indicates that iterative learning with token-level features could enhance entity-level representations for boosting downstream performance. In addition, for three vanilla baselines, LXMERT can get 2.28%, 6.38% and 15.36% higher than VisualBERT on \mathcal{D} , \mathcal{P} and \mathcal{H} , respectively, which may result from inputs of LXMERT containing positional information, but VisualBERT is pre-trained on the visual feature of input RoIs only. Subsequently, compared with VisualBERT, the non-pretrained M4C can increase by around 1.5%, 3% and 9% on \mathcal{D} , \mathcal{P} and \mathcal{H} , respectively. It may illustrate the significance of textual and positional features. However, due to the input feature differences between vanilla models, we conducted external experiments to explore the effects of multi-aspect features on all adopted models for a fair comparison.

M4C initially contains V , T , and P of input RoIs, and is a non-pretrained model with random initial weights. Thus, it may cause only slight improvements can be found in Table 5 after adding D and G into the model on three test sets. For VisualBERT, the vanilla model only contains RoIs visual features; after stacking more aspect features, the performance is gradually improved where the F1-score on \mathcal{H} increases from 65.25% (vanilla) to 72.65% (with V, T, S, D, G). There is an apparent increase after adding D and G into VisualBERT, which may contribute to the additional features making the input

Model	Input Features					Overall (F1-score)		
	V	T	P	D	G	Digital(\mathcal{D})	Printed(\mathcal{P})	Handwritten(\mathcal{H})
M4C	○	○	○	×	×	96.91	88.62	74.06
	○	○	○	○	×	96.32	89.52	74.98
	○	○	○	×	○	97.00	88.81	75.04
	○	○	○	○	○	97.22	89.78	74.89
	○	×	×	×	×	95.55	85.91	65.25
VisualBERT	○	×	○	×	×	95.84	85.93	67.25
	○	○	○	×	×	96.07	85.90	70.14
	○	○	○	×	○	96.61	87.43	70.79
	○	○	○	○	×	96.70	87.00	71.28
	○	○	○	○	○	96.73	87.18	72.65
LXMERT	○	×	×	×	×	97.83	92.29	80.51
	○	○	○	×	×	97.67	94.15	82.80
	○	○	○	×	○	97.70	94.59	82.60
	○	○	○	×	○	97.74	94.49	83.71
	○	○	○	○	○	97.74	95.07	84.43
Our Model	○	○	○	×	×	99.30	95.32	84.75
	○	○	○	×	○	99.12	95.50	85.77
	○	○	○	○	○	99.09	95.50	86.98
	○	○	○	○	○	99.09	95.50	86.98

Table 5: Different Background Colours. This is to show the results of vanilla model (first row of each model group), model with all aspect features (last row of each model group) on \mathcal{D} , \mathcal{P} , and \mathcal{H} Test Sets. RoI inputs contain visual (V), textual (T), positional (P), text density (D), gap distance (G) features, where \circ and \times represent using or not using the specific feature in that column.

PE Approach	Digital(\mathcal{D})	Printed(\mathcal{P})	Handwritten(\mathcal{H})
Without Positional Encoding	98.83	94.34	85.98
Linear Projection	99.09	95.50	86.98
XY-Positional Encoding	99.22	96.14	89.13

Table 6: Performance of Various Positional Encodings (PE)

representations more comprehensive. A similar trend also can be found in LXMERTs where T may contribute to a more positive effect (about 2% increasing) than VisualBERT on \mathcal{P} and \mathcal{H} . Regarding our model, although there is no noticeable improvement observed on \mathcal{D} , the positive effects on \mathcal{P} and \mathcal{H} can be found, especially for \mathcal{H} (from 84.56% to 86.98%). Generally, the proposed additional features could help the model better understand form structures to improve the model generalisation ability, notably when the appearance of input forms differs during the training and inference stage.

7.2.2 Positional Encoding Validation. To demonstrate the effectiveness of XY-pos, we conduct external experiments to compare the performance of the model with linear projection PE methods, also without any PE methods. Compared model without any PE in Table 6, linear projection and XY-pos can have an apparent increase in all three test sets. It illustrates that the positional information is significant to understand the layout structure for form understanding. Furthermore, from Table 6, we can find XY-pos reach a better performance than linear projection on \mathcal{D} , \mathcal{P} and \mathcal{H} . Significantly, the F1 increased from 86.98% to 89.13% on \mathcal{H} . It demonstrates XY-Pos may capture more positional information to understand various layout structures of input forms better.

7.2.3 Fine-grained Training Set Ratio. For exploring the training size influences for form understanding, we define fine-grained training set ratios (from 10% to 100%) to train our model and represent the evaluation performance on \mathcal{D} , \mathcal{P} and \mathcal{H} in Figure 9. With increasing training size, fluctuating increases can be observed for overall and breakdown performance on three test sets. The rapid

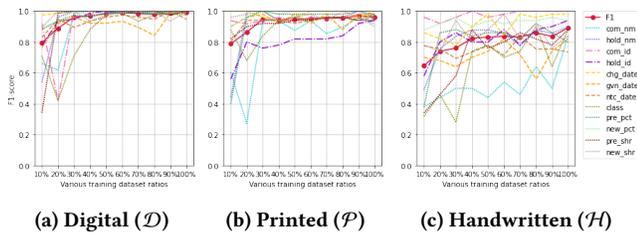


Figure 9: Performance of our Model with Fine-grained Training Set Ratios on Three Test Sets

Model	Acc	Breakdowns (Accuracy)											
		nm(name)		id		date		class	pct(percent)		shr(share)		
		com	hold	com	hold	chg	ntc		gvn	pre	new	pre	new
LXMERT	43.78	86.30	80.14	69.18	52.74	72.60	65.07	54.79	4.79	6.85	24.66	5.48	2.74
VisualBERT	57.48	88.36	78.08	63.01	59.59	39.04	70.55	37.67	47.95	50.00	28.77	53.42	73.29
M4C	69.35	83.56	80.82	63.70	73.29	65.75	57.53	82.88	61.64	65.75	58.90	59.59	78.77
Ours	72.72	91.78	81.51	71.92	71.23	80.82	80.82	84.93	71.92	70.55	80.14	50.00	57.53

Table 7: Overall performance and Breakdown Results of Key Information Extraction Task with PDFminer

increases can be observed on \mathcal{D} and \mathcal{P} before reaching 50% \mathcal{T} , following the stable trends after reaching this critical ratio. However, for handwritten set \mathcal{H} , we can find an apparent overall performance increase during the entire training size interval. It may demonstrate the wide variety of our training dataset \mathcal{H} can result in a more generic trained model to extract specific key information more effectively. Additionally, unlike the datasets providing general key-value annotations, we can show the training-size sensitivity of specific key-value pairs to analyse the robustness of the trained model with limited training data. For example, com_nm , com_id , and $hold_nm$ show more sensitive to training size, while date-related keys such as ntc_date are less sensitive.

7.3 Task B with PDF Parser

Real-world users, especially non-deep learning users, are challenged to get input RoIs through well pretrained layout analysing models because of lacking ground truth annotations and relevant background knowledge. Thus, using textlines extracted by specific PDF parsing tools is an alternative way to obtain the input of adopted well-trained models, such as the green rectangles in Figure 10. We use PDFminer (a widely used PDF parser¹⁴) to extract RoIs for replacing the ground truth RoIs in a digital set (\mathcal{D}) and feed them into three trained baseline models and our models. We use IoU = 0.5 as a threshold to calculate accuracy, as Table 7 shows.

Different from the trend of feeding ground truth RoIs during the testing stage, the non-pretrained model M4C can achieve much better accuracy (69.35%) compared with pre-trained LXMERT (35.24%) and VisualBERT (53.34%). The possible reason might be that too many noise RoIs detected by PDFminer feed into the large-scale pre-trained models, which decreases those heavy models' RoIs feature representation ability. Our model can achieve the highest overall performance among all tested models. It demonstrates that the proposed XY-pos enhanced dual-level model can improve generality to understand the input contents better. Regarding breakdown results, com_nm and $hold_nm$ could always perform better. Most of the

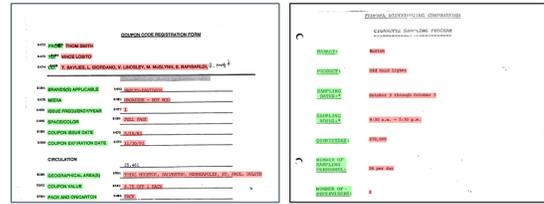


Figure 10: Two Document Samples with Selected Question and Answer Pairs from FUNSD Subset

key information extracted by our model could hold higher accuracy than other baselines, such as com_id and chg_date . Notably, LXMERT shows high sensitivity to the bounding box precision of input RoIs. However, because PDFminer extracted bounding boxes are inaccurate, especially in the table area, the keys located in the table are much lower than other models, such as $class$, pre_pct , pre_shr , etc.

8 CASE STUDY: TRANSFER LEARNING

In order to check the feasibility of our proposed data and model, we conducted a case study with transfer learning. Note that transfer learning aims to storing knowledge gained while solving our dataset and apply pre-trained models to a different but related problem. In this case study, we applied the best (our model) and second-best models (LXMERT), trained by the Form-NLU dataset, to the publicly available benchmark FUNSD. We use the selected key (question)-value (answer) pairs from the FUNSD [10] dataset like the samples shown in Figure 10¹⁵. Based on the testing results, even if the nature of FUNSD are entirely different from our Form-NLU dataset, the models trained on our dataset can still achieve sound performance. The result shows that LXMERT (the second best of in Task B) correctly predicted 80 (53.33%) samples out of the 150 FUNSD key(question)-value(answer) pairs while our model achieved 94 (62.67%) correct predictions. It demonstrates that our Form-NLU dataset can learn the general form layout and be applied in other benchmarks, in addition to this, our proposed model is efficient to extract a feature representation for form understanding.

9 CONCLUSION

We proposed Form-NLU, a new form structure understanding and key information extraction dataset. The proposed dataset covers the important point of view, enabling the interpretation of the form designer's specific intention and the alignment of user-written value on it. The dataset includes three form types: digital, printed, and handwritten, which cover diverse form appearances/layouts and deal with their noises. Moreover, we propose a new strong baseline for form structure understanding and key-value information extraction, which applies robust positional and logical relations. Our model outperformed all state-of-the-art models in key-value information extraction tasks. We do hope that our proposed dataset and model can be a great insight into form structure and information analysis, hence, we adopted with off-the-shelf pdf layout extraction tool and also provide its feasibility by conducting transfer learning.

¹⁴<https://pypi.org/project/pdfminer/>

¹⁵Please refer to https://github.com/adlnlp/form_nlu/blob/main/README.md#case-study-setup to check the setup detail

REFERENCES

- [1] Millicent Chang, Raymond da Silva Rosa, and Wilson Ng. 2016. The Informativeness of Substantial Shareholder Trading in the Lead up to a Takeover Bid. In *Asian Finance Association (AsianFA) 2016 Conference*.
- [2] Zewen Chi, Heyan Huang, Heng-Da Xu, Houjin Yu, Wanxuan Yin, and Xian-Ling Mao. 2019. Complicated table structure recognition. *arXiv preprint arXiv:1908.04729* (2019).
- [3] Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and psychological measurement* 20, 1 (1960), 37–46.
- [4] Yihao Ding, Zhe Huang, Runlin Wang, YanHang Zhang, Xianru Chen, Yuzhong Ma, Hyunsuk Chung, and Soyeon Caren Han. 2022. V-Doc: Visual questions answers with Documents. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 21492–21498.
- [5] Zhangxuan Gu, Changhua Meng, Ke Wang, Jun Lan, Weiqiang Wang, Ming Gu, and Liqing Zhang. 2022. Xylayoutlm: Towards layout-aware multimodal networks for visually-rich document understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4583–4592.
- [6] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2017. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*. 2961–2969.
- [7] Ronghang Hu, Amanpreet Singh, Trevor Darrell, and Marcus Rohrbach. 2020. Iterative answer prediction with pointer-augmented multimodal transformers for textvqa. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9992–10002.
- [8] Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and Furu Wei. 2022. LayoutLMv3: Pre-training for Document AI with Unified Text and Image Masking. *arXiv preprint arXiv:2204.08387* (2022).
- [9] Zheng Huang, Kai Chen, Jianhua He, Xiang Bai, Dimosthenis Karatzas, Shijian Lu, and CV Jawahar. 2019. Icdar 2019 robust reading challenge on scanned receipts ocr and information extraction. In *International Conference on Document Analysis and Recognition*.
- [10] Guillaume Jaume, Hazim Kemal Ekenel, and Jean-Philippe Thiran. 2019. Funsd: A dataset for form understanding in noisy scanned documents. In *2019 International Conference on Document Analysis and Recognition Workshops (ICDARW)*, Vol. 2. IEEE, 1–6.
- [11] Wonjae Kim, Bokyoung Son, and Ildoo Kim. 2021. Vilt: Vision-and-language transformer without convolution or region supervision. In *International Conference on Machine Learning*. PMLR, 5583–5594.
- [12] Liunan Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557* (2019).
- [13] Minghao Li, Yiheng Xu, Lei Cui, Shaohan Huang, Furu Wei, Zhounjun Li, and Ming Zhou. 2020. DocBank: A Benchmark Dataset for Document Layout Analysis. In *Proceedings of the 28th International Conference on Computational Linguistics*. 949–960.
- [14] Siwen Luo, Yihao Ding, Siqu Long, Josiah Poon, and Soyeon Caren Han. 2022. Doc-GCN: Heterogeneous Graph Convolutional Networks for Document Layout Analysis. In *Proceedings of the 29th International Conference on Computational Linguistics*. 2906–2916.
- [15] Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. 2021. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 2200–2209.
- [16] Seunghyun Park, Seung Shin, Bado Lee, Junyeop Lee, Jaeheung Surh, Minjoon Seo, and Hwalsuk Lee. 2019. CORN: a consolidated receipt dataset for post-OCR parsing. In *Workshop on Document Intelligence at NeurIPS 2019*.
- [17] Johannes Rausch, Octavio Martinez, Fabian Bissig, Ce Zhang, and Stefan Feuerriegel. 2021. Docparser: Hierarchical document structure parsing from renderings. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 4328–4338.
- [18] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems* 28 (2015).
- [19] Tomasz Stanisławek, Filip Graliński, Anna Wróblewska, Dawid Lipiński, Agnieszka Kaliska, Paulina Rosalska, Bartosz Topolski, and Przemysław Biecek. 2021. Kleister: key information extraction datasets involving long documents with complex layouts. In *International Conference on Document Analysis and Recognition*. Springer, 564–579.
- [20] Hao Tan and Mohit Bansal. 2019. Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490* (2019).
- [21] Ryota Tanaka, Kyosuke Nishida, and Sen Yoshida. 2021. VisualMRC: Machine Reading Comprehension on Document Images. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 13878–13888.
- [22] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [23] Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. 2015. Pointer networks. *Advances in neural information processing systems* 28 (2015).
- [24] Jiapeng Wang, Chongyu Liu, Lianwen Jin, Guozhi Tang, Jiaxin Zhang, Shuaitao Zhang, Qianying Wang, Yaqiang Wu, and Mingxiang Cai. 2021. Towards robust visual information extraction in real world: new dataset and novel solution. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 2738–2745.
- [25] Yiheng Xu, Minghao Li, Lei Cui, Shaohan Huang, Furu Wei, and Ming Zhou. 2020. 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*. 1192–1200.
- [26] Yiheng Xu, Tengchao Lv, Lei Cui, Guoxin Wang, Yijuan Lu, Dinei Florencio, Cha Zhang, and Furu Wei. 2021. Layoutlm: Multimodal pre-training for multilingual visually-rich document understanding. *arXiv preprint arXiv:2104.08836* (2021).
- [27] Yiheng Xu, Tengchao Lv, Lei Cui, Guoxin Wang, Yijuan Lu, Dinei Florencio, Cha Zhang, and Furu Wei. 2022. XFUND: A Benchmark Dataset for Multilingual Visually Rich Form Understanding. In *Findings of the Association for Computational Linguistics: ACL 2022*. 3214–3224.
- [28] Yang Xu, Yiheng Xu, Tengchao Lv, Lei Cui, Furu Wei, Guoxin Wang, Yijuan Lu, Dinei Florencio, Cha Zhang, Wanxiang Che, et al. 2021. LayoutLMv2: Multimodal Pre-training for Visually-rich Document Understanding. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. 2579–2591.
- [29] Yue Zhang, Bo Zhang, Rui Wang, Junjie Cao, Chen Li, and Zuyi Bao. 2021. Entity Relation Extraction as Dependency Parsing in Visually Rich Documents. *arXiv preprint arXiv:2110.09915* (2021).
- [30] Xinyi Zheng, Douglas Burdick, Lucian Popa, Xu Zhong, and Nancy Xin Ru Wang. 2021. Global table extractor (gte): A framework for joint table identification and cell structure recognition using visual context. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*. 697–706.
- [31] Xu Zhong, Elaheh ShafieiBavani, and Antonio Jimeno Yepes. 2020. Image-based table recognition: data, model, and evaluation. In *European Conference on Computer Vision*. Springer, 564–580.
- [32] Xu Zhong, Jianbin Tang, and Antonio Jimeno Yepes. 2019. Publaynet: largest dataset ever for document layout analysis. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*. IEEE, 1015–1022.