

StrucTeXT: Structured Text Understanding with Multi-Modal Transformers

Yulin Li*

Department of Computer Vision
Technology (VIS), Baidu Inc.
liyulin03@baidu.com

Xiameng Qin

Department of Computer Vision
Technology (VIS), Baidu Inc.
qinxiameng@baidu.com

Kun Yao, Junyu Han

Department of Computer Vision
Technology (VIS), Baidu Inc.
{yaokun01,hanjunyu}@baidu.com

Yuxi Qian*

Beijing University of Posts
and Telecommunications
qianyuxi@bupt.edu.cn

Chengquan Zhang†

Department of Computer Vision
Technology (VIS), Baidu Inc.
zhangchengquan@baidu.com

Yuechen Yu*

Department of Computer Vision
Technology (VIS), Baidu Inc.
yuyuechen@baidu.com

Yan Liu

Taikang Insurance Group
liuyan146@taikanglife.com

ABSTRACT

Structured text understanding on Visually Rich Documents (VRDs) is a crucial part of Document Intelligence. Due to the complexity of content and layout in VRDs, structured text understanding has been a challenging task. Most existing studies decoupled this problem into two sub-tasks: *entity labeling* and *entity linking*, which require an entire understanding of the context of documents at both token and segment levels. However, little work has been concerned with the solutions that efficiently extract the structured data from different levels. This paper proposes a unified framework named **StrucTeXT**, which is flexible and effective for handling both sub-tasks. Specifically, based on the transformer, we introduce a segment-token aligned encoder to deal with the entity labeling and entity linking tasks at different levels of granularity. Moreover, we design a novel pre-training strategy with three self-supervised tasks to learn a richer representation. StrucTeXT uses the existing Masked Visual Language Modeling task and the new Sentence Length Prediction and Paired Boxes Direction tasks to incorporate the multi-modal information across text, image, and layout. We evaluate our method for structured text understanding at segment-level and token-level and show it outperforms the state-of-the-art counterparts with significantly superior performance on the FUNSD, SROIE, and EPHOIE datasets.

*Equal contribution. This work is done when Yuxi Qian is an intern at Baidu Inc.

†Corresponding author.

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.

MM '21, October 20–24, 2021, Virtual Event, China

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8651-7/21/10...\$15.00

<https://doi.org/10.1145/3474085.3475345>

Jingtuo Liu, Errui Ding

Department of Computer Vision
Technology (VIS), Baidu Inc.
{liujingtuo,dingerrui}@baidu.com

CCS CONCEPTS

- Information systems → Document structure; Information extraction.

KEYWORDS

document understanding, document information extraction, pre-training

ACM Reference Format:

Yulin Li, Yuxi Qian, Yuechen Yu, Xiameng Qin, Chengquan Zhang, Yan Liu, Kun Yao, Junyu Han, and Jingtuo Liu, Errui Ding. 2021. StrucTeXT: Structured Text Understanding with Multi-Modal Transformers. In *Proceedings of the 29th ACM International Conference on Multimedia (MM '21)*, October 20–24, 2021, Virtual Event, China. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3474085.3475345>

1 INTRODUCTION

Understanding the structured document is a critical component of document intelligence that automatically explores the structured text information from the Visually Rich Documents (VRDs) such as forms, receipts, invoices, etc. Such task aims to extract the key information of text fields and the links among the semantic entities from VRDs, which named entity labeling and entity linking tasks [15] respectively. Structured text understanding has attracted increasing attention in both academia and industry. In reality, it plays a crucial role in developing digital transformation processes in office automation, accounting systems, and electronically archived. It offers businesses significant time savings on processing the million of forms and invoices every day.

Typical structure extraction methods rely on preliminary Optical Character Recognition (OCR) engines [19, 34, 39, 40, 47, 49] to understand the semantics of documents. As shown in Figure 1, the contents in a document can be located as several text segments (pink dotted boxes) by text detectors. The entity fields are presented in three forms: partial characters, an individual segment, and multiple segment lines. Traditional methods for entity labeling often formulated the task as a sequential labeling problem. In this setup,

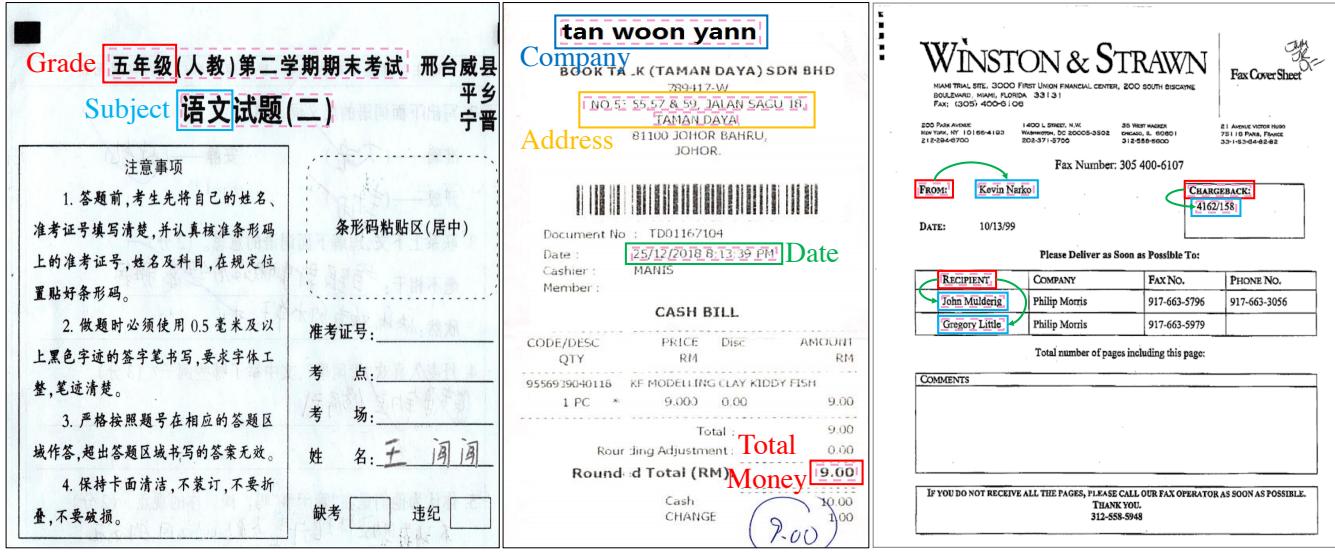


Figure 1: Examples of VRDs and their key extraction information. The dotted boxes are the text regions and the solid ones are the semantic entity regions. (a) The entity extraction in token-level characters. (b) The entity extraction in segment-level text lines. (c) The relationship extraction with key-value pairs at segment-level.

the text segments are serialized as a linear sequence with a pre-defined order. Then a **Named Entity Recognition (NER)** [17, 23] model is utilized to label each token such as word or character with an IOB (Inside, Outside, Beginning) tag. However, the capability is limited as the manner that is performed at token-level. As shown examples of Figure 1b and 1c, VRDs are usually organized in a number of text segments. The segment-level textual content presents richer geometric and semantic information, which is vital for structured text understanding. Several methods [1, 8, 14, 41] focus on a segment-level representation. On the contrary, they cannot cope with the entity composed of characters as shown in Figure 1a. Therefore, a comprehensive technique of structure extraction at both segment-level and token-level is worth considering.

Nowadays, accurate understanding of the structured text from VRDs remains a challenge. The key to success is the full use of multiple modal features from document images. Early solutions solve the entity tasks by only operating on plain texts, resulting in a semantic ambiguity. Noticing the rich visual information contained in VRDs, several methods [6, 16, 26, 32] exploit 2D layout information to provides complementation for textual content. Besides, for further improvement, mainstream researches [2, 21, 24, 30, 38, 48, 50] usually employ a shallow fusion of text, image, and layout to capture contextual dependencies. Recently, several pre-training models [28, 45, 46] have been proposed for joint learning the deep fusion of cross-modality on large-scale data and outperform counterparts on document understanding. Although these pre-training models consider all modalities of documents, they focus on the contribution related to the text side with less elaborate visual features.

To address the above limitations, in this paper, we propose a uniform framework named **StrucTexT** that incorporates the features from different levels and modalities to effectively improves

the understanding of various document structures. Inspired by recent developments in vision-language transformers [22, 35], we introduce a transformer encoder to learn cross-modal knowledge from both images of segments and tokens of words. In addition, we construct an extra segment ID embedding to associate visual and textual features at different granularity. Meanwhile, we attach a 2D position embedding to involve the layout clue. After that, a Hadamard product works on the encoded features between different levels and modalities for advanced feature fusion. Hence, StrucTexT can support segment-level and token-level tasks of structured text understanding in a single framework. Figure 2 shows the architecture of our proposed method.

To promote the representation capacity of multi-modality, we further introduce three self-supervised tasks for pre-training learning of text, image, and layout. Specifically, following the work of LayoutLM [45], the Masked Visual Language Modeling (MVLM) task is utilized to extract contextual information. In addition, we present two tasks named Sentence Length Prediction (SLP) and Paired Boxes Direction (PBD). SLP task predicts the segment length for enhancing the internal semantics of an entity candidate. PBD task is training to identify the relative direction within a sampled segment pair, which helps our framework discover the geometric structure topology. The three self-supervised tasks make full use of both textual and visual features of the documents. An unsupervised pre-training strategy with above all tasks is applied at first to get an enhanced feature encoder.

Major contributions of this paper are summarized as follows:

- (1) In this paper, we present a novel framework named StrucTexT to tackle the tasks of structured text understanding

- with a unified solution. It efficiently extracts semantic features from different levels and modalities to handle the **entity labeling and entity linking tasks**.
- (2) We introduce improved multi-task pre-training strategies to extract the multi-modal information from VRDs by self-supervised learning. In addition to the MVLM task that benefits the textual context, we proposed two new pre-training tasks of SLP and PBD to take advantage of image and layout features. We adopt the three tasks during the pre-training stage for a richer representation.
 - (3) Extensive experiments on real-world benchmarks show the superior performance of StrucTexT compared with the state-of-the-art methods. **Additional ablation studies demonstrate the effectiveness of our pre-training strategies.**

2 RELATED WORK

Structured Text Understanding The task of structured text understanding is to retrieve structured data from VRDs automatically. **It requires the model to extract the semantic structure of textual content robustly and effectively**, assigning the major purpose into two parts [15]: entity labeling and entity linking. **Generally speaking, the entity labeling task is to find named entities. The entity linking task is to extract the semantic relationships as key-value pairs between entities.** Most existing methods [6, 7, 13, 38, 45, 46, 48, 50] design a NER framework to perform entity labeling as a sequence labeling task at token-level. **However, traditional NER models organize text in one dimension depending on the reading order and are unsuitable for VRDs with complex layouts.** Recent studies [29, 38, 42, 45, 46, 48, 50] have realized the significance of segment-level features and incorporate a segment embedding to attach extra higher semantics. Although those methods, such as PICK [48] and TRIE [50], construct contextual features involving the segment clues, they revert to token-level labeling with NER-based schemes. Several works [1, 2, 14, 41] design their methods at segment-level to solve the tasks of entity labeling and entity linking. Cheng et al. [2] utilizes an attention-based network to explore one-shot learning for the text field labeling. DocStruct [41] predicts the key-value relations between the extracted text segments to establish a hierarchical document structure. With the graph-based paradigm, Carbonell et al. [1] and Hwang et al. [14] tackle the entity labeling and entity linking tasks simultaneously. However, they don't consider the situation where a text segment includes more than one category, which is difficult to identify the entity in token granularity.

In summary, the methods mentioned above can only handle one granularity representation. To this end, we propose a unified framework to support both token-level and segment-level structured extraction for VRDs. Our model is flexible to any granularity modeling tasks for structured text understanding.

Multi-Modal Feature Representations One of the most important modules of structured information extraction is to understand multi-modal semantic features. Previous works [3, 5, 7, 13, 17, 27, 31] usually adopt language models to extract entities from the plain text. These NLP-based approaches typically operate on text sequences and do not incorporate visual and layout information. Later studies [6, 16, 32] firstly tend to explore layout information

to aid entity extraction from VRDs. Post-OCR [13] reconstructs the text sequences based on their bounding boxes. VS2 [32] leverages the heterogeneous layout to perform the extraction in visual logical blocks. A range of other methods [6, 16, 51] represent a document as a 2D grid with text tokens to obtain the contextual embedding. After that, some researchers realize the necessity of multi-modal fusion and develop performance by integrating visual and layout information. GraphIE [29], PICK [48] and Liu et al. [21] design a graph-based decoder to improve the semantics of context information. Hwang et al. [14] and Wang et al. [41] leverage the relative coordinates and explore the link of each key-value pair. These methods only use simple early fusion strategies, such as addition or concatenation, without considering the semantic gap of different modalities. Recently, pre-training models [7, 22, 35, 36] show a strong feature representation using large-scale unlabeled training samples. Inspired by this, several works [28, 45, 46] combine pre-training techniques to improve multi-modal features. Pramanik et al. [28] introduces a multi-task learning-based framework to yield a generic document representation. LayoutLMv2 [46] uses 11 million scanned documents to obtain a pre-trained model, which shows the state-of-the-art performance in several downstream tasks of document understanding. However, these pre-training strategies mainly focus on the expressiveness of language but underuse the structured information from images. Hence, we propose a self-supervised pre-training strategy to better explore the potentials information from text, image, and layout. Compared with LayoutLMv2, the new strategy supports more useful features with less training data.

3 APPROACH

Figure 2 shows the overall illustration of StrucTexT. Given an input image with preconditioned OCR results, such as bounding boxes and content of text segments. We leverage various information from text, image, and layout aspects by a feature embedding stage. And then, the multi-modal embeddings are fed into the pre-trained transformer network to obtain rich semantic features. The transformer network has accomplished the cross-modality fusion by establishing interactions between the different modality inputs. At last, the Structured Text Understanding module receives the encoded features and carries out entity recognition for entity labeling and relation extraction for entity linking.

3.1 Multi-Modal Feature Embedding

Given a document image I with n text segments, we perform open source OCR algorithms [33, 52] to obtain the i -th segment region with the top-left and bottom-right bounding box $b_i = (x_0, y_0, x_1, y_1)$ and its corresponding text sentence $t_i = \{c_1^i, c_2^i, \dots, c_{l_i}^i\}$, where c is a word or character and l_i is the length of t_i .

Layout Embedding For every segment or word, we use the encoded bounding boxes as their layout information

$$L = \text{Emb}_l(x_0, y_0, x_1, y_1, w, h) \quad (1)$$

where Emb_l is a layout embedding layer and w, h is the shape of bounding box b . It is worth mentioning that we estimate the bounding box of a word by its belonging text segment in consideration of some OCR results without word-level information.

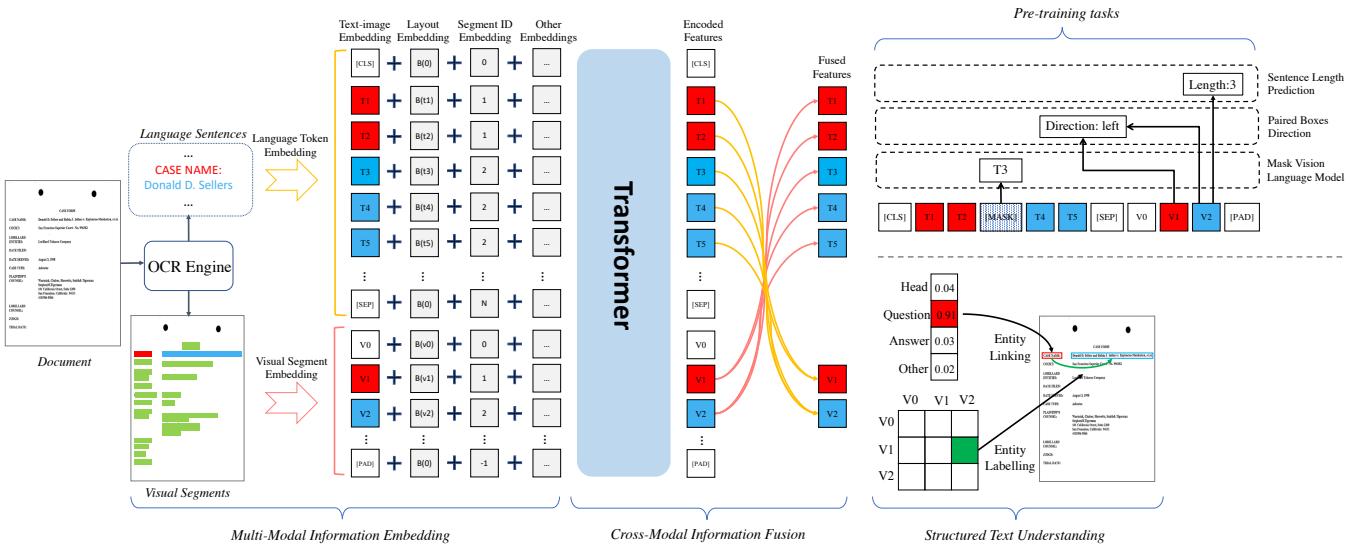


Figure 2: An overall illustration of the model framework and the inform extraction tasks for StrucText.

Language Token Embedding Following the common practice [7], we utilize the WordPiece [43] to tokenize text sentences. After that, all of text sentences are gathered as a sequence S by sorting the text segments from the top-left to bottom-right. Intuitively, a pair of special tags [CLS] and [SEP] are added at the beginning and end of the sequence, as $t_0 = \{[CLS]\}$, $t_{n+1} = \{[SEP]\}$. Thus, we can define the language sequence S as follows

$$\begin{aligned} S &= \{t_0, t_1, \dots, t_n, t_{n+1}\} \\ &= \{[CLS], c_1^1, \dots, c_{l_1}^1, \dots, c_1^n, \dots, c_{l_n}^n, [SEP]\} \end{aligned} \quad (2)$$

Then, we sum the embedded feature of S and layout **embedding L** to obtain the language embedding **T**

$$T = \text{Emb}_t(S) + L \quad (3)$$

where Emb_t is a text embedding layer.

Visual Segment Embedding In the model architecture, we use ResNet50 [44] with FPN [20] as the image feature extractor to generate feature maps of I . Then, the image feature of each text segment is extracted from the CNN maps by ROIAlign [10] according to b . The visual segment embedding V is computed as

$$V = \text{Emb}_v(\text{ROIAlign}(\text{CNN}(I), b)) + L \quad (4)$$

where Emb_v is the visual embedding layer. Furthermore, the entire feature maps of image I is embedded as V_0 to introduce the global information into image features.

Segment ID Embedding Compared with the vision-language tasks based on wild pictures, understanding the structured document requires higher semantics to identify the ambiguous entities. Thus, we propose a segment ID embedding S^{id} to allocate a unique number to each text segment with its image and text features, which makes an explicit alignment of cross-modality clues.

Other Embeddings In addition, we add two other embeddings [22, 35] into the input. The position embedding P^{id} encodes the indexes

from 1 to maximum sequence length, and the segment embedding M^{id} denotes the modality for each feature. All above embeddings have the same dimensions. In the end, the input of our model is represented as the combination of the embeddings.

$$\text{Input} = \text{Concat}(T, V) + S^{id} + P^{id} + M^{id} \quad (5)$$

Moreover, we append several [PAD] tags to fill the short input sequence to a fixed length. An empty bounding box with zeros is assigned to the special [CLS], [SEP], and [PAD] tags.

3.2 Multi-Modal Feature Enhance Module

StrucText collects multi-modal information from visual segments, text sentences, and position layouts to produce an embedding sequence. We support an image-text alignment between different granularities by leveraging the segment IDs mentioned above. At this stage, we perform a transformer network to encode the embedding sequence to establish deep fusion between modalities and granularities. Crucially, three self-supervised tasks encode the input features during the pre-training stage to learn task-agnostic joint representations. The details are introduced as follows, where patterns of all self-supervised tasks are as shown in Figure 3.

Task 1: Masked Visual Language Modeling

The MVLM task promotes capturing a contextual representation on the language side. Following the pattern of *masked multi-modal modeling* in ViLBERT [22], we select 15% tokens from the language sequences, mask 80% among them with a [MASK] token, replace 10% of that with random tokens, and keep 10% tokens unchanged. Then, the model is required to reconstruct the corresponding tokens. Rather than following the image region mask in ViLBERT, we retain all other information and encourage the model to hunt for the cross-modality clues at all possible.

Task 2: Segment Length Prediction

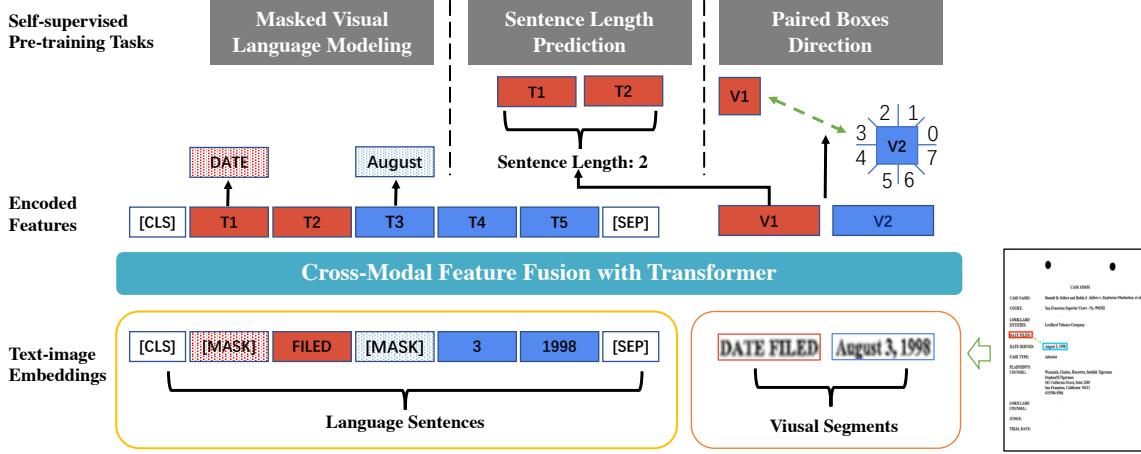


Figure 3: The illustration of cross-modal information fusion. Three self-supervised tasks MVLM, SLP, and PBD introduced in Section 3.1 are employed simultaneously on the visual and language embeddings in the pre-training stage.

Besides the MLVM, we introduce a new self-supervised task called Sequence Length Prediction (SLP) to excavate fine-grained semantic information on the image side. The SLP task asks the model to recognize the length of the text segment from each visual feature. In this way, we force the encoder to learn from the image feature, more importantly, the language sequence knowledge via the same segment ID. We argue that this information flow could accelerate the deep cross-modal fusion among textual, visual, and layout information.

Moreover, to avoid the disturbance of sub-words produced by WordPiece [43], we only count each first sub-word for keeping the same length between language sequences and image segments. Therefore, we build an extra alignment between two granularities, which is simple but effective.

Task 3: Paired Box Direction

Furthermore, our third self-supervised task, Paired Box Direction (PBD), is designed to exploit global layout information. The PBD task aims at learning a comprehensive geometric topology for document structures by predicting the pairwise spatial relationships of text segments. First of all, we divide the field of 360 degrees into eight identical buckets. Secondly, we compute the angle θ_{ij} between the text segment i and j and label it with one of the buckets. Next, we carry out the subtraction between two visual features on the image side and take the result $\Delta\hat{V}_{ij}$ as the input of the PBD

$$\Delta\hat{V}_{ij} = \hat{V}_i - \hat{V}_j \quad (6)$$

where we use the $\hat{\cdot}$ symbol to denote the features after transformer encoding. \hat{V}_i and \hat{V}_j express the visual features for i -th segment and j -th segment.

Finally, we define PBD as a classification task to estimate the relative positional direction with $\Delta\hat{V}_{ij}$.

3.3 Structural Text Understanding

Cross-granularity Labeling Module The cross-granularity labeling module supports both token-level entity labeling and segment-level entity labeling tasks. In this module, tokens with the same segment ID on the language side are aggregated into a segment-level textual feature through the arithmetic average

$$\hat{T}_i = \text{mean}(\hat{t}_i) = (\hat{c}_1 + \hat{c}_2 + \dots + \hat{c}_{l_i})/l_i \quad (7)$$

where \hat{t}_i means the features of i -th text sentence, \hat{c} is the feature of token, l_i is the sentence length. After that, a bilinear pooling layer is utilized to compute a Hadamard product to fuse the textual segment feature T_i and the visual segment feature V_i .

$$X_i = V_i * T_i \quad (8)$$

Finally, we apply a fully connected layer on the cross-modal features X_i to predict an entity label for segment i with the Cross-Entropy loss.

Segment Relationship Extraction Module The segment relationship extraction module is proposed for entity linking. Documents usually represent their structure as a set of hierarchical relations, such as key-value pair or table parsing. Inspired by DocStruct [41], we use an asymmetric parameter matrix M to extract the relationship from segments i to j in probability form

$$P_{i \rightarrow j} = \sigma(X_j M X_i^T) \quad (9)$$

where $P_{i \rightarrow j}$ is the probability of whether i links to j . M is a parameter matrix and σ is the sigmoid function.

We notice that most of the segment pairs in a document are not related. To alleviate the data sparsity and balance the number of related and unrelated pairs, we learn from the Negative Sampling method [25] and build a sampling set with non-fixed size. Our sampling set consists of the same number of positive and negative samples.

However, we also find the training process is unstable only using the above sampling strategy. To utmost handle the imbalanced

distribution of entity linking, we combine the Margin Ranking Loss and Binary Cross-Entropy to supervise the training simultaneously. Thus, the linking loss can be formulated as

$$\text{Loss} = \text{LossBCE} + \text{LossRank} \quad (10)$$

where the $\text{Loss}_{\text{Rank}}$ is computed as following

$$\text{LossRank}(P_i, P_j, y) = \max(0, -y * (P_i - P_j) + \text{Margin}), \quad (11)$$

Note that y equals 1 if (P_i, P_j) is the positive-negative samples pair or equals 0 for the negative-positive samples pair.

4 EXPERIMENTS

4.1 Datasets

In this section, we firstly introduce several datasets that are used for pre-training and evaluating our StrucTextT. The extensive experiments are conduct on three benchmark databases: FUNSD [15], SROIE [12], EPHOIE [38]. Moreover, we perform ablation studies to analyze the effects of each proposed component.

DOCBANK [18] contains 500K document pages (400K for training, 50K for validation and 50K for testing) for document layout analysis. We pre-train StrucTextT on the dataset.

RVL-CDIP [9] consists of 400,000 grayscale images in 16 classes, with 25,000 images per class. There are 320,000 training images, 40,000 validation images and 40,000 test images. We adopt RVL-CDIP for pre-training our model.

FUNSD [15] consists of 199 real, fully annotated, scanned form images. The dataset is split into 149 training samples and 50 testing samples. Three sub-tasks (word grouping, semantic entity labeling, and entity linking) are proposed to identify the semantic entity (i.e., questions, answers, headers, and other) and entity linking present in the form. We use the official OCR annotation and focus on the latter two tasks in this paper.

SROIE [12] is composed of 626 receipts for training and 347 receipts for testing. Every receipt contains four predefined values: company, date, address, and total. The segment-level text bounding box and the corresponding transcript are provided according to the annotations. We use the official OCR annotations and evaluate our model for receipt information extraction.

EPHOIE [38] is collected from actual Chinese examination papers with the diversity of text types and layout distribution. The 1,494 samples are divided into a training set with 1,183 images and a testing set with 311 images, respectively. Every character in the document is annotated with a label from ten predefined categories. The token-level entity labeling task is evaluated in this dataset.

4.2 Implementation

Following the typical pre-training and fine-tuning strategies, we train the model end-to-end. Across all pre-training and downstream tasks, we rescale the images and pad them to the size of 512×512 . The input sequence is set as a maximum length of 512.

4.2.1 Pre-training. We extract both token-level text features and segment-level visual features based on a unified joint model by the encoder. Due to time and computational resource restrictions, we

| Method | Prec. | Recall | F1 | Params. |
|-----------------------|--------------|--------------|--------------------------------|---------|
| LayoutLM_BASE [45] | 94.38 | 94.38 | 94.38 | 113M |
| LayoutLM_LARGE [45] | 95.24 | 95.24 | 95.24 | 343M |
| PICK [48] | 96.79 | 95.46 | 96.12 | - |
| VIES [38] | - | - | 96.12 | - |
| TRIE [50] | - | - | 96.18 | - |
| LayoutLMv2_BASE [46] | 96.25 | 96.25 | 96.25 | 200M |
| MatchVIE [37] | - | - | 96.57 | - |
| LayoutLMv2_LARGE [46] | 96.61 | 96.61 | 96.61 | 426M |
| Ours | 95.84 | 98.52 | 96.88 (± 0.15) | 107M |

Table 1: Model Performance (entity labeling) comparison on the SROIE dataset.

| Method | Prec. | Recall | F1 | Params. |
|-----------------------|--------------|--------------|-------------------------|---------|
| Carbonell et al. [1] | - | - | 64.0 | - |
| SPADE [14] | - | - | 70.5 | - |
| LayoutLM_BASE [45] | 75.97 | 81.55 | 78.66 | 113M |
| LayoutLM_LARGE [45] | 75.96 | 82.19 | 78.95 | 343M |
| MatchVIE [37] | - | - | 81.33 | - |
| LayoutLMv2_BASE [46] | 80.29 | 85.39 | 82.76 | 200M |
| LayoutLMv2_LARGE [46] | 83.24 | 85.19 | 84.20 | 426M |
| Ours | 85.68 | 80.97 | 83.09 (± 0.09) | 107M |

Table 2: Model Performance (entity labeling) comparison on the FUNSD dataset, We ignore entities belonging to the *other* category and use the mean performance of three classes (header, question, and answer) as our final results.

| Method | Reconstruction | | Detection | | | F1 |
|----------------------|----------------|-------------|--------------|--------------|--------------|-------------|
| | mAP | mRank | Hit@1 | Hit@2 | Hit@5 | |
| FUNSD [15] | - | - | - | - | - | 4.0 |
| Carbonell et al. [1] | - | - | - | - | - | 39.0 |
| LayoutLM* [45] | 47.61 | 7.11 | 32.43 | 45.56 | 66.41 | - |
| DocStruct [41] | 71.77 | 2.89 | 58.19 | 76.27 | 88.94 | - |
| SPADE [14] | - | - | - | - | - | 41.7 |
| Ours | 78.36 | 3.38 | 67.67 | 84.33 | 95.33 | 44.1 |

Table 3: Model Performance (entity linking) comparison on the FUNSD dataset. (LayoutLM* is implemented by [41])

choose the 12-layer transformer encoder with 768 hidden size and 12 attention heads. We initialize the transformer and the text embedding layer from the ERNIE_{BASE} [36]. The weights of the ResNet50 network is initialized using the ResNet_vd [11] pre-trained on the ImageNet [4]. The rest of the parameters are randomly initialized.

To obtain the pre-training OCR results, we apply the PaddleOCR¹ to extract the text segment in both DOCBANK and RVL-CDIP datasets. All three self-supervised tasks are trained for classification with the Cross-Entropy loss. The Adamax optimizer is used with an initial 5×10^{-5} learning rate for a warm-up. And then, we

¹<https://github.com/PaddlePaddle/PaddleOCR>

| Method | Entities | | | | | | | | | | |
|---------------|--------------|------------|--------------|--------------|--------------|--------------|--------------|----------------|--------------|--------------|--------------|
| | Subject | Test Time | Name | School | Exam Number | Seat Number | Class | Student Number | Grade | Score | Mean |
| TRIE [50] | 98.79 | 100 | 99.46 | 99.64 | 88.64 | 85.92 | 97.94 | 84.32 | 97.02 | 80.39 | 93.21 |
| VIES [38] | 99.39 | 100 | 99.67 | 99.28 | 91.81 | 88.73 | 99.29 | 89.47 | 98.35 | 86.27 | 95.23 |
| MatchVIE [37] | 99.78 | 100 | 99.88 | 98.57 | 94.21 | 93.48 | 99.54 | 92.44 | 98.35 | 92.45 | 96.87 |
| Ours | 99.25 | 100 | 99.47 | 99.83 | 97.98 | 95.43 | 98.29 | 97.33 | 99.25 | 93.73 | 97.95 |

Table 4: Model Performance (token-level entity labeling) comparison on the EPHOIE dataset.

keep 1×10^{-4} for 2~5 epochs and set a linear decay schedule in the rest of epochs. We pre-train our architecture in DOCBANK [18] and RVL-CDIP [9] dataset for 10 epochs with a batch size of 64 on 4 NVIDIA Tesla V100 32GB GPUs.

4.2.2 Fine-tuning. We fine-tune our StrucText on three information extraction tasks: entity labeling and entity linking at segment-level and entity labeling at token-level. For the segment-based entity labeling task, we aggregate token features of the text sentence via the arithmetic average and get the segment-level features by multiplying visual features and textual features. At last, a softmax layer is followed by the features for segment-level category prediction. The entity-level F1-score is used as the evaluation metric.

The entity linking task takes two segment features as input to obtain a pairwise relationship matrix. Then we pass the non-diagonal elements in the relationship matrix through a sigmoid layer to predict the binary classification of each relationship.

For the token-based entity labeling task, the output visual feature is expanded as a vector with the same length of its text sentence. Next, the extended visual features are element-wise multiplied with the corresponding textual features to obtain token-level features to predict the category of each token through a softmax layer.

We fine-tune our pre-trained model at all downstream tasks for 50 epochs with a batch size of 4 and a learning rate from 1×10^{-4} to 1×10^{-5} . We use the precision, recall, and F1-score as evaluation metrics for entity labeling. Following DocStruct [41] and SPADE [14], the performance of entity linking is estimated with Hit@1, Hit@2, Hit@5, mAp, mRank, and F1-score.

4.3 Comparison with the State-of-the-Arts

We evaluate our proposed StrucText in three publish benchmarks for both the entity labeling and entity linking tasks.

Segment-level Entity Labeling The comparison results are shown in Table 1. We can observe that StrucText exhibits a superior performance over baseline methods [37, 38, 45, 46, 48, 50] on SROIE. Specifically, our method obtains a precision of 95.84% and a Recall of 98.52% in SROIE, which surpass that of LayoutLMv2_LARGE [46] by 0.27% in F1-score.

As shown in Table 2, our method achieves competitive F1-score of 83.09% in FUNSD. Although LayoutLMv2_LARGE beats our F1-score by 1%, it is worth noting that LayoutLMv2_LARGE using a larger transformer consisting of 24 layers and 16 heads that contains 426M parameters. Further, our model using only 90K documents for pre-training compared to LayoutLMv2_LARGE which uses 11M documents. On the contrary, our model shows a better performance than LayoutLMv2_BASE under the same architecture settings. It

fully proves the superiority of our proposed framework. Moreover, to verify the performance gain is statistically stable and significant, we repeat our experiments five times to eliminate the random fluctuations and attach the standard deviation below the F1-score.

Segment-level Entity Linking As shown in Table 3, we compare our method with several state-of-the-arts on FUNSD for entity linking. The baseline method [15] gets the remarkably worst result with a simple binary classification for pairwise entities. The SPADE [14] shows a tremendous gain by leading a Graph into their model. Compared with the SPADE, our method has a 2.4% improvement and achieves 44.1% F1-score. Besides, we evaluate the performance in the mAp, mRank, and Hit metrics mentioned in DocStruct [41]. Our method attains 78.36% mAP, 79.19% Hit@1, 84.33% Hit@2, and 95.33% Hit@5 which outperforms DocStruct and obtains a competitive performance at the 3.38 mRank score.

Token-level Entity Labeling We further perform StrucText on EPHOIE. It is noticed that the entities annotated in this dataset are character-based. Therefore, we apply our StrucText to calculate the entities with the token-level prediction. Table 4 illustrates the overall performance of the EPHOIE dataset. Our StrucText contributes to a top-tier performance with 97.95%.

4.4 Ablation Study

We study the impact of individual components in StrucText and conduct ablation studies on the FUNSD and SROIE datasets.

| Dataset | Pre-training Tasks | Prec. | Recall | F1 |
|---------|--------------------|--------------|--------------|--------------|
| FUNSD | MVLM | 76.41 | 79.36 | 77.71 |
| | MVLM+PBD | 81.22 | 79.46 | 80.29 |
| | MVLM+SLP | 87.45 | 78.69 | 82.12 |
| | MVLM+PBD+SLP | 85.68 | 80.97 | 83.09 |
| SROIE | MVLM | 95.25 | 97.89 | 96.54 |
| | MVLM+PBD | 95.32 | 98.25 | 96.75 |
| | MVLM+SLP | 95.30 | 98.16 | 96.70 |
| | MVLM+PBD+SLP | 95.84 | 98.52 | 96.88 |

Table 5: Ablation Studies with entity labeling on the FUNSD and SROIE datasets.

Self-supervised Tasks in Pre-training In this study, we evaluate the impact of different pre-training tasks. As shown in Table 5, we can observe that the PBD and SLP tasks make better use of visual information. Specifically, compared with the model only trained

| Dataset | Modality | Prec. | Recall | F1 |
|---------|-------------------|--------------|--------------|--------------|
| FUNSD | Visual | 76.93 | 77.51 | 77.22 |
| | Language | 81.73 | 79.38 | 80.49 |
| | Visual + Language | 85.68 | 80.97 | 83.09 |
| SROIE | Visual | 90.14 | 92.11 | 91.11 |
| | Language | 94.54 | 97.91 | 96.18 |
| | Viusal + Language | 95.84 | 98.52 | 96.88 |

Table 6: Ablation studies with visual-only and language-only entity labeling on the FUNSD and SROIE datasets.

| Dataset | Granularity | Prec. | Recall | F1 |
|---------|-------------|--------------|--------------|--------------|
| FUNSD | Token | 81.20 | 82.10 | 81.59 |
| | Segment | 85.68 | 80.97 | 83.09 |
| SROIE | Token | 92.77 | 98.81 | 95.62 |
| | Segment | 95.84 | 98.52 | 96.88 |

Table 7: Ablation studies with the comparison of token-level and segment-level entity labeling on the FUNSD and SROIE datasets.

| | |
|--|--|
| Total (Excluding GST): 99.06 GST payable (6%): 5.94 Total (Inclusive of GST): 105.00 TOTAL: 105.00 | Item(s) Total (MYR) : 13.80 GST @ 6% : 0.83 Net Total (MYR) : 14.63 Rounding Adj. : other 0.02 Net Total Rounded (MYR) : (total) 14.63 |
| Other Total (incl GST) total(other) 109.05 Total Rounded 109.05 Cash Tendered 150.00 Change 40.95 VISA 50.90 | Sub-total] 50.92 Total Sales Incl GST total(other) 50.92 Rounding Adj. -0.02 Total After Adj Incl GST 50.90 VISA 50.90 |

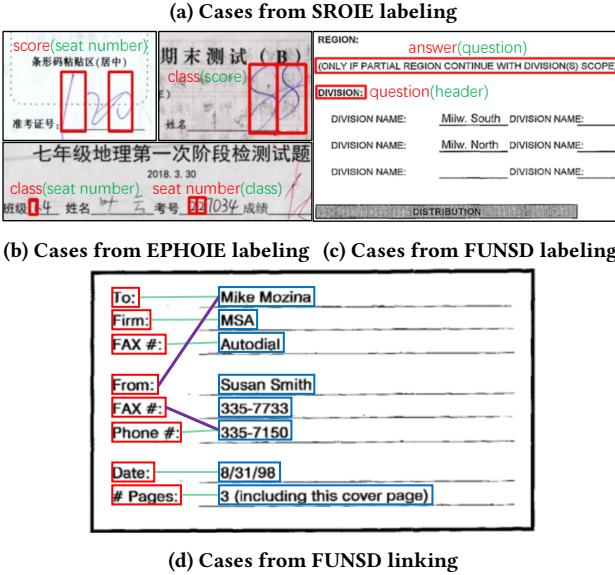


Figure 4: Visualization of badcases of StrucText. For the entity labeling task in (a), (b), and (c), the wrong cases are shown as the red boxes (the correct results are hidden) and their nearby text represents the predictions and ground-truths in red and green color, respectively. For the entity linking task in (d), the green/purple lines indicate the correct/error predicted linkings.

with the MVLM task, MVLM+PBD gains nearly 3% improvement in FUNSD and 0.2% improvement in SROIE. Meanwhile, the results turn out that the SLP task also improves the performance dramatically. Furthermore, the incorporation of the three tasks obtains the optimal performance compared with other combinations. It means that both the SLP and PBD tasks contribute to richer semantics and potential relationships between cross-modality.

Multi-Modal Features Profits As shown in Table 6, then we perform experiments in verifying the benefits of features in multiple modalities. The textual features perform better than visual ones, which we attribute more semantics to the textual content of documents. Moreover, combining visual and textual features can achieve higher performance with a notable gap, indicating complementarity between language and visual information. The results show that the multi-modal feature fusion in our model can get a richer semantic representation.

Granularity of Feature Fusion We also study the representations with different granularities towards the performance. In detail, we complete the experiments of entity labeling on SROIE and FUNSD in token-level supervision. As shown in Table 7, overall, the segment-based results perform better than token-based ones, which proves our opinion that the effectiveness of text segment.

4.5 Error Analysis

Although our work has achieved outstanding performance, we also observe some badcases of the proposed method. This section presents an error analysis on the qualitative results in SROIE, FUNSD and EPHOIE, respectively. In Figure 4a, our model makes the mistakes of wrong answers to the total in SROIE. We attribute the errors to the similar semantics of textual contents and the close distance of locations. In addition, as shown in Figure 4b, our model is confused by the similar style of digits, which demonstrates the relatively low performance in the numeral entities in EPHOIE, such as exam number, seat number, student number, and score in 4. However, these entities can certainly be distinguished by their keywords. To this end, a goal-directed information aggregation of key-value pair is well worth considering, and we would study it in future works. As shown in Figure 4c, the model is failed in recognizing the header and the question in FUNSD. We analyze the model is overfitting the layout position of training data. Then, according to Figure 4d, some links are assigned incorrectly. We attribute the errors to ambiguous semantics of relationships.

5 CONCLUSION

In this paper, we further explore improving the understanding of document text structure by using a unified framework. Our framework shows superior performance on three real-world benchmark datasets after applying novel pre-training strategies for the multi-modal and multi-granularity feature fusion. Moreover, we evaluate the influence of different modalities and granularities on the ability of entity extraction, thus providing a new perspective to study the problem of structured text understanding.

REFERENCES

- [1] Manuel Carbonell, Pau Riba, Mauricio Villegas, Alicia Fornés, and Josep Lladós. 2020. Named Entity Recognition and Relation Extraction with Graph Neural Networks in Semi Structured Documents. In *ICPR*. IEEE, 9622–9627.
- [2] Mengli Cheng, Minghui Qiu, Xing Shi, Jun Huang, and Wei Lin. 2020. One-shot Text Field labeling using Attention and Belief Propagation for Structure Information Extraction. In *ACM Multimedia*. ACM, 340–348.
- [3] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc Le, and Ruslan Salakhutdinov. 2019. Transformer-XL: Attentive Language Models beyond a Fixed-Length Context. In *ACL*. ACL, 2978–2988.
- [4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. 2009. ImageNet: A large-scale hierarchical image database. In *CVPR*. IEEE, 248–255.
- [5] Andreas Dengel and Bernt Klein. 2002. smartFIX: A Requirements-Driven System for Document Analysis and Understanding. In *DAS*. Springer, 433–444.
- [6] Timo I Denk and Christian Reisswig. 2019. Bertgrid: Contextualized embedding for 2d document representation and understanding. *arXiv preprint arXiv:1909.04948* (2019).
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [8] He Guo, Xiameng Qin, Jiaming Liu, Junyu Han, Jingtuo Liu, and Errui Ding. 2019. EATEN: Entity-Aware Attention for Single Shot Visual Text Extraction. In *ICDAR*. IEEE, 254–259.
- [9] Adam W. Harley, Alex Ufkes, and Konstantinos G. Derpanis. 2015. Evaluation of deep convolutional nets for document image classification and retrieval. In *ICDAR*. IEEE, 991–995.
- [10] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross B. Girshick. 2017. Mask R-CNN. In *ICCV*, 2961–2969.
- [11] Tong He, Zhi Zhang, Hang Zhang, Zhongyue Zhang, Junyuan Xie, and Mu Li. 2019. Bag of Tricks for Image Classification with Convolutional Neural Networks. In *CVPR*. IEEE, 558–567.
- [12] Zheng Huang, Kai Chen, Jianhua He, Xiang Bai, Dimosthenis Karatzas, Shijian Lu, and CV Jawahar. 2019. Icdar2019 competition on scanned receipt ocr and information extraction. In *ICDAR*. IEEE, 1516–1520.
- [13] Wonseok Hwang, Seonghyeon Kim, Minjoon Seo, Jinyeong Yim, Seunghyun Park, Sungrae Park, Junyeop Lee, Bado Lee, and Hwalsuk Lee. 2019. Post-OCR parsing: building simple and robust parser via BIO tagging. In *NeurIPS Workshop*.
- [14] Wonseok Hwang, Jinyeong Yim, Seunghyun Park, Sohee Yang, and Minjoon Seo. 2021. Spatial Dependency Parsing for Semi-Structured Document Information Extraction. In *ACL-IJCNLP*.
- [15] Guillaume Jaume, Hazim Kemal Ekenel, and Jean-Philippe Thiran. 2019. FUNSD: A Dataset for Form Understanding in Noisy Scanned Documents. In *ICDAR Workshop*. IEEE, 1–6.
- [16] Anoop R. Katti, Christian Reisswig, Cordula Guder, Sebastian Brarda, Steffen Bickel, Johannes Höhne, and Jean Baptiste Faddoul. 2018. Chargrid: Towards Understanding 2D Documents. In *EMNLP*. ACL, 4459–4469.
- [17] Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. 2016. Neural Architectures for Named Entity Recognition. In *ACL*. ACL, 260–270.
- [18] Minghao Li, Yiheng Xu, Lei Cui, Shaohan Huang, Furu Wei, Zhoujun Li, and Ming Zhou. 2020. DocBank: A Benchmark Dataset for Document Layout Analysis. In *COLING*. ICCL, 949–960.
- [19] Minghui Liao, Guan Pang, Jing Huang, Tal Hassner, and Xiang Bai. 2020. Mask textspotter v3: Segmentation proposal network for robust scene text spotting. In *ECCV*. Springer, 706–722.
- [20] Tsung-Yi Lin, Piotr Dollár, Ross B. Girshick, Kaiming He, Bharath Hariharan, and Serge J. Belongie. 2017. Feature Pyramid Networks for Object Detection. In *CVPR*. IEEE, 936–944.
- [21] Xiaojing Liu, Feiyu Gao, Qiong Zhang, and Huasha Zhao. 2019. Graph Convolution for Multimodal Information Extraction from Visually Rich Documents. In *NAACL-HLT*. ACL, 32–39.
- [22] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks. In *NeurIPS*. 13–23.
- [23] Xuezhe Ma and Eduard Hovy. 2016. End-to-end Sequence Labeling via Bidirectional LSTM-CNNs-CRF. In *ACL*. ACL, 1064–1074.
- [24] Bodhisattwa Prasad Majumder, Navneet Potti, Sandeep Tata, James Bradley Wendt, Qi Zhao, and Marc Najork. 2020. Representation Learning for Information Extraction from Form-like Documents. In *ACL*. ACL, 6495–6504.
- [25] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *NIPS*. 3111–3119.
- [26] Rasmus Berg Palm, Florian Laws, and Ole Winther. 2019. Attend, Copy, Parse End-to-end Information Extraction from Documents. In *ICDAR*. IEEE, 329–336.
- [27] Rasmus Berg Palm, Ole Winther, and Florian Laws. 2017. CloudScan - A Configuration-Free Invoice Analysis System Using Recurrent Neural Networks. In *ICDAR*. IEEE, 406–413.
- [28] Subhojeet Pramanik, Shashank Mujumdar, and Hima Patel. 2020. Towards a Multi-modal, Multi-task Learning based Pre-training Framework for Document Representation Learning. *CoRR abs/2009.14457* (2020). arXiv:2009.14457
- [29] Yujie Qian, Enrico Santus, Zhijing Jin, Jiang Guo, and Regina Barzilay. 2019. GraphIE: A Graph-Based Framework for Information Extraction. In *ACL*. ACL, 751–761.
- [30] Clément Sage, Alex Aussem, Véronique Eglin, Haytham Elghazel, and Jérémie Espinas. 2020. End-to-End Extraction of Structured Information from Business Documents with Pointer-Generator Networks. In *SPNLP*. ACL, 43–52.
- [31] Clément Sage, Alexandre Aussem, Haytham Elghazel, Véronique Eglin, and Jérémie Espinas. 2019. Recurrent Neural Network Approach for Table Field Extraction in Business Documents. In *ICDAR*. IEEE, 1308–1313.
- [32] Ritesh Sarkhel and Arnab Nandi. 2019. Visual Segmentation for Information Extraction from Heterogeneous Visually Rich Documents. In *SIGMOD*. ACM, 247–262.
- [33] Baoguang Shi, Xiang Bai, and Cong Yao. 2017. An End-to-End Trainable Neural Network for Image-Based Sequence Recognition and Its Application to Scene Text Recognition. *TPAMI* 39, 11 (2017), 2298–2304.
- [34] Bolan Su and Shijian Lu. 2014. Accurate Scene Text Recognition Based on Recurrent Neural Network. In *ACCV*. Springer, 35–48.
- [35] Weiji Su, Xizhou Zhu, Yu Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. 2020. VL-BERT: Pre-training of Generic Visual-Linguistic Representations. In *ICLR*.
- [36] Yi Sun, Shuhuan Wang, Yu-Kun Li, Shikun Feng, Hao Tian, Hua Wu, and Haifeng Wang. 2020. ERNIE 2.0: A Continual Pre-Training Framework for Language Understanding. In *AAAI AAAI*, 8968–8975.
- [37] Guozhi Tang, Lele Xie, Lianwen Jin, Jiapeng Wang, Jingdong Chen, Zhen Xu, Qianying Wang, Yaqiang Wu, and Hui Li. 2021. MatchVIE: Exploiting Match Relevancy between Entities for Visual Information Extraction. In *IJCAI ijcai.org*.
- [38] Jiapeng Wang, Chongyu Liu, Lianwen Jin, Guozhi Tang, Jiaxin Zhang, Shuaítiao Zhang, Qianying Wang, Yaqiang Wu, and Mingxiang Cai. 2021. Towards Robust Visual Information Extraction in Real World: New Dataset and Novel Solution. In *AAAI AAAI*, 2738–2745.
- [39] Pengfei Wang, Chengquan Zhang, Fei Qi, Zuming Huang, Mengyi En, Junyu Han, Jingtuo Liu, Errui Ding, and Guangming Shi. 2019. A Single-Shot Arbitrarily-Shaped Text Detector based on Context Attended Multi-Task Learning. In *ACM Multimedia*. ACM, 1277–1285.
- [40] Pengfei Wang, Chengquan Zhang, Fei Qi, Shanshan Liu, Xiaoqiang Zhang, Pengyuan Lyu, Junyu Han, Jingtuo Liu, Errui Ding, and Guangming Shi. 2021. PGNet: Real-time Arbitrarily-Shaped Text Spotting with Point Gathering Network. In *AAAI AAAI*, 2782–2790.
- [41] Zilong Wang, Mingjie Zhan, Xuebo Liu, and Ding Liang. 2020. DocStruct: A Multimodal Method to Extract Hierarchy Structure in Document for General Form Understanding. In *EMNLP*. ACL, 898–908.
- [42] Mengxi Wei, Yifan He, and Qiong Zhang. 2020. Robust Layout-aware IE for Visually Rich Documents with Pre-trained Language Models. In *SIGIR*. ACM, 2367–2376.
- [43] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144* (2016).
- [44] Saining Xie, Ross B. Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. 2017. Aggregated Residual Transformations for Deep Neural Networks. In *CVPR*. IEEE, 5987–5995.
- [45] 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 *KDD*. ACM, 1192–1200.
- [46] Yang Xu, Yiheng Xu, Tengchao Lv, Lei Cui, Furu Wei, Guoxin Wang, Yijuan Lu, Dinei Florencio, Cha Zhang, Wanxiang Che, et al. 2020. LayoutLMv2: Multi-modal pre-training for visually-rich document understanding. *arXiv preprint arXiv:2012.14740* (2020).
- [47] Deli Yu, Xuan Li, Chengquan Zhang, Tao Liu, Junyu Han, Jingtuo Liu, and Errui Ding. 2020. Towards accurate scene text recognition with semantic reasoning networks. In *CVPR*. 12113–12122.
- [48] Wenwen Yu, Ning Lu, Xianbiao Qi, Ping Gong, and Rong Xiao. 2020. PICK: Processing Key Information Extraction from Documents using Improved Graph Learning-Convolutional Networks. In *ICPR*. IEEE, 4363–4370.
- [49] Chengquan Zhang, Borong Liang, Zuming Huang, Mengyi En, Junyu Han, Errui Ding, and Xinghao Ding. 2019. Look More Than Once: An Accurate Detector for Text of Arbitrary Shapes. In *CVPR*. IEEE, 10552–10561.
- [50] Peng Zhang, Yunlu Xu, Zhanzhan Cheng, Shiliang Pu, Jing Lu, Liang Qiao, Yi Niu, and Fei Wu. 2020. TRIE: End-to-End Text Reading and Information Extraction for Document Understanding. In *ACM Multimedia*. ACM, 1413–1422.
- [51] Xiaohui Zhao, Endi Niu, Zhuo Wu, and Xiaoguang Wang. 2019. Cutie: Learning to understand documents with convolutional universal text information extractor. *arXiv preprint arXiv:1903.12363* (2019).
- [52] Xinyu Zhou, Cong Yao, He Wen, Yuzhi Wang, Shuchang Zhou, Weiran He, and Jiajun Liang. 2017. EAST: An Efficient and Accurate Scene Text Detector. In *CVPR*. IEEE, 2642–2651.