

LayoutLM: Pre-training of Text and Layout for Document Image Understanding

Yiheng Xu*

charlesyihengxu@gmail.com
Harbin Institute of Technology

Minghao Li*

liminghao1630@buaa.edu.cn
Beihang University

Lei Cui

lecu@microsoft.com
Microsoft Research Asia

Shaohan Huang

shaohanh@microsoft.com
Microsoft Research Asia

Furu Wei

fuwei@microsoft.com
Microsoft Research Asia

Ming Zhou

mingzhou@microsoft.com
Microsoft Research Asia

ABSTRACT

Pre-training techniques have been verified successfully in a variety of NLP tasks in recent years. Despite the widespread use of pre-training models for NLP applications, they almost **exclusively focus on text-level manipulation, while neglecting layout and style information that is vital for document image understanding**. In this paper, we propose the **LayoutLM** to jointly model interactions between text and layout information across scanned document images, which is beneficial for a great number of real-world document image understanding tasks such as information extraction from scanned documents. Furthermore, we also leverage image features to incorporate words' visual information into LayoutLM. To the best of our knowledge, this is the first time that text and layout are jointly learned in a single framework for document-level pre-training. It achieves new state-of-the-art results in several downstream tasks, including form understanding (from 70.72 to 79.27), receipt understanding (from 94.02 to 95.24) and document image classification (from 93.07 to 94.42). The code and pre-trained LayoutLM models are publicly available at <https://aka.ms/layoutlm>.

CCS CONCEPTS

• **Information systems** → **Business intelligence**; • **Computing methodologies** → **Information extraction**; **Transfer learning**; • **Applied computing** → **Document analysis**.

KEYWORDS

LayoutLM; pre-trained models; document image understanding

ACM Reference Format:

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 Conference on Knowledge Discovery and Data Mining (KDD '20)*, August 23–27, 2020, Virtual Event, CA, USA. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3394486.3403172>

*Equal contributions during internship at Microsoft Research Asia.

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.

KDD '20, August 23–27, 2020, Virtual Event, CA, USA

© 2020 Association for Computing Machinery.

ACM ISBN 978-1-4503-7998-4/20/08...\$15.00

<https://doi.org/10.1145/3394486.3403172>

1 INTRODUCTION

Document AI, or Document Intelligence¹, is a relatively new research topic that refers techniques for automatically reading, understanding, and analyzing business documents. Business documents are files that provide details related to a company's internal and external transactions, which are shown in Figure 1. They may be digital-born, occurring as electronic files, or they may be in scanned form that comes from written or printed on paper. Some common examples of business documents include purchase orders, financial reports, business emails, sales agreements, vendor contracts, letters, invoices, receipts, resumes, and many others. Business documents are critical to a company's efficiency and productivity. The exact format of a business document may vary, but the information is usually presented in natural language and can be organized in a variety of ways from plain text, multi-column layouts, and a wide variety of tables/forms/figures. Understanding business documents is a very challenging task due to the diversity of layouts and formats, poor quality of scanned document images as well as the complexity of template structures.

Nowadays, many companies extract data from business documents through manual efforts that are time-consuming and expensive, meanwhile requiring manual customization or configuration. Rules and workflows for each type of document often need to be hard-coded and updated with changes to the specific format or when dealing with multiple formats. To address these problems, document AI models and algorithms are designed to automatically classify, extract, and structuralize information from business documents, accelerating automated document processing workflows. Contemporary approaches for document AI are usually built upon deep neural networks from a computer vision perspective or a natural language processing perspective, or a combination of them. Early attempts usually focused on detecting and analyzing certain parts of a document, such as tabular areas. [7] were the first to propose a table detection method for PDF documents based on Convolutional Neural Networks (CNN). After that, [21, 24, 29] also leveraged more advanced Faster R-CNN model [19] or Mask R-CNN model [9] to further improve the accuracy of document layout analysis. In addition, [28] presented an end-to-end, multimodal, fully convolutional network for extracting semantic structures from document images, taking advantage of text embeddings from pre-trained NLP models. More recently, [15] introduced a Graph Convolutional Networks (GCN) based model to combine textual and visual information for

¹<https://sites.google.com/view/di2019>

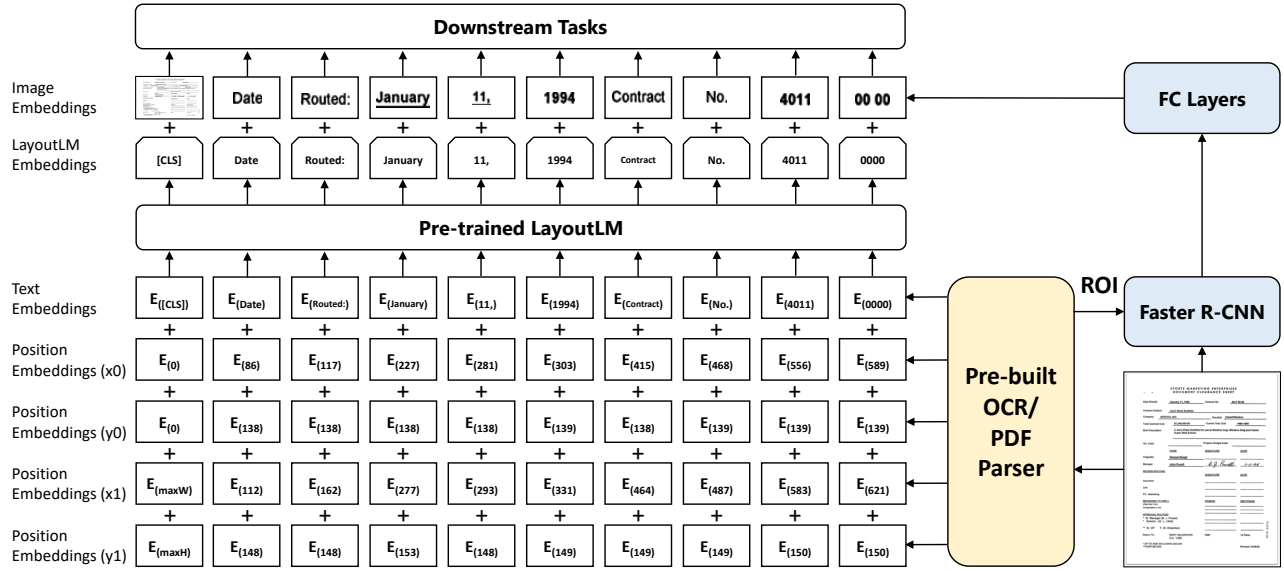


Figure 2: An example of LayoutLM, where 2-D layout and image embeddings are integrated into the original BERT architecture. The LayoutLM embeddings and image embeddings from Faster R-CNN work together for downstream tasks.

2.1 The BERT Model

The BERT model is an attention-based bidirectional language modeling approach. It has been verified that the BERT model shows effective knowledge transfer from the self-supervised task with large-scale training data. The architecture of BERT is basically a multi-layer bidirectional Transformer encoder. It accepts a sequence of tokens and stacks multiple layers to produce final representations. In detail, given a set of tokens processed using WordPiece, the input embeddings are computed by summing the corresponding word embeddings, position embeddings, and segment embeddings. Then, these input embeddings are passed through a multi-layer bidirectional Transformer that can generate contextualized representations with an adaptive attention mechanism.

There are two steps in the BERT framework: pre-training and fine-tuning. During the pre-training, the model uses two objectives to learn the language representation: Masked Language Modeling (MLM) and Next Sentence Prediction (NSP), where MLM randomly masks some input tokens and the objective is to recover these masked tokens, and NSP is a binary classification task taking a pair of sentences as inputs and classifying whether they are two consecutive sentences. In the fine-tuning, task-specific datasets are used to update all parameters in an end-to-end way. The BERT model has been successfully applied in a set of NLP tasks.

2.2 The LayoutLM Model

Although BERT-like models become the state-of-the-art techniques on several challenging NLP tasks, they usually leverage text information only for any kind of inputs. When it comes to visually rich documents, there is much more information that can be encoded into the pre-trained model. Therefore, we propose to utilize the visually rich information from document layouts and align them with the input texts. Basically, there are two types of features which

substantially improve the language representation in a visually rich document, which are:

Document Layout Information. It is evident that the relative positions of words in a document contribute a lot to the semantic representation. Taking form understanding as an example, given a key in a form (e.g., “Passport ID:”), its corresponding value is much more likely on its right or below instead of on the left or above. Therefore, we can embed these relative positions information as 2-D position representation. Based on the self-attention mechanism within the Transformer, embedding 2-D position features into the language representation will better align the layout information with the semantic representation.

Visual Information. Compared with the text information, the visual information is another significantly important feature in document representations. Typically, documents contain some visual signals to show the importance and priority of document segments. The visual information can be represented by image features and effectively utilized in document representations. For document-level visual features, the whole image can indicate the document layout, which is an essential feature for document image classification. For word-level visual features, styles such as bold, underline, and italic, are also significant hints for the sequence labeling tasks. Therefore, we believe that combining the image features with traditional text representations can bring richer semantic representations to documents.

2.3 Model Architecture

To take advantage of existing pre-trained models and adapt to document image understanding tasks, we use the BERT architecture as the backbone and add two new input embeddings: a 2-D position embedding and an image embedding.

2-D Position Embedding. Unlike the position embedding that models the word position in a sequence, 2-D position embedding aims to model the relative spatial position in a document. To represent the spatial position of elements in scanned document images, we consider a document page as a coordinate system with the top-left origin. In this setting, the bounding box can be precisely defined by (x_0, y_0, x_1, y_1) , where (x_0, y_0) corresponds to the position of the upper left in the bounding box, and (x_1, y_1) represents the position of the lower right. We add four position embedding layers with two embedding tables, where the embedding layers representing the same dimension share the same embedding table. This means that we look up the position embedding of x_0 and x_1 in the embedding table X and lookup y_0 and y_1 in table Y .

Image Embedding. To utilize the image feature of a document and align the image feature with the text, we add an image embedding layer to represent image features in language representation. In more detail, with the bounding box of each word from OCR results, we split the image into several pieces, and they have a one-to-one correspondence with the words. We generate the image region features with these pieces of images from the Faster R-CNN [19] model as the token image embeddings. For the [CLS] token, we also use the Faster R-CNN model to produce embeddings using the whole scanned document image as the Region of Interest (ROI) to benefit the downstream tasks which need the representation of the [CLS] token.

2.4 Pre-training LayoutLM

Task #1: Masked Visual-Language Model. Inspired by the masked language model, we propose the Masked Visual-language Model (MVLM) to learn the language representation with the clues of 2-D position embeddings and text embeddings. During the pre-training, we randomly mask some of the input tokens but keep the corresponding 2-D position embeddings, and then the model is trained to predict the masked tokens given the contexts. In this way, the **LayoutLM model not only understands the language contexts but also utilizes the corresponding 2-D position information, thereby bridging the gap between the visual and language modalities.**

Task #2: Multi-label Document Classification. For document image understanding, many tasks require the model to generate high-quality document-level representations. As the IIT-CDIP Test Collection includes multiple tags for each document image, we also use a Multi-label Document Classification (MDC) loss during the pre-training phase. Given a set of scanned documents, we use the document tags to supervise the pre-training process so that the model can cluster the knowledge from different domains and generate better document-level representation. Since the MDC loss needs the label for each document image that may not exist for larger datasets, it is optional during the pre-training and may not be used for pre-training larger models in the future. We will compare the performance of MVLM and MVLM+MDC in Section 3.

2.5 Fine-tuning LayoutLM

The pre-trained LayoutLM model is fine-tuned on three document image understanding tasks, including a form understanding task, a

receipt understanding task as well as a document image classification task. For the form and receipt understanding tasks, LayoutLM predicts {B, I, E, S, O} tags for each token and uses sequential labeling to detect each type of entity in the dataset. For the document image classification task, LayoutLM predicts the class labels using the representation of the [CLS] token.

3 EXPERIMENTS

3.1 Pre-training Dataset

The performance of pre-trained models is largely determined by the scale and quality of datasets. Therefore, we need a large-scale scanned document image dataset to pre-train the LayoutLM model. Our model is pre-trained on the IIT-CDIP Test Collection 1.0, which contains more than 6 million documents, with more than 11 million scanned document images. Moreover, each document has its corresponding text and metadata stored in XML files. The text is the content produced by applying OCR to document images. The metadata describes the properties of the document, such as the unique identity and document labels. Although the metadata contains erroneous and inconsistent tags, the scanned document images in this large-scale dataset are perfectly suitable for pre-training our model.

3.2 Fine-tuning Dataset

The FUNSD Dataset. We evaluate our approach on the FUNSD dataset for form understanding in noisy scanned documents. This dataset includes 199 real, fully annotated, scanned forms with 9,707 semantic entities and 31,485 words. These forms are organized as a list of semantic entities that are interlinked. Each semantic entity comprises a unique identifier, a label (i.e., question, answer, header, or other), a bounding box, a list of links with other entities, and a list of words. The dataset is split into 149 training samples and 50 testing samples. We adopt the word-level F1 score as the evaluation metric.

The SROIE Dataset. We also evaluate our model on the SROIE dataset for receipt information extraction (Task 3). The dataset contains 626 receipts for training and 347 receipts for testing. Each receipt is organized as a list of text lines with bounding boxes. Each receipt is labeled with four types of entities which are {company, date, address, total}. The evaluation metric is the exact match of the entity recognition results in the F1 score.

The RVL-CDIP Dataset. The RVL-CDIP dataset 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. The images are resized, so their largest dimension does not exceed 1,000 pixels. The 16 classes include {letter, form, email, handwritten, advertisement, scientific report, scientific publication, specification, file folder, news article, budget, invoice, presentation, questionnaire, resume, memo}. The evaluation metric is the overall classification accuracy.

3.3 Document Pre-processing

To utilize the layout information of each document, we need to obtain the location of each token. However, the pre-training dataset (IIT-CDIP Test Collection) only contains pure texts while missing

their corresponding bounding boxes. In this case, we re-process the scanned document images to obtain the necessary layout information. Like the original pre-processing in IIT-CDIP Test Collection, we similarly process the dataset by applying OCR to document images. **The difference is that we obtain both the recognized words and their corresponding locations in the document image. Thanks to Tesseract⁶, an open-source OCR engine, we can easily obtain the recognition as well as the 2-D positions.** We store the OCR results in hOCR format, a standard specification format which clearly defines the OCR results of one single document image using a hierarchical representation.

3.4 Model Pre-training

We initialize the weight of LayoutLM model with the pre-trained BERT base model. Specifically, our BASE model has the same architecture: a 12-layer Transformer with 768 hidden sizes, and 12 attention heads, which contains about 113M parameters. Therefore, we use the BERT base model to initialize all modules in our model except the 2-D position embedding layer. For the LARGE setting, our model has a 24-layer Transformer with 1,024 hidden sizes and 16 attention heads, which is initialized by the pre-trained BERT LARGE model and contains about 343M parameters. Following [4], we select 15% of the input tokens for prediction. We replace these masked tokens with the [MASK] token 80% of the time, a random token 10% of the time, and an unchanged token 10% of the time. Then, the model predicts the corresponding token with the cross-entropy loss.

In addition, we also add the 2-D position embedding layers with four embedding representations (x_0, y_0, x_1, y_1), where (x_0, y_0) corresponds to the position of the upper left in the bounding box, and (x_1, y_1) represents the position of the lower right. Considering that the document layout may vary in different page size, we scale the actual coordinate to a “virtual” coordinate: the actual coordinate is scaled to have a value from 0 to 1,000. Furthermore, we also use the ResNet-101 model as the backbone network in the Faster R-CNN model, which is pre-trained on the Visual Genome dataset [12].

We train our model on 8 NVIDIA Tesla V100 32GB GPUs with a total batch size of 80. The Adam optimizer is used with an initial learning rate of $5e-5$ and a linear decay learning rate schedule. The BASE model takes 80 hours to finish one epoch on 11M documents, while the LARGE model takes nearly 170 hours to finish one epoch.

3.5 Task-specific Fine-tuning

We evaluate the LayoutLM model on three document image understanding tasks: **Form Understanding**, **Receipt Understanding**, and **Document Image Classification**. We follow the typical fine-tuning strategy and update all parameters in an end-to-end way on task-specific datasets.

Form Understanding. This task requires extracting and structuring the textual content of forms. It aims to extract key-value pairs from the scanned form images. In more detail, this task includes two sub-tasks: semantic labeling and semantic linking. Semantic labeling is the task of aggregating words as semantic entities and assigning pre-defined labels to them. Semantic linking is the task

of predicting the relations between semantic entities. In this work, we focus on the semantic labeling task, while semantic linking is out of the scope. To fine-tune LayoutLM on this task, we treat semantic labeling as a sequence labeling problem. We pass the final representation into a linear layer followed by a softmax layer to predict the label of each token. The model is trained for 100 epochs with a batch size of 16 and a learning rate of $5e-5$.

Receipt Understanding. This task requires filling several pre-defined semantic slots according to the scanned receipt images. For instance, given a set of receipts, we need to fill specific slots (i.e., company, address, date, and total). Different from the form understanding task that requires labeling all matched entities and key-value pairs, the number of semantic slots is fixed with pre-defined keys. Therefore, the model only needs to predict the corresponding values using the sequence labeling method.

Document Image Classification. Given a visually rich document, this task aims to predict the corresponding category for each document image. Distinct from the existing image-based approaches, our model includes not only image representations but also text and layout information using the multimodal architecture in LayoutLM. Therefore, our model can combine the text, layout, and image information in a more effective way. To fine-tune our model on this task, we concatenate the output from the LayoutLM model and the whole image embedding, followed by a softmax layer for category prediction. We fine-tune the model for 30 epochs with a batch size of 40 and a learning rate of $2e-5$.

3.6 Results

Form Understanding. We evaluate the form understanding task on the FUNSD dataset. The experiment results are shown in Table 1. We compare the LayoutLM model with two SOTA pre-trained NLP models: BERT and RoBERTa [16]. The BERT BASE model achieves 0.603 and while the LARGE model achieves 0.656 in F1. Compared to BERT, the RoBERTa performs much better on this dataset as it is trained using larger data with more epochs. Due to the time limitation, we present 4 settings for LayoutLM, which are 500K document pages with 6 epochs, 1M with 6 epochs, 2M with 6 epochs as well as 11M with 2 epochs. It is observed that the LayoutLM model substantially outperforms existing SOTA pre-training baselines. With the BASE architecture, the LayoutLM model with 11M training data achieves 0.7866 in F1, which is much higher than BERT and RoBERTa with the similar size of parameters. In addition, we also add the MDC loss in the pre-training step and it does bring substantial improvements on the FUNSD dataset. Finally, the LayoutLM model achieves the best performance of 0.7927 when using the text, layout, and image information at the same time.

In addition, we also evaluate the LayoutLM model with different data and epochs on the FUNSD dataset, which is shown in Table 2. For different data settings, we can see that the overall accuracy is monotonically increased as more epochs are trained during the pre-training step. Furthermore, the accuracy is also improved as more data is fed into the LayoutLM model. As the FUNSD dataset contains only 149 images for fine-tuning, the results confirm that the pre-training of text and layout is effective for scanned document understanding especially with low resource settings.

⁶<https://github.com/tesseract-ocr/tesseract>

Modality	Model	Precision	Recall	F1	#Parameters
Text only	BERT _{BASE}	0.5469	0.671	0.6026	110M
	RoBERTa _{BASE}	0.6349	0.6975	0.6648	125M
	BERT _{LARGE}	0.6113	0.7085	0.6563	340M
	RoBERTa _{LARGE}	0.678	0.7391	0.7072	355M
Text + Layout MVLM	LayoutLM _{BASE} (500K, 6 epochs)	0.665	0.7355	0.6985	113M
	LayoutLM _{BASE} (1M, 6 epochs)	0.6909	0.7735	0.7299	113M
	LayoutLM _{BASE} (2M, 6 epochs)	0.7377	0.782	0.7592	113M
	LayoutLM _{BASE} (11M, 2 epochs)	0.7597	0.8155	0.7866	113M
Text + Layout MVLM+MDC	LayoutLM _{BASE} (1M, 6 epochs)	0.7076	0.7695	0.7372	113M
	LayoutLM _{BASE} (11M, 1 epoch)	0.7194	0.7780	0.7475	113M
Text + Layout MVLM	LayoutLM _{LARGE} (1M, 6 epochs)	0.7171	0.805	0.7585	343M
	LayoutLM _{LARGE} (11M, 1 epoch)	0.7536	0.806	0.7789	343M
Text + Layout + Image MVLM	LayoutLM _{BASE} (1M, 6 epochs)	0.7101	0.7815	0.7441	160M
	LayoutLM _{BASE} (11M, 2 epochs)	0.7677	0.8195	0.7927	160M

Table 1: Model accuracy (Precision, Recall, F1) on the FUNSD dataset

# Pre-training Data	# Pre-training Epochs	Precision	Recall	F1
500K	1 epoch	0.5779	0.6955	0.6313
	2 epochs	0.6217	0.705	0.6607
	3 epochs	0.6304	0.718	0.6713
	4 epochs	0.6383	0.7175	0.6756
	5 epochs	0.6568	0.734	0.6933
	6 epochs	0.665	0.7355	0.6985
1M	1 epoch	0.6156	0.7005	0.6552
	2 epochs	0.6545	0.737	0.6933
	3 epochs	0.6794	0.762	0.7184
	4 epochs	0.6812	0.766	0.7211
	5 epochs	0.6863	0.7625	0.7224
	6 epochs	0.6909	0.7735	0.7299
2M	1 epoch	0.6599	0.7355	0.6957
	2 epochs	0.6938	0.759	0.7249
	3 epochs	0.6915	0.7655	0.7266
	4 epochs	0.7081	0.781	0.7427
	5 epochs	0.7228	0.7875	0.7538
	6 epochs	0.7377	0.782	0.7592
11M	1 epoch	0.7464	0.7815	0.7636
	2 epochs	0.7597	0.8155	0.7866

Table 2: LayoutLM_{BASE} (Text + Layout, MVLM) accuracy with different data and epochs on the FUNSD dataset

Furthermore, we compare different initialization methods for the LayoutLM model including from scratch, BERT and RoBERTa. The results in Table 3 show that the LayoutLM_{BASE} model initialized with RoBERTa_{BASE} outperforms BERT_{BASE} by 2.1 points in F1. For the LARGE setting, the LayoutLM_{LARGE} model initialized with RoBERTa_{LARGE} further improve 1.3 points over the BERT_{LARGE} model. We will pre-train more models with RoBERTa as the initialization in the future, especially for the LARGE settings.

Receipt Understanding. We evaluate the receipt understanding task using the SROIE dataset. The results are shown in Table 4. As we only test the performance of the Key Information Extraction task in SROIE, we would like to eliminate the effect of incorrect OCR results. Therefore, we pre-process the training data by using the ground truth OCR and run a set of experiments using the baseline models (BERT & RoBERTa) as well as the LayoutLM model. The results show that the LayoutLM_{LARGE} model trained with 11M

Initialization	Model	Precision	Recall	F1
SCRATCH	LayoutLM _{BASE} (1M, 6 epochs)	0.5630	0.6728	0.6130
BERT _{BASE}	LayoutLM _{BASE} (1M, 6 epochs)	0.6909	0.7735	0.7299
RoBERTa _{BASE}	LayoutLM _{BASE} (1M, 6 epochs)	0.7173	0.7888	0.7514
SCRATCH	LayoutLM _{LARGE} (11M, 1 epoch)	0.6845	0.7804	0.7293
BERT _{LARGE}	LayoutLM _{LARGE} (11M, 1 epoch)	0.7536	0.8060	0.7789
RoBERTa _{LARGE}	LayoutLM _{LARGE} (11M, 1 epoch)	0.7681	0.8188	0.7926

Table 3: Different initialization methods for BASE and LARGE (Text + Layout, MVLM)

document images achieve an F1 score of 0.9524, which is significantly better than the first place in the competition leaderboard. This result also verifies that the pre-trained LayoutLM not only performs well on the in-domain dataset (FUNSD) but also outperforms several strong baselines on the out-of-domain dataset like SROIE.

Document Image Classification. Finally, we evaluate the document image classification task using the RVL-CDIP dataset. Document images are different from other natural images as most of the content in document images are texts in a variety of styles and layouts. Traditionally, image-based classification models with pre-training perform much better than the text-based models, which is shown in Table 5. We can see that either BERT or RoBERTa underperforms the image-based approaches, illustrating that text information is not sufficient for this task, and it still needs layout and image features. We address this issue by using the LayoutLM model for this task. Results show that, even without the image features, LayoutLM still outperforms the single model of the image-based approaches. After integrating the image embeddings, the LayoutLM achieves the accuracy of 94.42%, which is significantly better than several SOTA baselines for document image classification. It is observed that our model performs best in the "email" category while performs worst in the "form" category. We will further investigate how to take advantage of both pre-trained LayoutLM and image models, as well as involve image information in the pre-training step for the LayoutLM model.

4 RELATED WORK

The research of Document Analysis and Recognition (DAR) dates to the early 1990s. The mainstream approaches can be divided into three categories: rule-based approaches, conventional machine learning approaches, and deep learning approaches.

4.1 Rule-based Approaches

The rule-based approaches [6, 13, 18, 23] contain two types of analysis methods: bottom-up and top-down. Bottom-up methods [5, 13, 23] usually detect the connected components of black pixels as the basic computational units in document images, and the document segmentation process is to combine them into higher-level structures through different heuristics and label them according to different structural features. Docstrum algorithm [18] is among the earliest successful bottom-up algorithms that are based on the connected component analysis. It groups connected components on a polar structure to derive the final segmentation. [23] use a special distance-metric between different components to construct

a physical page structure. They further reduced the time complexity by using heuristics and path compression algorithms.

The top-down methods often recursively split a page into columns, blocks, text lines, and tokens. [6] propose replacing the basic unit with the black pixels from all the pixels, and the method decomposed the document using the recursive the X-Y cut algorithm to establish an X-Y tree, which makes complex documents decompose more easily. Although these methods perform well on some documents, they require extensive human efforts to figure out better rules, while sometimes failing to generalize to documents from other sources. Therefore, it is inevitable to leverage machine learning approaches in the DAR research.

4.2 Machine Learning Approaches

With the development of conventional machine learning, statistical machine learning approaches [17, 22] have become the mainstream for document segmentation tasks during the past decade. [22] consider the layout information of a document as a parsing problem, and globally search the optimal parsing tree based on a grammar-based loss function. They utilize a machine learning approach to select features and train all parameters during the parsing process. Meanwhile, artificial neural networks [17] have been extensively applied to document analysis and recognition. Most efforts have been devoted to the recognition of isolated handwritten and printed characters with widely recognized successful results. In addition to the ANN model, SVM and GMM [27] have been used in document layout analysis tasks. For machine learning approaches, they are usually time-consuming to design manually crafted features and difficult to obtain a highly abstract semantic context. In addition, these methods usually relied on visual cues but ignored textual information.

4.3 Deep Learning Approaches

Recently, deep learning methods have become the mainstream and de facto standard for many machine learning problems. Theoretically, they can fit any arbitrary functions through the stacking of multi-layer neural networks and have been verified to be effective in many research areas. [28] treat the document semantic structure extraction task as a pixel-by-pixel classification problem. They propose a multimodal neural network that considers visual and textual information, while the limitation of this work is that they only used the network to assist heuristic algorithms to classify candidate bounding boxes rather than an end-to-end approach. [26] propose a lightweight model of document layout analysis for mobile and cloud services. The model uses one-dimensional information of images for

Modality	Model	Precision	Recall	F1	#Parameters
Text only	BERT _{BASE}	0.9099	0.9099	0.9099	110M
	RoBERTa _{BASE}	0.9107	0.9107	0.9107	125M
	BERT _{LARGE}	0.9200	0.9200	0.9200	340M
	RoBERTa _{LARGE}	0.9280	0.9280	0.9280	355M
Text + Layout MVLM	LayoutLM _{BASE} (500K, 6 epochs)	0.9388	0.9388	0.9388	113M
	LayoutLM _{BASE} (1M, 6 epochs)	0.9380	0.9380	0.9380	113M
	LayoutLM _{BASE} (2M, 6 epochs)	0.9431	0.9431	0.9431	113M
	LayoutLM _{BASE} (11M, 2 epochs)	0.9438	0.9438	0.9438	113M
Text + Layout MVLM+MDC	LayoutLM _{BASE} (1M, 6 epochs)	0.9402	0.9402	0.9402	113M
	LayoutLM _{BASE} (11M, 1 epoch)	0.9460	0.9460	0.9460	113M
Text + Layout MVLM	LayoutLM _{LARGE} (1M, 6 epochs)	0.9416	0.9416	0.9416	343M
	LayoutLM _{LARGE} (11M, 1 epoch)	0.9524	0.9524	0.9524	343M
Text + Layout + Image MVLM	LayoutLM _{BASE} (1M, 6 epochs)	0.9416	0.9416	0.9416	160M
	LayoutLM _{BASE} (11M, 2 epochs)	0.9467	0.9467	0.9467	160M
Baseline	Ranking 1 st in SROIE	0.9402	0.9402	0.9402	-

Table 4: Model accuracy (Precision, Recall, F1) on the SROIE dataset

Modality	Model	Accuracy	#Parameters
Text only	BERT _{BASE}	89.81%	110M
	RoBERTa _{BASE}	90.06%	125M
	BERT _{LARGE}	89.92%	340M
	RoBERTa _{LARGE}	90.11%	355M
Text + Layout MVLM	LayoutLM _{BASE} (500K, 6 epochs)	91.25%	113M
	LayoutLM _{BASE} (1M, 6 epochs)	91.48%	113M
	LayoutLM _{BASE} (2M, 6 epochs)	91.65%	113M
	LayoutLM _{BASE} (11M, 2 epochs)	91.78%	113M
Text + Layout MVLM+MDC	LayoutLM _{BASE} (1M, 6 epochs)	91.74%	113M
	LayoutLM _{BASE} (11M, 1 epoch)	91.78%	113M
Text + Layout MVLM	LayoutLM _{LARGE} (1M, 6 epochs)	91.88%	343M
	LayoutLM _{LARGE} (11M, 1 epoch)	91.90%	343M
Text + Layout + Image MVLM	LayoutLM _{BASE} (1M, 6 epochs)	94.31%	160M
	LayoutLM _{BASE} (11M, 2 epochs)	94.42%	160M
Baselines	VGG-16 [1]	90.97%	-
	Stacked CNN Single [2]	91.11%	-
	Stacked CNN Ensemble [2]	92.21%	-
	InceptionResNetV2 [25]	92.63%	-
	LadderNet [20]	92.77%	-
	Multimodal Single [3]	93.03%	-
	Multimodal Ensemble [3]	93.07%	-

Table 5: Classification accuracy on the RVL-CDIP dataset

inference and compares it with the model using two-dimensional information, achieving comparable accuracy in the experiments. [11] make use of a fully convolutional encoder-decoder network that predicts a segmentation mask and bounding boxes, and the model significantly outperforms approaches based on sequential text or document images. [24] incorporate contextual information into the

Faster R-CNN model that involves the inherently localized nature of article contents to improve region detection performance.

Existing deep learning approaches for DAR usually confront two limitations: (1) The models often rely on limited labeled data while leaving a large amount of unlabeled data unused. (2) Current deep learning models usually leverage pre-trained CV models or

NLP models, but do not consider the joint pre-training of text and layout. LayoutLM addresses these two limitations and achieves much better performance compared with the previous baselines.

5 CONCLUSION AND FUTURE WORK

We present LayoutLM, a simple yet effective pre-training technique with text and layout information in a single framework. Based on the Transformer architecture as the backbone, LayoutLM takes advantage of multimodal inputs, including token embeddings, layout embeddings, and image embeddings. Meanwhile, the model can be easily trained in a self-supervised way based on large scale unlabeled scanned document images. We evaluate the LayoutLM model on three tasks: form understanding, receipt understanding, and scanned document image classification. Experiments show that LayoutLM substantially outperforms several SOTA pre-trained models in these tasks.

For future research, we will investigate pre-training models with more data and more computation resources. In addition, we will also train LayoutLM using the LARGE architecture with text and layout, as well as involving image embeddings in the pre-training step. Furthermore, we will explore new network architectures and other self-supervised training objectives that may further unlock the power of LayoutLM.

REFERENCES

- [1] Muhammad Zeshan Afzal, Andreas Kölsch, Sheraz Ahmed, and Marcus Liwicki. 2017. Cutting the Error by Half: Investigation of Very Deep CNN and Advanced Training Strategies for Document Image Classification. *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)* 01 (2017), 883–888.
- [2] Arindam Das, Saikat Roy, and Ujjwal Bhattacharya. 2018. Document Image Classification with Intra-Domain Transfer Learning and Stacked Generalization of Deep Convolutional Neural Networks. *2018 24th International Conference on Pattern Recognition (ICPR)* (2018), 3180–3185.
- [3] Tyler Dauphinee, Nikunj Patel, and Mohammad Mehdi Rashidi. 2019. Modular Multimodal Architecture for Document Classification. *ArXiv abs/1912.04376* (2019).
- [4] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186. <https://doi.org/10.18653/v1/N19-1423>
- [5] Jaekyu Ha, Robert M Haralick, and Ihsin T Phillips. 1995. Document page decomposition by the bounding-box project. In *Proceedings of 3rd International Conference on Document Analysis and Recognition*, Vol. 2. IEEE, 1119–1122.
- [6] Jaekyu Ha, Robert M Haralick, and Ihsin T Phillips. 1995. Recursive XY cut using bounding boxes of connected components. In *Proceedings of 3rd International Conference on Document Analysis and Recognition*, Vol. 2. IEEE, 952–955.
- [7] Leipeng Hao, Liangcai Gao, Xiaohan Yi, and Zhi Tang. 2016. A Table Detection Method for PDF Documents Based on Convolutional Neural Networks. *2016 12th IAPR Workshop on Document Analysis Systems (DAS)* (2016), 287–292.
- [8] Adam W. Harley, Alex Ufkes, and Konstantinos G. Derpanis. 2015. Evaluation of deep convolutional nets for document image classification and retrieval. *2015 13th International Conference on Document Analysis and Recognition (ICDAR)* (2015), 991–995.
- [9] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross B. Girshick. 2017. Mask R-CNN. *CoRR abs/1703.06870* (2017). [arXiv:1703.06870](http://arxiv.org/abs/1703.06870)
- [10] Guillaume Jaume, Hazim Kemal Ekenel, and Jean-Philippe Thiran. 2019. FUNSD: A Dataset for Form Understanding in Noisy Scanned Documents. *2019 International Conference on Document Analysis and Recognition Workshops (ICDARW)* 2 (2019), 1–6.
- [11] Anoop R Katti, Christian Reisswig, Cordula Guder, Sebastian Brarda, Steffen Bickel, Johannes Höhne, and Jean Baptiste Faddoul. 2018. Chargrid: Towards Understanding 2D Documents. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Brussels, Belgium, 4459–4469. <https://doi.org/10.18653/v1/D18-1476>
- [12] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, Michael Bernstein, and Li Fei-Fei. 2016. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations. <https://arxiv.org/abs/1602.07332>
- [13] Frank Lebourgeois, Z Bublinski, and H Emptoz. 1992. A fast and efficient method for extracting text paragraphs and graphics from unconstrained documents. In *Proceedings., 11th IAPR International Conference on Pattern Recognition. Vol. II. Conference B: Pattern Recognition Methodology and Systems*. IEEE, 272–276.
- [14] D. Lewis, G. Agam, S. Argamon, O. Frieder, D. Grossman, and J. Heard. 2006. Building a Test Collection for Complex Document Information Processing. In *Proceedings of the 29th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval* (Seattle, Washington, USA) (SIGIR '06). ACM, New York, NY, USA, 665–666. <https://doi.org/10.1145/1148170.1148307>
- [15] Xiaojing Liu, Feiyu Gao, Qiong Zhang, and Huasha Zhao. 2019. Graph Convolution for Multimodal Information Extraction from Visually Rich Documents. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Industry Papers)*. Association for Computational Linguistics, Minneapolis, Minnesota, 32–39. <https://doi.org/10.18653/v1/N19-2005>
- [16] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke S. Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *ArXiv abs/1907.11692* (2019).
- [17] S. Marinai, M. Gori, and G. Soda. 2005. Artificial neural networks for document analysis and recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 27, 1 (Jan 2005), 23–35. <https://doi.org/10.1109/TPAMI.2005.4>
- [18] L. O’Gorman. 1993. The document spectrum for page layout analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 15, 11 (Nov 1993), 1162–1173. <https://doi.org/10.1109/34.244677>
- [19] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. 2015. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 39 (2015), 1137–1149.
- [20] Ritesh Sarkhel and Arnab Nandi. 2019. Deterministic Routing between Layout Abstractions for Multi-Scale Classification of Visually Rich Documents. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*. International Joint Conferences on Artificial Intelligence Organization, 3360–3366. <https://doi.org/10.24963/ijcai.2019/466>
- [21] Sebastian Schreiber, Stefan Agne, Ivo Wolf, Andreas Dengel, and Sheraz Ahmed. 2017. DeepDeSRT: Deep Learning for Detection and Structure Recognition of Tables in Document Images. *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)* 01 (2017), 1162–1167.
- [22] Michael Shilman, Percy Liang, and Paul Viola. 2005. Learning nongenerative grammatical models for document analysis. In *Tenth IEEE International Conference on Computer Vision (ICCV’05) Volume 1*, Vol. 2. IEEE, 962–969.
- [23] Anikó Simon, J-C Pret, and A Peter Johnson. 1997. A fast algorithm for bottom-up document layout analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 19, 3 (1997), 273–277.
- [24] Carlos Soto and Shinjae Yoo. 2019. Visual Detection with Context for Document Layout Analysis. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 3462–3468. <https://doi.org/10.18653/v1/D19-1348>
- [25] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alex Alemi. 2016. Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning. In *AAAI*.
- [26] Matheus Palhares Viana and Dário Augusto Borges Oliveira. 2017. Fast CNN-Based Document Layout Analysis. *2017 IEEE International Conference on Computer Vision Workshops (ICCVW)* (2017), 1173–1180.
- [27] H. Wei, M. Baechler, F. Slimane, and R. Ingold. 2013. Evaluation of SVM, MLP and GMM Classifiers for Layout Analysis of Historical Documents. In *2013 12th International Conference on Document Analysis and Recognition*. 1220–1224. <https://doi.org/10.1109/ICDAR.2013.247>
- [28] Xiaowei Yang, Ersin Yumer, Paul Asente, Mike Kralej, Daniel Kifer, and C. Lee Giles. 2017. Learning to Extract Semantic Structure from Documents Using Multimodal Fully Convolutional Neural Networks. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2017), 4342–4351.
- [29] Xu Zhong, Jianbin Tang, and Antonio Jimeno-Yepes. 2019. PubLayNet: largest dataset ever for document layout analysis. *ArXiv abs/1908.07836* (2019).