

LAYOUTLMv2: MULTI-MODAL PRE-TRAINING FOR VISUALLY-RICH DOCUMENT UNDERSTANDING

Yang Xu^{1*}, Yiheng Xu^{2*}, Tengchao Lv^{2*}, Lei Cui², Furu Wei², Guoxin Wang³, Yijuan Lu³, Dinei Florencio³, Cha Zhang³, Wanxiang Che¹, Min Zhang⁴, Lidong Zhou²

¹Harbin Institute of Technology

²Microsoft Research Asia

³Microsoft Cloud&AI Team

⁴Soochow University

{yxu, car}@ir.hit.edu.cn

{v-yixu, v-telv, lecu, fuwei, lidongz}@microsoft.com

{guow, yijlu, dinei, chazhang}@microsoft.com

minzhang@suda.edu.cn

ABSTRACT

Pre-training of text and layout has proved effective in a variety of visually-rich document understanding tasks due to its effective model architecture and the advantage of large-scale unlabeled scanned/digital-born documents. In this paper, we present **LayoutLMv2** by pre-training text, layout and image in a multi-modal framework, where new model architectures and pre-training tasks are leveraged. Specifically, LayoutLMv2 not only uses the existing masked visual-language modeling task but also the new text-image alignment and text-image matching tasks in the pre-training stage, where cross-modality interaction is better learned. Meanwhile, it also integrates a spatial-aware self-attention mechanism into the Transformer architecture, so that the model can fully understand the relative positional relationship among different text blocks. Experiment results show that LayoutLMv2 outperforms strong baselines and achieves new state-of-the-art results on a wide variety of downstream visually-rich document understanding tasks, including FUNSD (0.7895 \rightarrow 0.8420), CORD (0.9493 \rightarrow 0.9601), SROIE (0.9524 \rightarrow 0.9781), Kleister-NDA (0.834 \rightarrow 0.852), RVL-CDIP (0.9443 \rightarrow 0.9564), and DocVQA (0.7295 \rightarrow 0.8672).

1 INTRODUCTION

Visually-rich Document Understanding (VrDU) aims to analyze scanned/digital-born business documents (images, PDFs, etc.) where structured information can be automatically extracted and organized for many business applications. Distinct from conventional information extraction tasks, the VrDU task not only relies on textual information, but also visual and layout information that is vital for visually-rich documents. For instance, the documents in Figure 1 include a variety of types such as digital forms, receipts, invoices and financial reports. Different types of documents indicate that the text fields of interest locate at different positions within the document, which is often determined by the style and format of each type as well as the document content. Therefore, to accurately recognize the text fields of interest, it is inevitable to take advantage of the cross-modality nature of visually-rich documents, where the textual, visual and layout information should be jointly modeled and learned end-to-end in a single framework.

The recent progress of VrDU lies primarily in two directions. The first direction is usually built on the shallow fusion between textual and visual/layout/style information (Yang et al., 2017a; Liu et al., 2019; Sarkhel & Nandi, 2019; Yu et al., 2020; Majumder et al., 2020; Wei et al., 2020; Zhang et al., 2020). These approaches leverage the pre-trained NLP and CV models individually and combine the information from multiple modalities for supervised learning. Although good performance has been achieved, these models often need to be re-trained from scratch once the document type is

*Equal contributions during internship at MSRA. Corresponding authors: Lei Cui and Furu Wei

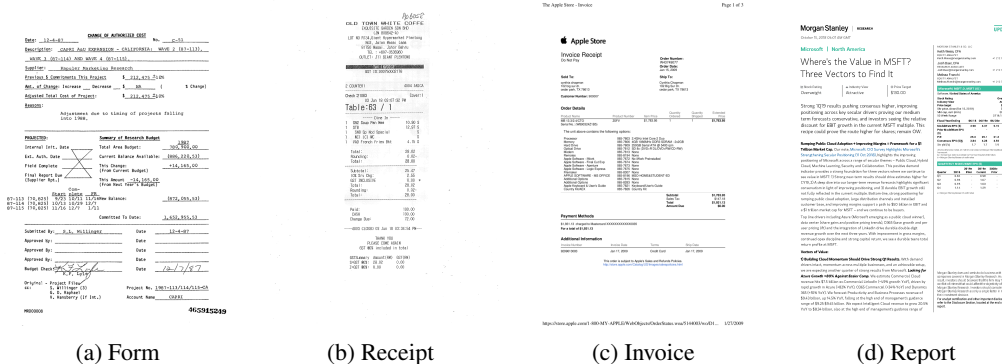


Figure 1: Visually-rich business documents with different layouts and formats

changed. In addition, the domain knowledge of one document type cannot be easily transferred into another document type, thereby the local invariance in general document layout (e.g. key-value pairs in a left-right layout, tables in a grid layout, etc.) cannot be fully exploited. To this end, the second direction relies on the deep fusion among textual, visual and layout information from a great number of unlabeled documents in different domains, where pre-training techniques play an important role in learning the cross-modality interaction in an end-to-end fashion (Lockard et al., 2020; Xu et al., 2020). In this way, the pre-trained models absorb cross-modal knowledge from different document types, where the local invariance among these layout and styles is preserved. Furthermore, when the model needs to be transferred into another domain with different document formats, only a few labeled samples would be sufficient to fine-tune the generic model in order to achieve state-of-the-art accuracy. Therefore, the proposed model in this paper follows the second direction, and we explore how to further improve the pre-training strategies for the VrDU task.

In this paper, we present an improved version of LayoutLM (Xu et al., 2020), aka **LayoutLMv2**. LayoutLM is a simple but effective pre-training method of text and layout for the VrDU task. Distinct from previous text-based pre-trained models, LayoutLM uses 2-D position embeddings and image embeddings in addition to the conventional text embeddings. During the pre-training stage, two training objectives are used, which are 1) a masked visual-language model and 2) multi-label document classification. The model is pre-trained with a great number of unlabeled scanned document images from the IIT-CDIP dataset (Lewis et al., 2006), and achieves very promising results on several downstream tasks. Extending the existing research work, we propose new model architectures and pre-training objectives in the LayoutLMv2 model. Different from the vanilla LayoutLM model where image embeddings are combined in the fine-tuning stage, we integrate the image information in the pre-training stage in LayoutLMv2 by taking advantage of the Transformer architecture to learn the cross-modality interaction between visual and textual information. In addition, inspired by the 1-D relative position representations (Shaw et al., 2018; Raffel et al., 2020; Bao et al., 2020), we propose the spatial-aware self-attention mechanism for the LayoutLMv2, which involves a 2-D relative position representation for token pairs. Different from the absolute 2-D position embeddings, the relative position embeddings explicitly provide a broader view for the contextual spatial modeling. For the pre-training strategies, we use two new training objectives for the LayoutLMv2 in addition to the masked visual-language model. The first is the proposed text-image alignment strategy, which covers text-lines in the image and makes predictions on the text-side to classify whether the token is covered or not on the image-side. The second is the text-image matching strategy that is popular in previous vision-language pre-training models (Tan & Bansal, 2019; Lu et al., 2019; Su et al., 2020; Chen et al., 2020; Sun et al., 2019), where some images in the text-image pairs are randomly replaced with another document image to make the model learn whether the image and OCR texts are correlated or not. In this way, LayoutLMv2 is more capable of learning contextual textual and visual information and the cross-modal correlation in a single framework, which leads to better VrDU performance. We select 6 publicly available benchmark datasets as the downstream tasks to evaluate the performance of the pre-trained LayoutLMv2 model, which are the FUNSD dataset (Jaume et al., 2019) for form understanding, the CORD dataset (Park et al., 2019) and the SROIE dataset (Huang et al., 2019) for receipt understanding, the Kleister-NDA dataset (Graliński

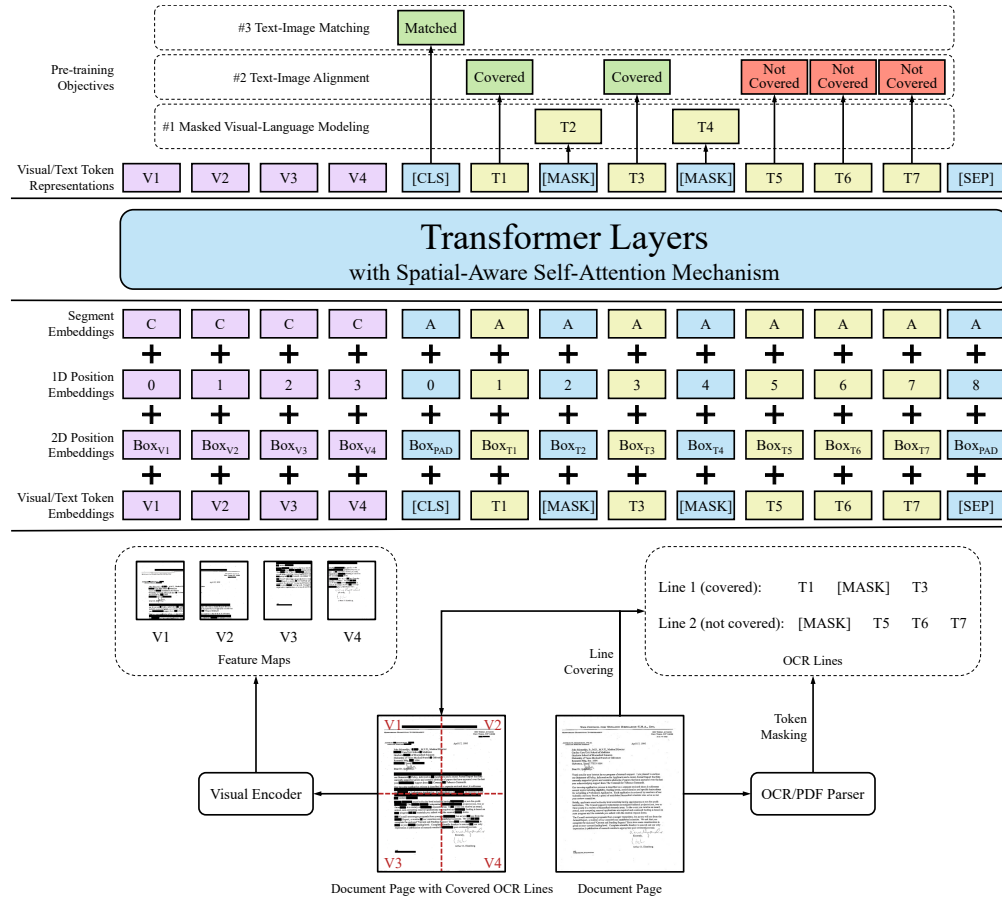


Figure 2: An illustration of the model architecture and pre-training strategies for LayoutLMv2

et al., 2020) for long document understanding with complex layout, the RVL-CDIP dataset (Harley et al., 2015) for document image classification, as well as the DocVQA dataset (Mathew et al., 2020) for visual question answering on document images. Experiment results show that the LayoutLMv2 model outperforms strong baselines including the vanilla LayoutLM and achieves new state-of-the-art results in these downstream VrDU tasks, which substantially benefits a great number of real-world document understanding tasks.

The contributions of this paper are summarized as follows:

- We propose a multi-modal Transformer model to integrate the document text, layout and image information in the pre-training stage, which learns the cross-modal interaction end-to-end in a single framework.
- In addition to the masked visual-language model, we also add text-image matching and text-image alignment as the new pre-training strategies to enforce the alignment among different modalities. Meanwhile, a spatial-aware self-attention mechanism is also integrated into the Transformer architecture.
- LayoutLMv2 not only outperforms the baseline models on the conventional VrDU tasks, but also achieves new SOTA results on the VQA task for document images, which demonstrates the great potential for the multi-modal pre-training for VrDU.

2 APPROACH

The overall illustration of the proposed LayoutLMv2 is shown in Figure 2. In this section, we will introduce the model architecture and pre-training tasks of the LayoutLMv2.

2.1 MODEL ARCHITECTURE

We build an enhanced Transformer architecture for the VrDU tasks, i.e. the multi-modal Transformer as the backbone of LayoutLMv2. The multi-modal Transformer accepts inputs of three modalities: text, image, and layout. The input of each modality is converted to an embedding sequence and fused by the encoder. The model establishes deep interactions within and between modalities by leveraging the powerful Transformer layers. The model details are introduced as follows, where some dropout and normalization layers are omitted.

Text Embedding We recognize text and serialize it in a reasonable reading order using off-the-shelf OCR tools and PDF parsers. Following the common practice, we use WordPiece (Wu et al., 2016) to tokenize the text sequence and assign each token to a certain segment $s_i \in \{[A], [B]\}$. Then, we add a [CLS] at the beginning of the token sequence and a [SEP] at the end of each text segment. The length of the text sequence is limited to ensure that the length of the final sequence is not greater than the maximum sequence length L . Extra [PAD] tokens are appended after the last [SEP] token to fill the gap if the token sequence is still shorter than L tokens. In this way, we get the input token sequence like

$$S = \{[\text{CLS}], w_1, w_2, \dots, [\text{SEP}], [\text{PAD}], [\text{PAD}], \dots\}, |S| = L$$

The final text embedding is the sum of three embeddings. Token embedding represents the token itself, 1D positional embedding represents the token index, and segment embedding is used to distinguish different text segments. Formally, we have the i -th text embedding

$$\mathbf{t}_i = \text{TokEmb}(w_i) + \text{PosEmb1D}(i) + \text{SegEmb}(s_i), 0 \leq i < L$$

Visual Embedding We use ResNeXt-FPN (Xie et al., 2016; Lin et al., 2017) architecture as the backbone of the visual encoder. Given a document page image I , it is resized to 224×224 then fed into the visual backbone. After that, the output feature map is average-pooled to a fixed size with the width being W and height being H . Next, it is flattened into a visual embedding sequence of length WH . A linear projection layer is then applied to each visual token embedding in order to unify the dimensions. Since the CNN-based visual backbone cannot capture the positional information, we also add a 1D positional embedding to these image token embeddings. The 1D positional embedding is shared with the text embedding layer. For the segment embedding, we attach all visual tokens to the visual segment [C]. The i -th visual embedding can be represented as

$$\mathbf{v}_i = \text{Proj}(\text{VisTokEmb}(I)_i) + \text{PosEmb1D}(i) + \text{SegEmb}([\text{C}]), 0 \leq i < WH$$

Layout Embedding The layout embedding layer aims to embed the spatial layout information represented by token bounding boxes in which corner coordinates and box shapes are identified explicitly. Following the vanilla LayoutLM, we normalize and discretize all coordinates to integers in the range $[0, 1000]$, and use two embedding layers to embed x-axis features and y-axis features separately. Given the normalized bounding box of the i -th text/visual token $\text{box}_i = (x_0, x_1, y_0, y_1, w, h)$, the layout embedding layer concatenates six bounding box features to construct a token-level layout embedding, aka the 2D positional embedding

$$\mathbf{l}_i = \text{Concat}(\text{PosEmb2D}_x(x_0, x_1, w), \text{PosEmb2D}_y(y_0, y_1, h)), 0 \leq i < WH + L$$

Note that CNNs perform local transformation, thus the visual token embeddings can be mapped back to image regions one by one with neither overlap nor omission. In the view of the layout embedding layer, the visual tokens can be treated as some evenly divided grids, so their bounding box coordinates are easy to calculate. An empty bounding box $\text{box}_{\text{PAD}} = (0, 0, 0, 0, 0, 0)$ is attached to special tokens [CLS], [SEP] and [PAD].

Multi-modal Encoder with Spatial-Aware Self-Attention Mechanism The encoder concatenates visual embeddings $\{\mathbf{v}_0, \dots, \mathbf{v}_{WH-1}\}$ and text embeddings $\{\mathbf{t}_0, \dots, \mathbf{t}_{L-1}\}$ to a unified sequence X and fuses spatial information by adding the layout embeddings to get the first layer input $\mathbf{x}^{(0)}$.

$$\mathbf{x}_i^{(0)} = X_i + \mathbf{l}_i, \text{ where } X = \{\mathbf{v}_0, \dots, \mathbf{v}_{WH-1}, \mathbf{t}_0, \dots, \mathbf{t}_{L-1}\}$$

Following the architecture of Transformer, we build our multi-modal encoder with a stack of multi-head self-attention layers followed by a feed-forward network. However, the original self-attention

mechanism can only implicitly capture the relationship between the input tokens with the absolute position hints. In order to efficiently model local invariance in the document layout, it is necessary to insert relative position information explicitly. Therefore, we introduce the spatial-aware self-attention mechanism into the self-attention layers. The original self-attention mechanism captures the correlation between query \mathbf{x}_i and key \mathbf{x}_j by projecting the two vectors and calculating the attention score

$$\alpha_{ij} = \frac{1}{\sqrt{d_{head}}} (\mathbf{x}_i \mathbf{W}^Q) (\mathbf{x}_j \mathbf{W}^K)^\top$$

We jointly model the semantic relative position and spatial relative position as bias terms and explicitly add them to the attention score. Let $\mathbf{b}^{(1D)}$, $\mathbf{b}^{(2D_x)}$ and $\mathbf{b}^{(2D_y)}$ denote the learnable 1D and 2D relative position biases respectively. The biases are different among attention heads but shared in all encoder layers. Assuming (x_i, y_i) anchors the top left corner coordinates of the i -th bounding box, we obtain the spatial-aware attention score

$$\alpha'_{ij} = \alpha_{ij} + \mathbf{b}_{j-i}^{(1D)} + \mathbf{b}_{x_j-x_i}^{(2D_x)} + \mathbf{b}_{y_j-y_i}^{(2D_y)}$$

Finally, the output vectors are represented as the weighted average of all the projected value vectors with respect to normalized spatial-aware attention scores

$$\mathbf{h}_i = \sum_j \frac{\exp(\alpha'_{ij})}{\sum_k \exp(\alpha'_{ik})} \mathbf{x}_j \mathbf{W}^V$$

2.2 PRE-TRAINING

We adopt three self-supervised tasks simultaneously during the pre-training stage, which are described as follows.

Masked Visual-Language Modeling Similar to the vanilla LayoutLM, we use the Masked Visual-Language Modeling (MVLM) to make the model learn better in the language side with the cross-modality clues. We randomly mask some text tokens and ask the model to recover the masked tokens. Meanwhile, the layout information remains unchanged, which means the model knows each masked token’s location on the page. The output representations of masked tokens from the encoder are fed into a classifier over the whole vocabulary, driven by a cross-entropy loss. To avoid visual clue leakage, we mask image regions corresponding to masked tokens on the raw page image input before feeding into the visual encoder. MVLM helps the model capture nearby tokens features. For instance, a masked blank in a table surrounded by lots of numbers is more likely to be a number. Moreover, given the spatial position of a blank, the model is capable of using visual information around to help predict the token.

Text-Image Alignment In addition to the MVLM, we propose the Text-Image Alignment (TIA) as a fine-grained cross-modality alignment task. In the TIA task, some text tokens are randomly selected, and their image regions are covered on the document image. We call this operation covering to avoid confusion with the masking operation in MVLM. During the pre-training, a classification layer is built above the encoder outputs. This layer predicts a label for each text token depending on whether it is covered, i.e., [Covered] or [Not Covered], and computes the binary cross-entropy loss. Considering the input image’s resolution is limited, the covering operation is performed at the line-level. When MVLM and TIA are performed simultaneously, TIA losses of the tokens masked in MVLM are not taken into account. This prevents the model from learning the useless but straightforward correspondence from [MASK] to [Covered].

Text-Image Matching Furthermore, a coarse-grained cross-modality alignment task, Text-Image Matching (TIM) is applied during the pre-training stage. We feed the output representation at [CLS] into a classifier to predict whether the image and text are from the same document page. Regular inputs are positive samples. To construct a negative sample, an image is either replaced by a page image from another document or dropped. To prevent the model from cheating by finding task features, we perform the same masking and covering operations to images in negative samples. The TIA target labels are all set to [Covered] in negative samples. We apply the binary cross-entropy loss in the optimization process.

2.3 FINE-TUNING

LayoutLMv2 produces representations with fused cross-modality information, which benefits a variety of VrDU tasks. Its output sequence provides representations at the token-level. Specifically, the output at [CLS] can be used as the global feature. For many downstream tasks, we only need to build a task specified head layer over the LayoutLMv2 outputs and fine-tune the whole model using an appropriate loss. In this way, LayoutLMv2 leads to much better VrDU performance by integrating the text, layout, and image information in a single multi-modal framework, which significantly improves the cross-modal correlation compared to the vanilla LayoutLM model.

3 EXPERIMENTS

3.1 DATA

In order to pre-train and evaluate LayoutLMv2 models, we select datasets in a wide range from the visually-rich document understanding area. Introduction to the dataset and task definitions along with the description of required data pre-processing are presented as follows.

Pre-training Dataset Following LayoutLM, we pre-train LayoutLMv2 on the IIT-CDIP Test Collection (Lewis et al., 2006), which contains over 11 million scanned document pages. We extract text and corresponding word-level bounding boxes from document page images with the Microsoft Read API.¹

FUNSD FUNSD (Jaume et al., 2019) is a dataset for form understanding in noisy scanned documents. It contains 199 real, fully annotated, scanned forms where 9,707 semantic entities are annotated above 31,485 words. The 199 samples are split into 149 for training and 50 for testing. The official OCR annotation is directly used with the layout information. The FUNSD dataset is suitable for a variety of tasks, where we focus on semantic entity labeling in this paper. Specifically, the task is assigning to each word a semantic entity label from a set of four predefined categories: question, answer, header or other. The entity-level F1 score is used as the evaluation metric.

CORD We also evaluate our model on the receipt key information extraction dataset, i.e. the public available subset of CORD (Park et al., 2019). The dataset includes 800 receipts for the training set, 100 for the validation set and 100 for the test set. A photo and a list of OCR annotations are equipped for each receipt. An ROI that encompasses the area of receipt region is provided along with each photo because there can be irrelevant things in the background. We only use the ROI as input instead of the raw photo. The dataset defines 30 fields under 4 categories and the task aims to label each word to the right field. The evaluation metric is entity-level F1. We use the official OCR annotations.

SROIE The SROIE dataset (Task 3) (Huang et al., 2019) aims to extract information from scanned receipts. There are 626 samples for training and 347 samples for testing in the dataset. The task is to extract values from each receipt of up to four predefined keys: company, date, address or total. The evaluation metric is entity-level F1. We use the official OCR annotations and results on the test set are provided by the official evaluation site.

Kleister-NDA Kleister-NDA (Graliński et al., 2020) contains non-disclosure agreements collected from the EDGAR database, including 254 documents for training, 83 documents for validation, and 203 documents for testing. This task is defined to extract the values of four fixed keys. We get the entity-level F1 score from the official evaluation tools.² Words and bounding boxes are extracted from the raw PDF file. We use heuristics to locate entity spans because the normalized standard answers may not appear in the utterance.

¹<https://docs.microsoft.com/en-us/azure/cognitive-services/computer-vision/concept-recognizing-text>

²<https://gitlab.com/filipg/geval>

RVL-CDIP RVL-CDIP (Harley et al., 2015) consists of 400,000 grayscale images, with 8:1:1 for the training set, validation set, and test set. A multi-class single-label classification task is defined on RVL-CDIP. The images are categorized into 16 classes, with 25,000 images per class. The evaluation metric is the overall classification accuracy. Text and layout information is extracted by Microsoft OCR.

DocVQA As a VQA dataset on the document understanding field, DocVQA (Mathew et al., 2020) consists of 50,000 questions defined on over 12,000 pages from a variety of documents. Pages are split into the training set, validation set and test set with a ratio of about 8:1:1. The dataset is organized as a set of triples $\langle \text{page image, questions, answers} \rangle$. Thus, we use Microsoft Read API to extract text and bounding boxes from images. Heuristics are used to find given answers in the extracted text. The task is evaluated using an edit distance based metric ANLS (aka average normalized Levenshtein similarity). Given that human performance is about 98% ANLS on the test set, it is reasonable to assume that the found ground truth which reaches over 97% ANLS on training and validation sets is good enough to train a model. Results on the test set are provided by the official evaluation site.

3.2 SETTINGS

Following the typical pre-training and fine-tuning strategy, we update all parameters and train whole models end-to-end for all the settings.

Pre-training LayoutLMv2 We train LayoutLMv2 models with two different parameter sizes. We set hidden size $d = 768$ in $\text{LayoutLMv2}_{\text{BASE}}$ and use a 12-layer 12-head Transformer encoder. While in the $\text{LayoutLMv2}_{\text{LARGE}}$, $d = 1024$ and its encoder has 24 Transformer layers with 16 heads. Visual backbones in the two models use the same ResNeXt101-FPN architecture. The numbers of parameters are 200M and 426M approximately for $\text{LayoutLMv2}_{\text{BASE}}$ and $\text{LayoutLMv2}_{\text{LARGE}}$, respectively.

The model is initialized from the existing pre-trained model checkpoints. For the encoder along with the text embedding layer, LayoutLMv2 uses the same architecture as UniLMv2 (Bao et al., 2020), thus it is initialized from UniLMv2. For the ResNeXt-FPN part in the visual embedding layer, the backbone of a Mask-RCNN (He et al., 2017) model trained on PubLayNet (Zhong et al., 2019) is leveraged.³ The rest of the parameters in the model are randomly initialized. We pre-train LayoutLMv2 models using Adam optimizer (Kingma & Ba, 2017; Loshchilov & Hutter, 2019), with the learning rate of 2×10^{-5} , weight decay of 1×10^{-2} , and $(\beta_1, \beta_2) = (0.9, 0.999)$. The learning rate is linearly warmed up over the first 10% steps then linearly decayed. $\text{LayoutLMv2}_{\text{BASE}}$ is trained with a batch size of 64 for 5 epochs, and $\text{LayoutLMv2}_{\text{LARGE}}$ is trained with a batch size of 2048 for 20 epochs on the IIT-CDIP dataset.

During the pre-training, we sample pages from the IIT-CDIP dataset and select a random sliding window of the text sequence if the sample is too long. We set the maximum sequence length $L = 512$ and assign all text tokens to the segment [A]. The output shape of the adaptive pooling layer is set to $W = H = 7$, so that it transforms the feature map into 49 image tokens. In MVLM, 15% text tokens are masked among which 80% are replaced by a special token [MASK], 10% are replaced by a random token sampled from the whole vocabulary, and 10% remains the same. In TIA, 15% of the lines are covered. In TIM, 15% images are replaced and 5% are dropped.

Fine-tuning LayoutLMv2 for Visual Question Answering We treat the DocVQA as an extractive QA task and build a token-level classifier on top of the text part of LayoutLMv2 output representations. Question tokens, context tokens and visual tokens are assigned to segment [A], [B] and [C], respectively. In the DocVQA paper, experiment results show that the BERT model fine-tuned on the SQuAD dataset (Rajpurkar et al., 2016) outperforms the original BERT model. Inspired by this, we add an extra setting, which is that we first fine-tune LayoutLMv2 on a Question Generation (QG) dataset followed by the DocVQA dataset. The QG dataset contains almost one million question-answer pairs generated by a generation model trained on the SQuAD dataset.

³“MaskRCNN ResNeXt101_32x8d FPN 3X” setting in <https://github.com/hpanwar08/detectron2>

Fine-tuning LayoutLMv2 for Document Image Classification This task depends on high-level visual information, thereby we leverage the image features explicitly in the fine-tuning. We pool the visual embeddings into a global pre-encoder feature, and pool the visual part of LayoutLMv2 output representations into a global post-encoder feature. The pre and post-encoder features along with the [CLS] output feature are concatenated and fed into the final classification layer.

Fine-tuning LayoutLMv2 for Sequence Labeling We formalize FUNSD, SROIE, CORD and Kleister-NDA as the sequence labeling tasks. To fine-tune LayoutLMv2 models on these tasks, we build a token-level classification layer above the text part of the output representations to predict the BIO tags for each entity field.

Baselines We select 3 baseline models in the experiments to compare LayoutLMv2 with the SOTA text-only pre-trained models as well as the vanilla LayoutLM model. Specifically, we compare LayoutLMv2 with BERT (Devlin et al., 2019), UniLMv2 (Bao et al., 2020) and LayoutLM (Xu et al., 2020) for all the experiment settings. We use the publicly available PyTorch models for BERT (Wolf et al., 2020) and LayoutLM,⁴ and use our in-house implementation for the UniLMv2 models. For each baseline approach, experiments are conducted using both the BASE and LARGE parameter settings.

3.3 RESULTS

FUNSD Table 1 shows the model accuracy on the FUNSD dataset which is evaluated using entity-level precision, recall and F1 score. For text-only models, the UniLMv2 models outperform the BERT models by a large margin in terms of the BASE and LARGE settings. For text+layout models, the LayoutLM family brings significant performance improvement over the text-only baselines, especially the LayoutLMv2 models. The best performance is achieved by the LayoutLMv2_{LARGE}, where an improvement of 3% F1 point is observed compared to the current SOTA results. This illustrates that the multi-modal pre-training in LayoutLMv2 learns better from the interactions from different modalities, thereby leading to the new SOTA on the form understanding task.

Model	Precision	Recall	F1	#Parameters
BERT _{BASE}	0.5469	0.6710	0.6026	110M
UniLMv2 _{BASE}	0.6349	0.6975	0.6648	125M
BERT _{LARGE}	0.6113	0.7085	0.6563	340M
UniLMv2 _{LARGE}	0.6780	0.7391	0.7072	355M
LayoutLM _{BASE}	0.7597	0.8155	0.7866	113M
LayoutLM _{LARGE}	0.7596	0.8219	0.7895	343M
LayoutLMv2 _{BASE}	0.8029	0.8539	0.8276	200M
LayoutLMv2 _{LARGE}	0.8324	0.8519	0.8420	426M
BROS (Anonymous, 2021)	0.8056	0.8188	0.8121	-

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

CORD Table 2 gives the entity-level precision, recall and F1 scores on the CORD dataset. The LayoutLM family significantly outperforms the text-only pre-trained models including BERT and UniLMv2, especially the LayoutLMv2 models. Compared to the baselines, the LayoutLMv2 models are also superior to the “SPADE” decoder method, as well as the “BROS” approach that is built on the “SPADE” decoder, which confirms the effectiveness of the pre-training for text, layout and image information.

SROIE Table 3 lists the entity-level precision, recall, and F1 score on Task 3 of the SROIE challenge. Compared to the text-only pre-trained language models, our LayoutLM family models have significant improvement by integrating cross-modal interactions. Moreover, with the same modal information, our LayoutLMv2 models also outperform existing multi-modal approaches (Anonymous,

⁴<https://github.com/microsoft/unilm/tree/master/layoutlm>

Model	Precision	Recall	F1	#Parameters
BERT _{BASE}	0.8833	0.9107	0.8968	110M
UniLMv2 _{BASE}	0.8987	0.9198	0.9092	125M
BERT _{LARGE}	0.8886	0.9168	0.9025	340M
UniLMv2 _{LARGE}	0.9123	0.9289	0.9205	355M
LayoutLM _{BASE}	0.9437	0.9508	0.9472	113M
LayoutLM _{LARGE}	0.9432	0.9554	0.9493	343M
LayoutLMv2 _{BASE}	0.9453	0.9539	0.9495	200M
LayoutLMv2 _{LARGE}	0.9565	0.9637	0.9601	426M
SPADE (Hwang et al., 2020)	-	-	0.9150	-
BROS (Anonymous, 2021)	0.9558	0.9514	0.9536	-

Table 2: Model accuracy (entity-level Precision, Recall, F1) on the CORD dataset

Model	Precision	Recall	F1	#Parameters
BERT _{BASE}	0.9099	0.9099	0.9099	110M
UniLMv2 _{BASE}	0.9459	0.9459	0.9459	125M
BERT _{LARGE}	0.9200	0.9200	0.9200	340M
UniLMv2 _{LARGE}	0.9488	0.9488	0.9488	355M
LayoutLM _{BASE}	0.9438	0.9438	0.9438	113M
LayoutLM _{LARGE}	0.9524	0.9524	0.9524	343M
LayoutLMv2 _{BASE}	0.9625	0.9625	0.9625	200M
LayoutLMv2 _{LARGE}	0.9661	0.9661	0.9661	426M
LayoutLMv2 _{LARGE} (Excluding OCR mismatch)	0.9904	0.9661	0.9781	426M
BROS (Anonymous, 2021)	0.9493	0.9603	0.9548	-
PICK (Yu et al., 2020)	0.9679	0.9546	0.9612	-
TRIE (Zhang et al., 2020)	-	-	0.9618	-
Top-1 on SROIE Leaderboard (Excluding OCR mismatch) ⁵	0.9889	0.9647	0.9767	-

Table 3: Model accuracy (entity-level Precision, Recall, F1) on the SROIE dataset (until 2020-12-24)

2021; Yu et al., 2020; Zhang et al., 2020), which demonstrates the model effectiveness. Eventually, the LayoutLMv2_{LARGE} single model can even beat the top-1 submission on the SROIE leaderboard.

Kleister-NDA Table 4 gives the entity-level F1 score of the Kleister-NDA dataset. As the labeled answers are normalized into a canonical form, we apply post-processing heuristics to convert the extracted date information into the “YYYY-MM-DD” format, and company names into the abbreviations such as “LLC” and “Inc.”. We report the evaluation results on the validation set because the ground-truth labels and the submission website for the test set are not available right now. The experiment results have shown that the LayoutLMv2 models improve the text-only and vanilla LayoutLM models by a large margin for the lengthy NDA documents, which also demonstrates that LayoutLMv2 can handle the complex layout information much better than previous models.

RVL-CDIP Table 5 shows the classification accuracy on the RVL-CDIP dataset, including text-only pre-trained models, the LayoutLM family as well as several image-based baseline models. As shown in the table, both the text and image information is important to the document image classification task because document images are text-intensive and represented by a variety of layouts and formats. Therefore, we observed that the LayoutLM family outperforms those text-only or image-only models as it leverages the multi-modal information within the documents. Specifically, the LayoutLMv2_{LARGE} model significantly improves the classification accuracy by more than 1.2% F1 point over the previous SOTA results, which achieves an accuracy of 95.64%. This also verifies that the pre-trained LayoutLMv2 model not only benefits the information extraction tasks in document

⁵Unpublished results, the leaderboard is available at <https://rrc.cvc.uab.es/?ch=13&com=evaluation&task=3>

Model	F1	#Parameters
BERT _{BASE}	0.779	110M
UniLMv2 _{BASE}	0.795	125M
BERT _{LARGE}	0.791	340M
UniLMv2 _{LARGE}	0.818	355M
LayoutLM _{BASE}	0.827	113M
LayoutLM _{LARGE}	0.834	343M
LayoutLMv2 _{BASE}	0.833	200M
LayoutLMv2 _{LARGE}	0.852	426M
RoBERTa _{BASE} in (Graliński et al., 2020)	0.793	125M

Table 4: Model accuracy (entity-level F1) on the validation set of the Kleister-NDA dataset using the official evaluation toolkit

Model	Accuracy	#Parameters
BERT _{BASE}	89.81%	110M
UniLMv2 _{BASE}	90.06%	125M
BERT _{LARGE}	89.92%	340M
UniLMv2 _{LARGE}	90.20%	355M
LayoutLM _{BASE} (w/ image)	94.42%	160M
LayoutLM _{LARGE} (w/ image)	94.43%	390M
LayoutLMv2 _{BASE}	95.25%	200M
LayoutLMv2 _{LARGE}	95.64%	426M
VGG-16 (Afzal et al., 2017)	90.97%	-
Single model (Das et al., 2018)	91.11%	-
Ensemble (Das et al., 2018)	92.21%	-
InceptionResNetV2 ⁶ (Szegedy et al., 2016)	92.63%	-
LadderNet (Sarkhel & Nandi, 2019)	92.77%	-
Single model (Dauphinee et al., 2019)	93.03%	-
Ensemble (Dauphinee et al., 2019)	93.07%	-

Table 5: Classification accuracy on the RVL-CDIP dataset

understanding but also the document image classification task through the effective model training across different modalities.

DocVQA Table 6 lists the Average Normalized Levenshtein Similarity (ANLS) scores on the DocVQA dataset of text-only baselines, LayoutLM family models and the previous top-1 on the leaderboard. With multi-modal pre-training, LayoutLMv2 models outperform LayoutLM models and text-only baselines by a large margin when fine-tuned on the train set. By using all data (train + dev) as the fine-tuning dataset, the LayoutLMv2_{LARGE} single model outperforms the previous top-1 on the leaderboard which ensembles 30 models. Under the setting of fine-tuning LayoutLMv2_{LARGE} on a question generation dataset (QG) and the DocVQA dataset successively, the single model performance increases by more than 1.6% ANLS and achieves the new SOTA.

3.4 ABLATION STUDY

To fully understand the underlying impact of different components, we conduct an ablation study to explore the effect of visual information, the pre-training tasks, spatial-aware self-attention mechanism, as well as different initialization models. Table 7 shows model performance on the DocVQA validation set. Under all the settings, we pre-train the models using all IIT-CDIP data for one epoch. The hyper-parameters are the same as those used to pre-train LayoutLMv2_{BASE} in Section 3.2.

⁶<https://medium.com/@jdegange85/benchmarking-modern-cnn-architectures-to-rvl-cdip-9dd0b7ec2955>

⁷Unpublished results, the leaderboard is available at <https://rrc.cvc.uab.es/?ch=17&com=evaluation&task=1>

Model	Fine-tuning set	ANLS	#Parameters
BERT _{BASE}	train	0.6354	110M
UniLMv2 _{BASE}	train	0.7134	125M
BERT _{LARGE}	train	0.6768	340M
UniLMv2 _{LARGE}	train	0.7709	355M
LayoutLM _{BASE}	train	0.6979	113M
LayoutLM _{LARGE}	train	0.7259	343M
LayoutLMv2 _{BASE}	train	0.7808	200M
LayoutLMv2 _{LARGE}	train	0.8348	426M
LayoutLMv2 _{LARGE}	train + dev	0.8529	426M
LayoutLMv2 _{LARGE} + QG	train + dev	0.8672	426M
Top-1 on DocVQA Leaderboard (30 models ensemble) ⁷	-	0.8506	-

Table 6: Average Normalized Levenshtein Similarity (ANLS) score on the DocVQA dataset (until 2020-12-24), “QG” denotes the data augmentation with the question generation dataset.

#	Model Architecture	Initialization	SASAM	MVLM	TIA	TIM	ANLS
1	LayoutLM _{BASE}	BERT _{BASE}		✓			0.6841
2a	LayoutLMv2 _{BASE}	BERT _{BASE} + X101-FPN		✓			0.6915
2b	LayoutLMv2 _{BASE}	BERT _{BASE} + X101-FPN		✓	✓		0.7061
2c	LayoutLMv2 _{BASE}	BERT _{BASE} + X101-FPN		✓		✓	0.6955
2d	LayoutLMv2 _{BASE}	BERT _{BASE} + X101-FPN		✓	✓	✓	0.7124
3	LayoutLMv2 _{BASE}	BERT _{BASE} + X101-FPN	✓	✓	✓	✓	0.7217
4	LayoutLMv2 _{BASE}	UniLMv2 _{BASE} + X101-FPN	✓	✓	✓	✓	0.7421

Table 7: Ablation study on the DocVQA dataset, where ANLS scores on the validation set are reported. “SASAM” means the spatial-aware self-attention mechanism. “MVLM”, “TIA” and “TIM” are the three proposed pre-training tasks. All the models are trained using all IIT-CDIP data for 1 epoch with the BASE model size.

“LayoutLM” denotes the vanilla LayoutLM architecture in (Xu et al., 2020), which can be regarded as a LayoutLMv2 architecture without visual module and spatial-aware self-attention mechanism. “X101-FPN” denotes the ResNeXt101-FPN visual backbone described in Section 3.2. We first evaluate the effect of introducing visual information. By comparing #1 and #2a, we find that LayoutLMv2 pre-trained with only MVLM can leverage visual information effectively. Then, we compare the two cross-modality alignment pre-training tasks TIA and TIM. According to the four results in #2, both tasks improve the model performance substantially, and the proposed TIA benefits the model more than the commonly used TIM. Using both tasks together is more effective than using either one alone. From the comparison result of #2d and #3, the spatial-aware self-attention mechanism can further improve the model accuracy. In the settings #3 and #4, we change the text-side initialization checkpoint from BERT to UniLMv2, and confirm that LayoutLMv2 benefits from the better initialization.

4 RELATED WORK

With the development of conventional machine learning, statistical machine learning approaches (Shilman et al., 2005; Marinai et al., 2005) have become the mainstream for document segmentation tasks during the past decade. Shilman et al. (2005) 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 (Marinai et al., 2005) 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 (Wei et al., 2013) 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.

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. Yang et al. (2017b) treat the document semantic structure extraction task as a pixel-by-pixel classification problem. They propose a multi-modal 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. Viana & Oliveira (2017) propose a lightweight model of document layout analysis for mobile and cloud services. The model uses one-dimensional information of images for inference and compares it with the model using two-dimensional information, achieving comparable accuracy in the experiments. Katti et al. (2018) 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. Soto & Yoo (2019) incorporate contextual information into the Faster R-CNN model that involves the inherently localized nature of article contents to improve region detection performance.

In recent years, pre-training techniques have become more and more popular in both NLP and CV areas, and have also been leveraged in the VrDU tasks. Devlin et al. (2019) introduced a new language representation model called BERT, which is designed to pre-train deep bidirectional representations from the unlabeled text by jointly conditioning on both left and right context in all layers. As a result, the pre-trained BERT model can be fine-tuned with just one additional output layer to create state-of-the-art models for a wide range of tasks. Bao et al. (2020) propose to pre-train a unified language model for both autoencoding and partially autoregressive language modeling tasks using a novel training procedure, referred to as a pseudo-masked language model. In addition, the two tasks pre-train a unified language model as a bidirectional encoder and a sequence-to-sequence decoder, respectively. Lu et al. (2019) proposed ViLBERT for learning task-agnostic joint representations of image content and natural language by extending the popular BERT architecture to a multi-modal two-stream model. Su et al. (2020) proposed VL-BERT that adopts the Transformer model as the backbone, and extends it to take both visual and linguistic embedded features as input. (Xu et al., 2020) proposed 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. This work is a natural extension of the vanilla LayoutLM, which takes advantage of textual, layout and visual information in a single multi-modal pre-training framework.

5 CONCLUSION

In this paper, we present a multi-modal pre-training approach for visually-rich document understanding tasks, aka LayoutLMv2. Distinct from existing methods for VrDU, the LayoutLMv2 model not only considers the text and layout information but also integrates the image information in the pre-training stage with a single multi-modal framework. Meanwhile, the spatial-aware self-attention mechanism is integrated into the Transformer architecture to capture the relative relationship among different bounding boxes. Furthermore, new pre-training objectives are also leveraged to enforce the learning of cross-modal interaction among different modalities. Experiment results on 6 different VrDU tasks have illustrated that the pre-trained LayoutLMv2 model has substantially outperformed the SOTA baselines in the document intelligence area, which greatly benefits a number of real-world document understanding tasks.

For future research, we will further explore the network architecture as well as the pre-training strategies for the LayoutLM family, so that we can push the SOTA results in VrDU to the new height. Meanwhile, we will also investigate the language expansion to make the multi-lingual LayoutLMv2 model available for different languages especially the non-English areas around the world.

REFERENCES

- Muhammad Zeshan Afzal, Andreas Kölsch, Sheraz Ahmed, and Marcus Liwicki. 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)*, 01:883–888, 2017.
- Anonymous. {BROS}: A pre-trained language model for understanding texts in document. In *Submitted to International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=punMXQEsPr0>. under review.
- Hangbo Bao, Li Dong, Furu Wei, Wenhui Wang, Nan Yang, Xiaodong Liu, Yu Wang, Songhao Piao, Jianfeng Gao, Ming Zhou, and Hsiao-Wuen Hon. Unilmv2: Pseudo-masked language models for unified language model pre-training, 2020.
- Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *ECCV*, 2020.
- Arindam Das, Saikat Roy, and Ujjwal Bhattacharya. Document image classification with intra-domain transfer learning and stacked generalization of deep convolutional neural networks. *2018 24th International Conference on Pattern Recognition (ICPR)*, pp. 3180–3185, 2018.
- Tyler Dauphinee, Nikunj Patel, and Mohammad Mehdi Rashidi. Modular multimodal architecture for document classification. *ArXiv*, abs/1912.04376, 2019.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 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)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://www.aclweb.org/anthology/N19-1423>.
- Filip Graliński, Tomasz Stanisławek, Anna Wróblewska, Dawid Lipiński, Agnieszka Kaliska, Paulina Rosalska, Bartosz Topolski, and Przemysław Biecek. Kleister: A novel task for information extraction involving long documents with complex layout, 2020.
- Adam W Harley, Alex Ufkes, and Konstantinos G Derpanis. Evaluation of deep convolutional nets for document image classification and retrieval. In *International Conference on Document Analysis and Recognition (ICDAR)*, 2015.
- Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- Z. Huang, K. Chen, J. He, X. Bai, D. Karatzas, S. Lu, and C. V. Jawahar. Icdar2019 competition on scanned receipt ocr and information extraction. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*, pp. 1516–1520, 2019. doi: 10.1109/ICDAR.2019.00244.
- Wonseok Hwang, Jinyeong Yim, Seunghyun Park, Sohee Yang, and Minjoon Seo. Spatial dependency parsing for semi-structured document information extraction, 2020.
- Guillaume Jaume, Hazim Kemal Ekenel, and Jean-Philippe Thiran. Funsd: A dataset for form understanding in noisy scanned documents. *2019 International Conference on Document Analysis and Recognition Workshops (ICDARW)*, Sep 2019. doi: 10.1109/icdarw.2019.10029. URL <http://dx.doi.org/10.1109/ICDARW.2019.10029>.
- Anoop R Katti, Christian Reisswig, Cordula Guder, Sebastian Brarda, Steffen Bickel, Johannes Höhne, and Jean Baptiste Faddoul. Chargrid: Towards understanding 2D documents. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 4459–4469, Brussels, Belgium, October-November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1476. URL <https://www.aclweb.org/anthology/D18-1476>.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017.

- D. Lewis, G. Agam, S. Argamon, O. Frieder, D. Grossman, and J. Heard. 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, SIGIR '06*, pp. 665–666, New York, NY, USA, 2006. Association for Computing Machinery. ISBN 1595933697. doi: 10.1145/1148170.1148307. URL <https://doi.org/10.1145/1148170.1148307>.
- Tsung-Yi Lin, Piotr Dollar, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature networks for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- Xiaojing Liu, Feiyu Gao, Qiong Zhang, and Huasha Zhao. 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)*, pp. 32–39, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-2005. URL <https://www.aclweb.org/anthology/N19-2005>.
- Colin Lockard, Prashant Shiralkar, Xin Luna Dong, and Hannaneh Hajishirzi. Zeroshotceres: Zero-shot relation extraction from semi-structured webpages. *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020. doi: 10.18653/v1/2020.acl-main.721. URL <http://dx.doi.org/10.18653/v1/2020.acl-main.721>.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks, 2019.
- Bodhisattwa Prasad Majumder, Navneet Potti, Sandeep Tata, James Bradley Wendt, Qi Zhao, and Marc Najork. Representation learning for information extraction from form-like documents. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 6495–6504, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.580. URL <https://www.aclweb.org/anthology/2020.acl-main.580>.
- S. Marinai, M. Gori, and G. Soda. Artificial neural networks for document analysis and recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(1):23–35, Jan 2005. ISSN 1939-3539. doi: 10.1109/TPAMI.2005.4.
- Minesh Mathew, Dimosthenis Karatzas, R. Manmatha, and C. V. Jawahar. Docvqa: A dataset for vqa on document images, 2020.
- Seunghyun Park, Seung Shin, Bado Lee, Junyeop Lee, Jaeheung Surh, Minjoon Seo, and Hwalsuk Lee. Cord: A consolidated receipt dataset for post-ocr parsing. 2019.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 2383–2392, Austin, Texas, November 2016. Association for Computational Linguistics. doi: 10.18653/v1/D16-1264. URL <https://www.aclweb.org/anthology/D16-1264>.
- Ritesh Sarkhel and Arnab Nandi. 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*, pp. 3360–3366. International Joint Conferences on Artificial Intelligence Organization, 7 2019. doi: 10.24963/ijcai.2019/466. URL <https://doi.org/10.24963/ijcai.2019/466>.

- Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, 2018. doi: 10.18653/v1/n18-2074. URL <http://dx.doi.org/10.18653/v1/N18-2074>.
- Michael Shilman, Percy Liang, and Paul Viola. Learning nongenerative grammatical models for document analysis. In *Tenth IEEE International Conference on Computer Vision (ICCV'05) Volume 1*, volume 2, pp. 962–969. IEEE, 2005.
- Carlos Soto and Shinjae Yoo. Visual detection with context for document layout analysis. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3462–3468, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1348. URL <https://www.aclweb.org/anthology/D19-1348>.
- Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. Vi-bert: Pre-training of generic visual-linguistic representations, 2020.
- Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. Videobert: A joint model for video and language representation learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alex Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *AAAI*, 2016.
- Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*, 2019.
- Matheus Palhares Viana and Dário Augusto Borges Oliveira. Fast cnn-based document layout analysis. *2017 IEEE International Conference on Computer Vision Workshops (ICCVW)*, pp. 1173–1180, 2017.
- H. Wei, M. Baechler, F. Slimane, and R. Ingold. Evaluation of svm, mlp and gmm classifiers for layout analysis of historical documents. In *2013 12th International Conference on Document Analysis and Recognition*, pp. 1220–1224, Aug 2013. doi: 10.1109/ICDAR.2013.247.
- Mengxi Wei, Yifan He, and Qiong Zhang. Robust layout-aware ie for visually rich documents with pre-trained language models. *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, Jul 2020. doi: 10.1145/3397271.3401442. URL <http://dx.doi.org/10.1145/3397271.3401442>.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 38–45, Online, October 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-demos.6. URL <https://www.aclweb.org/anthology/2020.emnlp-demos.6>.
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*, 2016.
- Saining Xie, Ross B. Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5987–5995, 2016.
- Yiheng Xu, Minghao Li, Lei Cui, Shaohan Huang, Furu Wei, and Ming Zhou. Layoutlm: Pre-training of text and layout for document image understanding. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD '20*,

pp. 1192–1200, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450379984. doi: 10.1145/3394486.3403172. URL <https://doi.org/10.1145/3394486.3403172>.

Xiao Yang, Ersin Yumer, Paul Asente, Mike Kraley, Daniel Kifer, and C. Lee Giles. Learning to extract semantic structure from documents using multimodal fully convolutional neural networks. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jul 2017a. doi: 10.1109/cvpr.2017.462. URL <http://dx.doi.org/10.1109/CVPR.2017.462>.

Xiaowei Yang, Ersin Yumer, Paul Asente, Mike Kraley, Daniel Kifer, and C. Lee Giles. Learning to extract semantic structure from documents using multimodal fully convolutional neural networks. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4342–4351, 2017b.

Wenwen Yu, Ning Lu, Xianbiao Qi, Ping Gong, and Rong Xiao. Pick: Processing key information extraction from documents using improved graph learning-convolutional networks, 2020.

Peng Zhang, Yunlu Xu, Zhanzhan Cheng, Shiliang Pu, Jing Lu, Liang Qiao, Yi Niu, and Fei Wu. Trie: End-to-end text reading and information extraction for document understanding, 2020.

Xu Zhong, Jianbin Tang, and Antonio Jimeno Yepes. Publaynet: largest dataset ever for document layout analysis. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*, pp. 1015–1022. IEEE, Sep. 2019. doi: 10.1109/ICDAR.2019.00166.