# Towards Efficient Resume Understanding: A Multi-Granularity Multi-Modal Pre-Training Approach

Feihu Jiang[1], Chuan Qin[2,3*], Jingshuai Zhang[2], Kaichun Yao[4],
Xi Chen[1], Dazhong Shen[5], Chen Zhu[2], Hengshu Zhu[2], Hui Xiong[6,7,8*]

[1]University of Science and Technology of China,
[2]Career Science Lab, BOSS Zhipin,
[3]PBC School of Finance, Tsinghua University,
[4]Institute of Software, Chinese Academy of Sciences,
[5]Shanghai Artificial Intelligence Laboratory,
[6]The Thrust of Artificial Intelligence, The Hong Kong University of Science and Technology (GuangZhou),
[7]The Department of Computer Science and Engineering, The Hong Kong University of Science and Technology,
[8]Guangzhou HKUST Fok Ying Tung Research Institute,
{jiangfeihu, chenxi0401}@mail.ustc.edu.cn, yaokaichun@outlook.com, xionghui@ust.hk
{chuanqin0426, zhangjingshuai0, dazh.shen, zc3930155, zhuhengshu}@gmail.com.

*Abstract*—In the contemporary era of widespread online recruitment, resume understanding has been widely acknowledged as a fundamental and crucial task, which aims to extract structured information from resume documents automatically. Compared to the traditional rule-based approaches, the utilization of recently proposed pre-trained document understanding models can greatly enhance the effectiveness of resume understanding. The present approaches have, however, disregarded the hierarchical relations within the structured information presented in resumes, and have difficulty parsing resumes in an efficient manner. To this end, in this paper, we propose a novel model, namely ERU, to achieve efficient resume understanding. Specifically, we first introduce a layout-aware multi-modal fusion transformer for encoding the segments in the resume with integrated textual, visual, and layout information. Then, we design three self-supervised tasks to pre-train this module via a large number of unlabeled resumes. Next, we fine-tune the model with a multi-granularity sequence labeling task to extract structured information from resumes. Finally, extensive experiments on a real-world dataset clearly demonstrate the effectiveness of ERU.

*Index Terms*—Resume understanding, pre-trained multi-modals

## I. INTRODUCTION

With the rapid growth and widespread adoption of online recruitment, employers receive over two hundred job applications for each available position [1]. Recently, a multitude of intelligent recruitment techniques have emerged to expedite the selection of the most suitable candidate [2]. Notably, the approach of resume understanding, also known as resume parsing, assumes a pivotal and indispensable role, as its ability to save recruiters from manually entering candidate information and instead automate the conversion of resumes into structured data. Such structured data can subsequently be utilized for various downstream applications, including person-job matching [3] and talent evaluation [4].

As illustrated in Figure 1, a resume conventionally incorporates the candidate's personal information (e.g., name), educational background (e.g., major, school), work experience, and various other pertinent details. The conversion of resumes into textual formats can be accomplished by utilizing optical character recognition (OCR), thereby rendering resume understanding as a text-mining task [5]. Traditional rule-based approaches extract structured information from resumes via various techniques, including keyword matching [6] and lexical analysis [7]. Furthermore, several studies regard it to be a sequence labeling task, wherein techniques such as hidden Markov models (HMM) [8] and conditional random fields (CRF) [9] are employed to forecast the category of each word within a textual sequence, such as *personal.name* and *education.major*, as shown in Figure 1.

In recent times, the utilization of pre-trained models for document understanding has become widespread [10]. These models have attained state-of-the-art performance on a variety of benchmark tasks by integrating multi-modal features, such as visual and layout information. Despite this, they still face several challenges in addressing the problem of resume understanding. To begin with, these transformer-based models are usually built on token-level granularity, rendering them difficult to apply when parsing long documents such as resumes. Furthermore, there exists a hierarchical relationship among

**Resume Example**

**SOPHIA**

**Details**
Phone Number: 3854234
Email: @gmail.com

**Education**
Columbus State University          SEPTEMBER 2016 — FEBRUARY 2019
Bachelor          Majors: Marketing & Business Management

**Work**
Company Name:          Software Limited Liability Company
Position:          Product Manager
Work Content:  led the Agile development of a Customer Relationship
Management (CRM) system tailored for small to medium-sized businesses.
The goal was to create an intuitive, scalable, and customizable ...

**Extracted Structured Data**

Personal.name: Sophia          Personal.email: @gmail.com
Personal.phone: 3854234
Education.school: Columbus State University          Education.degree: Bachelor
Education.date: February, 2019          Education.date: September, 2016
Education.majors: Marketing & Business Management
Work.company_name:          Software Limited Liability Company
Work.position: Product Manager

Fig. 1. An illustrative example of resume understanding.

the information that needs to be extracted from resumes, as exemplified by the personal information block depicted in Figure 1, which includes information like name, email, and phone number. However, prior document understanding models have not been capable of effectively modeling such hierarchical relationships.

To address the above challenges, we propose an efficient resume understanding model, namely ERU. Specifically, we first extract textual, visual, and layout inputs from each resume through an OCR tool and integrate them with a layout-aware multi-modal fusion transformer. Then, we pre-train the model on a large number of unlabeled resumes with three well-designed self-supervised objectives, including the masked language model, visual position alignment, and masked segment prediction. After that, we fine-tune the model with a multi-granularity sequence labeling task to extract the structured information from resumes. Finally, extensive experiments on a real-world dataset clearly demonstrates the effectiveness of our ERU.

## II. RELATED WORKS

**Resume Understanding.** Resume understanding aims to extract semantically structured information from resume documents, which can facilitate a wide range of intelligent recruitment applications, such as talent evaluation [4] and person-job matching [3]. Typically, resume understanding can be formulated as a text-mining task by using OCR tools to convert resumes into text format. Along this line, traditional studies predominantly rely on rule-based solutions like

keyword searching [6] and lexical analysis [7], or sequence labeling solutions, like HMM [8]. However, these methods usually suffer from high costs stemming from the need for expert-crafted feature engineering. Recently, researchers have attempted to investigate the problem using neural network-based techniques. For instance, *Chen et al.* [9] combined a bidirectional long short-term memory (BiLSTM) neural network with a CRF layer to parse resumes by leveraging Word2Vec features. Despite their success, the valuable visual and layout information in the resume document is often ignored, resulting in sub-optimal performance.

**Multi-modal Document Understanding.** Recently, the utilization of pre-trained models by integrating multi-modal features for document understanding has gained widespread adoption [10], [11]. In this direction, LayoutLM is the first to jointly learn textual and 2-D layout information in a unified model, with visual information extracted by Faster R-CNN also integrated into the token embeddings [10]. To facilitate computational efficiency, LayoutLMv3 replaces the visual feature extraction module with a simple linear projection head. On such basis, DocFormer [12], LayoutXLM [13], and LiLT [11] further use diverse techniques to combine image, text, and layout, and also investigate additional pre-training tasks to enhance the learning of multi-modal representations. However, these studies fail to achieve effective resume understanding, primarily due to their inability to handle lengthy resume documents that are often multi-paged and to consider the hierarchical relationship among the fields that require parsing.

## III. PRELIMINARY

In this study, we consider the multi-modal inputs of a resume to parse its structured information. Specifically, given a resume $R$, we leverage an OCR tool [1] to construct a sequence of segments $S = \{s_i\}_{i=1}^{|S|}$. Each segment $s_i$ contains: 1) textual input $text_i$ that includes a sequence of words, i.e., $text_i = \{c_j\}_{j=1}^{|s_i|}$; 2) layout inputs $b_i$ and $p_i$ that represent the bounding box $(b_i^0, b_i^1, b_i^2, b_i^3)$ [2] and page number of the segment $s_i$ respectively; and 3) visual input $v_i$ that denotes the visual information that cropped based on $b_i$. We regard resume understanding as a sequence labeling task, where the $C = \{e_i\}_{i=1}^{|C|}$ denotes all the possible fields, such as *personal.name* and *education.major* in Figure 1, and $L = \{l_i\}_{i=1}^{|S|}$ denotes the corresponding label sequence of $S$. The formal definition of the problem is as follows:

Give a set of resume documents $\mathcal{R}$, where each $R \in \mathcal{R}$ contains a sequence of segments $S$, the target of resume understanding is to learn a model $M$, which can predict the corresponding label $l_i$ for each segment $s_i \in S$ to achieve parsing the structure information in $R$.

---

[1] https://github.com/pymupdf/PyMuPDF

[2] In the bounding box, $(b_i^0, b_i^1)$ represents the position of the upper left and $(b_i^2, b_i^3)$ represents the position of the lower right.

## IV. METHOD

### A. Overview

Figure 2 presents a comprehensive overview of our ERU. To be specific, we first embed each segment from both textual and visual information. Then, a graph transformer is developed to model the interaction among different segments by representing each segment as two nodes, i.e., textual and visual nodes. Resume layout information is used to encode the spatial adjacency relationship among nodes. Next, we introduce the pre-training strategy of our model, where three self-supervised objectives are defined. Finally, the fine-tuning strategy on a relatively small-scale labeled dataset will be introduced based on the pre-trained model.

### B. Multi-Modal Embedding

We extract the textual and visual features from each segment.

**Textual Embedding.** To represent the text of one segment $s_i$, we used a 6-layer transformer which is initialized by BERT. We add two specific tokens [CLS] and [SEP] to the beginning and end of the token sequence and use [CLS] token as the text representation of the segment. As a result, the textual embedding vector can be derived by:

$$\text{text}_i = \left[[\text{CLS}], c_1, c_2, \cdots, c_{|s_i|}, [\text{SEP}]\right],$$
$$t_i = \text{BERT}\left(\text{text}_i\right). \tag{1}$$

**Visual Embedding.** To represent the visual information of each segment $s_i$, we enlarge the bounding box $b_i$ to $b_i'$ by a certain proportion to capture the visual information $v_i$ such as the style and color of the text. We use Faster R-CNN [14] to extract visual features of the region where the segment is located, i.e.,

$$v_i = \text{Linear}_v(\text{FasterRCNN}(I, b_i')). \tag{2}$$

### C. Layout-aware Multi-Modal Fusion Transformer

To model the interaction among different segments and different modalities, we revisit each resume $S$ as a complete graph $G_s$ with the node set $X = \{t_0, v_0, ..., t_{|S|}, v_{|S|}\}$. In other words, each segment $s_i$ is donated as two nodes $\{t_i, v_i\}$ corresponding to the textual and visual features. Then, we develop a 4-layers graph transformer on the graph $G_s$ by inducing both the absolute position bias of each node and the relative position among nodes with the guidance of the resume layout information.

**Absolute Position Bias.** For each segment $s_i$ ($t_i$ and $v_i$), there are three types of layout features, 1-D position, 2-D position, and the page number. The 1-D position is the absolute order $r_i$ obtained by parsing the resume from left to right and top to bottom. The 2-D position $b_i$ reflects the absolute position of the segment in the resume. We also encode the width and height of the segment since the same type of entities may have a similar size. We normalize the 2-D range to $(0, 1000)$. Thus the absolute position bias for each node $x_k = t_i$ or $v_i$ can be calculated by,

$$p_k = \text{Pos}_{2D}\left(b_i, width, height, page\right) + \text{Pos}_{1D}(r_i),$$
$$x_k^0 = x_k + p_k, \tag{3}$$

where $\text{Pos}_{2D}$ and $\text{Pos}_{1D}$ are implemented by MLP layers. By adding the absolute position bias to node input, the self-attention mechanism in the transformer can catch the absolute layout information.

**Relative Position Bias.** Indeed, the neighbor segments are often related, which indicates that the spatial adjacency relationship is also important to enhance the performance of resume understanding. Therefore, we follow Graphormer [15] and introduce the relative position bias to help the attention mechanism capture the neighbor information. Specifically, for each node pair $(x_m \in \{t_i, v_i\}, x_n \in \{t_j, s_j\})$, the relative position bias $b_{\phi(m,n)}$ is defined with the distance between segment $s_i$ and $s_j$:

$$b_{\phi(m,n)} = W_x\left(|b_i^0 - b_j^0|\right) + W_y\left(|b_i^1 - b_j^1|\right), \tag{4}$$

where $W_x$, $W_y$ are learnable parameters.

**Self-Attention.** Then, the self-attention mechanism can be defined as follows:

$$\alpha_{mn} = \frac{\left(x_m^0 W_Q\right)\left(x_n^0 W_K\right)^T}{\sqrt{d_f}} + b_{\phi(m,n)},$$
$$x_m' = \sum_n \frac{\exp\left(\alpha_{mn}\right)}{\sum_k \exp\left(\alpha_{mk}\right)} x_n^0 W_V, \tag{5}$$

where $W_Q$, $W_K$, and $W_V$ are the learnable query, key, and value metrics in the transformer models.

### D. Pre-training Objectives

Given the extensive corpus of unlabeled resumes, we design a pre-training strategy that includes three self-supervised tasks, each with its specific objectives, as detailed below.

**Maked Language Model.** The MLM task aims to make the text encoder adapt to the resume data. We initialize it with a pre-trained BERT model using 6 layers to leverage the existing knowledge. We mask the selected tokens in one segment text but retain the corresponding layout features. The pre-training objective is to reconstruct the masked tokens which promotes the model to capture a contextual representation with the layout clues. The detailed computation method is similar to the word prediction task in BERT. We denote this loss as $L_{MLM}$.

**Visual-Position Alignment.** To enhance the performance of extracting visual features, we propose to predict the relative position among segments with only their visual features, which encourages the visual features to capture the hidden relationship among neighbor segments. For simplicity, we divide the relative position into four directions $o = \{up, down, left, right\}$. For $v_i$, we randomly select a neighbor $v_i'$ and get a label $o_i$ according to their 2D position. We concatenate the two visual features and use a classification layer to predict the label. The $L_{VPA}$ is described as:

$$\mathcal{L}_{VPA} = -\sum_{i=1}^{|S|} \log p_{\theta_{vpa}}(o_i \mid [v_i : v_i']), \tag{6}$$

where $[:]$ denotes a concatenation operation, and the possibility function $p_\theta(b \mid a)$ is implemented by an MLP layer with parameter $\theta$ and the softmax function as the activation of the output layer, where $a$ is the input vector and the dimension of
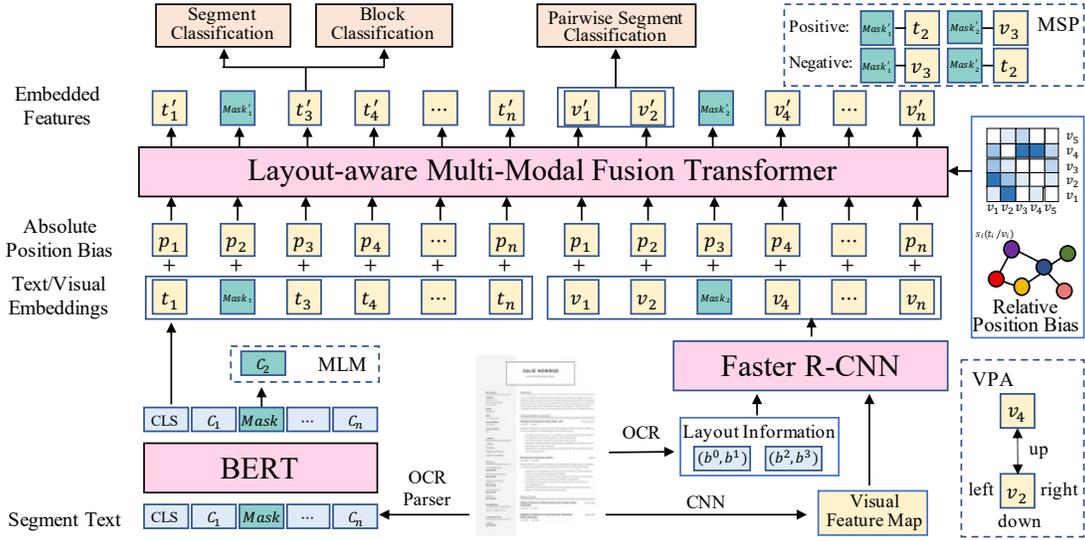
Fig. 2. An illustration of the proposed ERU.

the output is same as the dimension of $b$. To simplify, similar symbols will not be explained again in the future.

**Masked Segment Prediction.** To enhance the model's ability to integrate different modalities and learn contextual information, we drew inspiration from the MLM method of randomly masking tokens and applied it to the input sequence of segments. This involved randomly masking segments and allowing the model to learn the representation of the masked segments through context. Specifically, given the embedding set of textual and visual terms $X = \{t_0, ..., t_{|S|}, v_0, ..., v_{|S|}\}$, we randomly replace $q$ ($q < |S|$) term representations with a randomly initialized vector from both the textual and visual term sequence, ensuring that there is no overlap between the masked terms. This allows for the use of corresponding visual information when the text is masked, and vice versa. Thus we get a new input sequence with masked vectors $\hat{X} = \{t_0, \hat{t_1}, ..., t_{|S|}, \hat{v_0}, v_1, \hat{v_2}, ..., v_{|S|}\}$. After pulling $X$ and $\hat{X}$ into the transformer encoder defined in Section IV-C, we can get the output embedding set $X'$ and $\hat{X}'$, respectively. Then we use contrastive learning on the masked terms to minimize the distance between their representations in $X'$ and $\hat{X}'$. Specifically, for all masked term $x_m$, we follow [16] and use the following loss:

$$\mathcal{L}_{MSP} = -\sum_m \log \frac{\exp\left(\mathrm{Sim}\left(x'_m, \hat{x}'_m\right)/\tau\right)}{\sum_{k=1}^{N} \exp\left(\mathrm{Sim}\left(x'_m, \hat{x}'_k\right)/\tau\right)}, \quad (7)$$

where $\mathrm{Sim}(\cdot)$ measures the similarity between two vectors, such as Cosine similarity, $\hat{x}'_k$ is randomly selected from the set $\hat{X}'$ with the sampling number $N \ll |S|$, and $\tau$ is the temperature coefficient.

After the introduction of three self-supervised pre-training tasks, the final pre-training objective is defined as: $\mathcal{L}_{pre} = \lambda_{mlm}\mathcal{L}_{MLM} + \lambda_{vpa}\mathcal{L}_{VPA} + \lambda_{msp}\mathcal{L}_{MSP}$. where $\lambda_*$ are hyper-parameters to balance the different pre-training tasks.

### E. Multi-Granularity Sequence Labeling

After pre-training the model, we fine-tuned our model on a small-scale labeled dataset to classify the segments into corresponding labels. Specifically, for the output term representation sequence $X' = \{t'_0, t'_1, ..., t'_{|S|}, v'_0, v'_1, v'_2, ..., v'_{|S|}\}$, we use the representations of textual terms to predict both the class $l_i^{seg}$ of the corresponding segment and its corresponding block $l_i^{block}$, which provide a multi-granularity label prediction task. Moreover, we also introduce another task to provide better performance by predicting whether the given segment pair $g' = (t'_m, t'_n)$ belongs to the same blocks. If so, the label $l_i^{pair} = 1$, else 0. The final fine-tuning loss is as follows,

$$\mathcal{L}_f = -\sum_{i=1}^{|S|} \left( \log p_{\theta_{c1}}\left(l_i^{seg} \mid t'_i\right) + \log p_{\theta_{c2}}\left(l_i^{block} \mid t'_i\right) \right) + \sum_g \left(\log p_{\theta_{c3}}\left(l_i^{pair} \mid g'\right)\right). \quad (8)$$

## V. EXPERIMENTS

### A. Datasets

For the resume understanding experiments, we collected unlabeled resumes to pre-train our model. And we annotated the dataset for fine-tuning and test using an open-source PDF annotation tool named PAWLS[3]. The statistics of the dataset are summarized in Table I.

### B. Implementation

We set the maximum text length of a single segment to 32 and the maximum number of segments in a resume to 256. We used AdamW as the optimizer with the weight decay set to 0.01. The learning rates for the pre-trained model and other linear layers were 5e-5 and 1e-3, respectively. The temperature coefficient $\tau$ in MSP is set to 2. We set $\lambda_{mlm}$, $\lambda_{vpa}$ and $\lambda_{msp}$

[3]https://github.com/allenai/pawls

TABLE I
STATISTICS OF RESUME UNDERSTANDING DATASETS.

| Datasets | Pre-training Documents | Fine-tune Documents | | |
| --- | --- | --- | --- | --- |
| | | Training | Validation | Test |
| samples | 169,286 | 196 | 100 | 141 |
| avg of seg number | 88.90 | 95.49 | 98.34 | 102.42 |
| avg of seg length | 18.94 | 16.49 | 17.22 | 15.24 |
| avg of pages | 1.95 | 2.31 | 2.15 | 2.41 |

TABLE II
OVERALL PERFORMANCE OF RESUME UNDERSTANDING.

| Model | Precision | Recall | F1 |
| --- | --- | --- | --- |
| LLM IE | 63.37 ± 1.37 | 59.67 ± 2.03 | 61.24 ± 1.94 |
| BERT | 74.98 ± 2.21 | 79.36 ± 1.82 | 77.11 ± 0.54 |
| LiLT | 75.98 ± 1.16 | 89.47 ± 0.60 | 82.03 ± 0.50 |
| DocFormer | 83.81 ± 2.66 | 86.05 ± 4.50 | 83.45 ± 0.64 |
| LayoutXLM | 78.93 ± 1.66 | 89.14 ± 1.92 | 83.62 ± 0.28 |
| LayoutLMv3 | 82.96 ± 1.53 | 86.29 ± 2.66 | 84.36 ± 0.53 |
| ERU | **84.64 ± 0.32** | **91.08 ± 0.52** | **87.75 ± 0.22** |

TABLE III
ABLATION EXPERIMENTS OF DIFFERENT MODULES.

| Model | Precision | Recall | F1 |
| --- | --- | --- | --- |
| ERU | **84.64 ± 0.32** | **91.08 ± 0.52** | **87.75 ± 0.22** |
| -w/o VPA | 83.59 ± 0.15 | 90.47 ± 0.15 | 86.89 ± 0.12 |
| -w/o MSP | 81.20 ± 0.11 | 89.35 ± 0.15 | 85.07 ± 0.17 |
| -w/o MSP+VPA | 83.79 ± 0.11 | 85.53 ± 0.16 | 84.65 ± 0.11 |
| - w/o visual-emb | 84.53 ± 0.46 | 89.66 ± 0.45 | 87.02 ± 0.14 |
| - w/o relative-pos | 84.15 ± 0.42 | 90.16 ± 0.62 | 87.05 ± 0.18 |
| - w/o multi-gran | 84.41 ± 1.68 | 89.07 ± 0.26 | 87.20 ± 0.03 |

as 1,1 and 0.6 for the best performance. All experiments were conducted on 8 Tesla A800 80G GPUs.

### C. Baselines

We totally selected 6 models as baselines. **BERT** [17] is a typical sequence labeling method. **LiLT** [11] is trained on a single language and then directly fine-tuned on other languages. **DocFormer** [12] is a multi-modal pre-training model that adds spatial relative features. **LayoutXLM** [13] is a multi-modal pre-training model for multilingual document understanding. **LayoutLMv3** [18] is a general-purpose pre-training model for both text-centric and image-centric Document AI tasks with a unified architecture. **LLM IE** [19] is a generative information extraction method utilizing ChatGPT's advanced capabilities. To protect data privacy, our implementaion employs the open-source LLM Baichuan-13b which is enhanced by a supervised fine-tuning(sft) strategy.

### D. The Overall Performance

To evaluate the performance of the resume understanding task, following [9], we utilize precision, recall, and F1-score as the evaluation metrics. The comparison results for ERU and baseline methods are shown in Table II. According to the results, we observe that our model outperformed all the baselines, which demonstrates the effectiveness of ERU. In comparison to the top-performing LayoutLMv3, we find that ERU achieves above 1.68%, 4.79%, and 3.39% improvement on the precision, recall, and F1-score respectively. LayoutLMv3 operates on a token-level granularity and will divide long resumes into parts for individual processing, hindering its ability to fully comprehend a resume. In contrast, ERU is designed on a segment-level granularity, enabling it to understand a resume in its entirety. Additionally, ERU captures the hierarchical structure within resumes through a multi-

granularity sequence labeling task. As a result, our model outperforms all baseline models.

We have also conducted further experiments to compare generative information extraction with our discriminative information extraction approach. As shown in Tabel II, LLM based method still falls short of extractive models when dealing with complex and lengthy texts. In the experiments, we found that certain flaws in LLM, such as the hallucination phenomenon [20], [21], lead to the prediction of information not pertaining to the current resume. This is detrimental to the high accuracy required for resume parsing. Additionally, LLM's parsing efficiency is notably lower compared to discriminative models, highlighting the need for further research into LLM IE. In future work, we aim to focus on improving both the efficiency and accuracy of information extraction methods based on LLM.

### E. Ablation Studies

To validate the effectiveness of each component in ERU, we conducted a series of six ablation experiments. The first set of experiments involved comparing ERU with three of its variants which represent the removal of the corresponding pre-training tasks: ERU *-w/o VPA*, *-w/o MSP*, and *-w/o MSP+VPA*. The second set of experiments compared ERU with three other variants: 1) *-w/o visual-emb*, which omits the visual embedding in multi-modal embedding; 2) *-w/o relative-pos*, which excludes the relative position bias in the layout-aware multi-modal fusion transformer; and 3) *-w/o multi-gran*, which uses only the segment classification loss in multi-granularity sequence labeling.

The first three ablation studies, which focus on pre-training losses, provide significant insights, as detailed in Table III. Initially, the exclusion of the VPA pre-training task led to a noticeable decrease in ERU's performance, highlighting the role of high-quality visual feature extraction in enhancing the model's classification effectiveness. Furthermore, the absence of the MSP pre-training task resulted in an even more pronounced performance decline, emphasizing the importance of the model's ability to integrate features from different modalities during pre-training. By utilizing both MSP and VPA pre-training tasks, our model achieved nearly a 2.6% improvement in F1-score, underlining the effectiveness of these designed pre-training tasks.

And the remaining three ablation experiments also clearly demonstrate their contributions to the overall performance

as shown in Table III. The removal of the visual embedding (*-w/o visual-emb*) impacts the model's effectiveness in processing and integrating visual information. Likewise, excluding the relative position bias (*-w/o relative-pos*) affects the model's understanding of spatial relationships in the data. Lastly, employing only the segment classification loss in sequence labeling (*-w/o multi-gran*) restricts the model's capability to label sequences at different levels of granularity, which influences the overall performance.

## VI. CONCLUSION

In this paper we introduced a novel efficient resume understanding model, ERU, to automatically extract structured information from resume documents. To be more specific, we first designed a layout-aware multi-modal fusion transformer to encode the segments in the resume with text, visual and lauoyt features. Then we pre-trained the model on a large number of unlabeled resumes with three self-supervised tasks. After that, we fine-tuned the pre-trained model to extract structured information in resume with a multi-granularity sequence labeling task. Finally, extensive experiments have clearly demonstrated the effectiveness of ERU.

## REFERENCES

[1] Bart Turczynski, "2023 HR Statistics: Job Search, Hiring, Recruiting and Interviews," https://zety.com/blog/hr-statistics, 2023, [Online; accessed 26-April-2023].

[2] Harriet Rodney, Katarina Valaskova, and Pavol Durana, "The artificial intelligence recruitment process: How technological advancements have reshaped job application and selection practices," *Psychosociological Issues in Human Resource Management*, vol. 7, no. 1, pp. 42–47, 2019.

[3] Chuan Qin, Hengshu Zhu, Tong Xu, Chen Zhu, Liang Jiang, Enhong Chen, and Hui Xiong, "Enhancing person-job fit for talent recruitment: An ability-aware neural network approach," in *The 41st international ACM SIGIR conference on research & development in information retrieval*, 2018, pp. 25–34.

[4] Mingzhe Li, Xiuying Chen, Weiheng Liao, Yang Song, Tao Zhang, Dongyan Zhao, and Rui Yan, "Ezinterviewer: To improve job interview performance with mock interview generator," *arXiv preprint arXiv:2301.00972*, 2023.

[5] Divyanshu Chandola, Aditya Garg, Ankit Maurya, and Amit Kushwaha, "Online resume parsing system using text analytics," *Journal of Multi-Disciplinary Engineering Technologies*, vol. 9, 2015.

[6] Sunil Kumar Kopparapu, "Automatic extraction of usable information from unstructured resumes to aid search," in *2010 IEEE International Conference on Progress in Informatics and Computing*. IEEE, 2010, vol. 1, pp. 99–103.

[7] Christian Siefkes and Peter Siniakov, "An overview and classification of adaptive approaches to information extraction," *Journal on Data Semantics IV*, pp. 172–212, 2005.

[8] Yan Wentan and Qiao Yupeng, "Chinese resume information extraction based on semi-structured text," in *2017 36th Chinese Control Conference (CCC)*. IEEE, 2017, pp. 11177–11182.

[9] Jiaze Chen, Liangcai Gao, and Zhi Tang, "Information extraction from resume documents in pdf format," *Electronic Imaging*, vol. 28, pp. 1–8, 2016.

[10] 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*, 2020, pp. 1192–1200.

[11] Jiapeng Wang, Lianwen Jin, and Kai Ding, "Lilt: A simple yet effective language-independent layout transformer for structured document understanding," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2022, pp. 7747–7757.

[12] Srikar Appalaraju, Bhavan Jasani, Bhargava Urala Kota, Yusheng Xie, and R Manmatha, "Docformer: End-to-end transformer for document understanding," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 993–1003.

[13] Yiheng Xu, Tengchao Lv, Lei Cui, Guoxin Wang, Yijuan Lu, Dinei Florencio, Cha Zhang, and Furu Wei, "Layoutxlm: Multimodal pre-training for multilingual visually-rich document understanding," *arXiv preprint arXiv:2104.08836*, 2021.

[14] Ross Girshick, "Fast r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1440–1448.

[15] Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen, and Tie-Yan Liu, "Do transformers really perform badly for graph representation?," *Advances in Neural Information Processing Systems*, vol. 34, pp. 28877–28888, 2021.

[16] Aaron van den Oord, Yazhe Li, and Oriol Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.

[17] XiaoWei Li, Hui Shu, Yi Zhai, and ZhiQiang Lin, "A method for resume information extraction using bert-bilstm-crf," in *2021 IEEE 21st International Conference on Communication Technology (ICCT)*. IEEE, 2021, pp. 1437–1442.

[18] Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and Furu Wei, "Layoutlmv3: Pre-training for document ai with unified text and image masking," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 4083–4091.

[19] Tingyu Xie, Qi Li, Jian Zhang, Yan Zhang, Zuozhu Liu, and Hongwei Wang, "Empirical study of zero-shot ner with chatgpt," *arXiv preprint arXiv:2310.10035*, 2023.

[20] Feihu Jiang, Chuan Qin, Kaichun Yao, Chuyu Fang, Fuzhen Zhuang, Hengshu Zhu, and Hui Xiong, "Enhancing question answering for enterprise knowledge bases using large language models," in *International Conference on Database Systems for Advanced Applications*. Springer, 2024.

[21] Derong Xu, Wei Chen, Wenjun Peng, Chao Zhang, Tong Xu, Xiangyu Zhao, Xian Wu, Yefeng Zheng, and Enhong Chen, "Large language models for generative information extraction: A survey," *arXiv preprint arXiv:2312.17617*, 2023.