

# PPTSER: A Plug-and-Play Tag-guided Method for Few-shot Semantic Entity Recognition on Visually-rich Documents

Anonymous ACL submission

## Abstract

001 Visually-rich document information extraction  
002 (VIE) is a vital aspect of document under-  
003 standing, wherein Semantic Entity Recognition  
004 (SER) plays a significant role. However, few-  
005 shot SER on visually-rich documents remains  
006 relatively unexplored despite its considerable  
007 potential for practical applications. To address  
008 this issue, we propose a simple yet effective  
009 **Plug-and-Play Tag-guided** method for few-  
010 shot **Semantic Entity Recognition (PPTSER)**  
011 on visually-rich documents. PPTSER is built  
012 upon off-the-shelf multi-modal pre-trained  
013 models. It leverages the semantics of the tags  
014 to guide the SER task, reformulating SER into  
015 entity typing and span detection, handling both  
016 tasks simultaneously via cross-attention. Exper-  
017 imental results illustrate that PPTSER outper-  
018 forms existing fine-tuning and few-shot meth-  
019 ods, especially in low-data regimes. With  
020 full training data, PPTSER achieves compa-  
021 rable or superior performance to fine-tuning  
022 baseline. For instance, on the FUNSD bench-  
023 mark, our method improves the performance  
024 of LayoutLMv3-base in 1-shot, 3-shot and 5-  
025 shot scenarios by 15.61%, 2.13%, and 2.01%,  
026 respectively. Overall, PPTSER demonstrates  
027 promising generalizability, effectiveness, and  
028 plug-and-play nature for few-shot SER on  
029 visually-rich documents. The codes will be  
030 available at <https://anonymous.for.review>.

## 031 1 Introduction

032 Information extraction from visually-rich docu-  
033 ments (VIE) is a process that concentrates on ex-  
034 tracting pertinent information from various sources  
035 such as scanned images, documents, and PDF files.  
036 It effectively leverages layout and visual cues to de-  
037 code the content enclosed within these documents  
038 (Xu et al., 2020). As an important part of VIE,  
039 Semantic Entity Recognition (SER) aims to extract  
040 entity spans from the visually-rich document. SER  
041 has been hailed as a significant advancement in

the realm of document intelligence, and has found  
widespread applications in numerous sectors.

The advent of multi-modal pre-trained models  
(Xu et al., 2020; Li et al., 2021c; Gu et al., 2021;  
Huang et al., 2022b; Yu et al., 2023) has ushered  
in a rapid evolution in SER methodologies. These  
models, pre-trained on a large corpus of scanned  
documents in a self-supervised manner, have sig-  
nificantly enhanced the comprehension ability of  
SER. Despite the remarkable achievements of the  
multi-modal pre-trained models, they often rely on  
extensive data for fine-tuning. However, acquiring  
a large volume of well-annotated SER data poses  
significant challenges such as: (1) Acquiring such  
data necessitates substantial financial resources and  
time. Annotators are required to label a multitude  
of OCR detection boxes in the document, adhering  
to meticulously designed guidelines. Identification  
of content within a box and accurately assigning  
labels to them are also tedious tasks. (2) The avail-  
ability of data is often restricted due to privacy  
concerns. In scenarios involving sensitive infor-  
mation, such as invoices and insurance documents,  
data accessibility is severely limited due to the con-  
fidential nature of this information.

Despite the scarce research (Cheng et al., 2020;  
Yao et al., 2021; Wang and Shang, 2022) on few-  
shot Semantic Entity Recognition for visually-rich  
documents (few-shot SER), results have shown lim-  
itations in terms of generality and performance, and  
were limited to the specific application scenario.  
This paper, inspired by the comprehension capabil-  
ities of pre-trained models and the selective focus  
nature of the attention mechanism, introduces a  
novel approach called PPTSER, a **Plug-and-Play  
Tag-guided** method for few-shot **Semantic Entity  
Recognition** on visually-rich documents. The un-  
derlying principle of PPTSER consists of three  
main components: (1) **Semantic Understanding  
and Alignment:** Words related to SER tags are  
used as a prompt and are concatenated with the

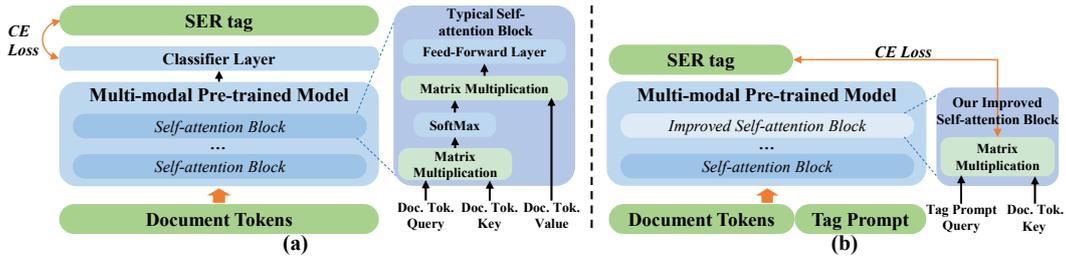


Figure 1: (a) Illustration of the traditional fine-tuning method. *Doc. Tok.* refers to *Document Tokens*. (b) Overview of our PPTSER method. PPTSER replaces the last self-attention block with an improved attention block and omits an extra classifier layer compared to traditional fine-tuning, which has less modules and parameters.

document’s text tokens. This combined input is then fed into a multi-modal pre-trained model. The motivation behind this is that the pre-trained model is expected to understand the semantics of both the document tokens and the tag-related prompt, thereby bringing the hidden states of the tokens and tag-related words for a specific entity type closer together. (2) **Decoupling of SER task:** SER task is segmented into Entity Classification and Entity Boundary Detection. This division aims to facilitate the resolution of boundary determination among adjacent entities of the same category within visually-rich documents. (3) **Efficient Usage of Multi-head Attention:** The attention weight obtained from the last attention block between the tag-related prompt and document tokens is directly used as the probability of tokens belonging to different tags. This mechanism, with different heads detecting various spans, is ideal for the SER task with numerous entity spans. By fully exploiting the weighted focus nature of the attention mechanism, the model eliminates the value transform layer, feed-forward layer in the last attention block, and omits a separate classifier layer compared to traditional fine-tuning methods (as depicted in Figure 1), As a result, the total parameter is reduced.

Extensive experiments are conducted to show the PPTSER’s effectiveness on commonly-used SER benchmarks, covering multiple languages, in few-shot to the full training set settings, and using different mainstream multi-modal pre-trained models.

The main contributions of this paper can be summarized as follows:

- We have demonstrated that the semantics of labels can effectively guide the SER task and have proposed a plug-and-play method ideal for few-shot SER on visually-rich documents. To the best of our knowledge, we are the first to propose a pluggable method that has shown effectiveness on various pre-trained models and languages.

- By innovatively leveraging the multi-head attention mechanism embedded in the pre-trained model, our method successfully extracts dense entities on visually-rich documents without adding any additional parameters.
- Experimental results show the superiority of our method over the traditional fine-tuning approaches in both few-shot and full-training-set scenarios. Moreover, PPTSER outperforms existing few-shot SER methods by significant margin, thereby underscoring its overall efficacy.

## 2 Method

### 2.1 Task Formulation

SER is usually formulated as a sequence labeling task. For given tokens from the document  $x = [x_i], i = 1, 2, \dots, n$ , SER aims to assign a label  $y_i \in \mathbb{C}$  for each token  $x_i$ , where  $\mathbb{C}$  is the SER label space. Subsequently, entity spans would be analyzed from the labeled tokens according to a specific scheme, such as BIO (Ramshaw and Marcus, 1995) and IO (Tjong Kim Sang and De Meulder, 2003).

In this paper, we primarily focus on the **In-Label-Space** setting for few-shot SER. Specifically, the pre-trained model is firstly fine-tuned on a small number of  $M$  annotated documents with label space  $\mathbb{C}$  and then directly evaluated on the test set with the same label space  $\mathbb{C}$ . This task presents a significant challenge as the model needs to learn the SER task with only limited training samples.

It is notable that in the context of few-shot SER, the few-shot setting of  $N$ -way  $K$ -shot indicates that *each of the  $N$  categories has  $K$  documents containing entities of that category as the support set*, as visually-rich documents are annotated at document level. Moreover, a document often contains entity spans of distinct types, causing potential overlaps

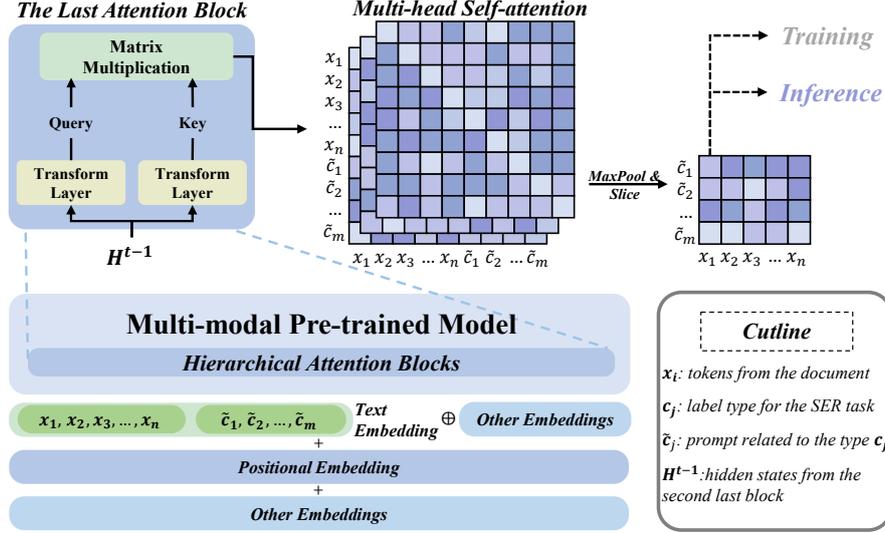


Figure 2: The overall architecture of PPTSER. *Other Embeddings* may include various embeddings such as *Visual Embedding*, *Type Embedding*, among others, with their presence and format dependent on the type of the pre-trained model used. In this architecture, the tokens extracted from documents and the tag-related prompt are concatenated and subsequently encoded with the pre-trained model. The attention weight, obtained from the last attention block between tokens and the prompt, is then used to ascertain whether the tokens correspond to the respective label type.

between the support sets for different entity types across  $N$  categories. Consequently, the overall number of annotated documents  $M < N \times K$ .

## 2.2 PPTSER

The fundamental concept and flow chart of PPTSER is shown in Figure 2. The method begins with the construction of a prompt based on SER tags. This prompt is then concatenated with the document tokens and jointly encoded using a unified pre-trained model. Within the transformer architecture of our model, attention weights between document tokens and the tag-related prompt are computed in hierarchical attention blocks. We use the attention weight between the tag-related prompt and document tokens, which can be considered as a form of cross-attention, obtained from the last attention block as the probability distribution of tokens belonging to different SER entity types.

### 2.2.1 Tag-related Prompt Construction and Target Generation

For an SER task with the label space  $\mathbb{C}$ , we need to construct tag-related words  $\tilde{c}_i$  for each  $c_i \in \mathbb{C}$ , and then the tag-related prompt  $\tilde{\mathbb{C}} = \{\tilde{c}_i\}$ ,  $i = 1, 2, \dots, m$  is built. In PPTSER, we simply use the tag names as the tag-related words.

To enable PPTSER to accurately identify the boundaries of entity spans, we employ BIO tagging scheme in our method. However, when dealing with an SER task involving entity types  $\mathbb{E} = \{e_i | e_0 = \text{Other}\}$ ,  $i = 0, 1, 2, \dots, m$

(where *Other* represents the entities that are not of interest), the label space would be  $\mathbb{C} = \{e_0, B_{e_i}, I_{e_i}\}$ , and the prompt would be  $\tilde{\mathbb{C}} = \{e_0, \text{beginning of } e_i, \text{inner of } e_i\}$ , where  $i = 1, 2, \dots, m$ . In such a scenario, the prompt  $\tilde{\mathbb{C}}$  becomes not only semantically redundant but also excessively long, potentially impeding the effective semantic learning of the document tokens.

Thus, we reframe the SER task with a BIO tagging scheme into two separate tasks: *entity typing* and *span detection*. Entity typing involves assigning an entity type for each document token, while span detection aims to identify whether tokens are at the beginning or interior of an entity span.

To further clarify, consider an SER task using BIO tagging scheme with a predefined entity type set  $\mathbb{E} = \{e_i | e_0 = \text{Other}\}$ ,  $i = 0, 1, 2, \dots, m$ . For entity typing, the label space and the tag-related prompt would be  $\mathbb{C}^{ent.} = \{c_i^{ent.} | c_i^{ent.} = e_i\}$  and  $\tilde{\mathbb{C}}^{ent.} = \{\tilde{c}_i^{ent.} | \tilde{c}_i^{ent.} = c_i^{ent.}\}$ , where  $i = 0, 1, 2, \dots, m$ ; And for span detection, the label space and the prompt would be  $\mathbb{C}^{det.} = \{c_1^{det.}, c_2^{det.}\}$  and  $\tilde{\mathbb{C}}^{det.} = \{\tilde{c}_1^{det.}, \tilde{c}_2^{det.}\}$ , where  $\mathbb{C}^{det.} = \tilde{\mathbb{C}}^{det.} = \{\text{beginning}, \text{inner}\}$ ; Then the full label space and prompt would be  $\mathbb{C} = \mathbb{C}^{ent.} \cup \mathbb{C}^{det.} = \{c_i^{ent.}, c_j^{det.}\}$  and  $\tilde{\mathbb{C}} = \tilde{\mathbb{C}}^{ent.} \cup \tilde{\mathbb{C}}^{det.} = \{\tilde{c}_i^{ent.}, \tilde{c}_j^{det.}\}$ , where  $i = 0, 1, 2, \dots, m$ ;  $j = 1, 2$ . For a token with an entity type of  $e_i (i \neq 0)$  located at the *beginning/inner* of an entity span, the corresponding labels would be  $c_i$  for entity typing and *beginning/inner* for span detection. However,

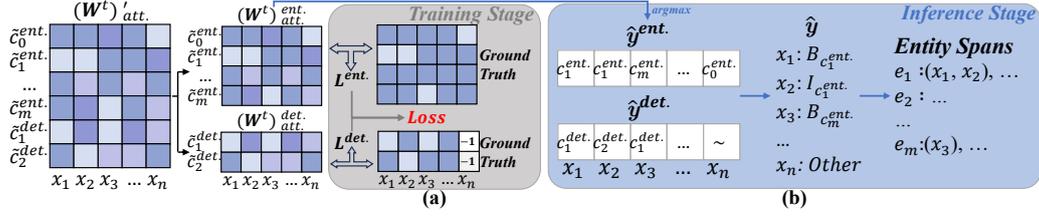


Figure 3: (a) PPTSER at training stage. Losses of entity typing and span detection are computed separately and then combined for the overall loss calculation. And  $-1$  signifies that the loss at those points is disregarded. (b) PPTSER at inference stage. Combined predictions of entity typing and span detection are utilized to analyze the entity spans.

for the token with an entity type of *Other*, the specific location of it within an entity span is irrelevant, and the loss for span detection is ignored here. Consequently, we can formulate the entity typing target  $\mathbf{y}^{ent.} = [y_i^{ent.}]$  and the span detection target  $\mathbf{y}^{det.} = [y_i^{det.}]$ , where  $i = 1, 2, \dots, n$ .

It is worth emphasizing that our PPTSER framework handles entity typing and span detection simultaneously. And prompts for them  $\tilde{\mathbf{C}}^{ent.}$  and  $\tilde{\mathbf{C}}^{det.}$  are encoded in parallel, allowing them to benefit from each other during the learning process.

### 2.2.2 Cross-attention within the Pre-trained Model

Once the tag-related prompt  $\tilde{\mathbf{C}}$  is constructed, it is concatenated with the document tokens  $\mathbf{x} = [x_i], i = 1, 2, \dots, n$ , forming a boosted input  $\mathbf{x}' = \mathbf{x} \oplus \tilde{\mathbf{C}} = [x_i, \tilde{c}_j^{ent.}, \tilde{c}_k^{det.}], i = 1, 2, \dots, n; j = 0, 1, 2, \dots, m; k = 1, 2, \dots, m$ . Then,  $\mathbf{x}'$  is used as the *Text Embedding* encoded in the pre-trained model. Notably, when  $\tilde{\mathbf{C}}$  involves other kinds of embedding, such as positional or visual embeddings, they are set to 0, since  $\tilde{\mathbf{C}}$  are hypothetical tokens not found in the document. Let's denote the hidden states from the second last block as  $\mathbf{H}^{t-1}$ :

$$\mathbf{H}^{t-1} = [\mathbf{h}_i^{t-1}, \tilde{\mathbf{h}}_j^{t-1}, \tilde{\mathbf{h}}_k^{t-1}] \quad (1)$$

where  $\mathbf{h}_i^{t-1}, \tilde{\mathbf{h}}_j^{t-1}, \tilde{\mathbf{h}}_k^{t-1}$  are the hidden states for  $\mathbf{x}, \tilde{\mathbf{C}}^{ent.}, \tilde{\mathbf{C}}^{det.}$ , correspondingly.

Then,  $\mathbf{H}^{t-1}$  is partitioned into multiple segments  $\mathbf{H}_i^{t-1}$  along the channel dimension, where queries  $\mathbf{Q}_i^t$  and keys  $\mathbf{K}_i^t$  of the  $i^{th}$  attention head are transformed as follows:

$$\mathbf{Q}_i^t = (\mathbf{W}_i^t)_q \mathbf{H}_i^{t-1} \quad (2)$$

$$\mathbf{K}_i^t = (\mathbf{W}_i^t)_k \mathbf{H}_i^{t-1} \quad (3)$$

where  $(\mathbf{W}_i^t)_q$  and  $(\mathbf{W}_i^t)_k$  are learnable weights embedded in the last attention block. And the self-attention weight of distinct heads is computed as below:

$$(\mathbf{W}_i^t)_{att.} = \mathbf{Q}_i^t (\mathbf{K}_i^t)^T \quad (4)$$

where  $(\mathbf{W}_i^t)_{att.}$  is a matrix with the shape of  $(n + m + 3) \times (n + m + 3)$ . From this matrix, we extract a sub-matrix  $(\mathbf{W}_i^t)'_{att.}$  that takes the prompt as queries and the document tokens as keys, which possesses the shape of  $(m+3) \times n$ .  $(\mathbf{W}_i^t)'_{att.}$  can be viewed as a form of cross-attention within the self-attention, which depicts the relationship between the tag-related prompt and document tokens.

We hypothesize that distinct heads of the attention mechanism enable the prompt to focus on distinct entity spans, which is suitable for the entity-rich scenario in visually-rich documents. We select the maximum weight across heads to get a summary relationship between the prompt and tokens:

$$(\mathbf{W}^t)'_{att.} = \max_{i \in \{1, 2, \dots, l\}} (\mathbf{W}_i^t)'_{att.} \quad (5)$$

Further,  $(\mathbf{W}^t)'_{att.}$  is partitioned into two components, namely  $(\mathbf{W}^t)_{att.}^{ent.}$  and  $(\mathbf{W}^t)_{att.}^{det.}$  as shown in Figure 3(a). These components use the hidden states of  $\tilde{\mathbf{C}}^{ent.}$  and  $\tilde{\mathbf{C}}^{det.}$  as queries, and possess the shape of  $(m+1) \times n$  and  $2 \times n$ , correspondingly.  $(\mathbf{W}^t)_{att.}^{ent.}$  and  $(\mathbf{W}^t)_{att.}^{det.}$  represent the probability distribution for document tokens belonging to distinct tags. The losses are then calculated as follows:

$$L^{ent.} = -\frac{1}{n} \sum_{i=1}^n \frac{\exp(w_{pi}^{ent.})}{\sum_{j=0}^m \exp(w_{ji}^{ent.})} \quad (6)$$

$$L^{det.} = -\frac{1}{n} \sum_{i=1}^n \frac{\exp(w_{qi}^{det.})}{\sum_{j=1}^2 \exp(w_{ji}^{det.})} \quad (7)$$

where  $L^{ent.}$  and  $L^{det.}$  are the losses for entity typing and span detection,  $w_{ij}^{ent.}$  and  $w_{ij}^{det.}$  are elements in  $(\mathbf{W}^t)_{att.}^{ent.}$  and  $(\mathbf{W}^t)_{att.}^{det.}$ , and  $y_i^{ent.} = c_p^{ent.}, y_i^{det.} = c_q^{det.}$ . And the total loss is formulated as follows:

$$Loss = L^{ent.} + \alpha L^{det.} \quad (8)$$

Here,  $\alpha$  is the ratio factor to balance the losses, and we set  $\alpha = 0.1$  for models with segment-level positional embeddings and  $\alpha = 1.5$  for models with word-level positional embeddings.

Modality		Text + Layout				Text + Layout + Image			
Methodology		BROS (AAAI 22)		LiLT (ACL 22)		LayoutLMv2 (ACL 21)		LayoutLMv3 (MM 22)	
		FT	Ours	FT	Ours	FT	Ours	FT	Ours
FUNSD	1-shot	48.08	<b>54.39</b> $\uparrow 6.31$	52.60	<b>55.64</b> $\uparrow 3.04$	48.22	<b>52.17</b> $\uparrow 3.95$	46.37	<b>61.98</b> $\uparrow 15.61$
	3-shot	64.34	<b>67.70</b> $\uparrow 3.36$	67.64	<b>69.17</b> $\uparrow 1.52$	61.66	<b>63.64</b> $\uparrow 1.98$	74.73	<b>76.86</b> $\uparrow 2.13$
	5-shot	67.77	<b>70.64</b> $\uparrow 2.87$	73.29	<b>75.26</b> $\uparrow 1.97$	65.86	<b>67.49</b> $\uparrow 1.63$	79.52	<b>81.53</b> $\uparrow 2.01$
	7-shot	68.21	<b>71.96</b> $\uparrow 3.75$	73.39	<b>75.71</b> $\uparrow 2.32$	66.55	<b>68.83</b> $\uparrow 2.28$	79.84	<b>81.60</b> $\uparrow 1.76$
	Full Data	83.83	<b>83.91</b> $\uparrow 0.08$	88.95	<b>89.07</b> $\uparrow 0.12$	83.52	<b>83.72</b> $\uparrow 0.20$	91.15	<b>92.01</b> $\uparrow 0.86$
CORD	1-shot	66.28	<b>68.48</b> $\uparrow 2.20$	70.04	<b>75.57</b> $\uparrow 5.54$	69.61	<b>69.97</b> $\uparrow 0.36$	70.35	<b>74.19</b> $\uparrow 3.84$
	3-shot	79.02	<b>79.61</b> $\uparrow 0.59$	81.64	<b>83.83</b> $\uparrow 2.19$	80.63	<b>81.66</b> $\uparrow 1.03$	82.05	<b>85.27</b> $\uparrow 3.22$
	5-shot	84.04	<b>84.37</b> $\uparrow 0.34$	85.52	<b>87.06</b> $\uparrow 1.54$	84.32	<b>84.53</b> $\uparrow 0.21$	85.83	<b>87.77</b> $\uparrow 1.94$
	7-shot	83.68	<b>84.09</b> $\uparrow 0.42$	85.35	<b>87.73</b> $\uparrow 2.38$	84.76	<b>85.31</b> $\uparrow 0.55$	86.94	<b>88.48</b> $\uparrow 1.54$
	Full Data	95.72	<b>95.75</b> $\uparrow 0.03$	95.80	<b>96.04</b> $\uparrow 0.25$	95.20	<b>95.63</b> $\uparrow 0.44$	96.34	<b>96.39</b> $\uparrow 0.05$
XFUND-zh	1-shot	-	-	60.10	<b>67.64</b> $\uparrow 7.54$	60.28	<b>68.26</b> $\uparrow 7.98$	52.92	<b>56.65</b> $\uparrow 3.73$
	3-shot	-	-	72.61	<b>74.17</b> $\uparrow 1.56$	74.37	<b>77.20</b> $\uparrow 2.83$	69.08	<b>75.24</b> $\uparrow 6.16$
	5-shot	-	-	77.40	<b>79.40</b> $\uparrow 2.00$	81.43	<b>82.34</b> $\uparrow 0.91$	75.25	<b>79.26</b> $\uparrow 4.01$
	7-shot	-	-	80.47	<b>81.38</b> $\uparrow 0.91$	82.25	<b>83.66</b> $\uparrow 1.41$	77.85	<b>80.97</b> $\uparrow 3.12$
	Full Data	-	-	90.47	<b>90.61</b> $\uparrow 0.14$	90.25	<b>90.79</b> $\uparrow 0.54$	91.61	<b>92.19</b> $\uparrow 0.58$

Table 1: F1 score (%) of PPTSER and traditional Fine-tuning methods. F1 score in **Bold** is better between our PPTSER and Fine-tuning. *FT* refers to *Fine-tuning* methods.

### 2.2.3 Decoding during the Inference Stage

The inference stage is shown in Figure 3(b). We first apply the  $\operatorname{argmax}$  operation on  $(\mathbf{W}^t)_{att}^{ent.}$  and  $(\mathbf{W}^t)_{att}^{det.}$  along distinct prompt words to get the predicted tag with the highest probability:

$$\hat{y}_i^{ent.} = \operatorname{argmax}_{j \in \{0,1,2,\dots,m\}} w_{ji}^{ent.} \quad (9)$$

$$\hat{y}_i^{det.} = \operatorname{argmax}_{j \in \{1,2\}} w_{ji}^{det.} \quad (10)$$

Then the prediction with BIO tagging scheme  $\hat{\mathbf{y}} = [\hat{y}_i], i = 1, 2, \dots, n$  is formulated as follows:

$$\hat{y}_i = \begin{cases} B_{\hat{y}_i^{ent.}} & , \hat{y}_i^{ent.} \neq Other, \hat{y}_i^{det.} = beginning \\ I_{\hat{y}_i^{ent.}} & , \hat{y}_i^{ent.} \neq Other, \hat{y}_i^{det.} = inner \\ Other & , \hat{y}_i^{ent.} = Other \end{cases} \quad (11)$$

Finally, the entity spans are analyzed from  $\hat{\mathbf{y}}$  using the BIO tagging scheme. Notably, spans not conform to BIO scheme, especially those starting with a token predicted as  $\hat{y}_i^{det.} = inner$ , are labeled as *Other*. This operation, aimed at enhancing predicting accuracy, is applied in both PPTSER and methods we compared to for a fair comparison.

To provide a more vivid demonstration of our method, we offer an example run with PPTSER, refer to Appendix C due to paper-length constraints.

## 3 Experiments

### 3.1 Experimental Settings

**Benchmarks.** We conducted experiments on several widely used SER benchmarks, including FUNSD (Jaume et al., 2019), CORD (Park et al., 2019) and XFUND (Xu et al., 2022). FUNSD

targets form understanding with 199 scanned documents related to market reports, commercials, and more. CORD, centered on receipt understanding, features both coarse (e.g., *menu*, *total*) and fine-grained (e.g., *menu.unitprice*, *menu.price*) annotations. This benchmark provides an official split of training, validation and test sets, and we strictly follow the procedure by selecting the model weight that achieved the best performance on the validation set for testing on the test set. XFUND focuses on document understanding covering multiple languages. In this article, our primary focus is on the Chinese subset of XFUND, denoted as XFUND-zh.

**Few-shot Settings.** PPTSER was evaluated on 1-shot, 3-shot, 5-shot, 7-shot and the full training set scenarios. With no official few-shot divisions in benchmarks mentioned above, we established our own following the process in Appendix A. We selected as few samples as possible while meeting the few-shot setting, which aligns with the real-world application. Due to the inherent instability of few-shot experiments, we randomly generated 5 different divisions for every scenario and tested each with 2 diverse random seeds. Hence, our experiment result is the average of 10 runs, ensuring the reliability and credibility of our findings.

### 3.2 Comparisons with Existing Fine-Tuning Methods

**Setup.** The foundation for our method is built upon several widely used multi-modal pre-trained models, incorporating different combinations of modalities as input. This includes **BROS** (Hong et al., 2022) and **LiLT** (Wang et al., 2022a) with textual and layout input, and **LayoutLMv2** (Xu et al.,

Modality		Text		Text + Layout			Text + Layout + Image	
Methodology		EntLM (NAACL 22)	COPNER (COLING 22)	LASER (ACL 22)	COPNER <sub>LiLT</sub> (COLING 22)	PPTSER <sub>LiLT</sub> (Ours)	COPNER <sub>LMv3</sub> (COLING 22)	PPTSER <sub>LMv3</sub> (Ours)
FUNSD	1-shot	24.32	19.37	38.47	55.15	<u>55.64</u>	51.19	<b>61.98</b>
	3-shot	34.94	31.21	44.88	68.66	69.17	<u>75.84</u>	<b>76.86</b>
	5-shot	39.55	35.13	49.31	73.43	75.26	<u>77.55</u>	<b>81.53</b>
	7-shot	41.41	37.31	52.56	73.35	75.71	<u>78.53</u>	<b>81.60</b>
	Full Data	67.42	64.58	69.23	87.74	89.07	<u>91.26</u>	<b>92.01</b>
CORD-Lv1	1-shot	74.29	68.61	66.80	86.97	<b>90.50</b>	86.98	<u>90.02</u>
	3-shot	83.68	82.25	76.09	94.16	<u>94.79</u>	94.03	<b>95.13</b>
	5-shot	87.11	86.08	82.23	94.86	<b>96.21</b>	<u>95.74</u>	<b>96.21</b>
	7-shot	87.31	86.74	83.61	95.04	<u>96.13</u>	96.06	<b>96.51</b>
	Full Data	95.93	95.90	96.56	99.21	<u>99.42</u>	<b>99.45</b>	<b>99.45</b>
CORD	1-shot	57.86	54.52	-	70.05	<b>75.57</b>	67.33	<u>74.19</u>
	3-shot	71.68	71.32	-	81.27	<u>83.83</u>	80.07	<b>85.27</b>
	5-shot	77.74	78.98	-	84.80	<u>87.06</u>	85.30	<b>87.77</b>
	7-shot	78.63	78.63	-	85.76	<u>87.73</u>	86.87	<b>88.48</b>
	Full Data	93.50	94.16	-	95.74	<u>96.04</u>	95.79	<b>96.39</b>
XFUND-zh	1-shot	26.38	23.29	-	48.76	<b>67.64</b>	54.26	<u>56.71</u>
	3-shot	37.22	37.49	-	64.59	<u>74.17</u>	71.27	<b>75.24</b>
	5-shot	43.54	44.36	-	69.03	<b>79.40</b>	76.37	<u>79.26</u>
	7-shot	46.62	46.90	-	74.44	<b>81.38</b>	79.29	<u>80.97</u>
	Full Data	66.20	67.11	-	89.17	90.61	<u>91.99</u>	<b>92.19</b>

Table 2: F1 score (%) of PPTSER and other Few-shot methods. F1 score in **Bold** is the best, and that with underline is the second best.

2021a) and **LayoutLMv3** (Huang et al., 2022b) with textual, layout and image input. Since BROS only supports English, we only tested it on FUNSD and CORD. For testing on XFUND-zh, we used **LayoutXLM** (Xu et al., 2021b), which is the multilingual version of LayoutLMv2. In our experiments, we utilized base-size pre-trained models.

**Results.** Table 1 showcases the performance of PPTSER against traditional fine-tuning methods. The results clearly demonstrate that our PPTSER outperforms traditional fine-tuning methods across all tested scenarios and benchmarks. This underscores the superior performance of PPTSER in diverse language contexts with various base models.

Overall, both PPTSER and the fine-tuning method demonstrate improved performance with increased training data. However, our PPTSER consistently outperforms previous fine-tuning methods in all few-shot settings, especially with exceptionally scarce data. In the 1-shot scenario on FUNSD, where only a single annotated document is available, PPTSER achieves gains of **+6.31%** with BROS, **+3.04%** with LiLT, **+3.95%** with LayoutLMv2 and the highest gain of **+15.62%** with LayoutLMv3, emphasizing its effectiveness in data-scarce situations. Notably, even when trained with the full training data, our PPTSER still achieves comparable performance to the fine-tuning method, and even outperforms it in certain scenarios. For example, we observe a gain of **+0.86%** on FUNSD with LayoutLMv3. This full data setting is often neglected in other few-shot research, further un-

derscoring the superiority of our approach when dealing with varying amounts of available data.

Our findings demonstrate that PPTSER is highly adaptive to different amounts of training data with distinct base models, making it an effective method for addressing the SER problem.

### 3.3 Comparisons with Existing Few-shot Methods

**Setup.** We selected the PPTSER models that performed better under different modality settings, denoted as **PPTSER<sub>LiLT</sub>** and **PPTSER<sub>LMv3</sub>**, and compared them with previous few-shot methods. For a comprehensive comparison, we re-implemented **LASER** (Wang and Shang, 2022) on our few-shot divisions. However, it can only handle the coarse-level typing for CORD (CORD-Lv1) and is limited to the English language. Since research on few-shot SER is rather limited, we selected two other few-shot NER methods for comparison. Specifically, We chose **COPNER** (Huang et al., 2022a) and **EntLM** (Ma et al., 2022b) due to their similar In-Label-Space setting with ours. Considering COPNER can also be used as a pluggable method, we also integrated it with LayoutLMv3 and LiLT, denoted as **COPNER<sub>LiLT</sub>** and **COPNER<sub>LMv3</sub>**.

**Results.** The overall experimental results are presented in Table 2. The results clearly show that PPTSER outshines existing few-shot NER and SER methods by a large margin. Interestingly, COPNER shows some degree of pluggability with various multi-modal pre-trained models, but PPTSER still

Model Designs	Baseline	Decoupling Strategies	Prompt Engineering		Aggregation Strategies	
	<i>default setting</i>	<i>plain BIO prompt</i>	<i>unrelated words</i>	<i>random embeddings</i>	<i>mean</i>	<i>single head</i>
1-shot	75.57	74.28	70.80	71.43	74.43	75.21
3-shot	83.83	82.41	82.69	83.07	83.05	83.92
5-shot	87.06	86.25	85.78	86.29	86.64	86.98
7-shot	87.73	86.49	86.42	86.95	86.74	87.24
Full Data	96.04	95.29	96.26	95.72	96.21	96.06
$\Delta$	-	<b>-1.10</b>	<b>-1.66</b>	<b>-1.35</b>	<b>-0.63</b>	<b>-0.17</b>

Table 3: F1 score (%) of PPTSER on CORD benchmark with different designs.  $\Delta$  denotes the average deviation of the F1 score relative to the *default setting*.

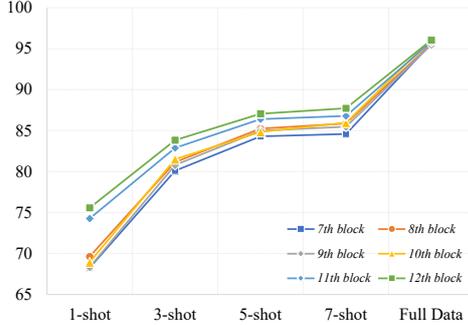


Figure 4: F1 score (%) of PPTSER on CORD benchmark with different settings when obtaining the attention weight from different blocks.

outperforms it across all settings and benchmarks.

In summary, our PPTSER surpasses existing few-shot NER and few-shot SER methods on various visually-rich documents, showcasing its effectiveness in handling few-shot SER challenge.

## 4 Ablation Study

We have conducted extensive analyses of our PPTSER to ensure its effectiveness and rationality. For convenience, experiments are conducted on the CORD dataset using PPTSER building upon LiLT.

**Origin of Attention Weights.** To pinpoint the source of superiority in PPTSER, we explored whether it stems from our meticulous design or the decreased over-fitting achieved through parameter reduction. We extracted attention weights from various blocks, including the default 12<sup>th</sup> block and shallower 7<sup>th</sup> ~ 11<sup>th</sup> block. And the experimental results shown in Figure 4 reveal that extracting attention weights from the last block is more effective than from other blocks, which has greatly assured the effectiveness of our design.

**Effectiveness of Decoupling Strategies.** Table 3 also shows comparisons with different frameworks of prompts. In this context, *default setting* refers to our design to decouple the SER task into entity typing and span detection then processing them concurrently, while *plain BIO prompt* refers to the direct usage of the aforementioned  $\tilde{C} = \{e_0, \text{beginning of } e_i, \text{inner of } e_i\}$  as the

prompt, without decoupling. The result shows that our decoupling avoids disrupting the language modeling of document tokens and performs better.

**Prompt Engineering.** We evaluated PPTSER using diverse prompt types. Beyond the *default setting* that uses tag names as prompts, we explored an *unrelated words* setting by replacing the whole prompts with irrelevant words like apple and orange. In the *random embeddings* setting, we replaced the whole prompt’s text embedding with random tensors. As Table 3 indicates, the default setting yields the highest score, suggesting the pre-trained model does grasp the prompt’s semantics, and tag semantics can direct the SER task. This indicates the careful selection and design of prompts can markedly influence model performance.

**Aggregation Strategies of Attention Weights.** Table 3 compares how various strategies to aggregate attention weights across different heads affect performance. While *default setting* and *mean* refer to obtain the maximum and average value across attention weights of distinct heads, *single head* uses just a single head of the attention weights to generate the final probability. The results indicate that the *max* operation outperforms others, which aligns with our hypothesis that different attention heads focus on entities with different semantics.

**Parameter Efficiency.** The parameter comparisons of our PPTSER methods and traditional fine-tuning are presented in Table 4. As the parameters might vary across diverse models and benchmarks, we offer a concise breakdown of the results from the methods we have tested. The results illustrate that our PPTSER has fewer parameters in comparison to traditional fine-tuning methods. For a more detailed analysis, please refer to Appendix E.

## 5 Related Works

**SER on Visually-rich Documents.** Although some early works of SER relied on heuristic algorithms (Simon et al., 1997; Schuster et al., 2013), the majority of research focused on neural network-based methods. Some of them leveraged textual

Methodology	BROS (AAAI 22)		LiLT (ACL 22)		LayoutLMv2 (ACL 21)		LayoutLMv3 (MM 22)	
	FT	Ours	FT	Ours	FT	Ours	FT	Ours
FUNSD	108.91M	<b>103.59M</b>	130.17M	<b>123.81M</b>	200.29M	<b>194.38M</b>	125.33M	<b>119.42M</b>
CORD	108.95M	<b>103.59M</b>	130.22M	<b>123.81M</b>	200.33M	<b>194.38M</b>	125.96M	<b>119.42M</b>
XFUND-zh	-	-	130.17M	<b>123.81M</b>	200.29M	<b>194.38M</b>	125.33M	<b>119.42M</b>

Table 4: Parameters of our PPTSER and traditional Fine-tuning methods. The metric in **Bold** indicates the method with fewer parameters. *FT* refers to *Fine-tuning* method.

features (Chiu and Nichols, 2016), visual features (Guo et al., 2019), or combined them with layout features (Yu et al., 2021; Wang et al., 2021a) to address this issue, but the emergence of multi-modal pre-trained models has revolutionized SER. These models are jointly pre-trained on a large-scale unlabeled document dataset with textual, layout, and even visual cues, so they have the potential to better understand a structured document. LayoutLM (Xu et al., 2020) was the first to combine textual and OCR positional features at the pre-training stage. Later, LayoutLMv2 (Xu et al., 2021a) and LayoutLMv3 (Huang et al., 2022b) further integrated visual features into the pre-training process with different architectures. Moreover, Wang et al. (2022a) advanced the model architecture with a language-agnostic layout transformer in their work, LiLT. Alongside the advancements in model structures, other works (Appalaraju et al., 2021; Li et al., 2021b,a; Hong et al., 2022; Luo et al., 2023) have focused on the diverse pre-training tasks to facilitate the fusion of diverse modalities at pre-training stage. While these advancements have improved SER capabilities to some extent, their few-shot learning abilities still require further examination.

**Few-shot SER on Visually-rich Documents.** Unlike SER, few-shot SER is not fully explored yet. Cheng et al. (2020) proposed to utilize graph-matching techniques (Zanfir and Sminchisescu, 2018), representing documents as graphs with nodes as OCR-scanned boxes. For an unseen document, the type of entities was determined by comparing the relationships in the graph of unseen document with those in the graphs of support documents. Yao et al. (2021) also adopted a graph-matching approach to address this issue, but the entity type was determined based on the relationships in different forms with more complex solvers. Taking a different way, Wang and Shang (2022) introduced a novel labeling scheme for SER. They reshaped SER as a generative task, and used LayoutReader (Wang et al., 2021b) for SER label generation. Although these studies preliminary explored few-shot SER, they lacked generality and pluggability, and their performances in various scenes

require further exploration and improvement.

**Few-shot NER in Plain Texts.** While few-shot SER on visually-rich documents has seen limited exploration, there has been extensive research on few-shot Named Entity Recognition (NER) in plain texts (Wang et al., 2022b; Das et al., 2022; Ma et al., 2022a; Cheng et al., 2023). However, only few of these studies have considered the scenario where only limited data in the target domain is available. Huang et al. (2022a) proposed using the NER tag as a prompt and employing contrastive learning to address this issue. On the other hand, Ma et al. (2022b) reformulated few-shot NER as a Language Modeling task and used the pre-trained Masked Language Model head to predict a word related to the entity type for each text token. However, since these methods are designed to address the NER problem with sparse entities in plain texts, they might not be suitable for entity-rich scenarios in visually-rich documents. Additionally, they do not emphasize the issue of detecting entity boundaries, without which adjacent entities of the same type might be erroneously merged into one. The comparison with two representative few-shot NER methods also show the effectiveness of our method.

## 6 Conclusion

In this paper, we present PPTSER, an innovative and efficient strategy for few-shot entity recognition on visually-rich documents using a plug-and-play, tag-guided approach. PPTSER redefines the SER task as a dual-function operation of entity typing and span detection, and utilizes the attention weight between document tokens and prompts related to SER tags as the target probability distributions. Our findings show that PPTSER is both effective and versatile in various data situations, from few-shot to full data scenarios. In the future, we plan to further investigate the capabilities of PPTSER across a range of VIE tasks like Entity Linking. In addition, we aim to explore PPTSER’s potential in other few-shot scenarios, particularly those outside of the In-Label-Space setting. It is our hope that our work will spark further research and advancements in the realm of few-shot SER.

582  
583  
584  
585  
586  
587  
588  
  
589  
  
590  
591  
592  
593  
594  
595  
596  
597  
  
598  
  
599  
600  
601  
602  
  
603  
604  
605  
606  
  
607  
608  
609  
610  
  
611  
612  
613  
614  
  
615  
616  
617  
618  
  
619  
620  
621  
622  
623  
  
624  
625  
626  
627  
  
628  
629  
630

## Limitations

Due to space constraints, our exploration of the few-shot SER setting is primarily limited to the In-Label-Space. Future research is essential to investigate the potential applications of our PPTSER in other few-shot settings and its adaptability to additional VIE tasks.

## Ethical Considerations

Our proposed PPTSER is a purely methodological innovation, which inherently avoids direct negative social impacts. By leveraging the self-attention mechanism within multi-modal pre-trained models without adding extra modules, it does not introduce additional ethical risks beyond those already present in the existing multi-modal pre-trained models.

## References

Srikar Appalaraju, Bhavan Jasani, Bhargava Urala Kota, Yusheng Xie, and R. Manmatha. 2021. Docformer: End-to-end transformer for document understanding. In *Proc. ICCV*, pages 993–1003.

Mengli Cheng, Minghui Qiu, Xing Shi, Jun Huang, and Wei Lin. 2020. One-shot text field labeling using attention and belief propagation for structure information extraction. In *Proc. ACM MM*, page 340–348.

Zifeng Cheng, Qingyu Zhou, Zhiwei Jiang, Xuemin Zhao, Yunbo Cao, and Qing Gu. 2023. Unifying token- and span-level supervisions for few-shot sequence labeling. *ACM Trans. Inf. Syst.*, 42(1).

Jason P.C. Chiu and Eric Nichols. 2016. Named Entity Recognition with Bidirectional LSTM-CNNs. *Transactions of the Association for Computational Linguistics*, 4:357–370.

Sarkar Snigdha Sarathi Das, Arzoo Katiyar, Rebecca Passonneau, and Rui Zhang. 2022. CONTaiNER: Few-shot named entity recognition via contrastive learning. In *Proc. ACL*, pages 6338–6353.

Jiuxiang Gu, Jason Kuen, Vlad I Morariu, Handong Zhao, Rajiv Jain, Nikolaos Barmpalios, Ani Nenkova, and Tong Sun. 2021. Unidoc: Unified pretraining framework for document understanding. In *Proc. NIPS*, volume 34, pages 39–50.

He Guo, Xiameng Qin, Jiaming Liu, Junyu Han, Jingtuo Liu, and Errui Ding. 2019. Eaten: Entity-aware attention for single shot visual text extraction. In *Proc. ICDAR*, pages 254–259.

Teakgyu Hong, DongHyun Kim, Mingi Ji, Wonseok Hwang, Daehyun Nam, and Sungrae Park. 2022. Bros: A pre-trained language model focusing on text

and layout for better key information extraction from documents. *Proc. AAAI*, 36(10):10767–10775. 631 632

Yucheng Huang, Kai He, Yige Wang, Xianli Zhang, Tieliang Gong, Rui Mao, and Chen Li. 2022a. COP-NER: Contrastive learning with prompt guiding for few-shot named entity recognition. In *Proc. COLING*, pages 2515–2527. 633 634 635 636 637

Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and Furu Wei. 2022b. Layoutlmv3: Pre-training for document ai with unified text and image masking. In *Proc. ACM MM*, page 4083–4091. 638 639 640 641

Guillaume Jaume, Hazim Kemal Ekenel, and Jean-Philippe Thiran. 2019. Funsd: A dataset for form understanding in noisy scanned documents. In *Workshop at ICDAR*, volume 2, pages 1–6. 642 643 644 645

Chenliang Li, Bin Bi, Ming Yan, Wei Wang, Songfang Huang, Fei Huang, and Luo Si. 2021a. StructuralLM: Structural pre-training for form understanding. In *Proc. ACL*, pages 6309–6318. 646 647 648 649

Peizhao Li, Jiuxiang Gu, Jason Kuen, Vlad I. Morariu, Handong Zhao, Rajiv Jain, Varun Manjunatha, and Hongfu Liu. 2021b. Selfdoc: Self-supervised document representation learning. In *Proc. CVPR*, pages 5652–5660. 650 651 652 653 654

Yulin Li, Yuxi Qian, Yuechen Yu, Xiameng Qin, Chengquan Zhang, Yan Liu, Kun Yao, Junyu Han, Jingtuo Liu, and Errui Ding. 2021c. Structext: Structured text understanding with multi-modal transformers. In *Proc. ACM MM*, page 1912–1920. 655 656 657 658 659

Chen Liang, Yue Yu, Haoming Jiang, Siawpeng Er, Ruijia Wang, Tuo Zhao, and Chao Zhang. 2020. Bond: Bert-assisted open-domain named entity recognition with distant supervision. In *Proc. ACM SIGKDD Int. Conf. Knowledge discovery & data mining*, page 1054–1064. 660 661 662 663 664 665

Chuwei Luo, Changxu Cheng, Qi Zheng, and Cong Yao. 2023. Geolayoutlm: Geometric pre-training for visual information extraction. In *Proc. CVPR*, pages 7092–7101. 666 667 668 669

Jie Ma, Miguel Ballesteros, Srikanth Doss, Rishita Anubhai, Sunil Mallya, Yaser Al-Onaizan, and Dan Roth. 2022a. Label semantics for few shot named entity recognition. In *Findings of ACL*, pages 1956–1971. 670 671 672 673 674

Ruotian Ma, Xin Zhou, Tao Gui, Yiding Tan, Linyang Li, Qi Zhang, and Xuanjing Huang. 2022b. Template-free prompt tuning for few-shot NER. In *Proc. NAACL-HLT*, pages 5721–5732. 675 676 677 678

Seunghyun Park, Seung Shin, Bado Lee, Junyeop Lee, Jaeheung Surh, Minjoon Seo, and Hwalsuk Lee. 2019. Cord: A consolidated receipt dataset for post-ocr parsing. In *Workshop on Document Intelligence at NeurIPS*. 679 680 681 682 683

684	Lance Ramshaw and Mitch Marcus. 1995. Text chunking using transformation-based learning. In <i>Workshop on Very Large Corpora</i> .	740
685		741
686		742
687	Daniel Schuster, Klemens Muthmann, Daniel Esser, Alexander Schill, Michael Berger, Christoph Weidling, Kamil Aliyev, and Andreas Hofmeier. 2013. Intellix – end-user trained information extraction for document archiving. In <i>Proc. ICDAR</i> , pages 101–105.	743
688		744
689		745
690		746
691		747
692		748
693	Anikó Simon, Jean-Christophe Pret, and A. Peter Johnson. 1997. A fast algorithm for bottom-up document layout analysis. <i>IEEE Trans. Pattern Anal. Mach. Intell.</i> , 19(3):273–277.	749
694		750
695		751
696		752
697	Erik F. Tjong Kim Sang and Fien De Meulder. 2003. Introduction to the CoNLL-2003 shared task: Language-independent named entity recognition. In <i>Proc. NAACL-HLT</i> , pages 142–147.	753
698		754
699		755
700		756
701	Jiapeng Wang, Lianwen Jin, and Kai Ding. 2022a. LiLT: A simple yet effective language-independent layout transformer for structured document understanding. In <i>Proc. ACL</i> , pages 7747–7757.	757
702		758
703		759
704		760
705	Jiapeng Wang, Chongyu Liu, Lianwen Jin, Guozhi Tang, Jiaxin Zhang, Shuaitao Zhang, Qianying Wang, Yaqiang Wu, and Mingxiang Cai. 2021a. Towards robust visual information extraction in real world: New dataset and novel solution. <i>Proc. AAAI</i> , 35(4):2738–2745.	761
706		762
707		
708		
709		
710		
711	Rui Wang, Tong Yu, Handong Zhao, Sungchul Kim, Subrata Mitra, Ruiyi Zhang, and Ricardo Henao. 2022b. Few-shot class-incremental learning for named entity recognition. In <i>Proc. ACL</i> , pages 571–582.	
712		
713		
714		
715		
716	Zilong Wang and Jingbo Shang. 2022. Towards few-shot entity recognition in document images: A label-aware sequence-to-sequence framework. In <i>Findings of ACL</i> , pages 4174–4186.	
717		
718		
719		
720	Zilong Wang, Yiheng Xu, Lei Cui, Jingbo Shang, and Furu Wei. 2021b. LayoutReader: Pre-training of text and layout for reading order detection. In <i>Proc. EMNLP</i> , pages 4735–4744.	
721		
722		
723		
724	Yang Xu, Yiheng Xu, Tengchao Lv, Lei Cui, Furu Wei, Guoxin Wang, Yijuan Lu, Dinei Florencio, Cha Zhang, Wanxiang Che, Min Zhang, and Lidong Zhou. 2021a. LayoutLMv2: Multi-modal pre-training for visually-rich document understanding. In <i>Proc. ACL</i> , pages 2579–2591.	
725		
726		
727		
728		
729		
730	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 <i>Proc. ACM SIGKDD Int. Conf. Knowledge discovery &amp; data mining</i> , page 1192–1200.	
731		
732		
733		
734		
735	Yiheng Xu, Tengchao Lv, Lei Cui, Guoxin Wang, Yijuan Lu, Dinei Florêncio, Cha Zhang, and Furu Wei. 2021b. Layoutxlm: Multimodal pre-training for multilingual visually-rich document understanding. <i>arXiv: Comp. Res. Repository</i> , abs/2104.08836.	
736		
737		
738		
739		
	Yiheng Xu, Tengchao Lv, Lei Cui, Guoxin Wang, Yijuan Lu, Dinei Florencio, Cha Zhang, and Furu Wei. 2022. XFUND: A benchmark dataset for multilingual visually rich form understanding. In <i>Findings of the ACL</i> , pages 3214–3224.	
	Minghong Yao, Zhiguang Liu, Liangwei Wang, Houqiang Li, and Liansheng Zhuang. 2021. One-shot key information extraction from document with deep partial graph matching. <i>arXiv: Comp. Res. Repository</i> , abs/2109.13967.	
	Wenwen Yu, Ning Lu, Xianbiao Qi, Ping Gong, and Rong Xiao. 2021. Pick: Processing key information extraction from documents using improved graph learning-convolutional networks. In <i>Proc. ICPR</i> , pages 4363–4370.	
	Yuechen Yu, Yulin Li, Chengquan Zhang, Xiaoqiang Zhang, Zengyuan Guo, Xiameng Qin, Kun Yao, Junyu Han, Errui Ding, and Jingdong Wang. 2023. Structextv2: Masked visual-textual prediction for document image pre-training. In <i>Proc. ICLR</i> .	
	Andrei Zanfir and Cristian Sminchisescu. 2018. Deep learning of graph matching. In <i>Proc. CVPR</i> , pages 2684–2693.	

## A Few-shot Divisions Generation

To cater to the real-world application scenarios, we have organized our few-shot divisions from the full training set with Algorithm 1. Our goal was to randomly select the minimum number of documents that satisfy the  $N$ -way  $K$ -shot requirement of *each of the  $N$  categories has  $K$  documents containing entities of that category as the support set*.

It is worth noticing that in the context of few-shot SER on visually-rich documents, the few-shot setting of  $N$ -way  $K$ -shot signifies that *each of the  $N$  categories has  $K$  documents containing entities of that category as the support set*, instead of *there are  $K$  entity spans for each of  $N$  entity types as the support set* for the setting of few-shot NER on plain texts.

## B Implementation Details

### B.1 Implementation Details of PPTSER

We used one NVIDIA 3090 to fine-tune our model with AdamW optimizer. The learning rate is  $5e - 5$  with a warm up ratio of 0.1, and we fine-tuned the model for 2000 iterations with a batch size of 8 by default. Besides the default augmentation strategies for images adopted in LayoutLMv2 and LayoutLMv3, we did not employ any additional augmentation strategies.

### B.2 Modification on few-shot NER method for visually-rich documents

In Section 3.3, we mentioned that we adapted two methods originally used for few-shot NER on plain text for few-shot SER on visually-rich documents. We will briefly introduce these modifications.

**COPNER (Huang et al., 2022a)**. The COPNER method employs contrastive learning, feeding both entity label semantics and sentences into a plain text pre-trained language model. This approach uses the hidden state output of the pre-trained model to calculate a contrastive loss between sentence tokens and label semantics, then determining the entity type of tokens. However, the original COPNER could only determine if a token belonged to an entity category, without recognizing boundaries between entities. Therefore, we also improved it with the *entity typing and span detection* framework introduced in our paper. That is, while determining the entity type of tokens, we also input the tokens *beginning* and *inner* into the pre-trained model to detect the entity boundary.

The model’s output hidden state is then used to calculate a contrastive loss between sentence tokens and these *beginning* and *inner* tokens.

Besides, we retained this core process but replaced the original language model pre-trained on pure text with a multi-modal pre-trained model. Experiments show that our use of multi-head cross-attention methods is more suitable for SER tasks on visually-rich documents, especially in Chinese contexts.

**EntLM (Ma et al., 2022b)**. EntLM treats NER as a task of Language Modeling. For testing a few-shot NER dataset on plain text, it first selects a related word for each entity type. Then, using the pre-trained Masked Language Modeling head of BERT, it predicts the probability distribution of each sentence token over these related words, thereby determining the probability distribution of tokens across different entity types.

The selection of related words relies on the distant data obtained from BOND (Liang et al., 2020), which uses BERT and the corpora from Wikipedia to create rough annotations for the NER test set. However, in the realm of visually-rich documents, such rough annotated data is not provided by BOND, and due to the relative abstract expression of SER tags from natural language expressions and the difference between structured documents and natural language expressions, it’s not feasible to obtain rough annotations using corpora from Wikipedia with BERT. Therefore, we directly use the ground truth annotations from the SER test set as distant data to find related words associated with entity types in the SER dataset. Although the experimental results on EntLM might be artificially high due to some exposure to the entity distribution in the test set, our proposed method significantly outperforms others that only accept text modality inputs, including EntLM.

In summary, methods for few-shot NER on plain text may not necessarily transition well to the task of few-shot SER on visually-rich documents. The notable performance of our proposed method in few-shot SER on visually-rich documents further highlights the innovation and contribution of our research.

---

**Algorithm 1:** Few-shot Divisions Generation

---

**Input:** Novel Dataset with the label space  $\mathbb{C} = \{c_1, c_2, \dots, c_N\}$ , full training set  $\mathcal{D}^{full}$

**Output:** N-way K-shot few-shot training set  $\mathcal{D}^{train}$

```
1  $\mathcal{D}^{train} = \{\}$ 
2 Number of documents that contain entities of  $c_i$  in  $\mathcal{D}^{full}$ :  $Q = \{c_1 : 0, c_2 : 0, \dots, c_N : 0\}$ 
3 Document set that contain entities of  $c_i$  in  $\mathcal{D}^{full}$ :  $R = \{c_1 : \{\}, c_2 : \{\}, \dots, c_N : \{\}\}$ 
4 for  $doc_i$  in  $\mathcal{D}^{full}$  do
5     for  $c_j$  in  $\mathbb{C}$  do
6         if  $doc_i$  contain entities of  $c_j$  then
7              $Q[c_j] += 1$ 
8              $R[c_j].append(doc_i)$ 
9         end
10    end
11 end
12  $Q' = sorted(Q, key = lambda x : x[1]) = \{c'_1 : n_1, c'_2 : n_2, \dots, c'_N : n_N\}$  ( $n_1 \leq n_2 \leq \dots \leq n_N$ )
13 Number of documents that contain entities of  $c_i$  in  $\mathcal{D}^{train}$ :  $S = \{c_1 : 0, c_2 : 0, \dots, c_N : 0\}$ 
14 for  $c'_i$  in keys of  $Q'$  do
15     for  $S[c'_i] < K$  do
16         if  $R[c'_i]$  is empty then
17             break
18         end
19         Randomly select a document  $doc_{candidate}$  from  $R[c'_i]$ 
20          $R[c'_i].pop(doc_{candidate})$ 
21         if  $doc_{candidate} \notin \mathcal{D}^{train}$  then
22              $\mathcal{D}^{train}.append(doc_{candidate})$ 
23             for  $c_j \in \mathbb{C}$  do
24                 if  $doc_{candidate}$  contain entities of  $c_j$  then
25                      $S[c_j] += 1$ 
26                 end
27             end
28         else
29             continue
30         end
31     end
32 end
```

---

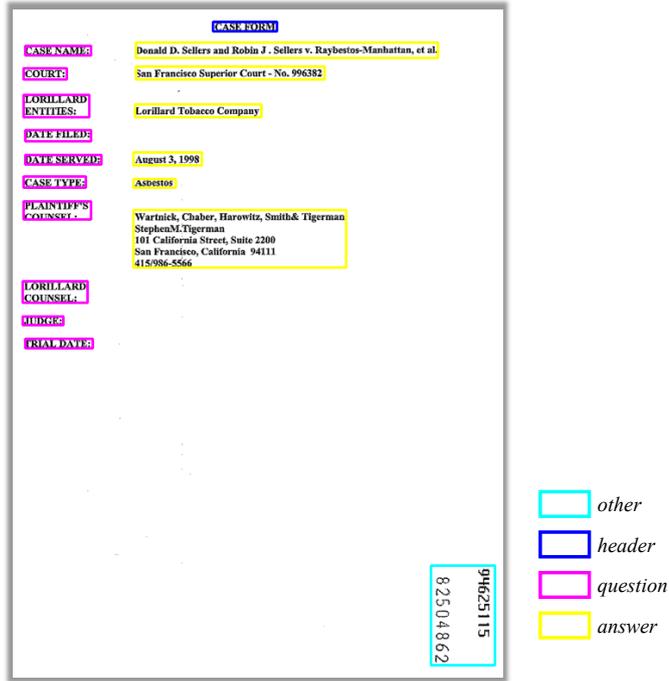


Figure 5: Illustration of a sample from FUNSD dataset. Different colored boxes represent the entities of distinct kinds. Zoom in for better view.

## C Examples when apply PPTSER

To more intuitively demonstrate our method, we provide an example from FUNSD dataset when apply our PPTSER method.

The FUNSD dataset includes three meaningful entity types: *header*, *question*, and *answer*, with all other uninteresting entities categorized as *other*. Hence, the entity type set is  $\mathbb{E} = \{other, header, question, answer\}$ . The label spaces for entity typing and span detection are  $\mathbb{C}^{ent.} = \{other, header, question, answer\}$  and  $\mathbb{C}^{det.} = \{beginning, inner\}$ , respectively. We directly use the label's names as tag-related prompts, with the prompts for entity typing and span detection being  $\tilde{\mathbb{C}}^{ent.} = "other\ header\ question\ answer"$  and  $\tilde{\mathbb{C}}^{det.} = "beginning\ inner"$ . These prompts are then concatenated to form the full tag-related prompt  $\tilde{\mathbb{C}} = "other\ header\ question\ answer\ beginning\ inner"$ .

Consider an example from the FUNSD dataset shown in Figure 5, with the document content "... CASE TYPE: Asbestos ... 82504862", where "..." indicates omitted parts. Here, "CASE TYPE:" belongs to the entity type of *question*, "Asbestos" to the entity type of *answer*, and "82504862" to *other*. Assuming the tokenizer splits the document into "CASE", "TYPE:", "Asbestos", and "82504862", their labels for entity

typing and span detection would be  $\mathbf{y}^{ent.} = [question, question, answer, other]$  and  $\mathbf{y}^{det.} = [beginning, inner, beginning, -1]$ .

Subsequently, "CASE", "TYPE:", "Asbestos", and "82504862" as document tokens are concatenated with the full tag-related prompt, and form the boosted input  $\mathbf{x}' = "CASE\ TYPE:\ Asbestos\ 82504862\ other\ header\ question\ answer\ beginning\ inner"$ . Then,  $\mathbf{x}'$  is input into the multi-modal pre-trained model as the *Text Embedding* to obtain the multi-head attention weight and the aggregated attention weight from the last block, as shown in Figure 6.

During training stage, as shown in Figure 7(a), the aggregated attention weight is split into attention weights between "other header question answer" and document tokens, as well as "beginning inner" and document tokens, which are then used to calculate the losses for entity typing and span detection, respectively, culminating in a combined total loss.

During inference stage, as shown in Figure 7(b), we select the document tokens with the highest probability for "other header question answer" and "beginning inner" as  $\hat{\mathbf{y}}^{ent.}$  and  $\hat{\mathbf{y}}^{det.}$ , then combine them to get the predictions under BIO tagging scheme  $\hat{\mathbf{y}}$  following the procedure in Section 2.2.3. And the entity spans are finally analyzed from  $\hat{\mathbf{y}}$ .

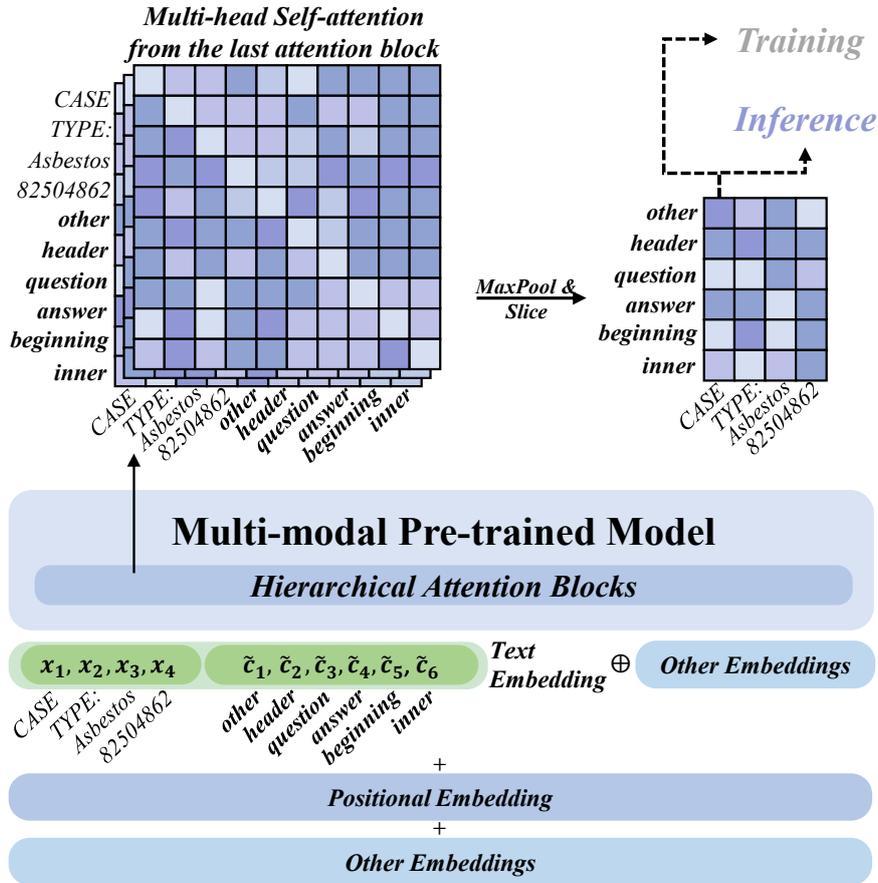


Figure 6: Example of a FUNSD sample run with PPTSER. In the image, **bolded words** in *Text Embeddings* indicate they are part of the prompt. In the Multi-head Self-attention section, the brightness of the color represents the magnitude of the value.

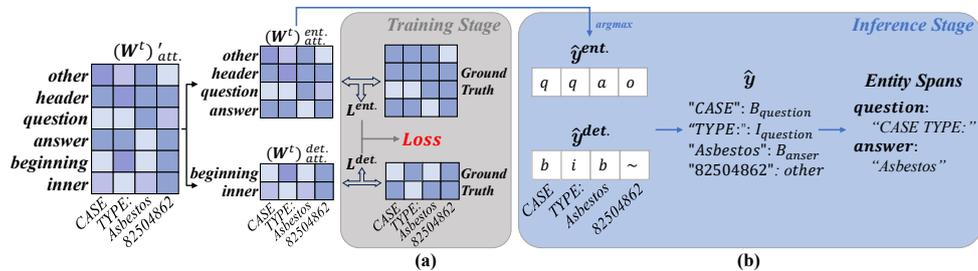


Figure 7: Example of a FUNSD sample running on PPTSER at *Training stage* and *Inference stage*. In the *Inference stage*,  $q$ ,  $a$ ,  $b$ ,  $i$  and  $o$  respectively stand for the *question*, *answer*, *beginning*, *inner* and *other* categories.

Modality		Text + Layout				Text + Layout + Image			
Methodology		BROS (AAAI 22)		LiLT (ACL 22)		LayoutLMv2 (ACL 21)		LayoutLMv3 (MM 22)	
		FT	Ours	FT	Ours	FT	Ours	FT	Ours
FUNSD	1-shot	49.17	<b>52.50</b>	<b>51.07</b>	50.91	44.15	<b>46.91</b>	42.67	<b>56.27</b>
	3-shot	63.07	<b>64.31</b>	65.65	<b>67.24</b>	<b>60.67</b>	60.65	72.66	<b>75.18</b>
	5-shot	66.31	<b>68.10</b>	71.11	<b>72.95</b>	63.45	<b>64.36</b>	77.29	<b>79.53</b>
	7-shot	69.24	<b>71.06</b>	72.49	<b>74.14</b>	65.94	<b>66.79</b>	79.22	<b>81.31</b>
	Full Data	83.42	<b>83.67</b>	88.62	<b>88.89</b>	83.54	<b>83.59</b>	91.41	<b>91.96</b>
CORD	1-shot	64.92	<b>68.04</b>	69.31	<b>75.35</b>	68.08	<b>69.73</b>	70.03	<b>74.02</b>
	3-shot	78.75	<b>79.26</b>	81.63	<b>83.85</b>	79.92	<b>81.35</b>	81.97	<b>85.09</b>
	5-shot	83.86	<b>84.23</b>	85.62	<b>87.03</b>	83.84	<b>84.42</b>	85.76	<b>87.71</b>
	7-shot	83.56	<b>83.76</b>	85.39	<b>87.74</b>	84.40	<b>85.09</b>	86.92	<b>88.37</b>
	Full Data	95.72	<b>95.88</b>	95.82	<b>96.06</b>	94.95	<b>95.64</b>	96.34	<b>96.39</b>
XFUND-zh	1-shot	-	-	59.03	<b>64.90</b>	59.17	<b>67.04</b>	46.64	<b>54.59</b>
	3-shot	-	-	71.33	<b>71.54</b>	73.17	<b>75.21</b>	62.73	<b>72.36</b>
	5-shot	-	-	75.09	<b>77.07</b>	79.21	<b>79.96</b>	69.39	<b>76.57</b>
	7-shot	-	-	77.34	<b>77.70</b>	79.68	<b>80.22</b>	72.03	<b>77.95</b>
	Full Data	-	-	87.92	<b>88.17</b>	88.60	<b>88.95</b>	89.09	<b>91.04</b>

(a) Precision (%) of our PPTSER and Traditional Fine-tuning methods.

Modality		Text + Layout				Text + Layout + Image			
Methodology		BROS (AAAI 22)		LiLT (ACL 22)		LayoutLMv2 (ACL 21)		LayoutLMv3 (MM 22)	
		FT	Ours	FT	Ours	FT	Ours	FT	Ours
FUNSD	1-shot	49.75	<b>57.89</b>	54.39	<b>61.77</b>	53.64	<b>59.15</b>	53.14	<b>70.07</b>
	3-shot	65.82	<b>71.50</b>	69.81	<b>71.35</b>	62.82	<b>67.01</b>	76.94	<b>78.64</b>
	5-shot	69.32	<b>73.43</b>	75.61	<b>77.76</b>	68.50	<b>70.97</b>	81.92	<b>83.65</b>
	7-shot	67.54	<b>72.94</b>	74.33	<b>77.38</b>	67.23	<b>71.04</b>	80.51	<b>81.91</b>
	Full Data	<b>84.26</b>	84.15	<b>89.30</b>	89.26	83.51	<b>83.85</b>	90.91	<b>92.05</b>
CORD	1-shot	67.70	<b>68.95</b>	70.79	<b>75.80</b>	<b>71.23</b>	70.22	70.67	<b>74.35</b>
	3-shot	79.30	<b>79.96</b>	81.65	<b>83.82</b>	81.36	<b>81.98</b>	82.14	<b>85.45</b>
	5-shot	84.21	<b>84.52</b>	85.43	<b>87.08</b>	<b>84.81</b>	84.63	85.91	<b>87.84</b>
	7-shot	83.79	<b>84.43</b>	85.32	<b>87.72</b>	85.12	<b>85.54</b>	86.97	<b>88.59</b>
	Full Data	<b>95.72</b>	95.61	95.78	<b>96.03</b>	95.45	<b>95.62</b>	96.34	<b>96.40</b>
XFUND-zh	1-shot	-	-	61.54	<b>70.68</b>	61.46	<b>69.68</b>	<b>61.21</b>	59.24
	3-shot	-	-	74.22	<b>77.22</b>	75.88	<b>79.46</b>	77.24	<b>78.58</b>
	5-shot	-	-	79.99	<b>82.02</b>	83.88	<b>85.02</b>	<b>82.27</b>	82.23
	7-shot	-	-	83.97	<b>85.44</b>	85.03	<b>87.43</b>	<b>84.76</b>	84.30
	Full Data	-	-	93.18	<b>93.20</b>	91.96	<b>92.72</b>	<b>94.27</b>	93.37

(b) Recall (%) of our PPTSER and Traditional Fine-tuning methods.

Table 5: Precision and Recall of PPTSER and Traditional Fine-tuning methods. Metrics in **Bold** is better between PPTSER and Fine-tuning and *FT* refers to *Fine-tuning* methods.

## D Further Analysis of Experimental Results

This section presents further analysis and additional performance metrics obtained from the main experiments.

### D.1 Further Analysis of Comparisons with Existing Fine-Tuning Methods

Table 5a and Table 5b present the precision and recall of PPTSER compared to the traditional fine-tuning method. The results demonstrate that PPTSER consistently outperforms the traditional fine-tuning method in most cases, leading to improved overall performance in terms of F1 scores.

Additionally, we observed from Table 1 and Table 5 that the improvement of PPTSER on the CORD dataset is generally less pronounced compared to its performance on other benchmarks. This prompted further investigation on our part.

As mentioned in Section 2.1, our  $N$ -way  $K$ -shot setting implies that *each of the  $N$  categories has  $K$  documents containing entities of that category as the support set*. Given that a single document in the CORD dataset cannot encompass entities from all categories, we selected more than  $K$  documents under the  $K$ -shot setup in the previous experiments. Consequently, we reselected 1 to 5 document samples as the training set for CORD and conducted additional experiments with LiLT, which showed the

Methodology	LayoutLMv2		LiLT		
	FT	Ours	FT	Ours	
Sample Number	1	33.28	<b>33.81</b>	40.46	<b>42.47</b>
	2	44.86	<b>47.09</b>	49.39	<b>54.18</b>
	3	54.24	<b>56.91</b>	56.31	<b>63.71</b>
	4	61.54	<b>62.23</b>	61.36	<b>69.19</b>
	5	63.89	<b>64.92</b>	64.09	<b>69.83</b>

(a) F1 score (%) of our PPTSER and Traditional Fine-tuning methods on CORD benchmark.

Methodology	LayoutLMv2		LiLT		
	FT	Ours	FT	Ours	
Sample Number	1	33.10	<b>34.06</b>	40.35	<b>40.08</b>
	2	42.31	<b>46.20</b>	47.26	<b>52.58</b>
	3	51.31	<b>57.50</b>	53.58	<b>62.92</b>
	4	59.58	<b>62.50</b>	60.17	<b>68.72</b>
	5	61.93	<b>64.65</b>	63.21	<b>69.47</b>

(b) Precision (%) of our PPTSER and Traditional Fine-tuning methods on CORD benchmark.

Methodology	LayoutLMv2		LiLT		
	FT	Ours	FT	Ours	
Sample Number	1	33.82	<b>33.75</b>	41.36	<b>45.28</b>
	2	48.01	<b>48.10</b>	51.81	<b>55.89</b>
	3	<b>57.60</b>	56.35	57.73	<b>64.52</b>
	4	63.64	<b>61.97</b>	62.63	<b>69.68</b>
	5	66.00	<b>65.22</b>	65.01	<b>70.19</b>

(c) Recall (%) of our PPTSER and Traditional Fine-tuning methods on CORD benchmark.

Table 6: Performances of our PPTSER and Traditional Fine-tuning methods on CORD benchmark with various numbers of sample as the support set. Metrics in **Bold** is the best and *FT* refers to *Fine-tuning* methods.

most significant improvement, and LayoutLMv2, which showed the least. The results, as illustrated in Table 6, led to several conclusions:

- **Sample Size:** The K-shot division on CORD often includes more than K samples. For example, in the 1-shot experiment, our training set averaged 7.6 samples. In contrast, on FUNSD, a single sample typically encompasses entities of all 4 types, resulting in only 1 sample in the training set for the 1-shot scenario. Our findings indicate that as the sample size increases, the performance gap between our method and conventional fine-tuning diminishes, yet our method retains its advantage. Therefore, under the same K-shot setting, CORD involves more samples than other datasets, which leads to a less pronounced improvement.
- **Label Complexity:** We believe that the labels in CORD are more complex and ab-

stract, making it harder for the model to grasp their semantic meanings compared to those in FUNSD. For instance, entity categories in FUNSD include *header*, *question*, and *answer*, whereas in CORD, they involve more abstract types like *menu.num* and *total.creditcardprice*. The experimental results indicate that when the sample size is extremely small, the improvement on CORD from our method is limited. However, this improvement increases rapidly with the sample size, suggesting our method can more accurately capture the relationship between document tokens and tag-related prompts with relatively more samples.

- **Pre-trained Models:** Different pre-trained models have varying degrees of understanding of labels. This could explain why some pre-trained models show weaker improvement on CORD. Our supplementary experiments reveal that the extent of improvement of our method is consistently better with the LiLT than the LayoutLMv2, suggesting that LiLT better understands the semantic information implied by the labels.

## D.2 Further Analysis of Comparisons with Existing Few-shot Methods

We also present a detailed comparison of PPTSER with other few-shot methods, including precision and recall metrics in Table 7a and Table 7b. Similar to the F1 score, models enhanced with PPTSER usually demonstrate superior performance compared to both few-shot NER and few-shot SER methods.

Modality		Text		Text + Layout			Text + Layout + Image	
Methodology		EntLM (NAACL 22)	COPNER (COLING 22)	LASER (ACL 22)	COPNER <sub>L<sub>1</sub>LT</sub> (COLING 22)	PPTSER <sub>L<sub>1</sub>LT</sub> (Ours)	COPNER <sub>LMv3</sub> (COLING 22)	PPTSER <sub>LMv3</sub> (Ours)
FUNSD	1-shot	22.85	18.67	36.61	<u>53.79</u>	50.91	49.01	<b>56.27</b>
	3-shot	33.39	30.97	46.71	67.26	67.24	<u>73.54</u>	<b>75.18</b>
	5-shot	37.23	32.57	46.80	71.40	72.95	<u>75.40</u>	<b>79.53</b>
	7-shot	40.29	35.43	51.17	72.81	74.14	<u>78.63</u>	<b>81.31</b>
	Full Data	67.53	63.39	69.08	87.45	88.89	<u>91.45</u>	<b>91.96</b>
CORD-Lv1	1-shot	73.23	67.42	65.56	86.93	<b>90.62</b>	86.80	<u>90.05</u>
	3-shot	82.85	81.37	75.43	94.39	94.97	93.99	<b>95.14</b>
	5-shot	86.87	85.87	82.07	94.94	<b>96.38</b>	95.65	<u>96.20</u>
	7-shot	86.65	86.54	83.54	95.12	96.25	96.05	<b>96.53</b>
	Full Data	95.93	95.83	96.50	99.23	<u>99.43</u>	<b>99.45</b>	<b>99.45</b>
CORD	1-shot	57.45	51.22	-	70.27	<b>75.35</b>	67.14	<u>74.02</u>
	3-shot	71.42	67.26	-	81.48	<u>83.85</u>	79.89	<b>85.09</b>
	5-shot	77.70	74.59	-	84.92	<u>87.03</u>	85.13	<b>87.71</b>
	7-shot	78.51	75.69	-	85.86	<u>87.74</u>	86.83	<b>88.37</b>
	Full Data	93.56	92.33	-	95.75	<u>96.06</u>	95.79	<b>96.39</b>
XFUND-zh	1-shot	27.02	23.98	-	49.49	<b>64.90</b>	53.37	<u>54.59</u>
	3-shot	35.94	35.72	-	64.69	<u>71.54</u>	69.98	<b>72.36</b>
	5-shot	43.18	42.93	-	68.21	<b>77.07</b>	73.48	<u>76.57</u>
	7-shot	45.13	44.66	-	72.61	<u>77.70</u>	76.42	<b>77.95</b>
	Full Data	64.75	65.59	-	87.23	88.17	<b>91.10</b>	<u>91.04</u>

(a) Precision (%) of PPTSER and other Few-shot methods.

Modality		Text		Text + Layout			Text + Layout + Image	
Methodology		EntLM (NAACL 22)	COPNER (COLING 22)	LASER (ACL 22)	COPNER <sub>L<sub>1</sub>LT</sub> (COLING 22)	PPTSER <sub>L<sub>1</sub>LT</sub> (Ours)	COPNER <sub>LMv3</sub> (COLING 22)	PPTSER <sub>LMv3</sub> (Ours)
FUNSD	1-shot	28.01	21.53	41.05	56.72	<u>61.77</u>	54.43	<b>70.07</b>
	3-shot	37.56	32.12	46.55	70.17	71.35	<u>78.32</u>	<b>78.64</b>
	5-shot	42.42	38.28	52.14	75.59	77.76	<u>79.83</u>	<b>83.65</b>
	7-shot	42.79	39.57	54.08	73.92	77.38	<u>78.45</u>	<b>81.91</b>
	Full Data	67.31	65.81	69.41	88.04	89.26	<u>91.06</u>	<b>92.05</b>
CORD-Lv1	1-shot	75.38	69.90	68.12	87.01	<b>90.38</b>	87.16	<u>89.99</u>
	3-shot	84.54	83.17	76.77	93.95	94.61	94.06	<b>95.12</b>
	5-shot	87.35	86.30	82.42	94.78	<u>96.04</u>	95.84	<b>96.23</b>
	7-shot	87.99	86.97	83.71	94.96	96.00	96.07	<b>96.50</b>
	Full Data	95.94	95.97	96.62	99.19	<u>99.42</u>	<b>99.45</b>	<b>99.45</b>
CORD	1-shot	58.29	53.79	-	69.83	<b>75.80</b>	67.51	<u>74.35</u>
	3-shot	71.95	68.35	-	81.08	<u>83.82</u>	80.24	<b>85.45</b>
	5-shot	77.78	74.69	-	84.68	<u>87.08</u>	85.47	<b>87.84</b>
	7-shot	78.76	76.51	-	85.66	<u>87.72</u>	86.91	<b>88.59</b>
	Full Data	93.45	92.68	-	95.72	<u>96.03</u>	95.79	<b>96.40</b>
XFUND-zh	1-shot	26.35	23.26	-	48.26	<b>70.68</b>	55.67	<u>59.24</u>
	3-shot	39.62	40.48	-	64.88	<u>77.22</u>	72.89	<b>78.58</b>
	5-shot	44.03	46.14	-	70.00	<u>82.02</u>	79.92	<b>82.23</b>
	7-shot	48.46	49.58	-	76.58	<b>85.44</b>	82.47	<u>84.30</u>
	Full Data	67.72	68.71	-	91.20	<u>93.20</u>	92.91	<b>93.37</b>

(b) Recall (%) of PPTSER and other Few-shot methods.

Table 7: Precision and Recall of PPTSER and other Few-shot methods. Metrics in **Bold** is the best, and that with underline is the second best.

CORD						
	<i>7<sup>th</sup> block</i>	<i>8<sup>th</sup> block</i>	<i>9<sup>th</sup> block</i>	<i>10<sup>th</sup> block</i>	<i>11<sup>th</sup> block</i>	<i>12<sup>th</sup> block</i>
1-shot	68.27	69.63	68.34	68.93	74.27	75.57
3-shot	80.09	81.16	80.81	81.51	82.89	83.83
5-shot	84.29	85.25	85.00	84.84	86.41	87.06
7-shot	84.58	85.86	85.48	85.97	86.79	87.73
Full Data	95.63	95.59	95.48	96.12	95.86	96.04
$\Delta$	<b>-3.47</b>	<b>-2.55</b>	<b>-3.03</b>	<b>-2.57</b>	<b>-0.80</b>	-

(a) F1 score (%) of PPTSER on CORD benchmark when obtaining the attention weight from different blocks.

CORD						
	<i>7<sup>th</sup> block</i>	<i>8<sup>th</sup> block</i>	<i>9<sup>th</sup> block</i>	<i>10<sup>th</sup> block</i>	<i>11<sup>th</sup> block</i>	<i>12<sup>th</sup> block</i>
1-shot	67.36	69.02	67.86	68.53	73.91	75.35
3-shot	79.76	81.19	80.78	81.46	82.76	83.85
5-shot	80.98	85.38	85.01	84.78	86.22	87.03
7-shot	84.41	85.92	85.52	85.97	86.68	87.74
Full Data	95.70	95.64	95.55	96.16	95.87	96.06
$\Delta$	<b>-4.36</b>	<b>-2.58</b>	<b>-3.06</b>	<b>-2.63</b>	<b>-0.92</b>	-

(b) Precision (%) of PPTSER on CORD benchmark when obtaining the attention weight from different blocks.

CORD						
	<i>7<sup>th</sup> block</i>	<i>8<sup>th</sup> block</i>	<i>9<sup>th</sup> block</i>	<i>10<sup>th</sup> block</i>	<i>11<sup>th</sup> block</i>	<i>12<sup>th</sup> block</i>
1-shot	69.21	70.25	68.83	69.34	74.64	75.80
3-shot	80.43	81.14	80.84	81.56	83.03	83.82
5-shot	80.41	85.13	84.99	84.89	86.60	87.08
7-shot	84.75	85.80	85.43	85.97	86.90	87.72
Full Data	95.57	95.54	95.41	96.09	95.85	96.03
$\Delta$	<b>-4.02</b>	<b>-2.52</b>	<b>-2.99</b>	<b>-2.52</b>	<b>-0.69</b>	-

(c) Recall (%) of PPTSER on CORD benchmark when obtaining the attention weight from different blocks.

Table 8: Performances of PPTSER on CORD benchmark when obtaining the attention weight from different blocks.  $\Delta$  denotes the average deviation of the corresponding metrics compared to that in the *12<sup>th</sup> block*.

### D.3 Further Analysis of Attention weights obtained from different blocks

We additionally provide the numerical metrics of distinct settings to obtain the attention weight from different blocks. Table 8 illustrates the experimental results, indicating that obtaining attention weights from the last block yields the best performance of F1 score, precision, and recall. Although the reduction of parameters alleviates over-fitting to some extent, since some shallower blocks outperform certain deeper ones in the 1-shot scenario, our default setting to obtain the attention weight from the last block significantly outperforms the alternative settings of obtaining the attention weight from shallower blocks. This finding strongly reinforces the effectiveness of our design.

### D.4 Further Analysis of Different Designs on PPTSER

Besides, we also offer the precision and recall of our PPTSER under different designs in Table 9a and Table 9b. The outcomes clearly indicate that our PPTSER, characterized by its meticulous design, outperforms other designs across all metrics evaluated, including the F1 score, Precision, and Recall. This superiority not only showcases the robustness of our design but also significantly substantiates the efficacy of our PPTSER framework in tackling the few-shot SER problem on visually-rich documents.

Model Designs	Baseline	Decoupling Strategies	Prompt Engineering		Aggregation Strategies	
	<i>default setting</i>	<i>plain BIO prompt</i>	<i>unrelated words</i>	<i>random embeddings</i>	<i>mean</i>	<i>single head</i>
1-shot	75.35	73.92	70.58	71.32	74.16	74.92
3-shot	83.85	82.44	82.83	83.24	83.04	84.07
5-shot	87.03	86.50	85.84	86.47	86.61	86.97
7-shot	87.74	86.64	86.44	87.09	86.76	87.30
Full Data	96.06	95.58	96.27	95.75	96.23	96.08
$\Delta$	-	<b>-0.99</b>	<b>-1.61</b>	<b>-1.23</b>	<b>-0.64</b>	<b>-0.14</b>

(a) Precision (%) of PPTSER on CORD benchmark with different designs.

Model Designs	Baseline	Decoupling Strategies	Prompt Engineering		Aggregation Strategies	
	<i>default setting</i>	<i>plain BIO prompt</i>	<i>unrelated words</i>	<i>random embeddings</i>	<i>mean</i>	<i>single head</i>
1-shot	75.80	74.66	71.03	71.54	74.69	75.51
3-shot	83.82	82.37	82.54	82.90	83.06	83.77
5-shot	87.08	86.01	85.72	86.11	86.67	86.99
7-shot	87.72	86.35	86.40	86.81	86.71	87.19
Full Data	96.03	95.00	96.24	95.70	96.18	96.03
$\Delta$	-	<b>-1.21</b>	<b>-1.70</b>	<b>-1.48</b>	<b>-0.63</b>	<b>-0.19</b>

(b) Recall (%) of PPTSER on CORD benchmark with different designs.

Table 9: Precision and Recall of PPTSER on CORD benchmark with different designs.  $\Delta$  denotes the average deviation of the corresponding metrics relative to the *default setting*.

## E In-depth Parameter Analysis of PPTSER over Traditional Fine-tuning

We provide a further analysis of the parameter counts in this section. As shown in Table 4, our PPTSER maintains consistent parameters across different benchmarks with the same pre-trained model. This is attributed to the fact that PPTSER does not necessitate an additional classifier layer, unlike the traditional fine-tuning method. Consequently, the parameter variance arises when employing the traditional fine-tuning method with the same pre-trained models on different benchmarks, owing to variations in the number of entity types present in those benchmarks. Furthermore, as PPTSER omits the value transform layer and the feed-forward layer in the final attention block, we achieve a reduction in the parameter count of the pre-trained model it is based on. Additionally, the extent of parameter reduction varies among different pre-trained models due to disparities in their architectural designs, resulting in slice differences in the eliminations of the modules.