

RE²: Region-Aware Relation Extraction from Visually Rich Documents

Anonymous EACL submission

Abstract

Current research in form understanding predominantly relies on large pre-trained language models, necessitating extensive data for pre-training. However, the importance of layout structure (i.e., the spatial relationship between the entity blocks in the visually rich document) to relation extraction has been overlooked. In this paper, we propose **RE²** that leverages region-level spatial structure among the entity blocks to improve their relation prediction. We design an edge-aware graph attention network to learn the interaction between entities while considering their spatial relationship defined by their region-level representations. We also introduce a constraint objective to regularize the model towards consistency with the inherent constraints of the relation extraction task. To support the research on relation extraction from visually rich documents and demonstrate the generalizability of **RE²**, we build a new benchmark dataset, DIVERSEFORM, that covers a wide range of domains. Extensive experiments on DIVERSEFORM and several public benchmark datasets demonstrate significant superiority and transferability of **RE²** across various domains and languages, with up to 18.88% absolute F-score gain over all high-performing baselines¹.

1 Introduction

Visually Rich Documents (VRDs) encompass various types such as *invoices*, *questionnaire forms*, *financial forms*, *legal documents*, and so on. These documents possess valuable layout information that aids in comprehending their content. Recent research (Liu et al., 2019; Jaume et al., 2019; Yu et al., 2020) has focused on extracting key information, such as entities and relations, from VRDs by leveraging their layout structures and Optical

¹We will make all the programs, model checkpoints and dataset publicly available once the paper is accepted.

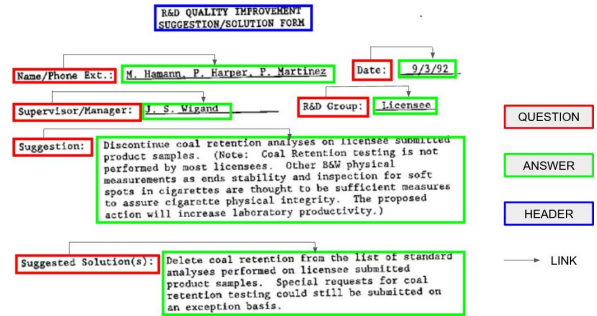


Figure 1: Example of entity and relation extraction from a visually rich document. The colored boxes represent three categories of semantic entities and the arrows represent relations between them.

Character Recognition (OCR) results². Figure 1 shows an example where entity recognition aims to identify blocks of text in certain categories, such as *Question(Q)*, *Answer(A)*, and *Header(H)*. Relation extraction further predicts the links among the entities, especially *Q-A* links indicating that the *A* block is the corresponding answer to the *Q* block.

Extracting key information, especially relations in VRDs is a challenging task. Though similar to traditional extraction tasks in text-only Natural Language Processing (NLP) (Grishman, 1997; Chen et al., 2022), inferring relations in VRDs poses additional challenges. They require not only understanding the semantic meaning of entities but also taking into account the layout information, e.g., the spatial structures among the entity blocks in original VRDs. Previous studies mainly focused on combining the text and layout with language model pre-training (Lu et al., 2019; Su et al., 2020; Chen et al., 2020; Powalski et al., 2021; Xu et al., 2022a; Wang et al., 2022a,b; Huang et al., 2022) or encoding the local layout information by constructing super-tokens (Qian et al., 2019; Liu et al., 2019; Yu et al., 2021; Lee et al., 2022, 2023). However, the

²Optical Character Recognition will recognize a set of bounding boxes and their corresponding text from VRDs where each bounding box can represent a single word or a cohesive group of words, both semantically and spatially.

063 layout of the VRDs, especially the relative spatial
064 relationship among the entity blocks, is still yet to
065 be effectively explored for relation extraction.

066 To this end, we propose **REgion-Aware Relation**
067 **Extraction (RE²)** that leverages region-level spa-
068 tial structures among the entities to reason about
069 their relations³. Specifically, given the question and
070 answer entities from each VRD, we define three
071 categories of region-level representations for each
072 entity block, through which we further characterize
073 the relative spatial relationship between each pair
074 of question and answer entities. We then employ
075 a layout-aware pre-trained language model (i.e.,
076 LayoutXLM (Xu et al., 2022a)) to encode the enti-
077 ties and an Edge-aware Graph Attention Network
078 (eGAT) to further learn the interaction between the
079 question and answer entities in a bipartite graph
080 while considering their spatial relationship. To en-
081 sure each answer is linked to at most one question,
082 we design a constraint-based learning objective to
083 guide the learning process, in combination with the
084 relation classification objective.

085 To validate the effectiveness of **RE²**, we con-
086 duct extensive experiments on various benchmark
087 datasets for a wide range of languages and do-
088 mains. We evaluate **RE²** on two public datasets
089 FUNSD (Jaume et al., 2019) and XFUND (Xu
090 et al., 2022b), under supervised, multitask transfer,
091 and zero-shot cross-lingual transfer settings. We
092 also create a new benchmark dataset **DIVERSE-**
093 **FORM** that covers diverse domains, such as Veter-
094 ans Affairs, visa applications, tax documents, air
095 transport and so on, and evaluate **RE²** for cross-
096 domain transfer. Experimental results show that
097 **RE²** outperforms the previous state-of-the-art ap-
098 proaches with a large margin on (almost) all lan-
099 guages and domains across all settings. Our abla-
100 tion studies also verify the significant benefit of the
101 region-level spatial structures of entity blocks for
102 relation extraction. The contributions of this work
103 are summarized as follows:

- 104 • We are the first to propose the region-level
105 entity representations and utilize them to char-
106 acterize the spatial structure among the entity
107 blocks, which have been proven to be signif-
108 icantly beneficial to relation extraction from
109 visually rich documents.
- 110 • We develop a new framework **RE²** that lever-
111 ages the spatial structures among the question

³This work mainly focuses on extracting *Q-A* relation given the gold *Question* and *Answer* entities.

and answer entities with an effective eGAT
network and regularizes model predictions
with a novel constraint objective. **RE²** demon-
strates superior performance across (almost)
all languages and domains under supervised,
cross-lingual, and cross-domain transfer set-
tings.

- We contribute **DIVERSEFORM**, a new bench-
mark dataset that covers a wide range of do-
mains to support the research on information
extraction from visually rich documents.

2 Related Work

Recent research on visually rich document infor-
mation extraction shows that incorporating 2D po-
sitional embedding and layout coordinates into the
pre-trained language models improves VRD under-
standing (Xu et al., 2020, 2022a; Huang et al., 2022;
Powalski et al., 2021; Wang et al., 2022b). To deal
with the variation of relation definitions, DocRel
(Li et al., 2022) proposes a contrastive learning
framework that utilizes the coherence of existing
relations in diverse enhanced positive views to gen-
erate relational representations. Zhang et al. (2021)
further explores entity relation extraction as depen-
dency parsing, incorporating minimum vertical and
horizontal distances between the entities as layout
heuristics. Compared with all these studies, our
approach is the first to propose and incorporate
multi-granular spatial structures among the entities,
which have shown to significantly improve relation
extraction from VRDs.

Graph Attention Networks (GAT) (Veličković
et al., 2018) have proven to be efficient for learn-
ing on graph-structured data (Zhang et al., 2022a).
This is exemplified by the work GraphDoc (Zhang
et al., 2022b), a multimodal graph attention-based
model that simultaneously utilizes text, layout, and
image information for visually rich document un-
derstanding. Though several studies (Liu et al.,
2019; Lee et al., 2022, 2023) have explored GNNs
for entity extraction from VRDs, we are the first to
design edge-aware GAT to improve relation extrac-
tion from VRDs, which presents additional chal-
lenges, encompassing spatial analysis to determine
entity layout on the page and semantics between
entities for identifying relationships. GNNs have
also been applied to relation extraction from tex-
tual documents (Zhu et al., 2019; Guo et al., 2019;
Zhang et al., 2018). However, these methods can-
not be directly adapted to relation extraction from
VRDs due to the fundamental differences in doc-

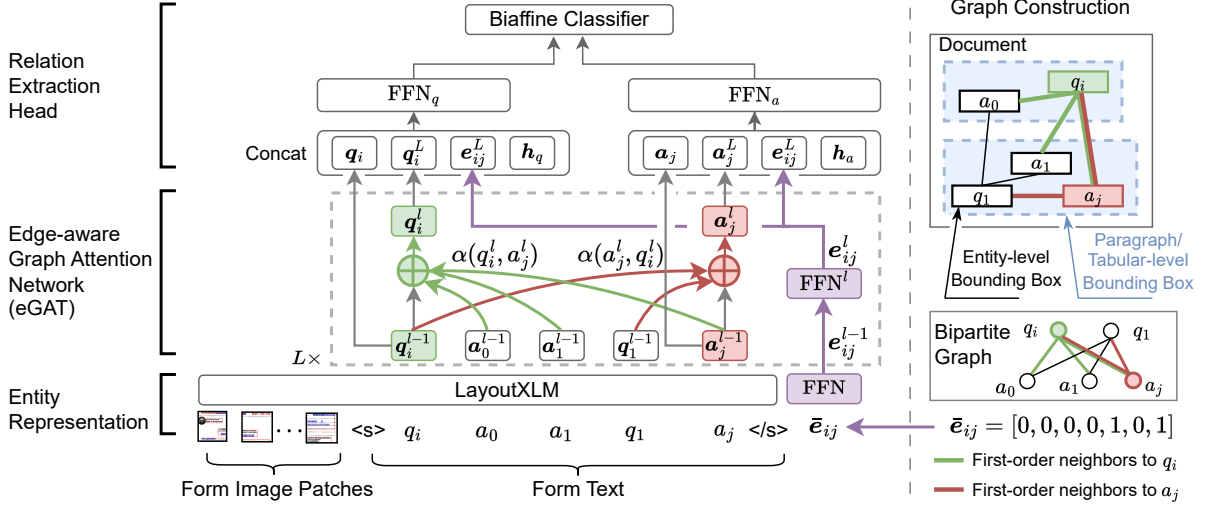


Figure 2: Overview of the **RE**gion-level **Relation Extraction (RE²)** framework. In the eGAT layer, the representation of each entity is updated based on the attention scores of its first-order neighbors.

ument formats, structures, and the key challenges encountered in relation extraction: text-only documents primarily rely on linguistic cues and phrases for relation extraction, whereas VRDs necessitate consideration of both semantics and spatial context. Given that, we innovatively incorporate a multi-granular layout heuristic into an edge-aware graph attention network, placing greater emphasis on capturing more fine-grained layout structures.

3 Approach

Given a visually rich document D , a set of question entities $Q = \{q_1, q_2, \dots, q_m\}$ and answers $A = \{a_1, a_2, \dots, a_n\}$, we aim to identify all the connected pairs (q, a) where $q \in Q$ and $a \in A$, indicating that a is the corresponding answer of q . Each q_i or a_j can be denoted as $\{[w_0, w_1, \dots, w_t], (x_0, y_0, x_1, y_1)\}$, where $[w_0, w_1, \dots, w_t]$ is the sequence of words denoting the entity span and (x_0, y_0, x_1, y_1) is the coordinates for the entity bounding box. Figure 2 illustrates our **RE²** framework that aims to leverage region-level spatial structures among the question and answer blocks to detect their association.

3.1 Entity Representation

We first learn the encoding of question and answer entities based on LayoutXML (Xu et al., 2022b), a layout-aware transformer-based model that has been extended to support multilingualism by pre-training on multilingual VRD datasets.

Given a set of question entities $Q = \{q_1, q_2, \dots, q_m\}$ and answers $A = \{a_1, a_2, \dots, a_n\}$ from document D , we obtain the entity embeddings $Q = \{q_1, q_2, \dots, q_m\}$, $A = \{a_1, a_2, \dots, a_n\}$, $q_i, a_i \in \mathbb{R}^{1 \times F}$, where F is the entity feature dimension⁴. For entities with multiple tokens, we use the embedding of their first tokens as their representations⁵.

3.2 Region-Aware Graph Construction

Based on the spatial structures of the input VRD, we define three distinct categories of regions (i.e., bounding box) for each entity: (1) an **entity-level bounding box** that refers to the bounding box encompassing the entire entity span and is obtained by merging the bounding boxes of all the words in a span obtained by OCR (Liu et al., 2019; Yu et al., 2020); (2) a **paragraph-level bounding box** that is defined as a visually distinct section for the paragraph where the entity occurs within a document and corresponds to the clustering of words that are located within a dense region. The paragraph-level bounding boxes are extracted by an existing tool, EasyOCR⁶, which takes the maximum horizontal and vertical distances between adjacent word-level bounding boxes as hyperparameters to merge them into paragraph-level bounding boxes; and (3) a **tabular-based bounding box** if the entity occurs in a tabular structure. We define a tabular-based

⁴We use bold symbols to denote vectors.

⁵Preliminary experiments showed use of first subtoken performed better than average embedding of all subtokens.

⁶<https://www.jaired.ai/easyocr/>

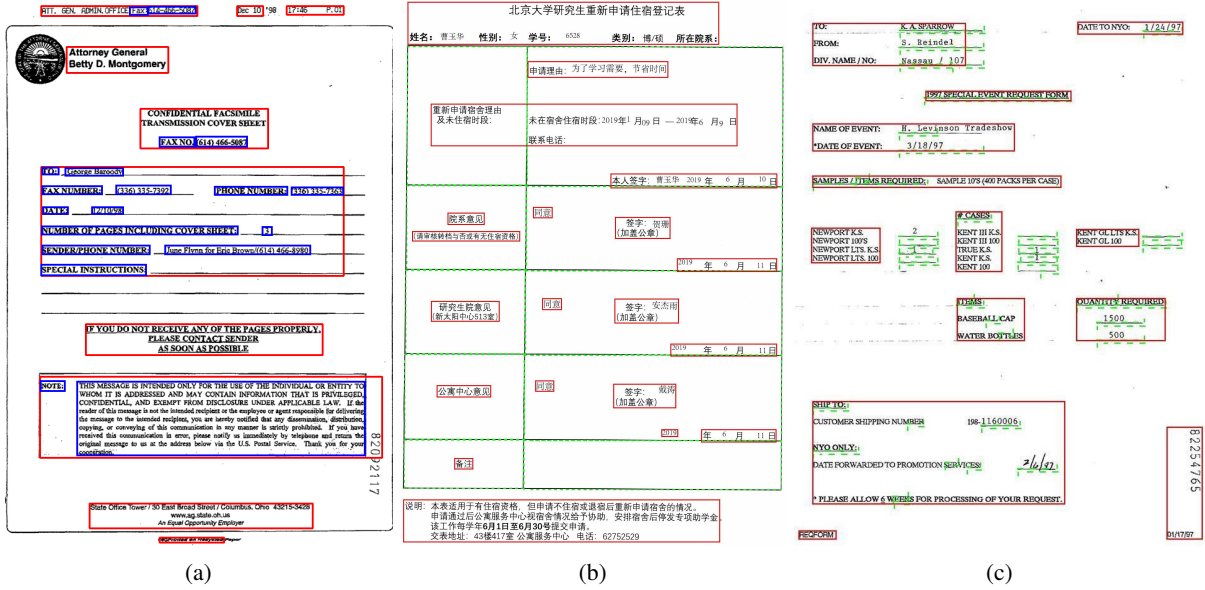


Figure 3: Entity level bounding box (for question and answer entities) are shown in blue, paragraph-level bounding box in red and tabular-based bounding box in green.

bounding box as the coordinates of a table cell. Note that each entity can only appear in either a paragraph or a table, so other than its entity-level bounding box, we always assign either a paragraph-level or tabular-based bounding box for each entity, instead of both. Our preliminary results show that a tabular-based bounding box is vital because tabular structures are usually not well-captured by existing OCR tools. Illustrations of the three types of regions are shown in Figure 3. The pseudocode for extracting paragraph/tabular regions is present in Appendix H.

To characterize the links between the question and answer entities, we further propose to construct a complete bipartite graph, $G = (Q, A, E)$, for each visually rich document, where the question entities $Q = \{q_1, q_2, \dots, q_m\}$ and answers $A = \{a_1, a_2, \dots, a_n\}$ are the nodes, and for each pair of q_i and a_j , there is an edge $e_{ij} \in E$ connecting them. Each entity is represented by the encoding learned from LayoutXLM as detailed in Section 3.1, and each edge is represented by a one-hot encoding vector based on the spatial relationship between the three categories of bounding boxes of the question and answer:

$$\bar{e}_{ij} = [I, E_{lr}^1, E_{tb}^1, E_{lr}^0, E_{tb}^0, R_{lr}, R_{tb}],$$

where each term is an indicator variable: I indicates whether the two entities are within the same paragraph/tabular region. If so, $I = 1$, otherwise, $I = 0$. When the two entities are from the same paragraph/tabular region, E_{lr}^1 and E_{tb}^1 further indicate the left-right (lr) and top-bottom (tb) spatial

relationship of their entity-level bounding boxes. For example, $E_{lr}^1 = 1$ indicates that the entity-level bounding boxes of the two entities have a left-right spatial relation, otherwise, $E_{lr}^1 = 0$. When the two entities are not from the same paragraph/tabular region, E_{lr}^0 and E_{tb}^0 indicate the left-right and top-bottom spatial relationship of their entity-level bounding boxes, while R_{lr} and R_{tb} indicate the left-right and top-bottom spatial relationship of their paragraph/tabular level bounding boxes. Note that when the two entities are from the same paragraph/tabular region, the indicators of E_{lr}^0 , E_{tb}^0 , R_{lr} , R_{tb} will be all zero.

To obtain a dense representation of each edge, we pass each one-hot encoding vector \bar{e}_{ij} to a feed-forward network, and the resulting vector $e_{ij} = \text{FFN}(\bar{e}_{ij})$ is assigned as the edge weight between q_i and a_j , where $e_{ij} \in \mathbb{R}^{1 \times F/2}$.

3.3 Edge-aware Graph Attention Network

We further propose an edge-aware graph attention network (eGAT), extended from the graph attention network (GAT) (Veličković et al., 2018) by incorporating the edge weights inferred by spatial information to learn the interaction between the question and answer nodes. In our experiments, eGAT consists of 2 encoding layers, while each layer updates the node embeddings based on the first-order neighbors with masked self-attention.

Specifically, given the node embeddings at layer l , $Q^l = \{q_1^l, q_2^l, \dots, q_m^l\}$, $A = \{a_1^l, a_2^l, \dots, a_n^l\}$, we first compute the attention weight between q_i

and a_j as follows

$$\begin{aligned} \text{att}(\mathbf{W}^l \mathbf{q}_i^l, \mathbf{W}^l \mathbf{a}_j^l) &= \mathbf{W}_{\text{att}}^\top (\mathbf{W}^l \mathbf{q}_i^l \parallel \mathbf{W}^l \mathbf{a}_j^l) \\ c(q_i^l, a_j^l) &= \text{LeakyReLU}(\text{att}(\mathbf{W}^l \mathbf{q}_i^l, \mathbf{W}^l \mathbf{a}_j^l)) \\ \alpha(q_i^l, a_j^l) &= \text{softmax}_j \left(\sum_j (e_{ij}^l \cdot c(q_i^l, a_j^l)) \right) \end{aligned}$$

where \cdot denotes scalar multiplication. $\mathbf{W}^l \in \mathbb{R}^{F' \times F}$ is a parameter matrix for shared linear transformation for \mathbf{q}_i^l and \mathbf{a}_j^l . $\mathbf{W}_{\text{att}} \in \mathbb{R}^{2F'}$ is a weight vector for the attention mechanism. \parallel denotes the catenation operation. $e_{ij}^l = \text{FFN}^l(e_{ij}^{l-1})$ where e_{ij}^0 is the initial dense representation e_{ij} of each edge.

The resulting edge-aware normalized attention scores are then used to update the hidden representations of the question and answer nodes, respectively, with residual connection:

$$\begin{aligned} \bar{\mathbf{q}}_i^{l+1} &= \mathbf{q}_i + \sum_{j \in \mathcal{N}_i} \alpha(q_i^l, a_j^l) \mathbf{W} \mathbf{a}_j^l \\ \bar{\mathbf{a}}_j^{l+1} &= \mathbf{a}_j + \sum_{i \in \mathcal{M}_j} \alpha(a_j^l, q_i^l) \mathbf{W} \mathbf{q}_i^l \end{aligned}$$

where \mathcal{N}_i and \mathcal{M}_j denotes the first order neighbors of q_i and a_j respectively.

For each layer of eGAT, we apply multi-head attention (Vaswani et al., 2017), where each attention head performs operations independently, and the mean of all attention heads is taken for aggregation. The updated representation of question node q_i and answer a_j is computed as follows:

$$\begin{aligned} \mathbf{q}_i^{l+1} &= \sigma \left(\frac{1}{K} \sum_{k=1}^K \left(\mathbf{q}_i + \sum_{j \in \mathcal{N}_i} \alpha(q_i^l, a_j^l)^k \mathbf{W}^k \mathbf{a}_j^l \right) \right) \\ \mathbf{a}_j^{l+1} &= \sigma \left(\frac{1}{K} \sum_{k=1}^K \left(\mathbf{a}_j + \sum_{i \in \mathcal{M}_j} \alpha(a_j^l, q_i^l)^k \mathbf{W}^k \mathbf{q}_i^l \right) \right) \end{aligned}$$

where K is the number of independent attention heads and \mathbf{W}^k denotes the weight matrix for the k^{th} attention head. $\sigma(\cdot)$ denotes a non-linear function (ELU is used for experiments). \mathbf{q}_i^{l+1} and \mathbf{a}_j^{l+1} are then used as input node embeddings for layer $l+1$.

3.4 Relation Extraction

Binary Relation Prediction For each pair of question q_i and answer a_j , we aim to predict a binary label, indicating whether the answer corresponds to the particular question or not. The representation of q_i or a_j is formed by concatenating the LayoutXLM embedding, the final node embedding

learned from eGAT, the edge representation of the pair of q_i and a_j , and an entity type representation (i.e., question or answer). The entity type representation is learned by an embedding layer during training. The inclusion of entity type embedding is important in determining the direction of the relation between the two entities. The resulting representations q_i and a_j are then passed through two feed-forward networks and a biaffine classifier (Dozat and Manning, 2017) to obtain a score $s_{i,j}$ for determining whether the pair is associated or not.

$$\begin{aligned} \mathbf{q}'_i &= \text{FFN}_q(\mathbf{q}_i \parallel \mathbf{q}_i^L \parallel \mathbf{e}_{ij}^L \parallel \mathbf{h}_q) \\ \mathbf{a}'_j &= \text{FFN}_a(\mathbf{a}_j \parallel \mathbf{a}_j^L \parallel \mathbf{e}_{ij}^L \parallel \mathbf{h}_a) \\ s_{ij} &= \mathbf{q}'_i \mathbf{U} \mathbf{a}'_j + \mathbf{V}(\mathbf{q}'_i \circ \mathbf{a}'_j) + \mathbf{b} \end{aligned}$$

where \mathbf{h}_q and \mathbf{h}_a are the type embeddings of question and answer entities. Note that \mathbf{h}_q and \mathbf{h}_a remain the same across all questions and answers, respectively. \mathbf{U} , \mathbf{V} and \mathbf{b} are trainable parameters. During training, the loss is computed following the cross-entropy loss

$$\mathcal{L}_b = - \sum y \cdot \log(p_{ij}).$$

where $y \in \{0, 1\}$ is the target binary label and $p_{ij} = \text{softmax}(s_{ij})$, indicating the probability of a relation between q_i and a_j .

Constraint Loss Our preliminary study shows that without any constraint, the model tends to predict multiple questions to be associated with one answer, which is against the definition of relation extraction for VRDs, where each answer is linked to at most one question. To address this issue, we incorporate the constraint into the learning process in the form of a constraint loss. Previous work (Li et al., 2019; Wang et al., 2020) demonstrated that declarative logical constraints can be converted into differentiable functions, and help regularize the model towards consistency with the logical constraints. We design a declarative logical constraint that holds true for relation extraction task from VRDs as follows, $\forall a_j \in A, \forall q_i \in Q$,

$$\text{rel}(q_i, a_j) \rightarrow \bigwedge_{q_k \in \mathbf{Q} \setminus \{q_i\}} \neg \text{rel}(q_k, a_j).$$

This means, for any $a_j \in A$, if there exists one relation link between a_j and any particular q_i among all questions, there cannot be another relation link

Model	EN	ZH	JA	ES	FR	IT	DE	PT	Avg.
XLM-RoBERTa _{BASE} (Conneau et al., 2020)	26.59	51.05	58.00	52.95	49.65	53.05	50.41	39.82	47.69
InfoXML _{BASE} (Chi et al., 2021)	29.20	52.14	60.00	55.16	49.13	52.81	52.62	41.70	49.10
LayoutXML _{BASE} (Xu et al., 2022b)	54.83	70.73	69.63	68.96	63.53	64.15	65.51	57.18	64.32
LiLT[InfoXML] _{BASE} (Wang et al., 2022a)	62.76	72.97	70.37	71.95	69.65	70.43	65.58	58.74	67.81
RE² (Our Approach)	71.76	79.60	75.36	75.59	76.38	77.45	75.86	59.76	73.98

Table 1: Language-specific fine-tuning results (F1%) on FUNSD(EN) and XFUND.

Model	EN	ZH	JA	ES	FR	IT	DE	PT	Avg.
XLM-RoBERTa _{BASE}	26.59	16.01	26.11	24.40	22.40	23.74	22.88	19.96	22.76
InfoXML _{BASE}	29.20	24.05	28.51	24.81	24.54	21.93	20.27	20.49	24.23
LayoutXML _{BASE}	54.83	44.94	44.08	47.08	44.16	40.90	38.20	36.85	43.88
LiLT[InfoXML] _{BASE}	62.76	47.64	50.81	49.68	52.09	46.97	41.69	42.72	49.30
RE² (Our Approach)	71.76	66.32	64.42	58.82	69.02	61.83	60.57	43.87	62.08

Table 2: Zero-shot cross-lingual results (F1%) (trained on EN (FUNSD) and tested on other languages)

for this answer a_j . We further define the following constraint loss derived from the logical constraints:

$$\mathcal{L}_c = y \cdot \left| \log(p_{ij}) - \frac{1}{|Q| - 1} \sum_{\substack{k=0 \\ k \neq i}}^{|Q|} \log(1 - p_{kj}) \right|$$

where Q denotes the whole set of questions in the document.

Overall Learning Objective The overall learning objective is a weighted combination of the binary cross entropy loss and the constraint loss:

$$\mathcal{L} = \beta \mathcal{L}_b + \delta \mathcal{L}_c$$

where β and δ are hyperparameters.

4 Experiment Settings

4.1 Datasets

The main challenge of relation extraction from visually rich documents lies in the varying layouts of form-like documents from different domains or languages. On the other hand, the region-level spatial structures defined in RE² are domain and language-independent. To validate the effectiveness of RE², we perform experiments on several benchmark datasets covering a wide range of languages and domains.

FUNSD The FUNSD dataset (Jaume et al., 2019) is a subset of the RVL-CDIP dataset (Harley et al., 2015) containing realistic scanned document images with ground truth OCR. It provides bounding boxes for entity spans and annotations for four types: *Question*, *Answer*, *Header*, and *Other*. The dataset includes relational links among these entities, with a particular focus on Question-Answer

links. We adopt the same data split and experimental settings as previous studies (Xu et al., 2022b; Wang et al., 2022a), using 149 documents for training, 50 documents for evaluation, and reporting the best performance on the evaluation set.

XFUND XFUND (Xu et al., 2022b) is a diverse multilingual dataset containing visually rich documents in seven languages: Portuguese, Chinese, Spanish, French, Japanese, Italian, and German. The dataset encompasses a wide range of document structures, reflecting variations in document conventions across languages and countries. With a total of 1,393 fully annotated forms, each language includes 149 forms in the train set and 50 forms in the test set. The dataset provides human-annotated ground truth OCR, entity annotations, and relation annotations. Notably, the XFUND dataset shares similarities with the FUNSD dataset in terms of its document format.

DIVERSEFORM To best demonstrate the performance of domain transfer of RE², we further create a new dataset, **DIVERSEFORM**, by curating government forms from Aggarwal et al. (2020) and Sarkar et al. (2020). These forms encompass a wide range of question types, including checkboxes, tables, multiple-choice questions (MCQs), and fill-in-the-blank fields. The domains of the forms cover various areas such as Veterans Affairs, visa applications, tax documents, air transport, legal forms, vehicle-related forms from the Department of Motor Vehicles (DMV), and miscellaneous forms from different government agencies. These forms are of single page and were originally empty and they are designed to collect confidential information such as health data and

Model	EN	ZH	JA	ES	FR	IT	DE	PT	Avg.
XLM-RoBERTa _{BASE}	36.38	67.97	68.29	68.28	67.27	69.37	68.87	60.82	63.41
InfoXML _{BASE}	36.99	64.93	64.73	68.28	68.31	66.90	63.84	57.63	61.45
LayoutXML _{BASE}	66.71	82.41	81.42	81.04	82.21	83.10	78.54	70.44	78.23
LiLT[InfoXML] _{BASE}	74.07	84.71	83.45	83.35	84.66	84.58	78.78	76.43	81.25
RE² (Our Approach)	74.11	88.25	82.27	83.23	86.83	84.02	81.89	71.04	81.46

Table 3: Multitask fine-tuning performance (F1%) on FUNSD(EN) and XFUND.

Model	DIVERSEFORM	FUNSD → DIVERSEFORM	DIVERSEFORM → FUNSD
LayoutXML _{BASE}	69.72	37.33	32.58
LiLT[InfoXML] _{BASE}	64.15	41.56	30.26
RE²	70.87	41.78	50.32

Table 4: Supervised results on **DIVERSEFORM** and cross-domain transfer results between **DIVERSEFORM** and FUNSD. (F1%)

tax details. To populate the forms, we employed two annotators who used synthetic data generated by The One Generator⁷ for fields such as names, addresses, and other necessary information. This approach ensures the privacy and security of individuals’ personal information while providing a realistic representation of the data typically found in these government forms. We then hire another annotator to label the *Question* and *Answer* entities as well as their relations for these documents using the annotation tool UBIAl⁸, which also offers its customized OCR model for extracting text from uploaded images. However, due to the serialized top-left to bottom-right text extraction approach of the OCR, the spans of entities are sometimes fragmented in complex layout forms. During the annotation process, these fragmented spans are identified and merged to achieve the correct serialization of spans. After labeling the entities and relations for these documents, we further hire three annotators to validate the annotations. All the annotators are senior undergraduate students majoring in Computer Science and are paid a rate of \$15/hour. We name the final annotated dataset as **DIVERSEFORM**, which comprises a total of 150 training documents and 50 testing documents. Details of **DIVERSEFORM** annotation and statistics is in Appendix C.

4.2 Experiment Results

Language-specific fine-tuning results are presented in Table 1, where each model is fine-tuned on language X and tested on language X. The experimental findings show that the proposed model outperforms all the baselines across all evaluated

languages. To evaluate the **cross-lingual zero-shot transfer** capability, the model is fine-tuned on the FUNSD dataset in English, followed by testing on multiple languages. The experimental results, as shown in Table 2, demonstrate the superiority of our model over the baseline approach in terms of zero-shot performance. This outcome provides compelling evidence that the incorporated region-level spatial structures and constraints for relation extraction exhibit effective transferability across different languages. We also conduct a significance test for both our approach and the best-performing baseline (i.e., LiLT[InfoXML]_{BASE} (Wang et al., 2022a)) under the settings of language-specific fine-tuning and cross-lingual zero-shot transfer. As shown in Table 7 in Appendix D, our approach significantly outperforms the baseline under both settings.

Furthermore, the **multitask fine-tuning** results are shown in Table 3, where the model is trained on the training sets of all languages and then tested on each individual language. The achieved results demonstrate superior performance, indicating that the model successfully learns the layout invariance present across different languages. By effectively capturing and leveraging the shared layout characteristics, the model exhibits improved generalization capabilities, leading to enhanced performance across diverse linguistic contexts. This highlights the significance of incorporating layout information in cross-lingual settings and underscores the model’s ability to adapt and transfer its learned knowledge to effectively process and understand documents in various languages.

Note that **RE²** shows less competitive performance on Portuguese (PT) due to more complex layout structures. Portuguese forms exhibit a com-

⁷<https://theonegenerator.com/>

⁸<https://ubiai.tools/>

Model	EN	ZH	JA	ES	FR	IT	DE	PT	Avg.
RE²	71.76	79.60	75.36	75.59	76.38	77.45	75.86	59.76	73.98
- node embedding	70.19	78.93	75.00	74.60	76.00	76.82	73.20	57.29	72.75
- edge embedding	57.42	69.37	67.93	72.01	73.73	69.67	63.48	55.61	66.15
- constraint loss	68.52	77.77	74.49	74.78	75.20	75.66	73.61	57.48	72.19
- entity level regions	44.69	76.89	66.71	73.11	62.44	70.63	62.10	44.30	62.61
- paragraph/tabular regions	71.57	79.5	74.17	72.05	74.98	76.79	74.55	57.49	72.64

Table 5: Ablation study results (F1%) on eGAT (node and edge embeddings), constraint loss, paragraph/tabular regions and entity level regions.

508 combination of mixed tables and paragraph structures, 548
509 making it challenging to determine the appropriate 549
510 usage for paragraph-level regions or tabular regions. 550
511 An example is shown in Appendix E. 551

512 We also assess the generalization of **RE²** and 552
513 two high-performing baselines based on **DIVERSE-** 553
514 **FORM** and FUNSD, which cover two sets of dis- 554
515 tinct domains. We conduct experiments under the 555
516 settings of both domain-specific fine-tuning and 556
517 cross-domain transfer where the models are trained 557
518 on one dataset and tested on the other. As shown in 558
519 Table 4, **RE²** significantly outperforms the two 559
520 strong baselines when fine-tuned on **DIVERSE-** 560
521 **FORM** and tested on **DIVERSEFORM** or FUNSD. 561
522 The improvement of **RE²** when it’s trained on 562
523 FUNSD and tested on **DIVERSEFORM** is marginal, 563
524 probably due to the greater diversity and complex- 564
525 ity in document layout of **DIVERSEFORM** com- 565
526 pared to FUNSD. 566

527 4.3 Ablation Study

528 **Effect of Node and Edge Embeddings from** 568
529 **eGAT** The node and edge embeddings from 569
530 eGAT are concatenated with the entity represen- 570
531 tations before being passed to the biaffine classifier. 571
532 A series of ablation studies are conducted to assess 572
533 the individual contributions of the layout informa- 573
534 tion. The results of these studies are presented in 574
535 Table 5. Figure 6 in Appendix F provides visual 575
536 evidence that solely relying on the updated node 576
537 embeddings from eGAT fails to adequately capture 577
538 the layout heuristics and results in the omission of 578
539 numerous relations. Conversely, employing only 579
540 the updated edge embeddings without considering 580
541 the node embeddings leads to an over-prediction 581
542 of relations with limited regard for the semantic 582
543 relevance of the entities involved. Optimal per- 583
544 formance is achieved through the joint utilization 584
545 of both node and edge embeddings, indicating the 585
546 importance of integrating both sources of informa- 586
547 tion to effectively capture the region-level spatial

548 structures and consider the semantic context of the 549
549 relations. 550

551 **Effect of Constraint Loss** The constraint loss 552
552 has been modeled to encourage each answer entity 553
553 to be linked to at most one question. Table 5 shows 554
554 that incorporating the constraint loss significantly 555
555 improves the F1 score of **RE²**, especially precision. 556
556 The detailed experimental results are evidenced in 557
557 Appendix G. 558

559 **Effect of Region Information** We also investi- 560
560 gate the impact of each category of regions on 561
561 characterizing the spatial relationship among the 562
562 entities and further affecting the performance of 563
563 **RE²**. As shown in Table 5, the inclusion of each 564
564 category of region information significantly im- 565
565 proves the performance of **RE²**. The absence of 566
566 entity-level regions resulted in a substantial de- 567
567 crease in performance, underscoring the vital role 568
568 of pairwise entity layout information, i.e., whether 569
569 the question and answer entities are arranged verti- 570
570 cally (top-bottom) or horizontally (left-right). Fig- 571
571 ure 7 in Appendix F shows an example to com- 572
572 pare the relation predictions with and without para- 573
573 graph/tabular regions, indicating that incorporating 574
574 paragraph/tabular regions helps prevent the model 575
575 from predicting relations across semantically differ- 576
576 ent regions. The result of this ablation study proves 577
577 the effectiveness of the multi-granular region infor- 578
578 mation. 579

577 5 Conclusion

578 In this work, we propose a novel entity relation 579
579 extraction model, **RE²**, that incorporates layout 580
580 heuristics and constraints that are generalizable 581
581 across different languages. Experimental results 582
582 on 8 different languages and our proposed dataset 583
583 **DIVERSEFORM** show the effectiveness of our 584
584 proposed method under four settings (language- 585
585 specific, cross-lingual zero-shot, multi-lingual fine- 586
586 tuning, and cross-domain transfer). 587

587 **Limitations**

588 In this work, we found the incorporation of layout
589 heuristics to be compelling and we are excited by
590 how leveraging region information improves per-
591 formance drastically. One of the limitations of our
592 model is its reliance on a relatively limited set of
593 heuristics and features. For instance, we have not
594 yet incorporated visual information and template-
595 based knowledge, which could potentially improve
596 the accuracy and robustness of the relation extrac-
597 tion task. Additionally, the current model employs
598 an exhaustive inference approach, considering all
599 possible relations during prediction. While this
600 ensures comprehensive coverage, it also results in
601 longer inference times for each relation type. These
602 limitations indicate avenues for further improve-
603 ment, such as exploring additional heuristics and
604 incorporating more efficient inference strategies,
605 to enhance the performance and efficiency of our
606 model.

607 **References**

608 Milan Aggarwal, Mausoom Sarkar, Hires Gupta, and
609 Balaji Krishnamurthy. 2020. Multi-modal associ-
610 ation based grouping for form structure extraction.
611 In *The IEEE Winter Conference on Applications of*
612 *Computer Vision*, pages 2075–2084.

613 Muhao Chen, Lifu Huang, Manling Li, Ben Zhou, Heng
614 Ji, and Dan Roth. 2022. New frontiers of information
615 extraction. In *Proceedings of the 2022 Conference*
616 *of the North American Chapter of the Association*
617 *for Computational Linguistics: Human Language*
618 *Technologies: Tutorials*.

619 Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El
620 Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and
621 Jingjing Liu. 2020. Uniter: Universal image-text
622 representation learning. In *ECCV*.

623 Zewen Chi, Li Dong, Furu Wei, Nan Yang, Saksham
624 Singhal, Wenhui Wang, Xia Song, Xian-Ling Mao,
625 Heyan Huang, and Ming Zhou. 2021. **InfoXLM: An**
626 **information-theoretic framework for cross-lingual**
627 **language model pre-training**. In *Proceedings of the*
628 *2021 Conference of the North American Chapter of*
629 *the Association for Computational Linguistics: Hu-*
630 *man Language Technologies*, pages 3576–3588, On-
631 line. Association for Computational Linguistics.

632 Alexis Conneau, Kartikay Khandelwal, Naman Goyal,
633 Vishrav Chaudhary, Guillaume Wenzek, Francisco
634 Guzmán, Edouard Grave, Myle Ott, Luke Zettle-
635 moyer, and Veselin Stoyanov. 2020. **Unsupervised**
636 **cross-lingual representation learning at scale**. In *Pro-*
637 *ceedings of the 58th Annual Meeting of the Asso-*
638 *ciation for Computational Linguistics*, pages 8440–

8451, Online. Association for Computational Lin- 639
guistics. 640

Timothy Dozat and Christopher D. Manning. 2017. 641
Deep biaffine attention for neural dependency pars- 642
ing. 643

Ralph Grishman. 1997. Information extraction: Tech- 644
niques and challenges. In *Information Extraction* 645
A Multidisciplinary Approach to an Emerging Infor- 646
mation Technology: International Summer School, 647
SCIE-97 Frascati, Italy, July 14–18, 1997, pages 10– 648
27. Springer. 649

Zhijiang Guo, Yan Zhang, and Wei Lu. 2019. **Atten-** 650
tion guided graph convolutional networks for relation 651
extraction. In *Proceedings of the 57th Annual Meet-* 652
ing of the Association for Computational Linguistics, 653
pages 241–251, Florence, Italy. Association for Com- 654
putational Linguistics. 655

Adam W. Harley, Alex Ufkes, and Konstantinos G. Der- 656
panis. 2015. **Evaluation of deep convolutional nets** 657
for document image classification and retrieval. 658

Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and 659
Furu Wei. 2022. **Layoutlmv3: Pre-training for docu-** 660
ment ai with unified text and image masking. In 661
Proceedings of the 30th ACM International Confer- 662
ence on Multimedia, MM ’22, page 4083–4091, New 663
York, NY, USA. Association for Computing Machin- 664
ery. 665

Guillaume Jaume, Hazim Kemal Ekenel, and Jean- 666
Philippe Thiran. 2019. **Funsd: A dataset for form** 667
understanding in noisy scanned documents. 668

Chen-Yu Lee, Chun-Liang Li, Timothy Dozat, Vin- 669
cent Perot, Guolong Su, Nan Hua, Joshua Ainslie, 670
Renshen Wang, Yasuhisa Fujii, and Tomas Pfister. 671
2022. **Formnet: Structural encoding beyond sequen-** 672
tial modeling in form document information extrac- 673
tion. 674

Chen-Yu Lee, Chun-Liang Li, Hao Zhang, Timothy 675
Dozat, Vincent Perot, Guolong Su, Xiang Zhang, Ki- 676
hyuk Sohn, Nikolai Glushnev, Renshen Wang, Joshua 677
Ainslie, Shangbang Long, Siyang Qin, Yasuhisa Fu- 678
jii, Nan Hua, and Tomas Pfister. 2023. **Formnetv2:** 679
Multimodal graph contrastive learning for form docu- 680
ment information extraction. 681

Tao Li, Vivek Gupta, Maitrey Mehta, and Vivek Sriku- 682
mar. 2019. **A logic-driven framework for consistency** 683
of neural models. In *Proceedings of the 2019 Confer-* 684
ence on Empirical Methods in Natural Language Pro- 685
cessing and the 9th International Joint Conference 686
on Natural Language Processing (EMNLP-IJCNLP), 687
pages 3924–3935, Hong Kong, China. Association 688
for Computational Linguistics. 689

Xin Li, Yan Zheng, Yiqing Hu, Haoyu Cao, Yunfei 690
Wu, Deqiang Jiang, Yinsong Liu, and Bo Ren. 2022. 691
Relational representation learning in visually-rich 692
documents. 693

694	Xiaojing Liu, Feiyu Gao, Qiong Zhang, and Huasha Zhao. 2019. Graph convolution for multimodal information extraction from visually rich documents. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Industry Papers)</i> , pages 32–39. Association for Computational Linguistics.	749
695		750
696		751
697		752
698		753
699		
700		754
701		755
		756
702	Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. <i>ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks</i> . Curran Associates Inc., Red Hook, NY, USA.	757
703		758
704		
705		
706	Rafał Powalski, Łukasz Borchmann, Dawid Jurkiewicz, Tomasz Dwojak, Michał Pietruszka, and Gabriela Pałka. 2021. Going full-tilt boogie on document understanding with text-image-layout transformer. In <i>Document Analysis and Recognition-ICDAR 2021: 16th International Conference, Lausanne, Switzerland, September 5–10, 2021, Proceedings, Part II 16</i> , pages 732–747. Springer.	759
707		760
708		761
709		762
710		763
711		764
712		765
713		
714	Yujie Qian, Enrico Santus, Zhijing Jin, Jiang Guo, and Regina Barzilay. 2019. GraphIE: A graph-based framework for information extraction . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 751–761, Minneapolis, Minnesota. Association for Computational Linguistics.	766
715		767
716		768
717		769
718		770
719		771
720		772
721		
722		
723	Mausoom Sarkar, Milan Aggarwal, Arneh Jain, Hires Gupta, and Balaji Krishnamurthy. 2020. Document structure extraction using prior based high resolution hierarchical semantic segmentation. In <i>European Conference on Computer Vision</i> , pages 649–666. Springer.	773
724		774
725		775
726		776
727		777
728		778
729		
730	Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. 2020. Vi-bert: Pre-training of generic visual-linguistic representations . In <i>International Conference on Learning Representations</i> .	779
731		780
732		781
733		782
734		783
735		
736		
737		784
738		785
739		786
740		787
741		788
742		789
743		
744		790
745		791
746		792
747		
748		
		793
		794
		795
		796
		797
		798
		799
		800
		801
		802
		803

804	Hao Zhu, Yankai Lin, Zhiyuan Liu, Jie Fu, Tat-Seng	– Annotate headers even if their questions	852
805	Chua, and Maosong Sun. 2019. Graph neural net-	haven't been answered	853
806	works with generated parameters for relation extrac-	– Headers do not have answers directly at-	854
807	tion . In <i>Proceedings of the 57th Annual Meeting of</i>	attached to them	855
808	<i>the Association for Computational Linguistics</i> , pages		
809	1331–1339, Florence, Italy. Association for Compu-		
810	tational Linguistics.		
811	Appendix		
812	A Data Preprocessing		
813	To accurately determine the layout heuristics, it	• Answer: A word, set of words, or sentence	856
814	is important to get the bounding box of the en-	written in response to a question.	857
815	tire entity span. If token-level bounding boxes are	– Responses in the form of checkbox op-	858
816	provided, the boxes can be merged to obtain a span-	tions count as answers.	859
817	level box. All the paragraph/tabular regions are de-	– In multiple choice type questions, all the	860
818	ected and their bounding boxes are obtained. We	options are annotated as answers (follow-	861
819	identify the region an entity belongs to by checking	ing FUNSD and XFUND)	862
820	the entity's Intersection over Union (IoU) with the		
821	regions and assign the region with the maximum	The guidelines of annotating relations are as fol-	863
822	IoU.	lows:	864
823	B Experiment Setting and	• Question-Answer: A link exists between a	865
824	Hyperparameters	question entity and an answering entity when	866
825	NVIDIA A40 GPU was used for all finetuning	the answer is a response to a particular ques-	867
826	tasks. To create paragraph-level regions using	tion.	868
827	EasyOCR, a horizontal merging of text boxes is	– When multiple answers exist for a ques-	869
828	carried out when their distance is within 2, and ver-	tion, there are multiple Question-Answer	870
829	tical merging is done when the distance is within	links from the same question entity.	871
830	1, while setting the paragraph flag to True. The	– Answers to a sub-question should only	872
831	model is trained end-to-end, with fine-tuning of the	be linked to the sub-question and not the	873
832	LayoutXLM base model. The eGAT layers and	parent question.	874
833	relation extraction head are trained from scratch,	• Question-Question: A link exists between a	875
834	using 2 eGAT layers for all experiments. The train-	question entity and another question entity	876
835	ing process entails 5000 steps with a batch size of	if one question is a sub-question of another	877
836	4, a learning rate of 5e-5, and a warm-up ratio of	question or one question is conditioned on the	878
837	0.1. The cross entropy loss is weighted at 1 and	answer of another question.	879
838	constraint loss is weighted at 0.02.	– For example, "If yes, . . ." type of ques-	880
839	C Annotation Details and Statistics of	tion has a Question-Question link with	881
840	DIVERSEFORM	the parent question.	882
841	The guidelines of annotating entities are as follows:	– A question that is split into multiple	883
842		fine-grained questions has a Question-	884
843	• Question: A word, set of words, or sentence	Question link between them. For ex-	885
844	worded or expressed so as to elicit information	ample, "Address" can have further ques-	886
845	from the person filling the form.	tions such as "Apt. No", "Street Name",	887
846	– Questions are annotated even if they	"City", "State", "Zip Code".	888
847	haven't been answered	• Header-Question: A link exists between a	889
848	– Questions and sub-questions are labeled	header entity and a question entity if the ques-	890
849	as the same type of entity.	tions are present under the section or subsec-	891
850	• Header: A word, set of words, or sentences	tion that is characterized by the header.	892
851	worded or expressed so as give context or en-	– If multiple questions exist under a header,	893
	capsulate a set of questions.	there are multiple Header-Question links	894
		from the same header entity.	895
		– Often confused with Question-Question	896
		links and can be differentiated based on	897

898 layout structure, font style, and other vi- 946
899 sual aspects of the questions from the 947
900 form. 948

901 The guidelines of annotation of tables are as 949
902 follows: We mainly deal with one dimensional 950
903 tables. For the case that each cell in the table is 951
904 related to both row and column questions, there will 952
905 be a Question-Question link between the questions 953
906 extracted from the row and column, indicating that 954
907 one question is a sub-question of another question 955
908 or one question is conditioned on the answer of 956
909 another question. This is part of the annotation 957
910 guidelines for FUNSD and our own dataset. Based
911 on these annotation rules, the constraint of one
912 answer having one question still holds.

913 Figure 4 shows the distribution of domains in **DI-**
914 **VERSEFORM**. Miscellaneous consists of forms for
915 voter registration, agriculture, scholarship, immi-
916 gration, property tax, etc. Veteran’s Affairs encom-
917 passes varying forms ranging from child support
918 payments to retirement funds. There is rich layout
919 variation within each domain shown in the chart.
920 The number of entities and relations of each type
921 in **DIVERSEFORM** are tabulated in Table 6.

922 **D Significance Test**

923 Table 7 shows the significance test results for both
924 our approach and the best performing baseline (i.e.,
925 LiLT[InfoXLM]_{BASE} (Wang et al., 2022a)) under
926 the settings of language-specific fine-tuning and
927 cross-lingual zero-shot transfer. The results for all
928 experiments reported were averaged across 3 runs.

929 **E Case Study**

930 Figure 5 visualizes paragraph-level regions, tabu-
931 lar regions, and predictions for a Portuguese form
932 in FUNSD. It shows that paragraph-level regions
933 are suitable for the top portion of the form, while
934 tabular regions specifically pertain to the bottom
935 table. In this particular form, the decision was
936 made to adopt paragraph-level regions, resulting
937 in the exclusion of the tabular layout despite its
938 ability to convey more information. We acknowl-
939 edge that there are instances where our proposed
940 approach may struggle to accurately distinguish be-
941 tween paragraph-level and tabular regions, leading
942 to a performance decrease.

943 **F Visualizations of Ablation Results**

944 Figure 6 shows the visualization of predictions
945 of the ablation study of node and edge embed-

946 dings. Figure 7 shows the visualization of pre-
947 dictions of the ablation study of incorporating para-
948 graph/tabular regions.

949 **G Ablation Results of Constraint Loss**

950 The constraint loss has been modeled to encour-
951 age each answer entity to be linked to at most one
952 question. Table 8 shows that incorporating the con-
953 straint loss significantly improves the F1 score of
954 **RE**², especially precision.

955 **H Pseudocode**

956 The following pseudocode extracts the tabular and
957 paragraph-level regions from a VRD.

Algorithm 1 IdentifyHorizontalAndVerticalLines(image)

- 1: Apply horizontal kernel to the image
 - 2: Apply vertical kernel to the image
 - 3: Find horizontal lines
 - 4: Find vertical lines
 - 5: **return** Combined horizontal and vertical lines
-

Algorithm 2 FindBoundingBoxes(lines)

- 1: TabularBoxList = []
 - 2: Find contours in lines
 - 3: **for** each contour **do**
 - 4: Compute the bounding box
 - 5: Append the box to TabularBoxList
 - 6: **end for**
 - 7: **return** TabularBoxList
-

Algorithm 3 SortBoxesByArea(boundingBoxes)

- 1: Sort the bounding boxes by area in increasing order
 - 2: **return** boundingBoxes
-

Algorithm 4 AppendBoxToList(boundingBoxes, text)

- 1: FinalBoxList = []
 - 2: **for** each box in boundingBoxes **do**
 - 3: **if** the box contains any text and has no intersection with existing boxes in FinalBoxList **then**
 - 4: Append the box to FinalBoxList
 - 5: **end if**
 - 6: **end for**
 - 7: **return** FinalBoxList
-

Algorithm 5 CheckAllTextPresent(FinalBoxList, text)

- 1: **if** all the text in the document is present in the boxes in FinalBoxList **then**
 - 2: **return** True
 - 3: **else**
 - 4: **return** False
 - 5: **end if**
-

Algorithm 6 GetMissingText(FinalBoxList, text)

- 1: missingText = []
 - 2: **if** the text is not present inside the bounding boxes of any of the FinalBoxList **then**
 - 3: Append text to missingText
 - 4: **end if**
 - 5: **return** missingText
-

Algorithm 7 AppendMissingTextBoxes(FinalBoxList, missingText, ParagraphRegions)

- 1: **for** each missing text in missingText **do**
 - 2: **if** missing text is present in any paragraph region in ParagraphRegions **then**
 - 3: Append paragraph region to FinalBoxList
 - 4: **end if**
 - 5: **end for**
 - 6: **return** FinalBoxList
-

Split	Entities			Relations		
	Question	Answer	Header	Question-Answer	Question-Question	Header-Question
Training	3,087	3,585	230	1,172	594	546
Test	956	1,048	57	520	270	164

Table 6: Statistics of entities and relations in DIVERSEFORM

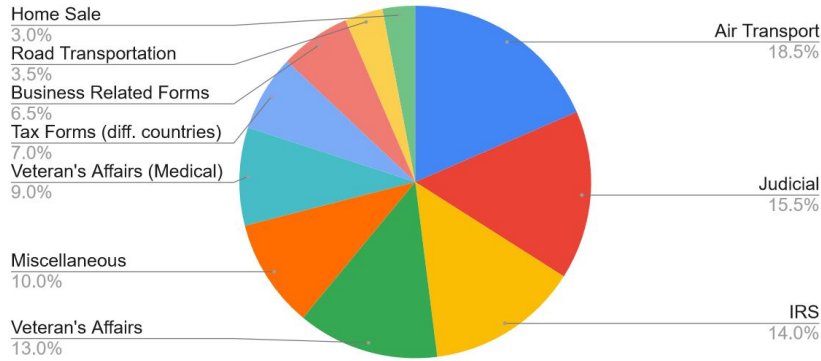


Figure 4: Domain distribution of DIVERSEFORM.

Nº	Nome dos Membros da Família	Idade	Grav de Parentesco	Atividade Laboral	Renda Bruta Mensal (RS)
01	Ana Costa	22	candidato(a)	professor	100
02	Jose Dias	34	candidato	professor	200
03	Pedro Dias	54	candidato	médico	200
04	Fernanda Braga	22	candidata	trabalhador	150
05	Luis Braga	27	candidato	professor	150
06	Filipe Da Cama	37	candidato	médico	180
TOTAL DA RENDA BRUTA MENSAL FAMILIAR (RS)					

(a)

(b)

Nº	Nome dos Membros da Família	Idade	Grav de Parentesco	Atividade Laboral	Renda Bruta Mensal (RS)
01	Ana Costa	22	candidato(a)	professor	100
02	Jose Dias	34	candidato	professor	200
03	Pedro Dias	54	candidato	médico	200
04	Fernanda Braga	22	candidata	trabalhador	150
05	Luis Braga	27	candidato	professor	150
06	Filipe Da Cama	37	candidato	médico	180
TOTAL DA RENDA BRUTA MENSAL FAMILIAR (RS)					

(c)

Figure 5: Visualization of paragraph-level regions (a), tabular regions (b) and predictions (c) for a Portuguese form in XFUND.

一级指标	二级指标	三级指标	三级指标目标值		
			(2020年)	(2021年)	
产出	数量指标	普通高职(专科)院校	增加	增加	
		中等职业教育	增加	550万左右	
	质量指标	举办全省职业院校技能大赛	1100左右	每年1100左右	
		“中国职业教育”高职院校和专业建设任务	启动	基础完成	
	过程指标	“一带一路”沿线国家职业教育合作	启动	基础完成	
		建设程序	启动	启动	
	时效指标	职业院校办学满意度	逐年改善	≥85%	
		工程造价符合率	符合	符合	
	效益	社会效益指标	公办高职院校生均财政拨款水平	≥2.2万	≥2.2万
			全省职业院校技能大赛	顺利举办	顺利举办, 以赛促教, 以赛促学, 以赛促改
经济效益指标		普通高职院校初次就业率	≥92%	≥92%	
		普通高职院校就业率 (%)	≥93%	≥93%	
环境效益指标		中等职业院校在校生规模	≥2000人	≥2000人	
		高中阶段教育毛入学率	≥90%	≥90%	
可持续发展指标		经济欠发达地区公办高职院校办学水平和综合实力提升数 (所)	18	18	
		服务人员满意度	≥85%	≥85%	

(a) Ground Truth

(b) Concatenating only node embeddings

一级指标	二级指标	三级指标	三级指标目标值		
			(2020年)	(2021年)	
产出	数量指标	普通高职(专科)院校	增加	增加	
		中等职业教育	增加	550万左右	
	质量指标	举办全省职业院校技能大赛	1100左右	每年1100左右	
		“中国职业教育”高职院校和专业建设任务	启动	基础完成	
	过程指标	“一带一路”沿线国家职业教育合作	启动	基础完成	
		建设程序	启动	启动	
	时效指标	职业院校办学满意度	逐年改善	≥85%	
		工程造价符合率	符合	符合	
	效益	社会效益指标	公办高职院校生均财政拨款水平	≥2.2万	≥2.2万
			全省职业院校技能大赛	顺利举办	顺利举办, 以赛促教, 以赛促学, 以赛促改
经济效益指标		普通高职院校初次就业率	≥92%	≥92%	
		普通高职院校就业率 (%)	≥93%	≥93%	
环境效益指标		中等职业院校在校生规模	≥2000人	≥2000人	
		高中阶段教育毛入学率	≥90%	≥90%	
可持续发展指标		经济欠发达地区公办高职院校办学水平和综合实力提升数 (所)	18	18	
		服务人员满意度	≥85%	≥85%	

(c) Concatenating only edge embeddings

一级指标	二级指标	三级指标	三级指标目标值		
			(2020年)	(2021年)	
产出	数量指标	普通高职(专科)院校	增加	增加	
		中等职业教育	增加	550万左右	
	质量指标	举办全省职业院校技能大赛	1100左右	每年1100左右	
		“中国职业教育”高职院校和专业建设任务	启动	基础完成	
	过程指标	“一带一路”沿线国家职业教育合作	启动	基础完成	
		建设程序	启动	启动	
	时效指标	职业院校办学满意度	逐年改善	≥85%	
		工程造价符合率	符合	符合	
	效益	社会效益指标	公办高职院校生均财政拨款水平	≥2.2万	≥2.2万
			全省职业院校技能大赛	顺利举办	顺利举办, 以赛促教, 以赛促学, 以赛促改
经济效益指标		普通高职院校初次就业率	≥92%	≥92%	
		普通高职院校就业率 (%)	≥93%	≥93%	
环境效益指标		中等职业院校在校生规模	≥2000人	≥2000人	
		高中阶段教育毛入学率	≥90%	≥90%	
可持续发展指标		经济欠发达地区公办高职院校办学水平和综合实力提升数 (所)	18	18	
		服务人员满意度	≥85%	≥85%	

(d) Concatenating node & edge embeddings

Figure 6: Visualization of predictions of the ablation study of node and edge embeddings, where red lines denote the question span, green lines denote the answer span, and blue lines denote the question answer relation predictions.

Algorithm 8 GetParagraphTabularRegions(imageFile)

```

1: image = LoadImage(imageFile)
2: text = OCR(imageFile)
3: lines = IdentifyHorizontalAndVerticalLines(image)
4: boundingBoxes = FindBoundingBoxes(lines)
5: boundingBoxes = SortBoxesByArea(boundingBoxes)
6: FinalBoxList = AppendBoxToList(boundingBoxes, text)
7: if CheckAllTextPresent(FinalBoxList, text) then
8:   OutputResult(FinalBoxList)
9: else
10:  ParagraphRegions = GetParagraphRegionsFromEasyOCR(image)
11:  missingText = GetMissingText(FinalBoxList, text)
12:  FinalBoxList = AppendMissingTextBoxes(FinalBoxList, missingText, ParagraphRegions)
13:  OutputResult(FinalBoxList)
14: end if

```

DIVISION NAME: Detroit North # REPS: 13
 DIVISION NAME: Detroit East # REPS: 12
 DIVISION NAME: Grand Rapids # REPS: 13
 DIVISION NAME: Detroit South # REPS: 13
 DIVISION NAME: Detroit West # REPS: 12
 DIVISION NAME: Flint # REPS: 13

DISTRIBUTION

Direct Accounts and Chains Headquartered within the Region
(15+ Stores) Stocking No Old Gold Light Box 100's

Name of Account	Ind/Lor Volume	Number of Stores	Name of Account	Ind/Lor Volume	Number of Stores
Quality	15115	8	Speedy Q	19613	16
Bay State	819	8			
Schwartz Oil	14113	8			
Wilbur	14113	8			
Quaker	11799	8			
Spartan	1188	9			
Indy	16615	9			
Arby's	11313	209			
Imoco/Oil	16311	31			

Direct Accounts and Chains Headquartered Outside the Region
(15+ Stores) Stocking No Old Gold Light Box 100's

Name of Account	Ind/Lor Volume	Number of Stores	Name of Account	Ind/Lor Volume	Number of Stores
Clark		12			
Fox		12			
Fa-Eye		12			
Way		12			
Miss Jewel		14			
Dairy Queen		15			
Moe's		16			
ACA America		31			

(a) Ground Truth

DIVISION NAME: Detroit North # REPS: 13
 DIVISION NAME: Detroit East # REPS: 12
 DIVISION NAME: Grand Rapids # REPS: 13
 DIVISION NAME: Detroit South # REPS: 13
 DIVISION NAME: Detroit West # REPS: 12
 DIVISION NAME: Flint # REPS: 13

DISTRIBUTION

Direct Accounts and Chains Headquartered within the Region
(15+ Stores) Stocking No Old Gold Light Box 100's

Name of Account	Ind/Lor Volume	Number of Stores	Name of Account	Ind/Lor Volume	Number of Stores
Quality	15115	8	Speedy Q	19613	16
Bay State	819	8			
Schwartz Oil	14113	8			
Wilbur	14013	8			
Quaker	11799	8			
Spartan	1188	9			
Indy	16615	9			
Arby's	11313	209			
Imoco/Oil	16311	31			

Direct Accounts and Chains Headquartered Outside the Region
(15+ Stores) Stocking No Old Gold Light Box 100's

Name of Account	Ind/Lor Volume	Number of Stores	Name of Account	Ind/Lor Volume	Number of Stores
Clark		12			
Fox		12			
Fa-Eye		12			
Way		12			
Miss Jewel		14			
Dairy Queen		15			
Moe's		16			
ACA America		31			

(b) Predictions without paragraph/tabular region information

DIVISION NAME: Detroit North # REPS: 13
 DIVISION NAME: Detroit East # REPS: 12
 DIVISION NAME: Grand Rapids # REPS: 13
 DIVISION NAME: Detroit South # REPS: 13
 DIVISION NAME: Detroit West # REPS: 12
 DIVISION NAME: Flint # REPS: 13

DISTRIBUTION

Direct Accounts and Chains Headquartered within the Region
(15+ Stores) Stocking No Old Gold Light Box 100's

Name of Account	Ind/Lor Volume	Number of Stores	Name of Account	Ind/Lor Volume	Number of Stores
Quality	15115	8	Speedy Q	19613	16
Bay State	819	8			
Schwartz Oil	14113	8			
Wilbur	14113	8			
Quaker	11799	8			
Spartan	1188	9			
Indy	16615	9			
Arby's	11313	209			
Imoco/Oil	16311	31			

Direct Accounts and Chains Headquartered Outside the Region
(15+ Stores) Stocking No Old Gold Light Box 100's

Name of Account	Ind/Lor Volume	Number of Stores	Name of Account	Ind/Lor Volume	Number of Stores
Clark		12			
Fox		12			
Fa-Eye		12			
Way		12			
Miss Jewel		14			
Dairy Queen		15			
Moe's		16			
ACA America		31			

(c) Predictions with paragraph/tabular region information

Figure 7: Visualization of predictions of the ablation study of incorporating paragraph/tabular regions.

Setting	RE²		Baseline		P-value
	Mean	SD	Mean	SD	
Language-Specific Fine-Tuning	73.98	6.15	67.81	4.98	0.0447
Zero-Shot Cross-Lingual	62.08	8.53	49.30	6.55	0.0047

Table 7: Significance Test Results

Model	EN			ZH			JA			ES		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1
RE²	69.71	73.74	71.67	76.80	82.62	79.60	70.16	81.39	75.36	70.26	81.78	75.59
RE²- constraint loss	58.76	82.16	68.52	74.77	81.01	77.77	69.13	80.75	74.49	69.41	81.05	74.78
Model	FR			IT			DE			PT		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1
RE²	70.05	83.97	76.38	74.34	80.83	77.45	73.01	78.94	75.86	48.06	78.99	59.76
RE²- constraint loss	71.54	80.57	75.79	72.32	79.32	75.66	71.74	75.59	73.61	46.98	74.02	57.48

Table 8: Precision, Recall and F1 score of ablation study of Constraint Loss on **RE²**