



Enhancing Visually-Rich Document Understanding via Layout Structure Modeling

Qiwei Li[†], Zuchao Li^{†*}, Xiantao Cai*, Bo Du and Hai Zhao

Limitations of Previous Work

1. Mismatch between the raw order and proper understanding order.

GSH TRAUMA REGISTRY DATA FORM	
[Patient sticker]	BP: _____ Respiration rate: _____
	Pulse: _____
	GCS: Eyes: _____ Verbal: _____ Motor: _____ Total: _____
	Neurological status: _____
Race: <input type="checkbox"/> Black <input type="checkbox"/> White <input type="checkbox"/> Coloured <input type="checkbox"/> Other _____	
<input type="checkbox"/> Alert <input type="checkbox"/> Responds to verbal stimuli <input type="checkbox"/> Unresponsive <input type="checkbox"/> Responds to painful stimuli	

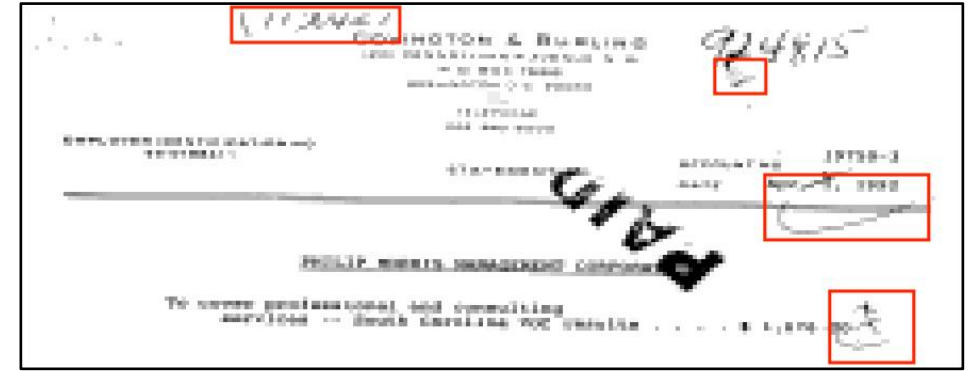
(a) **Raw Order:** GSH TRAUMA REGISTRY DATA FORM BP: Respiration rate: [Patient sticker] Pulse: GCS: Eyes: Verbal: Motor: Total: Neurological status: Race: Black White Coloured Other Alert Responds to verbal stimuli Unresponsive Responds to painful stimuli

(b) **Proper Order:** GSH TRAUMA REGISTRY DATA FORM [Patient sticker] Race: Black White Coloured Other BP: Respiration rate: Pulse: GCS: Eyes: Verbal: Motor: Total: Neurological status: Alert Responds to verbal stimuli Unresponsive Responds to painful stimuli

2. Inadequate disclosure of document structure by visual information.

1	BBQ Chicken	41,000
1	Sedang	0
ITEMS: 1		41,000
Total: 50,000		41,000
Pay Cash		Change 9,000

(a) Faulty Text Box



(b) Interference Information in Resized Image

GraphLayoutLM: Layout Language Model Enhanced by Layout Struture Graph

1. Layout graph modeling.

• Hierarchical relationship

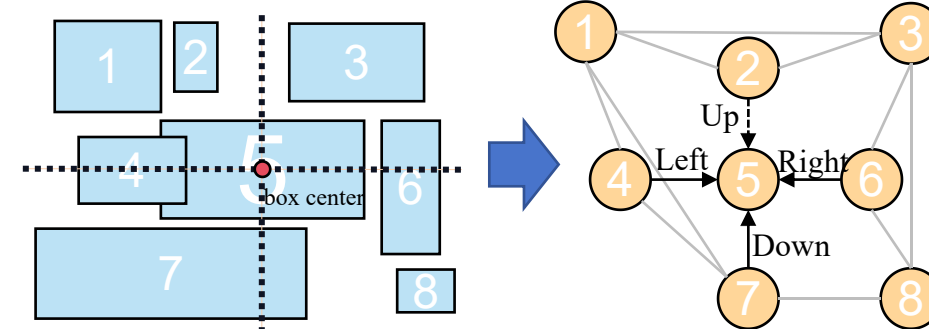
Parent node identification rule of region P :

$$L(n_i) = \begin{cases} 1, & \text{box}_i = \text{Top}(P) \text{ and } \text{box}_i = \text{Left}(P) \\ 0, & \text{otherwise} \end{cases}$$

Layout Hierarchical subtree construction of region P :

$$T_p = \langle n_p, C_p n_p, \text{Parent} - \text{Child} \rangle$$

• Positional relationship



• Layout structure graph construction

Hierarchical relationships are used to build the layout structure tree T_g .

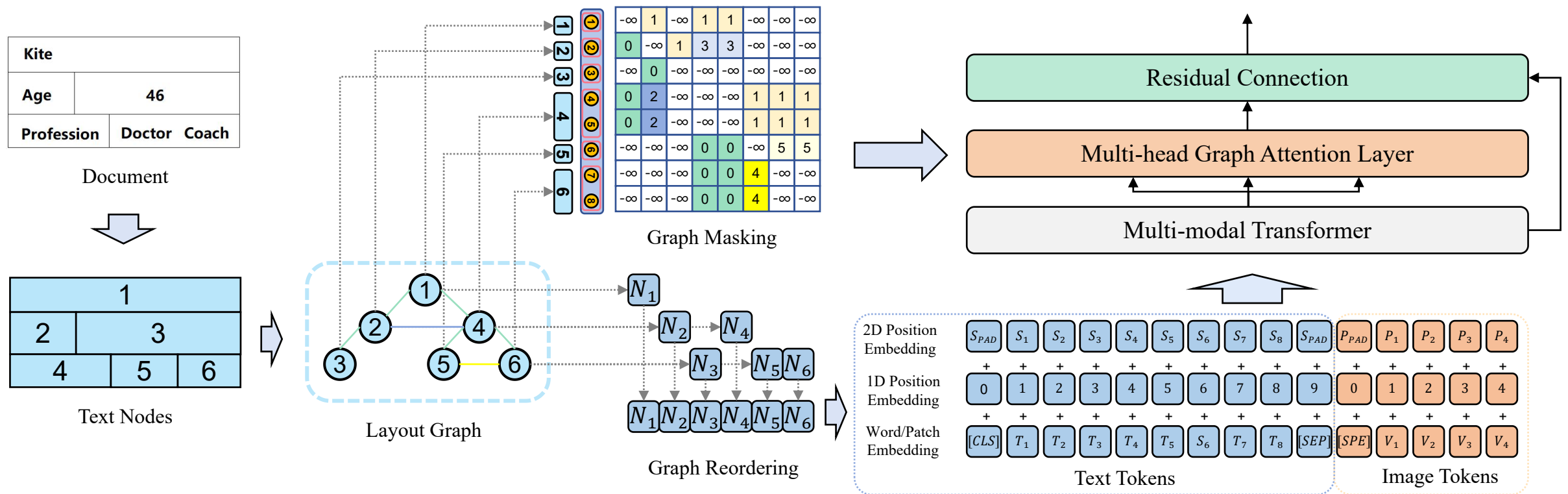
Positional relationships are used to enrich the tree into a layout structure graph G .

Representation of G :

$$G = T_g \cup \{ \langle n_i, n_j, \text{rel}(n_i, n_j) \rangle \mid n_i, n_j \in T_g, \text{sibling} \}$$

2. The architecture of GraphLayoutLM.

• GraphLayoutLM utilizes document layout structures graph through Graph Reordering and Graph Masking strategies.



Experiment Results

1. Experiment results of SER task on English datasets

Model	Parameters	Modality	FUNSD			CORD		
			Precision	Recall	F1	Precision	Recall	F1
BERT _{BASE}	110M	T	54.69	67.10	60.26	88.33	91.07	89.68
RoBERTa _{BASE}	125M	T	63.49	69.75	66.48	-	-	93.54
BROS _{BASE}	110M	T+L	81.16	85.02	83.05	95.58	95.14	95.36
LayoutLM _{BASE}	160M	T+L+I	76.77	81.95	79.27	94.37	95.08	94.72
XYLayoutLM _{BASE}	-	T+L+I	-	-	83.35	-	-	-
LayoutLMv2 _{BASE}	200M	T+L+I	80.29	85.37	82.76	94.53	95.39	94.95
DocFormer _{BASE}	183M	T+L+I	80.76	86.09	83.34	96.52	96.14	96.33
ERNIE-Layout _{BASE}	-	T+L+I	-	-	90.28	-	-	96.61
LayoutLMv3 _{BASE}	133M	T+L+I	-	-	90.29	-	-	96.56
LayoutLMv3 _{BASE} [†]	133M	T+L+I	90.82	91.55	91.19	96.35	96.71	96.53
GraphLayoutLM_{BASE} (Ours)	135M	T+L+I+G	92.46	93.85	93.15	97.02	97.53	97.28
BERT _{LARGE}	340M	T	61.13	70.85	65.63	88.86	91.68	90.25
RoBERTa _{LARGE}	355M	T	67.80	73.91	70.72	-	-	93.80
LayoutLM _{LARGE}	343M	T+L	75.96	82.19	78.95	94.32	95.54	94.93
BROS _{LARGE}	340M	T+L	82.81	86.31	84.52	-	-	97.28
StructuralLM _{LARGE}	355M	T+L	83.52	86.81	85.14	-	-	-
LayoutLMv2 _{LARGE}	426M	T+L+I	83.24	85.19	84.20	95.65	96.37	96.01
DocFormer _{LARGE}	536M	T+L+I	82.29	86.94	84.55	97.25	96.74	96.99
ERNIE-Layout _{LARGE}	-	T+L+I	-	-	93.12	-	-	97.21
LayoutLMv3 _{LARGE}	368M	T+L+I	-	-	92.08	-	-	97.46
LayoutLMv3 _{LARGE} [†]	368M	T+L+I	91.51	92.70	92.10	97.45	97.52	97.49
GraphLayoutLM_{LARGE} (Ours)	372M	T+L+I+G	94.49	94.30	94.39	97.75	97.75	97.75

2. Experiment results of SER task on Chinese datasets

Model	Modality	XFUND	
		F1	
XLM-RoBERTa _{BASE}	T	87.74	
XLM-RoBERTa _{LARGE}	T	89.25	
LayoutXLM _{BASE}	T+L+I	89.24	
LayoutXLM _{LARGE}	T+L+I	91.61	
XYLayoutLM _{BASE}	T+L+I	91.76	
ERNIE-LayoutX _{BASE} [‡]	T+L+I	88.58	
LayoutLMv3-Chinese _{BASE} [‡]	T+L+I	92.02	
LayoutLMv3-Chinese _{BASE} [†]	T+L+I	91.82	
GraphLayoutLM-Chinese_{BASE} (Ours)	T+L+I+G	93.56	

3. Ablation Study

Dataset	Graph Reorder	Graph Mask	Accuracy	Precision	Recall	F1
FUNSD	✗	✗	84.76	90.82	91.55	91.19
	✓	✗	85.70	92.36	93.15	92.75
	✗	✓	86.75	91.73	92.05	91.89
	✓	✓	88.39	92.46	93.85	93.15
CORD	✗	✗	97.11	96.35	96.71	96.53
	✓	✗	97.33	96.79	97.16	96.97
	✗	✓	97.88	96.94	97.23	97.09
	✓	✓	98.01	97.02	97.53	97.28
XFUND	✗	✗	85.87	89.79	93.94	91.82
	✓	✗	85.61	89.98	94.37	92.12
	✗	✓	90.88	91.58	94.43	92.99
	✓	✓	91.19	91.80	95.38	93.56