

# Neural Collaborative Graph Machines

for Table Structure Recognition

CVPR 2022 Hao Liu, et al. (Tencent). Reporter-TongKe Nee, 2023/04/01

# 1. Background

TSR (table structure recognition) task is recognizing the physical and logical structure of tables which are usually represented by image.

The conventional digital table image recognition has been relatively mature. Recent years, some researchers have focused on more challenging table structure recognition tasks, such as **recognizing the distortion table**.

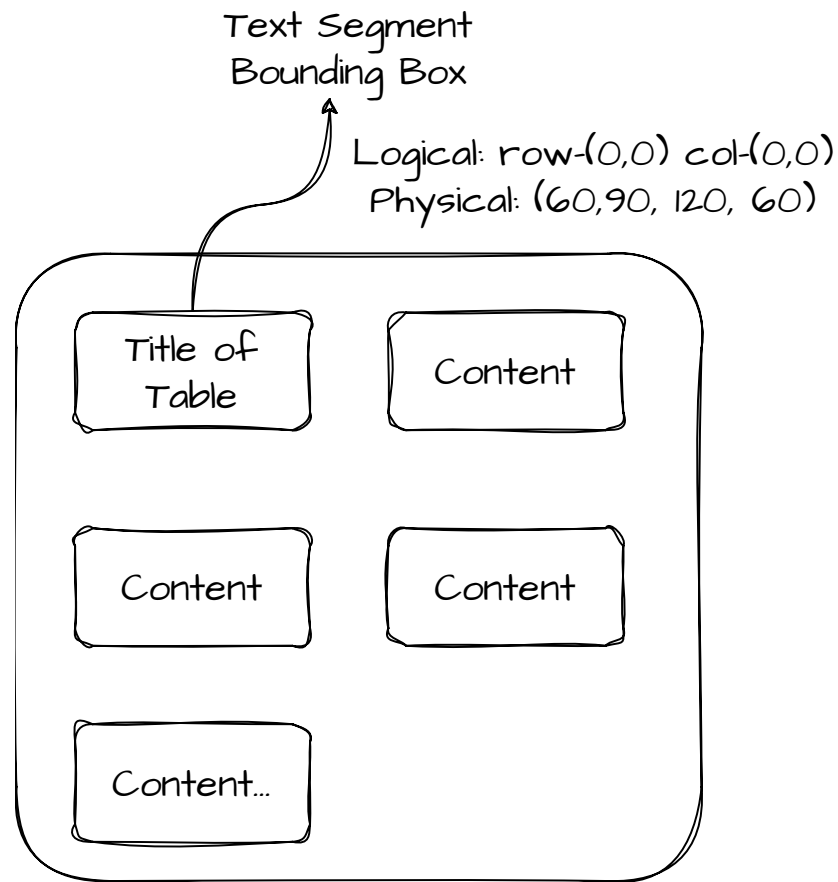
	POS tagging information				
	adj	verb	idiom	noun	other
pos	1,230	734	1,026	266	642
neg	785	904	746	165	797
neu	918	7,569	2,016	12,668	10,214
sum	2,933	9,207	3,788	13,099	11,653

An example of table structure recognition<sup>1</sup>

<sup>1</sup>Neural Collaborative Graph Machines for Table Structure Recognition

# Logical and Physical Position

The smallest unit in table recognition is text segment bounding box, and the task is to predict the position of each bounding box, in which the logical position refers to the row position (start (end) row, start (end) column), and the physical position refers to the coordinate position  $(x, y, w, h)$ .



An example of physical structure and logical structure

# Distorted Table and Regular Table

This image shows a distorted version of a table. The table is curved and wavy, making it difficult to read. It contains numerical data in a grid format. The columns are labeled 1 through 20, and the rows are labeled 1 through 20. The data is presented in a way that is not standard for a regular table.

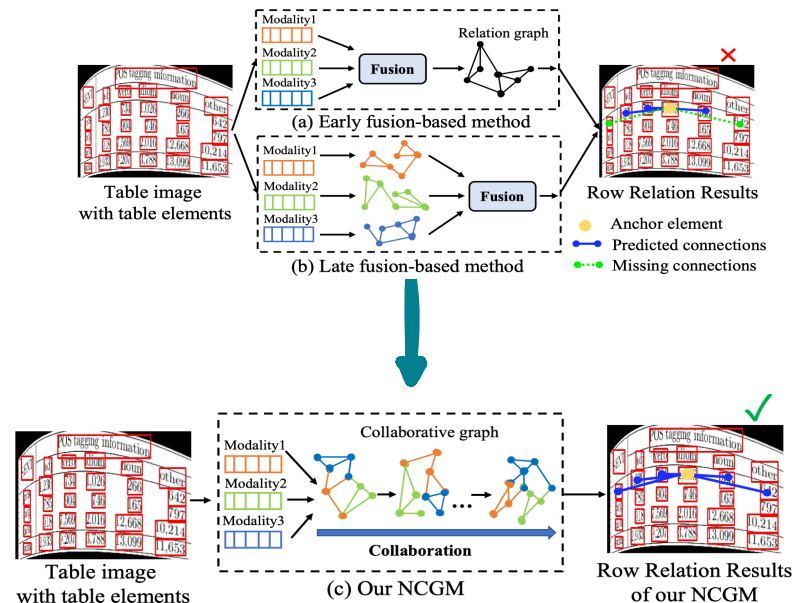
Preparation	Active agent(s) and concentration(s)	No neutralizers in sampling fluid		P-value (vs. reference) <sup>a</sup>	Neutralizers in sampling fluid		P-value (vs. reference) <sup>a</sup>
		Samples without detectable test bacteria	Mean log <sub>10</sub> reduction ± SD		Samples without detectable test bacteria	Mean log <sub>10</sub> reduction ± SD	
Product A	Ethanol (61%, w/w) chlorhexidine gluconate (1%, w/w)	9 / 15	4.8 ± 1.5	0.02	0 / 15	2.7 ± 0.4	0.033
Reference alcohol	Iso-propanol (60%, v/v)	0 / 15	2.7 ± 0.8	n.a.	0 / 15	2.5 ± 0.9	n.a.
Product B	Ethanol (83%, w/w)	0 / 15	3.3 ± 0.7	0.44	0 / 15	3.3 ± 0.7	0.11

Distorted vs Regular Table<sup>1</sup>

# Comparison

Traditional multi-modal fusion methods lack of attention to modal collaboration, so they do not perform well on distorted tables.

NCGM utilize **graph** to construct different modality relationship, and achieves good performance on the **distorted tables**.



Early (Late) fusion-based method with NCGM.

# Related Work

After entering the era of deep learning, there are several main ways to identify table structure:

## Boundary extraction-based

Based on the boundary extraction technology, the table structure is identified by detecting the **horizontal and vertical separator** in the table.

## Generative model-based

This method uses the **encoder-decoder** schema to generate HTML text to represent the table structure.

# Related Work

After entering the era of deep learning, there are several main ways to identify table structure:

## Graph-based

This method is mainly based on the **graph** model, and describes the table structure by representing the table elements as nodes and the relationship between them as edges.

## Transformer-based multimodal

Use early fused embedding as input (VLBERT, LayoutLM etc.), and the method of using two modalities to do co-attentional fusion (ViLBERT).

# Motivation

In table recognition, inductive biases of different modalities may affect the performance of the modality.

For example, when identifying regular tables, the coordinates of tables will dominate, but when dealing with distorted tables, they will become unreliable.

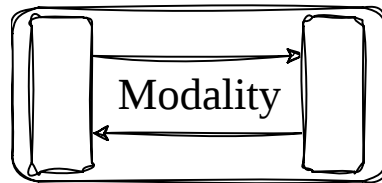


# Hetero-TSR Problem

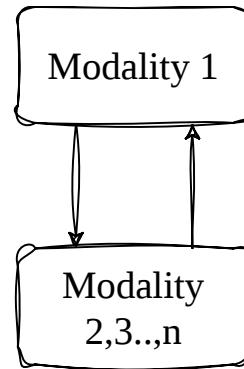
Can different modalities **collaborate** with each other rather than **interfering** with each other?

Two aspects to utilize the modalities.

- **Intra-Modality** (Among the modality)
- **Inter-Modality** (Between the modalities)



Intra

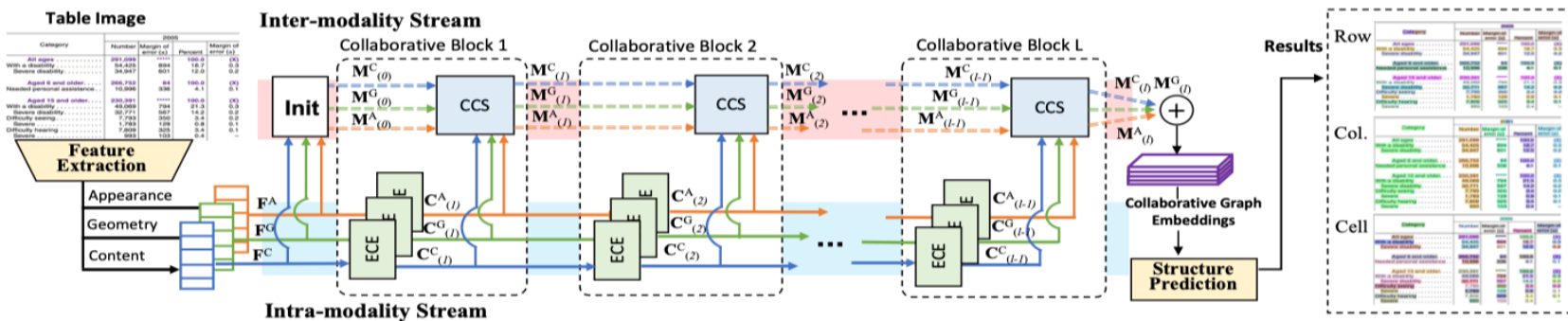


Inter

Intra-modality and Inter-modality.

## 2. Model Architecture

1. Three modalities: **Geometry**, **Appearance** and **Content**
2. Collaborative block: **ECE** (Intra modality) and **CCS** (Inter modality)
3. Structure prediction: **Fully connected layers**



The architecture of NCGM.<sup>1</sup>

<sup>1</sup>Neural Collaborative Graph Machines for Table Structure Recognition

# Feature Extraction

Three different modalities are used in the NCGM model, namely **Geometry**, **Appearance** and **Content**.

## Geometry

$\mathbf{F}^G \in \mathbb{R}^{N \times d}$ , Geometry modality mainly uses the **spatial information** of the cell.

## Appearance

$\mathbf{F}^A \in \mathbb{R}^{N \times d}$ , Appearance modality mainly uses cell **visual information**, background color, pixel information and so on.

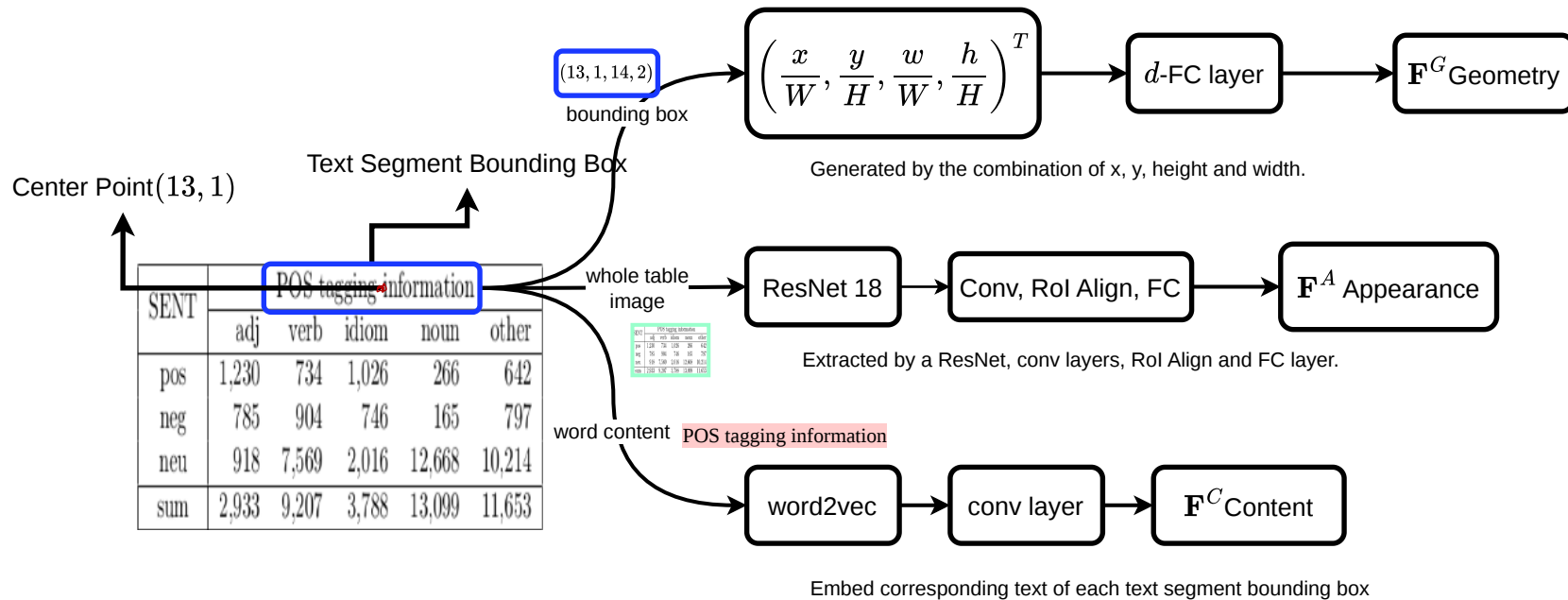
# Feature Extraction

Three different modalities are used in the NCGM model, namely **Geometry**, **Appearance** and **Content**.

## Content

$\mathbf{F}^C \in \mathbb{R}^{N \times d}$ , The Content modality mainly uses the **text content** of the cell.

# Feature Extraction

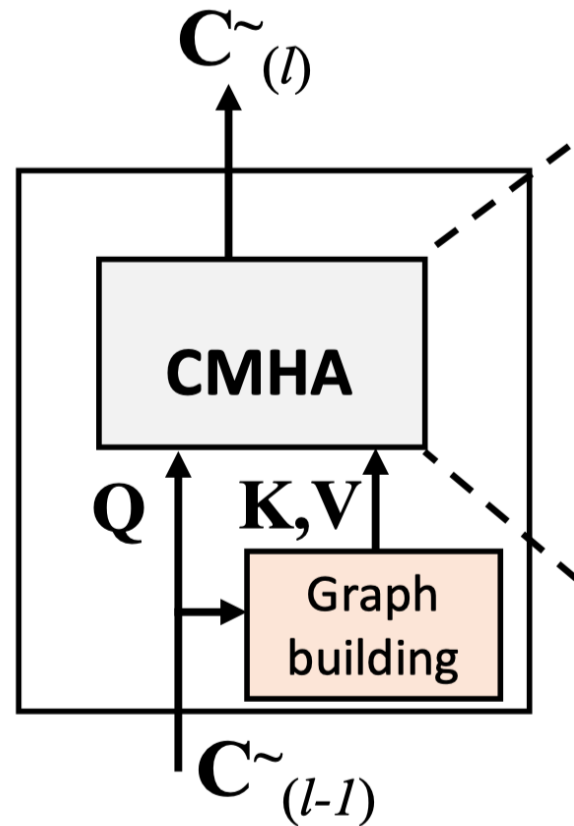


Three modalities

# ECE Module

Ego Context Extractor module (ECE) is used to extract the context information of intra modality. The feature of modality that input ECE will be constructed as a graph.

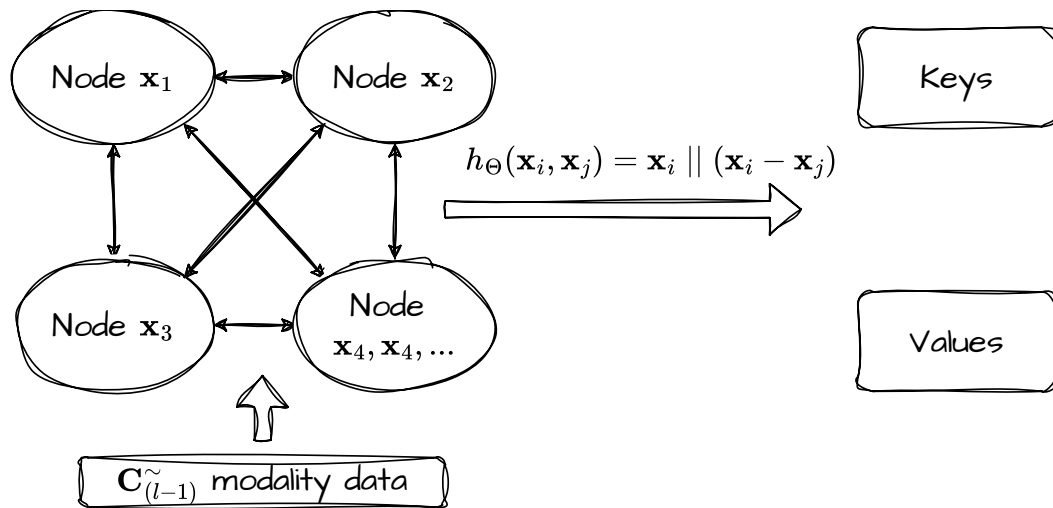
This input  $\mathbf{C}_{(l-1)}^{\sim}$  is copied into two copies, one as the **Query** input to CMHA and the other into the Graph building module to model the intra-modality context which is taken as Keys **K** and values **V**.



ECE Module Architecture

# ECE Module

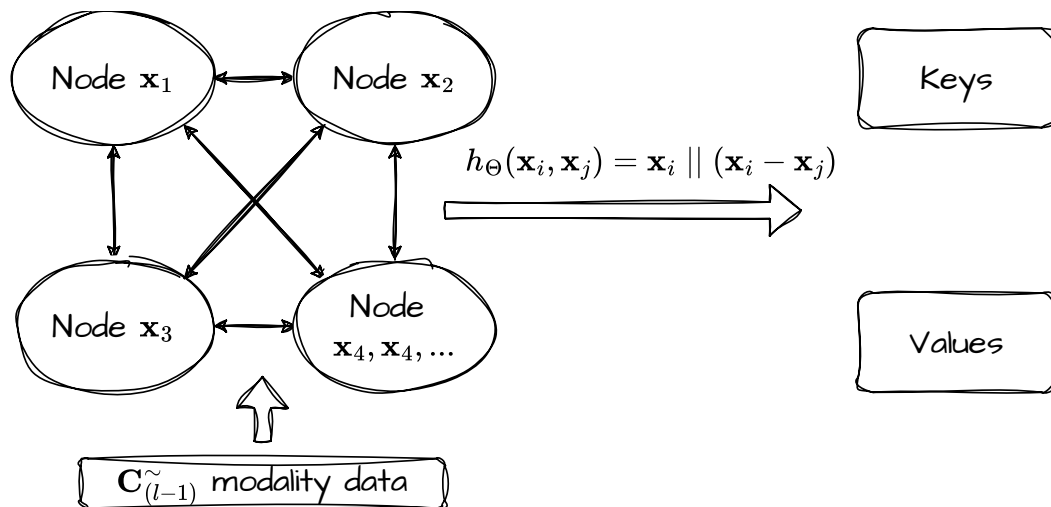
$\mathbf{G}^{\sim} = \{\mathcal{V}, \mathcal{E}\} \in \{\mathbf{G}^G, \mathbf{G}^A, \mathbf{G}^C\}$ . The nodes in  $\mathcal{V}$  is the node (cell) feature set, and  $\mathcal{E}$  is the full connected graph edge set.



Graph Building Module

# ECE Module

Authors adopt the following asymmetric edge function to combine graph edge features to each node.

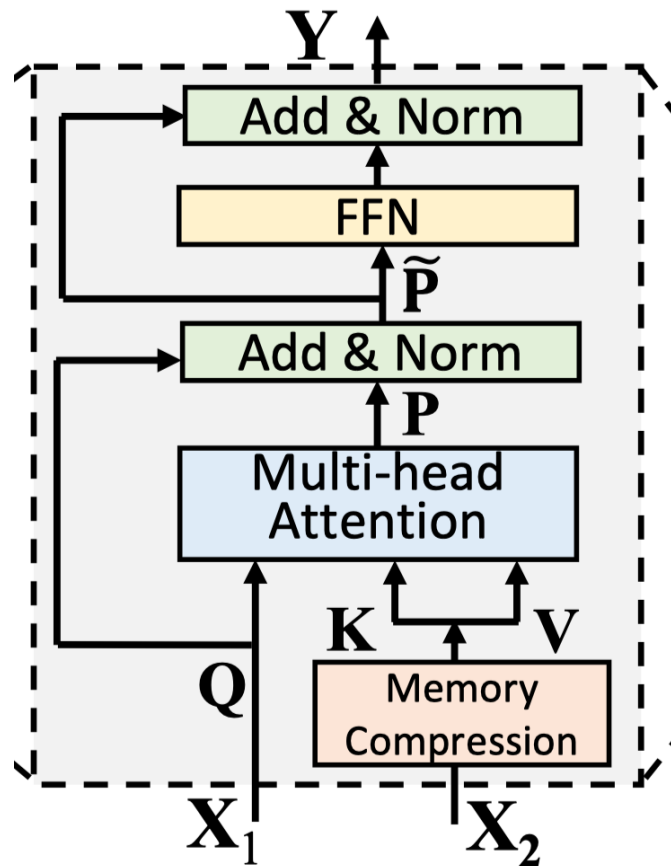




# CMHA Module

In order to reduce the amount of computation caused by too many dimensions, we introduce a memory compression module. In addition, we also introduce residual links to make the query information flow unimpeded.

$$\begin{aligned} \mathbf{Y} &= \text{Add\&Norm}(\text{FFN}(\tilde{\mathbf{P}}), \tilde{\mathbf{P}}) \\ \tilde{\mathbf{P}} &= \text{Add\&Norm}(\mathbf{Q}, \mathbf{P}), \\ \mathbf{P} &= \text{MHA}(\mathbf{Q}, \text{MC}(\mathbf{K}), \text{MC}(\mathbf{V})), \end{aligned}$$

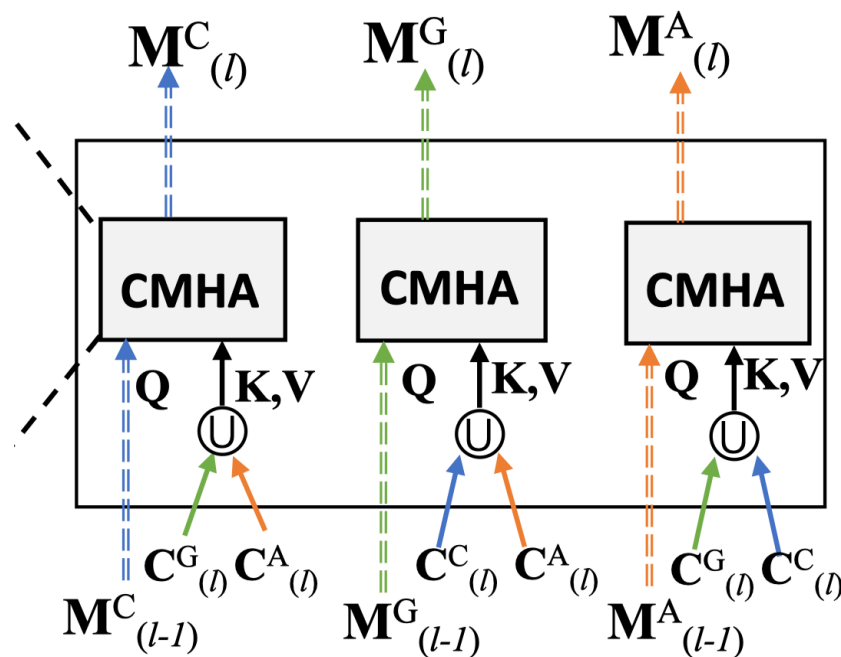


Compressed Multi-head Attention Module

# CCS Module

The Cross Context Synthesis module (CCS) fuses heterogeneous context graph embeddings collaboratively and learns collaborative patterns between different modalities.

Specifically, it takes three modality inputs as queries, and the union of the other two modalities as keys and values inputs to enable query modality to fully learn information from the other two modalities.



CCS Module Architecture

# Structure Prediction

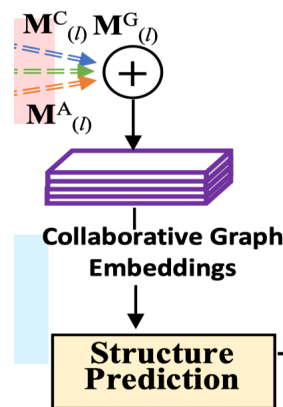
Based on the output embeddings

$$\mathbf{E} = \{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_N\} \in \mathbb{R}^{N \times d_e}$$

a set of node pairs is constructed, where each element is a vector formed by two node vectors concatenate, and then predicted by the full connection layer.

$$\mathbf{U} = \{\mathbf{u}_{1,1}, \mathbf{u}_{1,2}, \dots, \mathbf{u}_{i,j}, \dots, \mathbf{u}_{N,N}\} \in \mathbb{R}^{N^2 \times 2d_e}$$

Is to predict whether two nodes (cell) are in the same row, the same column, and to restore the table structure.



Results

Row

Col.

Cell

2005				
Category	Number	Margin of error (s)	Percent	Margin of error (s)
All ages .....	291,099	-----	100.0	(X)
With a disability .....	54,425	894	18.7	0.3
Severe disability .....	34,947	601	12.0	0.2
Aged 6 and older .....	266,782	84	100.0	(X)
Needed personal assistance .....	10,996	336	4.1	0.1
Aged 15 and older .....	230,391	-----	100.0	(X)
With a disability .....	49,069	794	21.3	0.3
Severe disability .....	32,771	567	14.2	0.2
Difficulty seeing .....	7,793	350	3.4	0.2
Severe .....	1,783	129	0.6	0.1
Difficulty hearing .....	7,809	325	3.4	0.1
Severe .....	993	103	0.4	—

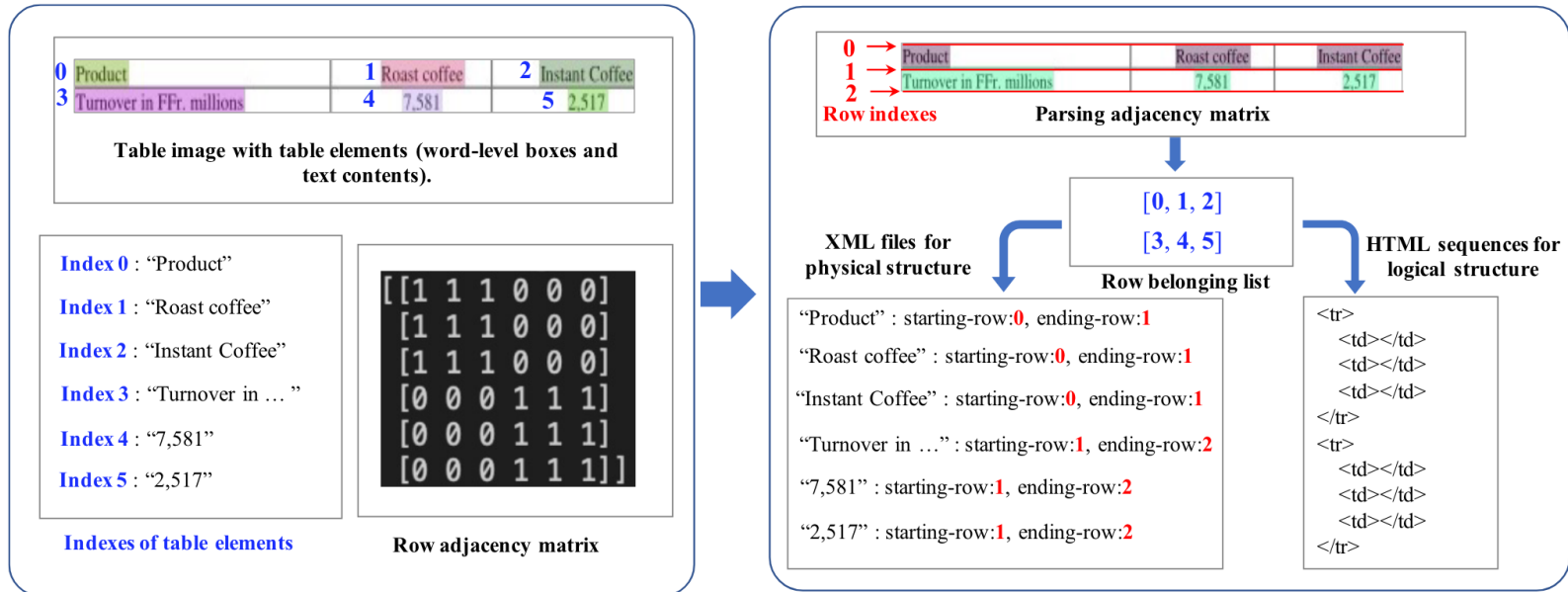
2005				
Category	Number	Margin of error (s)	Percent	Margin of error (s)
All ages .....	291,099	-----	100.0	(X)
With a disability .....	54,425	894	18.7	0.3
Severe disability .....	34,947	601	12.0	0.2
Aged 6 and older .....	266,782	84	100.0	(X)
Needed personal assistance .....	10,996	336	4.1	0.1
Aged 15 and older .....	230,391	-----	100.0	(X)
With a disability .....	49,069	794	21.3	0.3
Severe disability .....	32,771	567	14.2	0.2
Difficulty seeing .....	7,793	350	3.4	0.2
Severe .....	1,783	129	0.6	0.1
Difficulty hearing .....	7,809	325	3.4	0.1
Severe .....	993	103	0.4	—

2005				
Category	Number	Margin of error (s)	Percent	Margin of error (s)
All ages .....	291,099	-----	100.0	(X)
With a disability .....	54,425	894	18.7	0.3
Severe disability .....	34,947	601	12.0	0.2
Aged 6 and older .....	266,782	84	100.0	(X)
Needed personal assistance .....	10,996	336	4.1	0.1
Aged 15 and older .....	230,391	-----	100.0	(X)
With a disability .....	49,069	794	21.3	0.3
Severe disability .....	32,771	567	14.2	0.2
Difficulty seeing .....	7,793	350	3.4	0.2
Severe .....	1,783	129	0.6	0.1
Difficulty hearing .....	7,809	325	3.4	0.1
Severe .....	993	103	0.4	—

Predicting the structure of the table.

# Structure Prediction



Post processing. Convert adjacency matrix containing relationships to spanning information.

### 3. Result

Evaluation setting: Different methods utilize different information. Some methods use the cell/text bounding box, while others do not. Therefore, they design two different steps:

- **Step A:** Only with table image
- **Step B:** Along with the cell/text segment bounding box and text content

ICDAR-2013-P							
Method	Train Dataset	Setup-A			Setup-B		
		P	R	F1	P	R	F1
DGCNN [34]	Sci. + IC13-P	-	-	-	98.6	99.0	98.8
TabStr. [38]	Sci. + IC13-P	93.0	90.8	91.9	99.1	99.3	99.2
GTE [49]	Pub. + IC13-P	94.4	92.7	93.5	-	-	-
LGPMA [35]	Sci. + IC13-P	96.7	99.1	97.9	-	-	-
C-CTRNet [26]	WTW + IC19	95.5	88.3	91.7	-	-	-
FLAG-Net [24]	Sci. + IC13-P	97.9	<b>99.3</b>	98.6	99.2	99.5	99.3
<b>NCGM</b>	Sci. + IC13-P	<b>98.4</b>	<b>99.3</b>	<b>98.8</b>	<b>99.3</b>	<b>99.9</b>	<b>99.6</b>

ICDAR-2019							
DGCNN [34]	Sci. + IC19	80.3	77.8	79.0	-	-	-
TabStr. [38]	Sci. + IC19	82.2	78.7	80.4	97.5	95.8	96.6
C-CTRNet [26]	WTW	-	-	80.8	-	-	-
FLAG-Net [24]	Sci. + IC19	<b>85.2</b>	83.8	84.5	96.1	96.3	96.2
<b>NCGM</b>	Sci. + IC19	84.6	<b>86.1</b>	<b>85.3</b>	<b>98.9</b>	<b>98.8</b>	<b>98.8</b>

ICDAR 2013 Partial and ICDAR 2019

### 3. Result

SciTSR						
DGCNN [34]	Sci.	-	-	-	97.0	98.1 97.6
TabStr. [38]	Sci.	92.7	91.3	92.0	98.9	99.3 99.1
LGPMA [35]	Sci.	98.2	99.3	98.8	-	- -
FLAG-Net [24]	Sci.	<b>99.7</b>	99.3	99.5	<b>99.8</b>	99.5 99.6
<b>NCGM</b>	Sci.	<b>99.7</b>	<b>99.6</b>	<b>99.6</b>	99.7	<b>99.8 99.7</b>
SciTSR-COMP						
DGCNN [34]	Sci.	-	-	-	96.3	97.4 96.9
TabStr. [38]	Sci.	90.9	88.2	89.5	98.1	98.7 98.4
LGPMA [35]	Sci.	97.3	98.7	98.0	-	- -
FLAG-Net [24]	Sci.	98.4	98.6	98.5	98.6	99.0 98.8
<b>NCGM</b>	Sci.	<b>98.7</b>	<b>98.9</b>	<b>98.8</b>	<b>98.8</b>	<b>99.3 99.0</b>

SciTSR and SciTSR-COMP

# Result

$$\text{TEDS}(T_a, T_b) = 1 - \frac{\text{EditDist}(T_a, T_b)}{\max(|T_a|, |T_b|)}$$

where *EditDist* denotes tree-edit distance, and  $|T|$  is the number of nodes in  $T$ . The table recognition performance of a method on a set of test samples is defined as the mean of the TEDS score between the recognition result and ground truth of each sample.<sup>[1]</sup>

TableBank		
Method	Train Dataset	Setup-A
		BLEU
Image-to-Text [22]	TableBank	73.8
TabStruct-Net [38]	SciTSR	91.6
FLAG-Net [24]	SciTSR	93.9
<b>NCGM</b>	SciTSR	<b>94.6</b>
PubTabNet		
Method	Train Dataset	Setup-A
		TEDS
EDD [50]	PubTabNet	88.3
TabStruct-Net [38]	SciTSR	90.1
GTE [49]	PubTabNet	93.0
LGPMA [35]	PubTabNet	94.6
FLAG-Net [24]	SciTSR	95.1
<b>NCGM</b>	SciTSR	<b>95.4</b>

Logical structure recognition

<sup>1</sup>Image-based table recognition: data, model, and evaluation

# Ablation Study of Modality Fusion

Fusion Method	Input		Intra.			Inter.		Setup-B		
	Mix.	Ind.	DG.	Tr.	ECE	Con.	CCS	P	R	F1
Early Fusion	✓	✗	✓	✗	✗	✗	✗	96.3	97.4	96.8
	✓	✗	✗	✓	✗	✗	✗	95.1	95.6	95.3
	✓	✗	✗	✗	✓	✗	✗	97.8	98.3	98.0
Late Fusion	✗	✓	✓	✗	✗	✓	✗	96.9	98.2	97.5
	✗	✓	✗	✓	✗	✓	✗	94.9	96.1	95.5
	✗	✓	✗	✗	✓	✓	✗	98.4	98.2	98.3
<b>NCGM</b>	✗	✓	✗	✗	✓	✗	✓	<b>98.8</b>	<b>99.3</b>	<b>99.0</b>

Modality Fusion Abalation Study



# Ablation Study of Multi-Modality

Input Modality			Setup-B		
A	G	C	P	R	F1
✓	✗	✗	89.8	47.9	62.5
✗	✓	✗	97.9	97.7	97.8
✗	✗	✓	70.5	39.0	50.2
✓	✓	✗	98.6	98.3	98.4
✗	✓	✓	98.0	95.0	96.5
✓	✗	✓	87.6	89.3	88.4
✓	✓	✓	<b>98.8</b>	<b>99.3</b>	<b>99.0</b>

Modality Fusion Abalation Study

# Thinking about modalities collaboration

- What does ECE learn from the intra-modality?

Separate attention heads may learn to look for various relationships between inputs and introducing more sparsity and diversity for attention may improve performance and interpretability<sup>1</sup>.

- How do different modalities collaborate with each other?

Multi-head Attention.

---

<sup>1</sup>Sparse and constrained attention for neural machine translation.

# Conclusion

Authors proposed a novel graph-based method for heterogeneous table structure recognition through **intra-modality and inter-modality collaboration**.

Tests on various open data sets show the **effectiveness** of the method, and the importance of multi-modal cooperation for table structure identification.

**Limitations** include increased computational complexity and potential training collapse with deeper blocks. Future work can address these through refining the attention model.



# Thank You!

Any questions?

2023/04/01