# CaseGNN: Graph Neural Networks for Legal Case Retrieval with Text-Attributed Graphs

Yanran Tang, Ruihong Qiu, Yilun Liu, Xue Li, and Zi Huang

The University of Queensland
{yanran.tang, r.qiu, yilun.liu, helen.huang}@uq.edu.au,
xueli@eecs.uq.edu.au

**Abstract.** Legal case retrieval is an information retrieval task in the legal domain, which aims to retrieve relevant cases with a given query case. Recent research of legal case retrieval mainly relies on traditional bag-of-words models and language models. Although these methods have achieved significant improvement in retrieval accuracy, there are still two challenges: (1) **Legal structural information neglect**. Previous neural legal case retrieval models mostly encode the unstructured raw text of case into a case representation, which causes the lack of important legal structural information in a case and leads to poor case representation; (2) **Lengthy legal text limitation**. When using the powerful BERT-based models, there is a limit of input text lengths, which inevitably requires to shorten the input via truncation or division with a loss of legal context information. In this paper, a graph neural networks-based legal case retrieval model, CaseGNN, is developed to tackle these challenges. To effectively utilise the legal structural information during encoding, a case is firstly converted into a Text-Attributed Case Graph (TACG), followed by a designed Edge Graph Attention Layer and a readout function to obtain the case graph representation. The CaseGNN model is optimised with a carefully designed contrastive loss with easy and hard negative sampling. Since the text attributes in the case graph come from individual sentences, the restriction of using language models is further avoided without losing the legal context. Extensive experiments have been conducted on two benchmarks from COLIEE 2022 and COLIEE 2023, which demonstrate that CaseGNN outperforms other state-of-the-art legal case retrieval methods. The code has been released on https://github.com/yanran-tang/CaseGNN.

**Keywords:** Legal Case Retrieval · Graph Neural Networks.

## 1 Introduction

Legal case retrieval (LCR) is a specialised and indispensable retrieval task that focuses on retrieving relevant cases given a query case. For legal practitioners such as judges and lawyers, using retrieval models is more efficient than manually finding relevant cases by looking into thousands of legal documents. It is said that 59% of lawyers in the US are using web-based software to get technical services

and solution suggestions[1]. LCR also greatly helps a broader community who has legal questions but does not want to spend money on expensive consultation fees.

Existing LCR models can be categorised into two streams: statistical models and language models (LM). Statistical models [16,32,40] focus on measuring term frequency as the case similarity while LMs [1,2,7–9,20,22,26,37,41,48,49,51,55] conduct nearest neighbour search with case representations from language models. Among the LMs for LCR task, there are different case text matching strategies that focus on sentence [54], paragraph [41] and whole-case [20] levels.

Although the powerful LMs have achieved higher accuracy performance compared to traditional statistical models in LCR, two critical challenges still remain unsolved. (1) **Legal structural information neglect**. Under the context of legal domain, the case structural information typically refers to the relationship among different elements in a legal case, such as parties, crime activities and evidences. Recent LMs for LCR are trying to encode the unstructured raw text of a legal case into a high dimensional representation to measure the similarity with different text matching strategies [20,41,54]. However, only using bag-of-words statistical retrieval models or sequence LMs will restrict the interactions between different elements of cases, which will cause a significant loss of the useful structural information of legal cases. (2) **Lengthy legal text limitation**. As in the study of the LCR benchmark, COLIEE2023 [13], the average length of a case is 5,566 tokens [44], which exceeds most input limitations of LMs, such as 512-token limit of BERT-based models [11]. Therefore, most existing researches rely on truncation [20] or division [41] to shorten the input text. These pre-processing methods for adapting the LMs' input length will lead to incomplete legal case text and finally cause the loss of legal information from a global view.

To address the above two challenges, a novel CaseGNN framework is proposed in this paper. Firstly, for each case, the informative and useful case structural information that refers to the parties, crime activities or evidences of cases will be extracted by Named Entity Recognition and Relation Extraction tools to construct a case graph. Furthermore, to collaboratively utilise the textual and structural information in legal cases, the extracted structural information will be represented by LMs to transform the case graph into a Text-Attributed Case Graph (TACG). Secondly, to effectively obtain a case representation from the TACG, a CaseGNN framework is proposed by utilising the Edge Graph Attention Layer (EdgeGAT) and a readout function to obtain a graph level representation for retrieval. Finally, to train CaseGNN, a contrastive loss is designed to incorporate effective easy and hard negative samples. Empirical experiments are conducted on two benchmark datasets COLIEE 2022 [12] and COLIEE 2023 [13], which demonstrates that the proposed case structural information and CaseGNN can achieve state-of-the-art performance on LCR by effectively leveraging the structural information. The main contributions of this paper are summarised as follows:

– A CaseGNN framework is proposed for LCR to tackle the challenges on incorporating legal structural information and avoiding overlong input text.

---

[1] https://www.clio.com/blog/lawyer-statistics/

- A Text-Attributed Case Graph (TACG) is developed to transform the format of unstructured case text into structural and textual format.
- A GNN layer called edge graph attention layer (EdgeGAT) is designed to learn the representation in the TACG.
- Extensive experiments conducted on two benchmark datasets demonstrate the state-of-the-art performance of CaseGNN.

## 2  Related Work

### 2.1  Legal Case Retrieval Models

As a specialised information retrieval (IR) task, the methods of LCR task can be categorised into two streams as IR task: statistical retrieval models [16, 32, 40] and language models [17, 31, 33, 39]. In LCR task, TF-IDF [16], BM25 [40] and LMIR [32] are the statistical models that also frequently used, which are all based on calculating the text matching score by utilising term frequency and inverse document frequency of words in legal case. General LMs [10, 11, 25, 30] are highly used for LCR task by using LMs to encode the case into representative embeddings for the powerful language understanding ability of LMs [1, 2, 4–9, 22, 23, 26, 37, 43, 45, 48, 49, 51, 53–56]. In the state-of-the-art research, to tackle the long text problem in legal domain, BERT-PLI [41] divides cases into paragraphs and calculates the similarity between two paragraphs while SAILER [20] directly truncates the case text to cope with the input limit of LM.

### 2.2  Graph Neural Networks

GNN models can effectively capture the structural information from graph data [14, 19, 24, 34–36, 47]. To further utilise the edge information, SCENE [29] proposed a GNN layer that can deal with the edge weights. Recently, text-attributed graph are widely used for the capacity of combining both the text understanding ability of LMs and the structural information of graphs, such as TAPE [15], G2P2 [50] and TAG [21].

There are two existing graph-based legal understanding methods, LegalGNN for legal case recommendation [52] and SLR for LCR [28]. Both methods utilise an external legal knowledge database, such as legal concepts and charges, to construct a knowledge graph with human knowledge while encoding legal cases with general LMs. Our proposed CaseGNN is different from these two methods that there is no external knowledge and the encoding of a case actually uses the structural information from the case itself.

## 3  Preliminary

In the following, a bold lowercase letter denotes a vector, a bold uppercase letter denotes a matrix, a lowercase letter denotes a scalar or a sequence of words, and a scripted uppercase letter denotes a set.

### 3.1  Task Definition

In legal case retrieval, given the query case $q$, and the set of $n$ candidate cases $D = \{d_1, d_2, ..., d_n\}$, our task is to retrieve a set of relevant cases $D^* = \{d_i^* | d_i^* \in D \land relevant(d_i^*, q)\}$, where $relevant(d_i^*, q)$ denotes that $d_i^*$ is a relevant case of the query case $q$. The relevant cases are called precedents in legal domain, which refer to the historical cases that can support the judgement of the query case.

### 3.2  Graph Neural Networks

*Graph:* A graph is denoted as $G = (V, E)$, where a node $v$ with feature $\mathbf{x}_v \in \mathbb{R}^d$ for $v \in V$, and an edge with feature $\mathbf{e}_{uv} \in \mathbb{R}^d$ for $e \in E$ between node $u$ and $v$.

*Graph Neural Networks:* GNNs utilise node features, edge features, and the graph structure to learn representations for nodes, edges and the graph. Most GNNs use iterative neighbourhood aggregation to calculate the representations. After $l - 1$ iterations of aggregation, the output features of a node $v$ after $l$-th layer is:

$$\mathbf{h}_v^l = \mathrm{Map}^l(\mathbf{h}_v^{l-1}, \mathrm{Agg}\ (\mathbf{h}_v^{l-1}, \mathbf{h}_u^{l-1}, \mathbf{h}_{e_{uv}}) : u \in N(v))), \tag{1}$$

where $\mathbf{h}_v^l \in \mathbb{R}^d$ is the node representation of $v$ at $l$th layer, $\mathbf{h}_{e_{uv}} \in \mathbb{R}^d$ is the edge representation between node $u$ and $v$, and $N(v)$ is the neighbour node set of node $v$. Specially, the input of the first layer is initialised as $\mathbf{h}_v^0 = \mathbf{x}_v$. Agg and Map are two functions that can be formed in different ways, where Agg performs aggregation to the neighbour node features and edge features while Map utilises the node self features and the neighbour features together for mapping node $v$ to a new feature vector. To generate a graph representation $\mathbf{h}_G$, a Readout function is used to transform all the node features:

$$\mathbf{h}_G = \mathrm{Readout}(G). \tag{2}$$

## 4  Method

### 4.1  Text-Attributed Case Graph

Text-Attributed Case Graph (TACG) aims to convert the unstructured case text into a graph. To construct a TACG, the structure and the features of the graph will be obtained by using information extraction tools and language models.

**Information Extraction.** To leverage the legal structural information for graph construction, named entity recognition tool and relation extraction tool are used for information extraction. From the legal perspective, the determining factor of relevant cases is the alignment of *legal fact* and *legal issue* [44]. Specifically, legal fact is a basic part of a case that describe "who, when, what, where and why" while legal issue is the legal disputes between parties of a case and need to be settled by judges [44]. The details of generating legal fact and legal issue
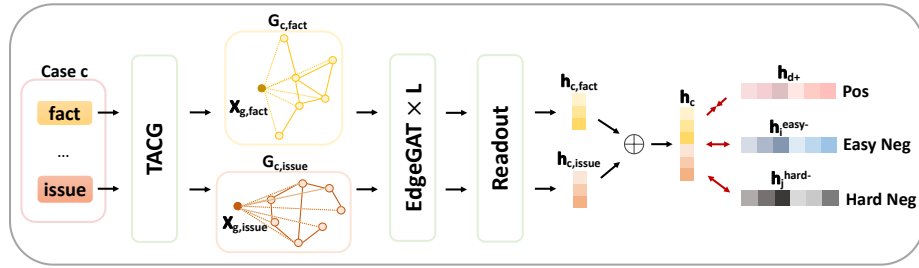
Fig. 1: The framework of CaseGNN. Given legal case $c$, the legal fact and the legal issue sections are converted into TACG based on information extraction and text encodings. The TACG is processed by L layers of EdgeGAT and a Readout function to obtain an overall case graph representation. The whole framework is trained with the contrastive loss with positive and negative samples.

can be found in PromptCase [44]. Therefore, in this paper, the important legal structural information refers to the relation triplets that generated from legal fact and legal issue, which are extracted from a case text using the PromptCase framework [44]. For example, in COLIEE datasets collected from the federal court of Canada [12, 13], a triplet example is extracted as $(applicant, is, Canadian)$ from a sentence "The applicant is a Canadian." in legal fact of a case, where $applicant$ denotes the "who" and $is, Canadian$ refers to the "what" in the case. After conducting information extraction, a set of triplets $R = \{(h, r, t)_{i=1:n}\}$ can be obtained, where $h$ is the head entity, $t$ is the tail entity, $r$ is the relation between $h$ and $t$, and $n$ is the number of triplets in a case.

**Graph Construction.** With the extracted set of triplets, the graph construction is to convert these triplets into a case graph. For a legal case $c$ and its triplet set $R_{c,\text{fact}}$ and $R_{c,\text{issue}}$ specifically for its *legal fact* and *legal issue*, the TACG is con-



Fig. 2: The constructed legal fact graph, $G_{c,\text{fact}}$ and legal issue graph, $G_{c,\text{issue}}$ of case $c$.

structed as $G_{c,\text{fact}} = (V_{c,\text{fact}}, E_{c,\text{fact}})$ and $G_{c,\text{issue}} = (V_{c,\text{issue}}, E_{c,\text{issue}})$. For $G_{c,\text{fact}}$, $V_{c,\text{fact}}$ includes the set of nodes of all head and tail entities $h$ and $t$ in $R_{c,\text{fact}}$ as $v_h$ and $v_t$, and $E_{c,\text{fact}}$ includes the set of edges corresponding to the relations $r$ from head entity $h$ to tail entity $t$ in $R_{c,\text{fact}}$ as $e_{v_h v_t}$. The same construction process is applied to the $G_{c,\text{issue}}$ for the legal issue. Additionally, for both the legal fact graph and the legal issue graph of a case, two virtual node $v_{\text{g, fact}}$ and $v_{\text{g, issue}}$ that representing the global textual semantics are added to fact graph and issue graph respectively. To help propagate the global information in the
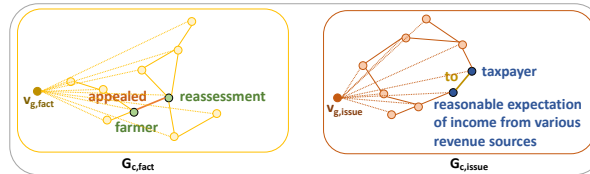
graph representation learning, this virtual global node is connected to every node in the graph. The detailed illustration of the TACG is demonstrated in Fig. 2.

**Text Attribute.** In the case graph above with the extracted entities as nodes and relations as edges, the node and edge features are obtained by using language models encoding of the text in nodes and edges. For node $u$, node $v$, and edge $e_{uv}$ in the case, the text attribute encoding is processed as :

$$\mathbf{x}_u = \mathrm{LM}(t_u); \quad \mathbf{x}_v = \mathrm{LM}(t_v); \quad \mathbf{x}_{e_{uv}} = \mathrm{LM}(t_{e_{uv}}), \tag{3}$$

where $t_u, t_v, t_{e_{uv}}$ is the text of node $u$, node $v$ and edge $e_{uv}$ respectively, and LM is a pre-trained language model, such as BERT [11], SAILER [20] or PromptCase [44]. $\mathbf{x}_u \in \mathbb{R}^d$, $\mathbf{x}_v \in \mathbb{R}^d$, and $\mathbf{x}_{e_{uv}} \in \mathbb{R}^d$ are the output of text-attributed encoding and serve as the feature vector of node $u$, $v$ and edge $e_{uv}$. For the virtual global nodes $v_{\mathrm{g,fact}}$ and $v_{\mathrm{g,issue}}$, the node feature is whole text encoding of the legal fact and legal issue extracted by using PromptCase [44] from a case respectively. The feature of the edge between the virtual gloabal node and other nodes such as entity $u$ will directly reuse the feature of other nodes $u$ in the TACG to simplify the feature extraction:

$$
\begin{aligned}
\mathbf{x}_{v_{\mathrm{g,fact}}} &= \mathrm{LM}(t_{\mathrm{fact}}); \quad \mathbf{x}_{e_{uv_{\mathrm{g,fact}}}} = \mathbf{x}_u; \\
\mathbf{x}_{v_{\mathrm{g,issue}}} &= \mathrm{LM}(t_{\mathrm{issue}}); \quad \mathbf{x}_{e_{uv_{\mathrm{g,issue}}}} = \mathbf{x}_u.
\end{aligned} \tag{4}
$$

### 4.2   Edge Graph Attention Layer

After obtaining the text-attributed features of nodes and edges, a self-attention module will be used to aggregate the nodes and its neighbour nodes and edges information to an informative representation. Moreover, to avoid over-smoothing, a residual connection is added. As previous study shows, multi-head attention has better performance than original attention [46]. According to multi-head attention mechanism, the update of node $v$ with $K$ attention heads is defined as:

$$\mathbf{h}_v^{'} = \mathbf{W}_s \cdot \mathbf{h}_v + \underset{k=1:K}{\mathrm{Avg}} \Big( \sum_{u \in N(v)} \alpha^k (\mathbf{W}_n^k \cdot \mathbf{h}_u + \mathbf{W}_e^k \cdot \mathbf{h}_{e_{uv}}), \tag{5}$$

where $\mathbf{h}_v^{'} \in \mathbb{R}^{d'}$ is the updated node feature, and Avg means the average of the output vectors of K heads. Specially, the input of the first EdgeGAT layer is initialised as $\mathbf{h}_v = \mathbf{x}_v$, $\mathbf{h}_u = \mathbf{x}_u$ and $\mathbf{h}_{e_{uv}} = \mathbf{x}_{e_{uv}}$. All the weight matrices $\mathbf{W}$ are in $\mathbb{R}^{d \times d'}$. Specifically, $\mathbf{W}_s$ is the node self update weight matrix, $\mathbf{W}_n$ is the neighbour node update weight matrix and $\mathbf{W}_e$ is the edge update weight matrix respectively. $\alpha^k$ is the attention weight in the attention layer as:

$$\alpha^k = \mathrm{Softmax}(\mathrm{LeakyReLU}({\mathbf{w}_{\mathrm{att}}^k}^T [\mathbf{W}_n^k \cdot \mathbf{h}_v \parallel \mathbf{W}_n^k \cdot \mathbf{h}_u \parallel \mathbf{W}_e^k \cdot \mathbf{h}_{e_{uv}}])), \tag{6}$$

where Softmax is the softmax function, LeakyReLU is the non-linear function, $\mathbf{w}_{\mathrm{att}}^k \in \mathbb{R}^{3d'}$ is the weight vector of attention layer and $\parallel$ denotes concatenation of vectors. Specifically, the same edge features are reused to make the EdgeGAT simpler. Further model development of updating edge can be designed.

### 4.3   Readout Function

With the updated node and edge representations, a graph readout function is designed to obtain the case graph representation for the case $c$:

$$\mathbf{h}_{c,\text{fact}} = \text{Readout}(G_{c,\text{fact}}), \quad \mathbf{h}_{c,\text{issue}} = \text{Readout}(G_{c,\text{issue}}), \tag{7}$$

where Readout is an aggregation function to output an overall representation in the graph level. One example is the average pooling of all the node embeddings:

$$\mathbf{h}_{c,\text{fact}} = \text{Avg}(\mathbf{h}_{v_i}|v_i \in V_{c,\text{fact}}); \quad \mathbf{h}_{c,\text{issue}} = \text{Avg}(\mathbf{h}_{v_i}|v_i \in V_{c,\text{issue}}). \tag{8}$$

In addition, since the virtual global node has already been in TACG, the updated virtual global node vector $\mathbf{h}_g$ can be considered as the final representation of the case graph:

$$\mathbf{h}_{c,\text{fact}} = \mathbf{h}_{g,\text{fact}}; \quad \mathbf{h}_{c,\text{issue}} = \mathbf{h}_{g,\text{issue}}. \tag{9}$$

For the experiments in this paper, the final virtual global node vector is used as the graph representation.

For case $c$, the fact graph feature $\mathbf{h}_{c,\text{fact}}$ and issue graph feature $\mathbf{h}_{c,\text{issue}}$ are generated by using $\mathbf{h}_{c,\text{fact}}$ and $\mathbf{h}_{c,\text{issue}}$ respectively. Therefore, the case graph representation $\mathbf{h}_c \in \mathbb{R}^{2d'}$ is the concatenation of $\mathbf{h}_{c,\text{fact}}$ and $\mathbf{h}_{c,\text{issue}}$:

$$\mathbf{h}_c = \mathbf{h}_{c,\text{fact}} \parallel \mathbf{h}_{c,\text{issue}}. \tag{10}$$

### 4.4   Objective Function

To train the CaseGNN model for the LCR task, it is required to distinguish the relevant cases from irrelevant cases given a large candidate case pool. To provide the training signal, contrastive learning is a tool that aims at pulling the positive samples closer while pushing the negative samples far away used in retrieval tasks [20, 49, 55]. In this paper, given a query case $q$ and a set of candidate case $D$ that includes both relevant cases $d^+$ and irrelevant cases $d^-$, the objective function is defined as a contrastive loss:

$$\ell = -\log \frac{e^{(s(\mathbf{h}_q, \mathbf{h}_{d^+}))/\tau}}{e^{(s(\mathbf{h}_q, \mathbf{h}_{d^+}))/\tau} + \sum_{i=1}^n e^{(s(\mathbf{h}_q, \mathbf{h}_{d_i^{easy-}}))/\tau} + \sum_{j=1}^m e^{(s(\mathbf{h}_q, \mathbf{h}_{d_j^{hard-}}))/\tau}}, \tag{11}$$

where $s$ is the similarity metric such as dot product or cosine similarity, $n$ is the number of easy negative samples, $m$ is the number of hard negative samples, and $\tau$ is the temperature coefficient. During training, the positive samples are given by the ground truth from the dataset. The easy negative samples are randomly sampled from the whole candidate pool as well as using the in-batch samples from other queries. For the hard negative samples, it is designed to make use of harder samples to effectively guide the training. Therefore, hard negative samples are sampled based on the BM25 relevance score. If a candidate case has a high score from BM25 yet it is not a positive case, such a case is considered as a hard negative case because it has a high textual similarity to the query case while it is still not a positive case. The overall pipeline is detailed in Fig. 1.

## 5   Experiments

### 5.1   Setup

**Datasets.**   To evaluate the proposed CaseGNN, the experiments are conducted on two benchmark LCR datasets, COLIEE2022 [12] and COLIEE2023 [13] from the Competition on Legal Information Extraction/Entailment (COLIEE), where the cases are collected from the federal court of Canada. Given a query case,

Table 1: Statistics of datasets.

| Datasets | COLIEE2022 | | COLIEE2023 | |
|---|---|---|---|---|
| | train | test | train | test |
| # Query | 898 | 300 | 959 | 319 |
| # Candidates | 4415 | 1563 | 4400 | 1335 |
| # Avg. relevant cases | 4.68 | 4.21 | 4.68 | 2.69 |
| Avg. length (# token) | 6724 | 6785 | 6532 | 5566 |
| Largest length (# token) | 127934 | 85136 | 127934 | 61965 |

relevant cases are retrieved from the entire candidate pool. The difference between two datasets are: (1) Although the training sets have overlap, the test sets are totally different; (2) As shown in Table 1, the average relevant cases numbers per query are different, leading to different difficulties in finding relevant cases. These datasets focus on the most widely used English legal case retrieval benchmarks and CaseGNN can be easily extended to different languages with the corresponding information extraction tools and LMs.

**Metrics.**   In this experiment, the metric of precision (P), recall (R), Micro F1 (Mi-F1), Macro F1 (Ma-F1), Mean Reciprocal Rank (MRR), Mean Average Precision (MAP) and normalized discounted cumulative gain (NDCG) are used for evaluation. According to the previous LCR works [20, 27, 44], top 5 ranking results are evaluated. All metrics are the higher the better.

**Baselines.**   According to the recent research [20, 44], 5 popular and state-of-the-art methods are compared as well as the competition winners:

- **BM25** [40]: a strong retrieval benchmark that leverages both term frequency and inverse document frequency for retrieval tasks.
- **LEGAL-BERT** [8]: a legal LM that is pre-trained on large English corpus.
- **MonoT5** [30]: a pre-trained LM that utilises T5 [38] architecture for document ranking tasks.
- **SAILER** [20]: a pre-trained legal structure-aware LM that obtains competitive performance on both datasets.
- **PromptCase** [44]: an input reformulation method that works on LM for LCR, which achieves sate-of-the-art performance on COLIEE2023 [13] dataset. Two-stage usage of PromptCase with BM25 is evaluated as well.

**Implementation.**   The French text in both datasets are removed. The spaCy[2], Stanford OpenIE [3] and LexNLP[3] packages are used for information extraction.

---

[2] https://spacy.io/

[3] https://github.com/LexPredict/lexpredict-lexnlp

Table 2: Overall performance on COLIEE2022 and COLIEE2023 (%). Underlined numbers indicate the best baselines. Bold numbers indicate the best performance of all methods. Both one-stage and two-stage results are reported.

| Methods | COLIEE2022 | | | | | | | COLIEE2023 | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | P@5 | R@5 | Mi-F1 | Ma-F1 | MRR@5 | MAP | NDCG@5 | P@5 | R@5 | Mi-F1 | Ma-F1 | MRR@5 | MAP | NDCG@5 |
| **One-stage** | | | | | | | | | | | | | | |
| BM25 | 17.9 | 21.2 | 19.4 | 21.4 | 23.6 | 25.4 | 33.6 | 16.5 | 30.6 | 21.4 | 22.2 | 23.1 | 20.4 | 23.7 |
| LEGAL-BERT | 4.47 | 5.30 | 4.85 | 5.38 | 7.42 | 7.47 | 10.9 | 4.64 | 8.61 | 6.03 | 6.03 | 11.4 | 11.3 | 13.6 |
| MonoT5 | 0.71 | 0.65 | 0.60 | 0.79 | 1.39 | 1.41 | 1.73 | 0.38 | 0.70 | 0.49 | 0.47 | 1.17 | 1.33 | 0.61 |
| SAILER | 16.6 | 15.2 | 14.0 | 16.8 | 17.2 | 18.5 | 25.1 | 12.8 | 23.7 | 16.6 | 17.0 | 25.9 | 25.3 | 29.3 |
| PromptCase | 17.1 | 20.3 | 18.5 | 20.5 | 35.1 | 33.9 | 38.7 | 16.0 | 29.7 | 20.8 | 21.5 | 32.7 | 32.0 | 36.2 |
| CaseGNN (Ours) | **35.5±0.2** | **42.1±0.2** | **38.4±0.3** | **42.4±0.1** | **66.8±0.8** | **64.4±0.9** | **69.3±0.8** | **17.7±0.7** | **32.8±0.7** | **23.0±0.5** | **23.6±0.5** | **38.9±1.1** | **37.7±0.8** | **42.8±0.7** |
| **Two-stage** | | | | | | | | | | | | | | |
| SAILER | **23.8** | 25.7 | 24.7 | 25.2 | 43.9 | 42.7 | 48.4 | 19.6 | 32.6 | 24.5 | 23.5 | 37.3 | 36.1 | 40.8 |
| PromptCase | 23.5 | 25.3 | 24.4 | **30.3** | 41.2 | 39.6 | 45.1 | **21.8** | 36.3 | **27.2** | 26.5 | 39.9 | 38.7 | 44.0 |
| CaseGNN (Ours) | 22.9±0.1 | **27.2±0.1** | **24.9±0.1** | 27.0±0.1 | **54.9±0.4** | **54.0±0.5** | **57.3±0.6** | 20.2±0.2 | **37.6±0.5** | 26.3±0.3 | **27.3±0.2** | **45.8±0.9** | **44.4±0.8** | **49.6±0.8** |

Two-stage experiment uses the top 10 retrieved cases by BM25 as the first stage result. The embedding size are set to 768 based on BERT. The number of EdgeGAT layers are set to 2 and the number of EdgeGAT heads are chosen from {1, 2, 4}. The training batch sizes are chosen from {16, 32, 64, 128}. The Dropout [42] rate of every layer's representation is chosen from {0.1, 0.2, 0.3, 0.4, 0.5}. Adam [18] is applied as optimiser with the learning rate chosen from {0.00001, 0.00005, 0.0001, 0.0005, 0.000005} and weight decay from {0.00001, 0.0001, 0.001, 0.01}. For every query during training, the number of positive sample is set to 1; the number of randomly chosen easy negative sample is set to 1; the number of hard negative samples is chosen from {1, 5, 10, 30}. The in-batch samples from other queries are also employed as easy negative samples. SAILER [20] is chosen as the LM model to generate the text attribute of nodes and edges, which is a BERT-based model that pre-trained and fine-tuned on large corpus of legal cases.

## 5.2 Overall Performance

In this experiments, the overall performance of CaseGNN is evaluated on COL-IEE2022 and COLIEE2023 by comparing with state-of-the-art models, as shown in Table 2. According to the results, CaseGNN achieves the best performance compared with all the baseline models by a large margin in both one-stage and two-stage settings. In COLIEE2022, one-stage CaseGNN has a much higher performance than other two-stage methods, and in COLIEE2023, one-stage CaseGNN has a comparable results with two-stage methods.

For one-stage retrieval setting, compared to the state-of-the-art performance on COLIEE2022 and COLIEE2023, CaseGNN significantly improved the LCR performance by utilising the important legal structural information with graph neural network. Compared with the strong baseline of traditional retrieval model BM25, CaseGNN achieves outstanding performance. The reason of inferior performance of BM25 is because only using the term frequency will ignore the important legal semantics of a case. The performance of CaseGNN also outperforms LEGAL-BERT, a legal corpus pre-trained LCR model, which indicates that

Table 3: Ablation study. (%)

| Methods | COLIEE2022 | | | | | | | COLIEE2023 | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | P@5 | R@5 | Mi-F1 | Ma-F1 | MRR@5 | MAP | NDCG@5 | P@5 | R@5 | Mi-F1 | Ma-F1 | MRR@5 | MAP | NDCG@5 |
| PromptCase | 17.1 | 20.3 | 18.5 | 20.5 | 35.1 | 33.9 | 38.7 | 16.0 | 29.7 | 20.8 | 21.5 | 32.7 | 32.0 | 36.2 |
| w/o $v_g$ | 1.6±0.1 | 2.9±0.1 | 2.1±0.1 | 2.2±0.1 | 4±0.1 | 4.0±0.1 | 4.8±0.2 | 1.6±0.1 | 2.9±0.1 | 2.1±0.1 | 2.2±0.1 | 4.0±0.1 | 2.9±0.1 | 4.8±0.1 |
| Avg Readout | 30.5±0.5 | 36.2±0.6 | 33.1±0.5 | 36.8±0.5 | 61.3±0.4 | 59.0±0.1 | 64.6±0.5 | 17.6±0.4 | 32.6±0.7 | 22.8±0.5 | 23.6±0.4 | 37.7±1.0 | 36.6±0.7 | 41.7±0.5 |
| CaseGNN | 35.5±0.2 | 42.1±0.2 | 38.4±0.3 | 42.4±0.1 | 66.8±0.8 | 64.4±0.9 | 69.3±0.8 | 17.7±0.7 | 32.8±0.7 | 23.0±0.5 | 23.6±0.5 | 38.9±1.1 | 37.7±0.8 | 42.8±0.7 |

only using legal corpus to simply pre-trained on BERT-based model is not enough for the difficult LCR task. The performances of MonoT5 model on two datasets are the poorest in the experiment, which may for the reason that MonoT5 is pre-trained for text-to-text tasks instead of information retrieval tasks. Although SAILER model uses the structure-aware architecture, the performances are not that good as CaseGNN, which shows that the combining of both TACG and GNN model can largely improve the learning and understanding ability of model. Since CaseGNN model uses the fact and issue format to construct the TACG to encode a graph representation, the worse performances of PromptCase indicates the importance of transforming fact and issue into TACG and applying to GNN to learn an expressive case graph representation for LCR.

For two-stage retrieval setting, all methods use a BM25 top10 results as the first stage retrieval and conduct the re-ranking based on these ten retrieved cases. SAILER and PromptCase are compared since these two methods have a comparable one-stage retrieval performance. CaseGNN outperforms both methods in most metrics, especially those related to ranking results, such as MRR@5, MRR and NDCG@5. In COLIEE2022, although CaseGNN has a much higher performance than other baselines, CaseGNN actually cannot benefit from a two-stage retrieval since BM25 cannot provide a higher and useful first stage ranking result compared with CaseGNN itself. In COLIEE2023, all methods can benefit from the two-stage retrieval and CaseGNN can further improve the performance over the one-stage setting.

### 5.3   Ablation Study

The ablation study is conducted to verify the effectiveness of the graph components of CaseGNN: (1) not using any graph, which is equivalent to PromptCase; (2) using TACG without the virtual global node (w/o $v_g$); and (3) using the average of updated node features as case graph representation (Avg Readout). The experiments are conducted on both datasets and measured under all metrics. Results are reported in Table 3. Only one-stage experiments are considered.

As shown in Table 3, the CaseGNN framework with all the proposed components can significantly outperform the other variants for both datasets. PromptCase utilises the text encodings of the legal fact and the legal issue to obtain the case representation, which serves as a strong baseline for LCR. The virtual global node in TACG comes from the encoding of the legal fact and the legal issue. For the w/o $v_g$ variant, the performance is the worst that the model almost cannot learn any useful information because without the proper text encodings,

Table 4: Effectiveness of GNNs. (%)

| Methods | COLIEE2022 | | | | | | | COLIEE2023 | | | | | | |
|---------|------|------|-------|-------|-------|------|--------|------|------|-------|-------|-------|------|--------|
| | P@5 | R@5 | Mi-F1 | Ma-F1 | MRR@5 | MAP | NDCG@5 | P@5 | R@5 | Mi-F1 | Ma-F1 | MRR@5 | MAP | NDCG@5 |
| GCN | 21.3±0.3 | 25.3±0.4 | 23.2±0.2 | 26.0±0.4 | 46.0±0.2 | 44.4±0.1 | 49.7±0.3 | 12.8±0.2 | 23.7±0.3 | 16.5±0.1 | 16.9±0.2 | 27.8±0.8 | 26.8±0.6 | 31.6±0.7 |
| GAT | 29.3±0.1 | 34.8±0.3 | 31.8±0.1 | 35.3±0.2 | 59.2±0.5 | 56.9±0.3 | 62.3±0.7 | 17.4±0.3 | 32.2±0.5 | 22.5±0.3 | 23.1±0.5 | 37.5±0.4 | 36.4±0.4 | 41.4±0.4 |
| EdgeGAT | 35.5±0.2 | 42.1±0.2 | 38.4±0.3 | 42.4±0.1 | 66.8±0.8 | 64.4±0.9 | 69.3±0.8 | 17.7±0.7 | 32.8±0.7 | 23.0±0.5 | 23.6±0.5 | 38.9±1.1 | 37.7±0.8 | 42.8±0.7 |

the overall semantics are not effectively encoded. For the AvG Readout variant, the readout function of CaseGNN is set to using the average node embeddings as the case graph representation. This variant is outperformed by CaseGNN because in this variant, the readout function ignores the information in the edge features. Nevertheless, Avg Readout has a better result compared with PromptCase, which verifies that the graph structure in the case can provide useful information.

### 5.4    Effectiveness of GNNs

To validate the effectiveness of the EdgeGAT layer, CaseGNN is compared with variants of substituting EdgeGAT with GCN [19] and GAT [47]. The experiments are conducted on both datasets and evaluated with all metrics. The results are shown in Table 4. Only one-stage experiments are considered. For a fair comparison, all variants will be trained with the proposed TACG.

As shown in the experimental results, EdgeGAT has the highest performance compared with the widely used GNN models GCN and GAT. The outstanding results of EdgeGAT on LCR tasks is because in the proposed CaseGNN method, there is a TACG module including both node and edge features in the case graph. Correspondingly, EdgeGAT has a novel design to incorporate the edge features into the case representation calculation. These edge features are important in terms of the legal information contained in the relations between different entities extracted from the legal case. For both GCN and GAT, since these general GNN layers do not have the capability to utilise the edge information, only the node information encoded from the entities are used in the calculation, which leads to information loss of the legal case. More specifically, GAT has a better performance compared with GCN. This phenomenon is aligned with the performance gap between GAT and GCN on other general graph learning tasks because of the graph learning ability difference between GAT and GCN.

### 5.5    Parameter Sensitivity

In this experiment, the temperature coefficient $\tau$ and the number of hard negative samples in the contrastive loss in Equation (11) are investigated for their parameter sensitivity. The results are presented in Fig. 3 and Fig. 4 respectively.

*Temperature coefficient.* As shown in Fig. 3, temperature is chosen from {0.01, 0.1, 0.5, 1, 3}. For both datasets, $\tau$ set to 0.1 achieves the best performance. When the temperature is too large, the similarity score will be flatten and it
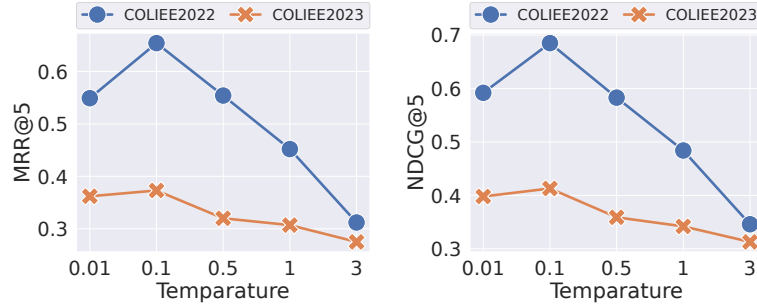
Fig. 3: Parameter sensitivity for the temperature $\tau$ in the contrastive loss.

cannot provide sufficient training signal to the model via the contrastive loss. In the contrast, when $\tau$ is too small, it will extremely sharpen the similarity distribution, which will make the objective function neglect the less significant prediction in the output.

*Number of hard negative samples.* According to Fig. 4, the choice of number of hard negative samples are from $\{0, 1, 5, 10\}$. The hard negative samples are sampled from highly rank irrelevant cases by BM25. Different numbers of hard negative samples in Equation (11) have different impacts to the final model. When there is no hard negative samples in the training objective, the model will
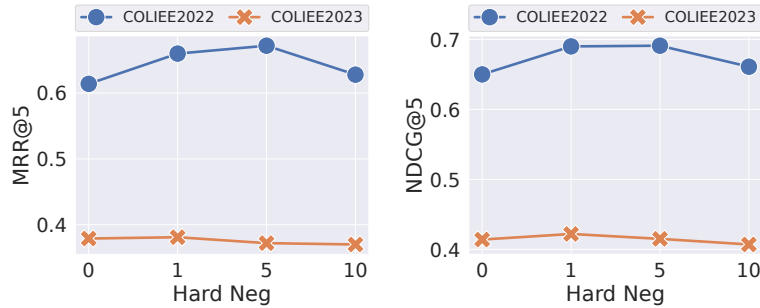


Fig. 4: Parameter sensitivity for the number of hard negative samples in the contrastive loss.

be only trained with easy negative samples by random sampling. Model trained without hard negative samples will have an inferior performance compared with a proper selected number of hard negative samples. This is because hard negative samples can provide a more strict supervision signal to train the CaseGNN. The performance decreases when there are too many hard negative samples, which is

because the training task becomes extremely difficult and the model can barely obtain useful information from the training signal.
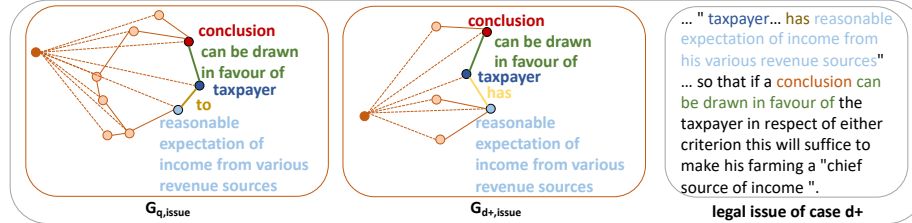


Fig. 5: TACGs of cases successfully retrieved by CaseGNN but not by PromptCase.

### 5.6   Case Study

To further demonstrate how the graph structure helps with the LCR task in CaseGNN, the constructed TACG is visualised in Fig. 5, where CaseGNN successfully performs retrieval while LM-based PromptCase fails. In this visualisation, the constructed TACGs of the legal issue in both the query case $q$ and the ground truth candidate case $d+$, $G_{q,\text{issue}}$ and $G_{d+,\text{issue}}$ are presented. From this visualisation, it it clear that the graph structure can bring multiple entities together to create a candidate graph that is similar to the query graph. On the contrary, from the case text of the candidate $d+$ on the right of Fig. 5, the corresponding language of these entities and relationships are far from each other, which leads to the unsuccessful retrieval of PromptCase for only using sequential LM without case structural information.

## 6   Conclusion

This paper identifies two challenges remaining in recent LM-based LCR models about the legal structural information neglect and the lengthy legal text limitation. To overcome the problems from these challenges, this paper proposes a novel framework, CaseGNN, with specific design to utilise the graph neural networks. A Text-Attributed Case Graph (TACG) is developed to transform the unstructured case text into structural graph data. To learn an effective case representation, an Edge Graph Attention Layer (EdgeGAT) is developed. Extensive experiments conducted on two benchmark datasets verify the state-of-the-art performance and the effectiveness of CaseGNN.

# References

1. Abolghasemi, A., Verberne, S., Azzopardi, L.: Improving bert-based query-by-document retrieval with multi-task optimization. In: ECIR (2022)
2. Althammer, S., Askari, A., Verberne, S., Hanbury, A.: Dossier@coliee 2021: Leveraging dense retrieval and summarization-based re-ranking for case law retrieval. CoRR **abs/2108.03937** (2021)
3. Angeli, G., Premkumar, M.J.J., Manning, C.D.: Leveraging linguistic structure for open domain information extraction. In: ACL (2015)
4. Askari, A., Abolghasemi, A., Pasi, G., Kraaij, W., Verberne, S.: Injecting the BM25 score as text improves bert-based re-rankers. In: ECIR (2023)
5. Askari, A., Peikos, G., Pasi, G., Verberne, S.: Leibi@coliee 2022: Aggregating tuned lexical models with a cluster-driven bert-based model for case law retrieval. CoRR **abs/2205.13351** (2022)
6. Askari, A., Verberne, S.: Combining lexical and neural retrieval with longformer-based summarization for effective case law retrieval. In: DESIRES. CEUR (2021)
7. Askari, A., Verberne, S., Abolghasemi, A., Kraaij, W., Pasi, G.: Retrieval for extremely long queries and documents with RPRS: a highly efficient and effective transformer-based re-ranker. CoRR **abs/2303.01200** (2023)
8. Chalkidis, I., Fergadiotis, M., Malakasiotis, P., Aletras, N., Androutsopoulos, I.: LEGAL-BERT: the muppets straight out of law school. CoRR **abs/2010.02559** (2020)
9. Chalkidis, I., Kampas, D.: Deep learning in law: early adaptation and legal word embeddings trained on large corpora. Artif. Intell. Law **27**(2), 171–198 (2019)
10. Dai, Z., Callan, J.: Context-aware sentence/passage term importance estimation for first stage retrieval. CoRR **abs/1910.10687** (2019)
11. Devlin, J., Chang, M., Lee, K., Toutanova, K.: BERT: pre-training of deep bidirectional transformers for language understanding. In: NAACL-HLT (2019)
12. Goebel, R., Kano, Y., Kim, M.Y., Loro, M.N., Minh, N.L., Rabelo, J., Rossi, J., Satoh, K., Savelka, J., Shao, Y., Shimazu, A., Tojo, S., Tran, V., Valvoda, J., Westermann, H., Yamada, H., Yoshioka, M., Wehnert, S.: Competition on legal information extraction/entailment (COLIEE) (2022)
13. Goebel, R., Kano, Y., Kim, M.Y., Loro, M.N., Minh, N.L., Rabelo, J., Rossi, J., Satoh, K., Savelka, J., Shao, Y., Shimazu, A., Tojo, S., Tran, V., Valvoda, J., Westermann, H., Yamada, H., Yoshioka, M., Wehnert, S.: Competition on legal information extraction/entailment (COLIEE) (2023)
14. Hamilton, W.L., Ying, Z., Leskovec, J.: Inductive representation learning on large graphs. In: NeurIPS (2017)
15. He, X., Bresson, X., Laurent, T., Hooi, B.: Explanations as features: Llm-based features for text-attributed graphs. CoRR **abs/2305.19523** (2023)
16. Jones, K.S.: A statistical interpretation of term specificity and its application in retrieval. J. Documentation **60**(5), 493–502 (2004)
17. Khattab, O., Zaharia, M.: Colbert: Efficient and effective passage search via contextualized late interaction over BERT. In: SIGIR (2020)
18. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: ICLR (2015)
19. Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks. In: ICLR (2017)
20. Li, H., Ai, Q., Chen, J., Dong, Q., Wu, Y., Liu, Y., Chen, C., Tian, Q.: SAILER: structure-aware pre-trained language model for legal case retrieval. CoRR **abs/2304.11370** (2023)

21. Li, Y., Hooi, B.: Prompt-based zero- and few-shot node classification: A multimodal approach. CoRR **abs/2307.11572** (2023)
22. Liu, B., Hu, Y., Wu, Y., Liu, Y., Zhang, F., Li, C., Zhang, M., Ma, S., Shen, W.: Investigating conversational agent action in legal case retrieval. In: ECIR (2023)
23. Liu, B., Wu, Y., Zhang, F., Liu, Y., Wang, Z., Li, C., Zhang, M., Ma, S.: Query generation and buffer mechanism: Towards a better conversational agent for legal case retrieval. Inf. Process. Manag. (2022)
24. Liu, Y., Qiu, R., Huang, Z.: Cat: Balanced continual graph learning with graph condensation. CoRR **abs/2309.09455** (2023)
25. Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., Stoyanov, V.: Roberta: A robustly optimized BERT pretraining approach. CoRR **abs/1907.11692** (2019)
26. Ma, Y., Ai, Q., Wu, Y., Shao, Y., Liu, Y., Zhang, M., Ma, S.: Incorporating retrieval information into the truncation of ranking lists for better legal search. In: SIGIR (2022)
27. Ma, Y., Shao, Y., Wu, Y., Liu, Y., Zhang, R., Zhang, M., Ma, S.: Lecard: A legal case retrieval dataset for chinese law system. In: SIGIR (2021)
28. Ma, Y., Wu, Y., Ai, Q., Liu, Y., Shao, Y., Zhang, M., Ma, S.: Incorporating structural information into legal case retrieval. ACM Trans. Inf. Syst. (2023)
29. Monninger, T., Schmidt, J., Rupprecht, J., Raba, D., Jordan, J., Frank, D., Staab, S., Dietmayer, K.: SCENE: reasoning about traffic scenes using heterogeneous graph neural networks. IEEE Robotics Autom. Lett. (2023)
30. Nogueira, R., Jiang, Z., Pradeep, R., Lin, J.: Document ranking with a pretrained sequence-to-sequence model. In: EMNLP (2020)
31. Nogueira, R.F., Yang, W., Lin, J., Cho, K.: Document expansion by query prediction. CoRR **abs/1904.08375** (2019)
32. Ponte, J.M., Croft, W.B.: A language modeling approach to information retrieval. SIGIR (2017)
33. Qiao, Y., Xiong, C., Liu, Z., Liu, Z.: Understanding the behaviors of BERT in ranking. CoRR **abs/1904.07531** (2019)
34. Qiu, R., Huang, Z., Li, J., Yin, H.: Exploiting cross-session information for session-based recommendation with graph neural networks. ACM Trans. Inf. Syst. (2020)
35. Qiu, R., Li, J., Huang, Z., Yin, H.: Rethinking the item order in session-based recommendation with graph neural networks. In: CIKM (2019)
36. Qiu, R., Yin, H., Huang, Z., Chen, T.: GAG: global attributed graph neural network for streaming session-based recommendation. In: SIGIR (2020)
37. Rabelo, J., Kim, M., Goebel, R.: Semantic-based classification of relevant case law. In: JURISIN (2022)
38. Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J.: Exploring the limits of transfer learning with a unified text-to-text transformer. J. Mach. Learn. Res. (2020)
39. Reimers, N., Gurevych, I.: Sentence-bert: Sentence embeddings using siamese bert-networks. In: EMNLP-IJCNLP (2019)
40. Robertson, S.E., Walker, S.: Some simple effective approximations to the 2-poisson model for probabilistic weighted retrieval. In: SIGIR (1994)
41. Shao, Y., Mao, J., Liu, Y., Ma, W., Satoh, K., Zhang, M., Ma, S.: BERT-PLI: modeling paragraph-level interactions for legal case retrieval. In: IJCAI (2020)
42. Srivastava, N., Hinton, G.E., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: a simple way to prevent neural networks from overfitting. J. Mach. Learn. Res. (2014)

43. Sun, Z., Xu, J., Zhang, X., Dong, Z., Wen, J.: Law article-enhanced legal case matching: a model-agnostic causal learning approach. CoRR **abs/2210.11012** (2022)
44. Tang, Y., Qiu, R., Li, X.: Prompt-based effective input reformulation for legal case retrieval. CoRR **abs/2309.02962** (2023)
45. Tran, V.D., Nguyen, M.L., Satoh, K.: Building legal case retrieval systems with lexical matching and summarization using A pre-trained phrase scoring model. In: ICAIL (2019)
46. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: NeurIPS. pp. 5998–6008 (2017)
47. Velickovic, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., Bengio, Y.: Graph attention networks. In: ICLR (2018)
48. Vuong, T., Nguyen, H., Nguyen, T., Nguyen, H., Nguyen, T., Nguyen, H.: NOWJ at COLIEE 2023 - multi-task and ensemble approaches in legal information processing. CoRR **abs/2306.04903** (2023)
49. Wang, Z.: Legal element-oriented modeling with multi-view contrastive learning for legal case retrieval. In: IJCNN (2022)
50. Wen, Z., Fang, Y.: Augmenting low-resource text classification with graph-grounded pre-training and prompting. In: SIGIR (2023)
51. Xiao, C., Hu, X., Liu, Z., Tu, C., Sun, M.: Lawformer: A pre-trained language model for chinese legal long documents. AI Open **2**, 79–84 (2021)
52. Yang, J., Ma, W., Zhang, M., Zhou, X., Liu, Y., Ma, S.: Legalgnn: Legal information enhanced graph neural network for recommendation. ACM Trans. Inf. Syst. (2022)
53. Yao, F., Xiao, C., Wang, X., Liu, Z., Hou, L., Tu, C., Li, J., Liu, Y., Shen, W., Sun, M.: LEVEN: A large-scale chinese legal event detection dataset. In: ACL (2022)
54. Yu, W., Sun, Z., Xu, J., Dong, Z., Chen, X., Xu, H., Wen, J.: Explainable legal case matching via inverse optimal transport-based rationale extraction. In: SIGIR (2022)
55. Zhang, H., Dou, Z., Zhu, Y., Wen, J.R.: Contrastive learning for legal judgment prediction. ACM Trans. Inf. Syst. **41**(4),  25 (2023)
56. Zhong, H., Wang, Y., Tu, C., Zhang, T., Liu, Z., Sun, M.: Iteratively questioning and answering for interpretable legal judgment prediction. In: AAAI (2020)