# PDRS: Precedent Decision Retrieval System for Legal Documents

MURAT SAHIN, Bilkent University, Turkey

MEHMET CAN SAKIROGLU, Bilkent University, Turkey

We propose Precedent Decision Retrieval System (PDRS) to address the difficulties lawyers face when searching for and retrieving court judgments' precedent decisions. The system aims to reduce the extensive human resources and effort required for lawyers to manually seek relevant precedent decisions while working on their cases. It provides an effective solution for legal experts by assisting in the information retrieval process. The system includes a custom dataset of almost 300,000 effective documents scraped from the Turkish National Judicial Network Information System (UYAP). Its purpose is to retrieve the most relevant document(s) from this dataset. The methodology involves combining classical and learning-based approaches, specifically BM25 and BERT. To evaluate the system's performance, we used an unconventional method. The dataset's documents are divided into pseudo query and document pairs to form an evaluation set, which is then assessed. During assessment, we saw that over 20% of the queries we extracted was able to return its own document pair in Top-10 results. Additionally, a user-friendly interface is provided for ease of use. This work ultimately enables more precise document retrieval and time-efficiency.

CCS Concepts: • **Information systems** → **Query representation**; *Top-k retrieval in databases*; *Relevance assessment*; *Retrieval effectiveness.*

Additional Key Words and Phrases: Computational Law, Precedent Judgment Retrieval, Information Retrieval System, Search Engine

## 1 INTRODUCTION

In the field of legal practice, lawyers often need to look at the precedent decisions when working on their cases. This task requires a significant amount of time and effort, as legal professionals search through extensive databases to locate relevant court judgments. Recognizing the challenges presented by this laborious and time-consuming task, our aim is to create a system that simplifies and accelerates the process of retrieving relevant precedent decisions, ultimately improving the effectiveness of legal research.

The innovative dataset and distinctive evaluation method employed in our Precedent Decision Retrieval System for legal documents are the contribution of our work. Instead of relying on a classical information retrieval collection, we utilize a custom database comprising over 400,000 Turkish court decisions sourced from the Turkish National Judicial Network Information System (UYAP) [1]. This novel dataset provides a realistic and diverse representation of the legal scenarios of Turkey, which enriches the data behind our work. Our evaluation approach addresses the absence of predefined queries in the dataset. Rather than relying on established queries, we extract relevant information by taking the sentence following the phrase 'In summary to the petition;...' as a query substitute. The subsequent removal of these queries from the original documents constitutes an effective methodology for evaluating the system's performance.

Alongside this methodology, the synergy between classical methods and learning based methods are utilized. By combining these approaches, our system aims to effectively match queries with precedent judgments, thereby offering legal professionals a better precedent decision retrieval experience.

The sections to be followed can be summarized as follows:

Firstly, the *Related Work* section is a brief review of legal information retrieval literature and methodologies, offering insights into challenges and solutions, and establishing context for our approach. Then, the *Method* section outlines our approach to implementing the Precedent Decision Retrieval System, covering query and document representation,

similarity measurement, and the retrieval process. We also outline our unique evaluation approach, addressing the absence of predefined queries. After that, the *Dataset* section introduces our dataset, which consists of over 400,000 court decisions from the Turkish National Judicial Network Information System (UYAP). Then, in the *Experimental Results* section, we present outcomes from our system evaluation. This section discusses the effectiveness of the Precedent Decision Retrieval System and demonstrates its potential to enhance legal research efficiency. Moreover, in the *Discussion* section, we present the reasoning behind differences in the behaviors of the different models, followed by a brief explanation of the future plans. Lastly, in the *Conclusion* section, we summarize our contributions, the system's impact, and its potential benefits for legal practitioners.

## 2  RELATED WORK

The University of Alberta's COLIEE competition [3] has been of great interest to researchers and has been referenced in several papers, particularly in the area of natural language processing applications in the legal field. The competition consists of many tasks alongside the legal case retrieval task. Teams in the 2023 edition had different approaches [2], most of them applied a form of classical IR approaches such as BM25 or transformer-based methods such as BERT. There were teams combining both, the best performing team was utilizing a pre-trained language model. In that domain, it seems that classical methods still stand a chance, therefore we are utilizing them for our task.

The studies on the retrieval of previous cases have been increasing globally in the last few years, but there are still not many studies in the field of Turkish law. There have been only two works which are written by the same research group. In their first work [4], authors introduced the first prior case retrieval task for Turkish courts, and they analyzed different IR methods using this task. They collected an unlabeled dataset consisting over 300,000 case decision texts from the Court of Cassation. Also, they created a retrieval dataset by selecting 26 query sentences and 10 relevant documents. They first ranked decision texts using BM25, and after obtaining the top 20 documents, they leveraged an RNN autoencoder structure to extract encoding vectors and then they reranked the documents by calculating the similarity with the query vector. They saw significant improvement with the combination of those two approaches compared to only using document vectors, they connect it to the fact that many sentences are repeated in this domain, and BM25 algorithms rely on exact matching.

In their present research [5], they redirect their attention to utilizing BERTurk, which is an extension of the BERT model trained on Turkish corpora. To fine-tune BERTurk, they use the same dataset of over 300,000 collected decisions. The resulting BERTurk-Legal is then utilized to vectorize documents and evaluate similarity. They measured their performance by the retrieval dataset from their previous paper. The switch from traditional methods such as BM25 and LSTM to the advanced BERTurk-Legal boosted their metrics' performance. This indicates that we could also achieve greater success in our experiments using transformers.

Given the abundance of legal sources, the collection of data is generally not difficult, but the labeling of these documents is the real challenge. Previous studies, in general, have attempted to address this problem as a priority, especially when no dataset exists. Our work also addressed this problem early on.

## 3 METHOD

Our approach to the Precedent Decision Retrieval Task (See **Figure 1**) consists of the following steps:

(1) Measure the similarity between the query and document vectors and obtain the scores for each document
   (a) using BM25
   (b) using BERT
(2) Combine BM25 and BERT scores
(3) Rank documents according to the combined scores
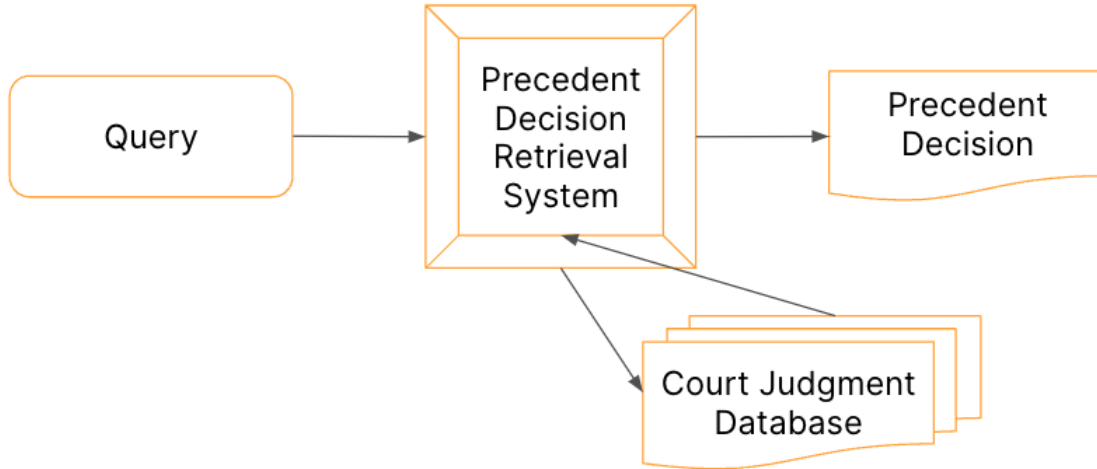


Fig. 1. A visualization of the Precedent Decision Retrieval Task

In the construction of information retrieval systems such as ours, it is a common approach to use different methods such as TF-IDF, BM25, word embeddings, BERT, and their weighted combinations for the purpose of ranking. Our Precedent Decision Retrieval System utilizes a hybrid approach that combines the best of both worlds. Our approach aims to retrieve documents that share both structural/statistical properties and semantic aspects with the given query.

To implement this idea, we calculate document scores for a given query using both BM25 and BERT for ranking purposes. The combination of these scores is used to create a final ranking of documents for retrieval. This approach enables the retrieval of semantically and structurally similar documents, making it suitable for precedent decision retrieval case. A visualization of this approach for the Precedent Decision Retrieval System can be observed in **Figure 2**.

As for the details, we combined the scores from BM25 and BERT by using a weighted average between their scores. We experimented with different weights and selected the best one, as explained thoroughly in the Experimental Results section. Moreover, the specific BERT model we use is a BERTurk-based model[1] as our data is in the language of Turkish.

We have implemented our methodology without the use of existing search engine tools, as this gave us more flexibility to implement our approach and extend our choices to achieve a more accurate retrieval system during inference.

After implementing the system, evaluating it is still non-trivial due to the lack of a classical information retrieval dataset with query-document pairs that encode relevant information. To address this, we use the sentence following the

---

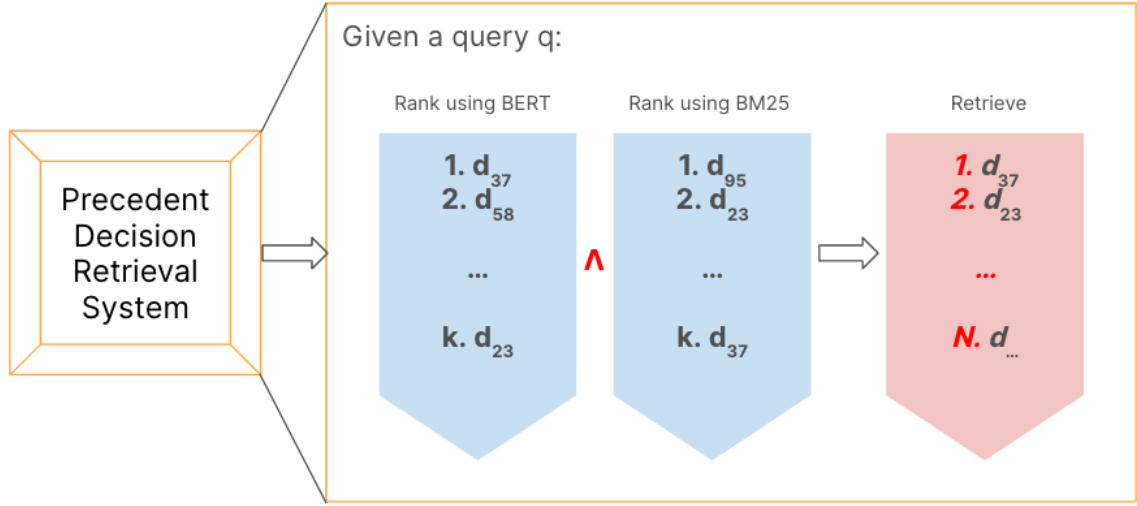[1]https://huggingface.co/emrecan/bert-base-turkish-cased-mean-nli-stsb-tr

Fig. 2. A visualization of the Precedent Decision Retrieval System

phrase 'In summary to the petition;...' as a substitute for the query and we then extract these queries from the original documents. At the end of this process, we have query-document pairs, similar to those found in a classical information retrieval collection. Therefore, we evaluate the Precedent Judgment Retrieval System using this constructed collection.

After completing these steps, a metric such as the rank of the corresponding document from which the query originated can be used. The goal is to minimize this metric. For instance, if the document d is the pair-document of query q, we would ideally like to see document d retrieved as the first document. As the rank of its retrieval increases, the metric indicates a worse evaluation of the system. For this purpose, we calculated the Top-N accuracies, where N equals 10, 100, and 1000, as discussed further in the Experimental Results section.

In addition, we offer a user-friendly graphical interface for our retrieval system. The interface is easy to use as it only requires the user to input a query, and then provides them with a ranked retrieval of the documents as precedent decisions. Examples of the use of this interface are shown in **Figure 3a** and **Figure 3b**.



(a) UI: Example Search Results

(b) UI: Example Retrieved Precedent Decision

Fig. 3. UI: Example Usage

## 4 DATASET

A dataset of over 400,000 local court decisions was collected from the Turkish National Judicial Network Information System (UYAP).

To extract queries from each document, we locate the text segment containing the phrase 'In summary to the petition;...' using HTML tags, and extract the text from there. This provides us with our query-document pairs. Note that not all documents fit this template, so our total number of decisions decreased to nearly 300,000.

We then performed additional preprocessing steps on both queries and documents, removing HTML tags and obtaining raw text. The data still contains punctuation, numbers, and capitalization, which is suitable for transformer-based approaches. However, for classical term frequency-based approaches, the data requires cleaning and tokenization. Thus, we implement the subsequent preprocessing steps to prepare the data for classical methods:

- Removing non-alphabetic characters
- Converting to lowercase
- Removing stopwords
- Tokenizing text

## 5 EXPERIMENTAL RESULTS

We have evaluated the performance of our Precedent Decision Retrieval System using the processed version of our data, which contains pairs of queries and their corresponding documents, as explained in our methodology.

To assess the effectiveness of our approach, we used the Top-N accuracy metrics with N values of *10*, *100*, and *1000*. We believe that these values provide a solid understanding of the system's performance.

As a brief definition, if the total number of queries we experiment with is *T*, and the number of document-pairs our retrieval system was able to correctly retrieve in the *top-N* documents after ranking is *Q*; then *Top-N Accuracy* can be explained using the **Equation 1**.

$$\text{Top-N Accuracy} = \frac{Q}{T} \tag{1}$$

Mainly; the *Top-10 Accuracy* of our retrieval system is **0.2136**, the *Top-100 Accuracy* of our retrieval system is **0.3168**, and the *Top-1000 Accuracy* of our retrieval system is **0.4654**. Moreover, we also evaluated our system for all the N values up to these thresholds separately to see the change of the performance of our system with respect to the value of N. The related figures of these experiments can be seen in **Figure 4a**, **Figure 4b**, and **Figure 4c** respectively for the Top-N accuracies on *N<10*, *N<100*, and *N<1000*.

Moreover, we experimented with our approach to weight selection for the combination of BM25 and BERT scores. We evaluated our systems performance using the Top-100 scores for all 11 weight distributions ranging from *0.0-to-1.0* from *1.0-to-0.0* for BERT and BM25, respectively. **Figure 5** shows a heatmap of the results obtained using different weight combinations. The diagonal contains the informative part of the matrix, while the other values are zero for visualization purposes. The highest performing weight combination is **0.5-to-0.5** for BERT and BM25, as shown in **Figure 5**.
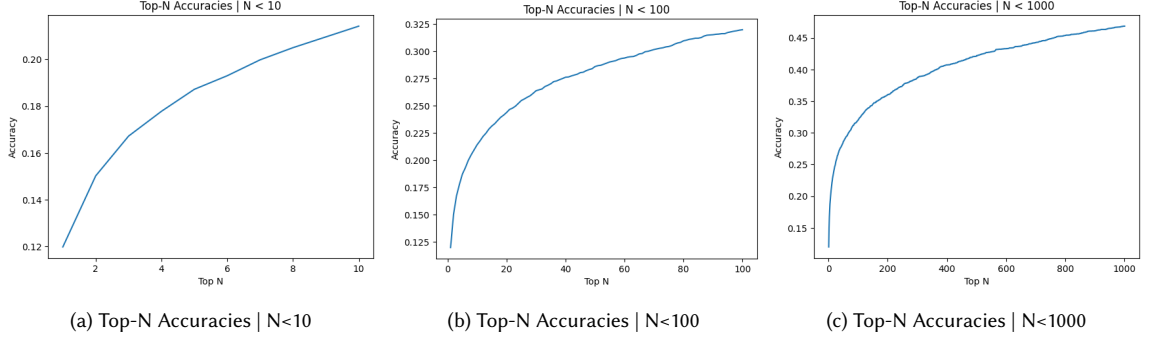
(a) Top-N Accuracies | N<10

(b) Top-N Accuracies | N<100

(c) Top-N Accuracies | N<1000

Fig. 4. Top-N Accuracies



Fig. 5. Top-100 Accuracies of Different Weight Combinations of BM25 and BERT

## 6 DISCUSSION

We propose a unique evaluation metric that we use in our work, as our data is not a conventional information retrieval collection. Our methodology for performance evaluation has not been previously implemented in this area, to our knowledge. To measure our effectiveness, we used Top-N accuracies obtained after running the system on our split

dataset. We observed a logarithmic trend in our effectiveness with respect to N values. These results are interpreted as follows:

If our retrieval system can retrieve successfully the document related to the query, it is usually able to do so within the first 10 or 100 documents. However, if it cannot retrieve the related document, even looking at the first 1000 documents or more will likely not yield results. This indicates our confidence in the retrieval results, as successful retrievals are ranked at the top. It is important to note that these rankings are based on the entire dataset. This means that we are searching for the Top-10 documents in the ranking, while the ranking itself includes hundreds of thousands of documents. Additionally, even in cases where the related document of a query is not retrieved in the Top-10 or Top-100 documents, we have observed that the highly ranked documents are still good suggestions for the given query. Therefore, our rankings are highly competent. Additional human evaluations could further demonstrate this competency in future work.

Moreover, as we briefly mentioned in the previous sections, we used BM25 and BERT in combination with equal weights to actualize our Precedent Decision Retrieval System. Our motivation behind this decision was to achieve a retrieval system that is not only based on structural similarity nor only based on semantic similarity, but actually based on both. Hence, we believe that this combination allows us to retrieve most effectively. Furthermore, it is important to note that our experimental results are also supportive of this methodology as we have observed that the most effective system we have is based on the *0.5-to-0.5* weighted usage of *BERT and BM25 in combination*. In our experiments with the various queries we gave to the system, we also observed that using an equal combination of BERT and BM25 was returning much more useful results compared to using only BM25 or using only BERT. Hence, using the advantages of both approaches gives the best results.

## 7 CONCLUSION

In conclusion, we introduce the Precedent Decision Retrieval System for legal documents. Our system blends traditional methods, such as BM25, with advanced transformer-based models, such as BERT. Our work is distinguished by the unique dataset sourced from UYAP and our evaluation approach. Our experimental results suggest that the combination of BM25 and BERT is highly effective, and the overall results are promising for this task. We anticipate that our system will be helpful for legal research by efficiently retrieving relevant precedent decisions. The purpose of this work is to connect legal practitioners with the vast amount of information contained in court judgments through the proposed system.

## REFERENCES

[1] [n. d.]. Emsal Karar UYAP. https://emsal.uyap.gov.tr/. Accessed: 2023-12-01.

[2] Randy Goebel, Yoshinobu Kano, Mi-Young Kim, Juliano Rabelo, Ken Satoh, and Masaharu Yoshioka. 2023. Summary of the Competition on Legal Information, Extraction/Entailment (COLIEE) 2023. *Proceedings of the Nineteenth International Conference on Artificial Intelligence and Law* (2023). https://api.semanticscholar.org/CorpusID:261583286

[3] Yoshinobu Kano, Masaharu Yoshioka, Yasuhiro Aoki, Randy Goebel, Ken Satoh, Diana Inkpen, Michel Custeau, Hai-Thanh Nguyen, Thi-Hai-Yen Vuong, Rohan Debbarma, et al. 2023. COLIEE-2023: Competition on Legal Information Extraction and Entailment. In *Proceedings of the 18th International Conference on Artificial Intelligence and Law*. ACM.

[4] Ceyhun E. Öztürk, Ömer Köksal, S. Baris Özçelik, and Aykut Koç. 2022. Prior Case Retrieval for the Court of Cassation of Turkey. *2022 16th International Conference on Signal-Image Technology & Internet-Based Systems (SITIS)* (2022), 539–544. https://api.semanticscholar.org/CorpusID:258072447

[5] Ceyhun E. Öztürk, Elektrik ve Elektronik, Mühendisli⌣gi Bölümü, and Hukuk Fakültesi. 2023. A Transformer-Based Prior Legal Case Retrieval Method. *2023 31st Signal Processing and Communications Applications Conference (SIU)* (2023), 1–4. https://api.semanticscholar.org/CorpusID:261313717