

---

# Encoded Summarization: Summarizing Documents into Continuous Vector Space for Legal Case Retrieval

Vu Tran<sup>†</sup>, Minh Le Nguyen<sup>†</sup>, Satoshi Tojo<sup>‡</sup>,  
and Ken Satoh<sup>‡</sup>

Published 2020-01-25 in AI and Law. DOI: <https://doi.org/10.1007/s10506-020-09262-4>

**Abstract** We present our method for tackling a legal case retrieval task by introducing our method of encoding documents by summarizing them into continuous vector space via our phrase scoring framework utilizing deep neural networks. On the other hand, we explore the benefits from combining lexical features and latent features generated with neural networks. Our experiments show that lexical features and latent features generated with neural networks complement each other to improve the retrieval system performance. Furthermore, our experimental results suggest the importance of case summarization in different aspects: using provided summaries and performing encoded summarization. Our approach achieved F1 of 65.6% and 57.6% on the experimental datasets of legal case retrieval tasks.

**Keywords** legal case · document retrieval · document summarization · deep learning · document representation

## 1 Introduction

Automatic legal document processing systems can speed up significantly the work of experts, which, otherwise, requires significant time and efforts. One crucial kind of such systems, automatic information retrieval whose systems, in place of experts, process over enormous amount of documents, for example, legal case reports, which are accumulated rapidly over time (the number of filings in the U.S. district courts for civil cases and criminal defendants is 344,787 in 2017 <sup>1</sup>). A case document contains a large volume of contents as the case may last days or even years. This one problem challenges the construction of an effective automatic legal case retrieval system.

In 2018, the Competition on Legal Information Extraction and Entailment (COLIEE) introduced two new tasks involving processing legal case documents: a legal case retrieval task and a legal case entailment task together with the previously introduced other two tasks about statute law Kano et al. (2018). The legal

---

<sup>†</sup>Japan Advanced Institute of Science and Technology; <sup>‡</sup>The National Institute of Informatics, Japan

<sup>1</sup> <http://www.uscourts.gov/statistics-reports/judicial-business-2017>

**Table 1** Statistics of candidate case documents in COLIEE 2018 and 2019 training data. (\*) Only count documents having an expert summary.

Property	2018		2019	
	Max	Avg.	Max	Avg.
#words/doc	85,551	5,690	9,666	2,665
#paragraphs/doc	1,117	43	119	22
#summary-words/doc*	8,827	589	3,085	242

case retrieval task is “to explore and evaluate case law retrieval technologies that are both effective and reliable”. The legal case entailment task is to “identify which paragraph in the noticed case entails the decision” given the noticed cases (assumed to be correctly retrieved) and a decision. A legal case document processing system fulfilling the two tasks would benefit lawyers in finding relevant information to construct arguments for their own objectives. This work tackles the first task: legal case retrieval.

In the format of COLIEE 2018 and COLIEE 2019, the legal case retrieval task involves reading a new case  $q$ , and extracting supporting cases  $c_1^*$ ,  $c_2^*$ , ...,  $c_n^*$  for the decision of  $q$  from a given list of candidate cases. The candidate cases that support for the decision of a new case are called ‘noticed cases’.

We tackle the task of finding the cases having supporting relationship with a new case indirectly through similarity measure. Our system extracts various kinds of features indicating the similarity of a new case and a previous case. The system is trained to score the relevance of the two cases by weighting the extracted similarity features. The similarity features are computed by comparing the different kinds of representations of the two cases including textual and vector representations. While the system does not provide definite answer to the supporting relationship, it learns from data to predict relevant cases by learning how the similarity features and the supporting relationship are related.

A legal case document contains case details and may contain other information such as citing cases, noticed cases, notices statutes, or editor drafted summaries. The case details are presented in form of paragraphs which can be fact statements, discussed legal points or the case decision (Fig. 1). The summary, if present, contains court decision, decisive facts, decisive legal points, and several key phrases, which is drafted by an editor.

Legal case documents usually contain huge amount of contents. As in Table 1, in COLIEE 2018, a legal case document contains  $\approx 5.7$ K words and 43 paragraphs in average, and could goes over 80K words and 1K paragraphs. This challenges the efficiency of not only human experts but also automatic retrieval systems. Editor summarization condensates contents by  $\approx 90\%$  which results in  $\approx 10\%$  key contents.

In COLIEE 2019, we observed the similar and different challenges. First, the candidate cases are  $\approx 2.7$ K-token long in average (Table 1). The difficulty of reading too long texts still emerges. We may pursue the idea that using summary as the main source of information. However, the dataset of COLIEE 2019 is different from the one of COLIEE 2018. While in COLIEE 2018, most of the candidate cases have a summary, in COLIEE 2019, more than  $\approx 47$ K in a total of 57K candidate cases are confirmed to have no summary (indicated with the note “This case is unedited, therefore contains no summary”). This means that summarization over candidate

**Summary:**

A human rights complaint alleged the federal government’s underfunding of welfare services for on-reserve First Nations children resulted in a lower level of services for those children than for other Canadian children whose welfare services were provincially funded. /\* ... \*/  
 The Federal Court held that, while the Tribunal had the power to decide this issue in advance of a full hearing on the merits, the process followed was not fair. /\* ... \*/

Administrative Law - Topic 547

The hearing and decision - Decisions of the tribunal - Reasons for decision - When required - [See second Civil Rights - Topic 7046].

Administrative Law - Topic 2608

Natural justice - Evidence and proof - Extraneous or irrelevant considerations - [See first Civil Rights - Topic 7046].

/\* ... \*/

**Paragraphs:**

[1] Mactavish, J. : The Government of Canada funds child welfare services for First Nations children living on reserves. The provinces fund child welfare services for all other Aboriginal and non-Aboriginal children.

[2] The First Nations Child and Family Caring Society and the Assembly of First Nations filed a human rights complaint with the Canadian Human Rights Commission in which they allege that the Government of Canada under-funds child welfare services for on-reserve First Nations children. /\* ... \*/

/\* ... \*/

[254] In my view, the ordinary meaning of the phrase “differentiate adversely in relation to any individual” on a prohibited ground of discrimination is to treat someone differently than you might otherwise have done because of the individual’s membership in a protected group. /\* ... \*/

/\* ... \*/

[395] As a result, the three applications for judicial review are granted.

/\* ... \*/

[396] THIS COURT ORDERS AND ADJUDGES that /\* ... \*/

**Fig. 1** Illustration of a legal case document from Federal Court of Canada. “/\* ... \*/”: omitted contents. Other information about citing cases, noticed cases, notices statutes, etc. are omitted.

case requires additional effort so that we can compare a query’s summary with a candidate’s summary.

We develop our system with representing a legal case document from its highlight contents. One way is to look at the editor drafted summary or catchphrases if these are available. The summary concisely states the decision of the case with the main arguments supporting the decision. While “catchphrases have an indicative function rather than informative, they present all the legal point considered instead that just summarizing the key points of a decision” (Galgani et al., 2012b). Catchphrases give a quick impression on what the case is about: “the function of catchwords is to give a summary classification of the matters dealt with in a case. [...] Their purpose is to tell the researcher whether there is likely to be anything in the case relevant to the research topic” (Olsson and of Judicial Administration, 1999). On one hand, catchphrases help lawyers/researchers quickly grasp the points of a case, without having to read the entire document, which saves significant time and effort for finding/studying relevant cases. On the other hand, catchphrases help improves the performance of automatic case retrieval systems.

**Table 2** Example of catchphrases found in legal case reports.

MIGRATION - partner visa - appellant sought to prove domestic violence by the provision of statutory declarations made under State legislation - “statutory declaration” defined by the Migration Regulations 1994 (Cth) to mean a declaration “under” the Statutory Declarations Act 1959 (Cth) in Div 1.5 - contrary intention in reg 1.21 as to the inclusion of State declarations under s 27 of the Acts Interpretation Act - statutory declaration made under State legislation is not a statutory declaration “under” the Commonwealth Act - appeal dismissed

Despite of the benefits, catchphrases are not always available in legal case documents, and are drafted by legal experts, which requires huge efforts when considering the enormous number of legal case documents. It is, therefore, beneficial to build automatic catchphrase generation systems for both old documents not having drafted catchphrases and new documents. Developing such systems, however, is challenging as the complexity of catchphrases shown in Table 2.

Approaches for generating catchphrases are based on phrase scoring derived from common model for retrieval: lexical matching with term frequency-inverse document frequency (Galgani et al., 2012a,b; Mandal et al., 2017b). The approaches are bounded by the limit of lexical matching, and corpus-wide statistical information. The limit of lexical matching can be lifted by moving to distributed vector space, for instance, distributed word embeddings in which common models are *Word2Vec* (Mikolov et al., 2013) and *GloVe* (Pennington et al., 2014). Corpus-wide statistical information has limit capability to identify catchphrases which are not really specific to some document but commonly used in several others.

In the COLIEE datasets, the legal documents may or may not have a drafted summary. Even using the drafted summary only may still result in limited performance as we observed in the datasets that the summary of a query case may not be similar to some of its noticed cases. We would like to build a system that is able to extract more informative features or key contents from a legal case document than just the summary.

We present our work on developing a legal case summarization system and on top of its core component - phrase scoring framework, building a legal case retrieval system.

First, we build a learning model to extract catchphrases for new documents with the knowledge from previously seen documents and the expert drafted catchphrases thereof. Our system utilizes deep neural networks which have been widely used in natural language processing (Liu and Zhang, 2018) to learn the direct relationship between gold catchphrases and document phrases. This results in our phrase scoring framework which is used to identify important phrases from a given legal case document.

On top of the phrase scoring framework, we develop our legal case document representation method which summarizes the document into continuous vector space. The representation is used as latent features for constructing case relevance ranking model, the core component of the retrieval system.

We also explore the benefits of employing various types of similarity measurement belonging to lexical similarity (keyword matching) and semantic similarity (meaning matching).

On one hand, the lexical similarity and semantic similarity differ from each other and can potentially complement each other as well. The lexical similarity is

obtained with approaches where the texts are compared by the direct surface forms with probably some transformations such as stemming, lemmatization, stopword removal, etc. High lexical similarity can present high matching, but low lexical similarity does not say much.

On the other hand, semantic similarity can provide the measurement where the surface forms are mismatched, for example, by paraphrasing. Semantic similarity can be learned in unsupervised fashion where common approaches are using statistical methods and benefits from huge available corpora (e.g. Wikipedia, GoogleNews, etc.) (Le and Mikolov, 2014; Levy and Goldberg, 2014; Mikolov et al., 2013; Pennington et al., 2014). Those methods treat a document as bag/sequence of words equally. Other information in the documents such as important words or phrases, or the document hierarchy when considered may provide significant information.

## 2 Encoded Summarization: Composing Document Vector from Phrase Scoring via Summary

In this section, we describe the method to compose document representations from phrase scoring via summary. When dealing with the legal case retrieval task, we observed several obstacles. First, the candidate cases are 5.7K-token long in average. This poses the problem of understanding the reason of selecting the cases as supporting cases. We, then, chose another approach which is comparing the summaries of each query and its candidate cases. We, however, found that the summary of the query is not necessarily lexically similar to the summary of the candidate cases. Moreover, some candidate cases do not have summary at all. We would like to obtain the summary for each and every candidate cases, and furthermore, the summary should be comparable with the summary of the corresponding query. One approach is to map the summaries into vector space with word embeddings (*word2vec* or *GloVe*) or document-embeddings (*doc2vec*). We come with another approach of document-embeddings which is to weight the document contents and perform weighted composition. For weighting the document contents, we build our phrase scoring framework to learn a scoring model based on the document summary.

### 2.1 The phrase scoring model

In this phase, we present our scoring model and how to train it using documents and their corresponding drafted summary.

#### 2.1.1 Constructing our scoring model architecture

We score each phrase in a document based on its contexts: its words, enclosing sentence, and document. Our approach takes advantage of the core property of word embedding techniques by Google *word2vec*, *GloVe*, etc.: contextual similarity, the similarity of two words is measured as the amount of common contexts where they appear. The phrase scoring model architecture is illustrated in Fig. 2.

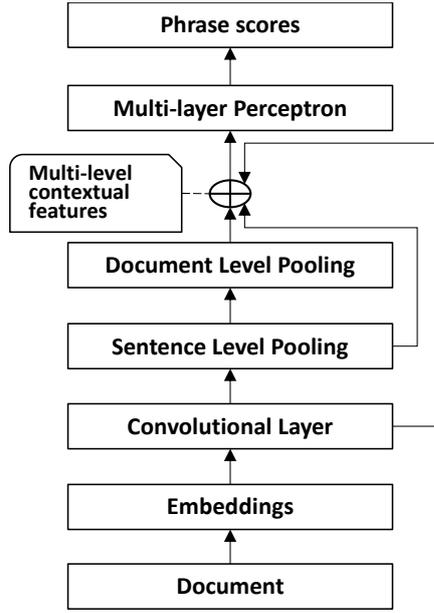


Fig. 2 Scoring model pipeline.

We adapt convolutional neural networks (CNNs), which are successfully used in text modeling (Kim, 2014; Severyn and Moschitti, 2015; Johnson and Zhang, 2015; Kalchbrenner et al., 2014), to encode each local context into latent feature space. Specifically, document phrase (summary phrase) features are captured by applying convolutional operations with window size  $2k + 1$  covering the word,  $k$  left and  $k$  right neighbors.

Given a document, we denote  $w_j^{s_i}$  as word  $j^{th}$  of sentence  $i^{th}$ . The features of an  $n$ -gram phrase  $p_j = \{w_j, w_{j+1}, \dots, w_{j+l-1}\}$  of a sentence are captured using convolutional neural layer as follows:

$$\mathbf{f}_{p_j} = ReLU \left( \mathbf{W}^c \begin{bmatrix} \mathbf{v}(w_j) \\ \mathbf{v}(w_{j+1}) \\ \dots \\ \mathbf{v}(w_{j+l-1}) \end{bmatrix} \right) \quad (1)$$

where,  $\mathbf{v}(\cdot) : \mapsto \mathbb{R}^d$ : word embedding vector lookup map,  $l$ : corresponding to the window size containing  $l$  contiguous words,  $[\cdot] \in \mathbb{R}^{dl}$ : concatenated embedding vector,  $\mathbf{W}^c \in \mathbb{R}^{c \times dl}$ : convolution kernel matrix with  $c$  filters,  $\mathbf{f}_{p_j} \in \mathbb{R}^c$ : phrase feature vector,  $ReLU$ : rectified linear unit activation.

Sentence (catchphrase) features are, then, captured by applying max pooling over the whole sentence (catchphrase).

$$\mathbf{f}_{s_i} = \text{max-pooling}_j(\mathbf{f}_{p_j^{s_i}}) \quad (2)$$

$$\mathbf{f}_{c_i} = \text{max-pooling}_j(\mathbf{f}_{p_j^{c_i}}) \quad (3)$$

where max pooling are operated over each dimension of vectors  $\mathbf{f}_{\mathbf{p}_{i,j}^s}$  ( $\mathbf{f}_{\mathbf{p}_{i,j}^c}$ ).

Document features are captured by applying max pooling over the document (not including summary). With the same max pooling operation as above, we compute document features as:

$$\mathbf{f}_d = \text{max-pooling}_i(\mathbf{f}_{s_i}) \quad (4)$$

The document features depend on only the document sentence, thereby, independent from the gold summary which are obviously not available for new documents.

Finally, we apply a multilayer perceptron (MLP) with one hidden and one output layer

$$MLP(\mathbf{x}) = \text{sigmoid}(\mathbf{W}_2 \cdot \tanh(\mathbf{W}_1 \cdot \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2) \quad (5)$$

to compute the score of each phrase  $p_j^{s_i}$  ( $p_j^{c_i}$ ) as

$$P(p_s, s, d) = MLP \left( \begin{bmatrix} f_{p_j^{s_i}} \\ f_{s_i} \\ f_d \end{bmatrix} \right) \quad (6)$$

$$P(p_c, c, d) = MLP \left( \begin{bmatrix} f_{p_j^{c_i}} \\ f_{c_i} \\ f_d \end{bmatrix} \right) \quad (7)$$

where the hidden layer computes the phrase representative features respecting to its local use, its enclosing sentence, and its document. The word representative features are feed to the output layer to compute word score (ranging from 0.0 to 1.0).

### 2.1.2 Training our scoring model

Main objective: given a document, summary phrases are “expected” to have higher score than document phrases.

First, we denote mean  $E$  and standard deviation  $std$  of word scores  $P$  for each document  $d$  in the following equations, which we will use to describe our objective as set of constraints, then formulated into loss function to be optimized.

$$E_c = E[P(p_c, c, d)] \text{ where } p_c \in c, c \in d \quad (8)$$

$$std_c = std[P(p_c, c, d)] \text{ where } p_c \in c, c \in d \quad (9)$$

$$E_s = E[P(p_s, s, d)] \text{ where } p_s \in s, s \in d \quad (10)$$

$$std_s = std[P(p_s, s, d)] \text{ where } p_s \in s, s \in d \quad (11)$$

$$E_{c,d'} = E[P(p_c, c, d')] \text{ where } p_c \in c, c \notin d' \quad (12)$$

Where  $p, c, s, d$  stand for phrase, summary sentence, document sentence, and the whole document respectively.  $c \notin d'$  means  $c$  is not a summary of document  $d'$ .

The main objective is realized by comparing the mean scores of summary phrases and document phrases:

- (o1) The mean score of summary phrases is higher than the mean score of document phrases:  $E_c > E_s$ .
- (o2) The mean score of summary phrases is lower than document phrases when comparing a summary with a document that the summary does not belong to:  $E_{c,d'} < E_{s'}$ . This is the negative constraint as opposed to the constraint **o1**.

The above two constraints are straightforward as the positive and negative factors of the objective. However, the comparison of the mean values does not guarantee to obtain to good scoring model as the score boundaries are not considered yet.

- (o3) The maximum score of summary phrases is higher than the maximum score of document phrases. It is expected that there exist concise summary phrases which is typical and representative for the document but could not found in the document. Such summary phrases should get higher scores than document phrases. The estimation  $E + std$  is used for representing max instead of hard max, whereby the constraint is realized as  $(E_c + std_c) > (E_s + std_s)$ .
- (o4) The minimum score of summary phrases is higher than the mean score of document phrases. Once again, to emphasize the importance of summary phrases, all summary phrases should get higher score than the average score of document phrases. The estimation  $E - std$  is used for representing min instead of hard min, whereby the constraint is realized as  $(E_c - std_c) > E_s$ .

We also add the following additional constraint to keep the scores from collapsing, which acts as regularization.

- (o5) Scores should not have small variance:  $std_c \neq 0, std_s \neq 0$ .

The loss function, hence, is composed from the constraints (**o1-5**) as follows.

$$\begin{aligned}
\mathfrak{L} = \sum_d \max(0, m - (a_1(E_c - E_s) & \\
& + a_2(\frac{1}{|\{d'\}|} \sum_{d' \neq d} E_{s'} - E_{c,d'}) \\
& + b_1((E_c + std_c) - (E_s + std_s)) & (13) \\
& + b_2((E_c - std_c) - E_s) \\
& - b_3(std_c) - b_4(std_s) \\
& ))
\end{aligned}$$

Note that rather imposing hard constraints, we compose the loss function with soft constraints. This means that some constraints may not be strictly satisfied after the training process. However, the violations of such constraints still incur certain losses and benefit the learning process.

## 2.2 Document Vector Composition

We present our method of composing document vectors from the phrase scoring model.

Given a document, we obtain its phrase scores and internal representations at three levels: phrase level, sentence level and document level. Then, we compose the document vector as:

$$\mathbf{g}(d) = \frac{\sum_{i,j} P(p_j^{s_i}, s_i, d) \times [\mathbf{f}_d; \mathbf{f}_{s_i}; \mathbf{f}_{p_j^{s_i}}]}{\sum_{i,j} P(p_j^{s_i}, s_i, d)} \quad (14)$$

Given a document, the composition weights the document contents based on their scores obtained from the phrase scoring framework. Important contents should have high contribution or affection to the final document vector. The component representations are multi-level contextual features which are the internal representations of the phrase scoring model. These internal representations contain the features which are learned to be used as base for scoring the surface contents. By using the multi-level contexts, the final document vector embeds the weighted multi-level contextual information including phrase level and sentence level contexts.

This composition resembles summarization where we weight the document internal representations by its summary. Thus, we call this composition encoded summarization.

### 2.3 Generating Text Summary

In this phase, we generate a summary for given a document by selecting and joining document phrases scored by the phrase scoring model. The process is as follows.

- Rank document phrases by their phrasal scores.
- Select phrases with scores from high to low.
- Join overlapping phrases into a longer phrase.
- Stop when the summary length exceeds length-threshold  $t$ .

The result summary is a list of phrases. The shortest phrases contain  $l$  words ( $l$  is the window size of the convolutional neural layer). The longest phrases are the sentences themselves.

## 3 Document Encoding and Relevance Modeling

### 3.1 Lexical Features

We estimate the lexical features by performing lexical matching between a query and a candidate case in different types of n-grams, skip-grams, longest common subsequence to measure various degrees of lexical similarity.

- N-gram matching: measuring n-gram overlapping between a query and a candidate case. We employ unigram and bigram models.
- Skip-bigram matching: measuring the co-occurrence of all word pairs in their sentence order. This allows the same non-continuous word pairs could be found in both query and candidate.

- We also employ the unigram+skip-gram model which balances the unigram matching and skip-gram matching.
- Longest common subsequence: measuring the strictly ordered overlapping scattering over the texts. We employ two variants: standard version and distance-weighted version. The distance-weighted version favors subsequences with shorter distances among words.

For each matching formula, we compute the matching scores by 3 different factors:

- Recall: normalized by query, measuring the percentage of the query contents found in the candidate.
- Precision: normalized by candidate, measuring the percentage of the candidate contents found in the query.
- F-measure: harmony score of the previous two.

$$f\text{-measure} = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

To have more precise comparison between a query and a candidate, we apply the following 4 matching options:

- Summary vs. Summary: we compute the matching of the query’s summary with the candidate’s summary. This matching represents the comparison of the highlights between the query and the candidate.
- Paragraphs vs. Summary: we compute the matching of the query’s paragraphs with the candidate’s summary. This matching represents the ratio of the candidate summary mentioning relevant details.
- Summary vs. Paragraphs: we compute the matching of the query’s summary with the candidate’s paragraphs. This matching represents the ratio of the query’s highlights mentioned in the candidate’s details.
- Paragraphs vs. Paragraphs: we compute the matching of the query’s paragraphs with the candidate’s paragraphs. This matching represents the ratio of the query’s details also occurred in the candidate’s details.

For COLIEE 2019 dataset, since most of the candidate cases do not have a summary, we perform summary generation in two ways: using the lead sentence of each paragraph and the generated summary described in Section 2.3. This results in 6 matching options for COLIEE 2019 dataset.

The coding for lexical features is in the form of q-c described as follows.

- q is a subset of query components including its expert summary (s) and paragraphs (p).
- c is a subset of candidate components including its expert summary (s) and paragraphs (p). As the case of COLIEE 2019 dataset, we use the lead sentences (l) and the generated summary (e) instead of unavailable expert summary (s).
- Each component of q is compared with each component of c.

For example, the lexical method sp-sp (q=sp, c=sp) means we perform 4 matching options: Summary vs. Summary, Summary vs. Paragraphs, Paragraphs vs. Summary, Paragraphs vs. Paragraphs, and the lexical method s-p (q=s, c=p) means we only perform Summary vs. Paragraphs matching. We use this naming for presenting lexical features’ impact analysis in our experiments.

In total, we collect lexical features from 6 matching formulas and 3 matching factors and 4 matching options, which results in 72 lexical features for measuring lexical matching between a query and each of its candidates. For COLIEE 2019 dataset, since most of the candidate cases do not have a summary, we perform summary generation in two ways: using the lead sentence of each paragraph and the generated summary described in Section 2.3. with the two additional matching options, we obtain 108 lexical features for COLIEE 2019 dataset.

### 3.2 Latent Features in Continuous Vector Space

We utilize several approaches for encoding documents into continuous vector space as follows.

- *word-embeddings*: From word vectors, we apply three kinds of vector compositions for producing document vectors: max pooling, average pooling, hierarchical pooling. The word vectors can be obtained from word embedding models, for example, Google word2vec (Mikolov et al., 2013) or GloVe (Pennington et al., 2014). In this work, we use the pre-trained word embeddings published by Stanford University<sup>2</sup>. Average pooling and max pooling are used to extract a fixed size feature vector from a sequence of vectors (Kim, 2014; Severyn and Moschitti, 2015; Chen et al., 2017). Hierarchical pooling composes the document vector from the average pooling and the max pooling of the entire document, together with the average pooling of sentence level max pooling which adds sentence-boundary dependent features.
- *doc2vec*(Le and Mikolov, 2014): This is a method for mapping text blocks into vector space. The method considers texts as sequences of tokens regardless of presented structures.
- Encoded summarization: We apply our method described in Section 2. The phrase scoring model is trained only on COLIEE 2018 dataset where the over 50K candidate cases have a summary. The pre-trained model is applied directly to COLIEE 2019 dataset without re-training. To compare with this method, we also derive encoding methods based on the above *word-embeddings*, and *doc2vec* compositions, but apply them only on the expert summary part of each document. The derived methods are noted with “(summary)”.

### 3.3 Query-Candidate Relevance Vector

The relevance vector consists of the features indicating the relevance of a candidate given a query. We compose this vector from lexical features and latent features.

The lexical features are computed by lexical matching which by themselves present the relevance measurement.

For the latent features which are encoded information in continuous vector space, by comparing each dimension independently, we can estimate the compatibility of a query and a candidate over the dimension. Thus, we compute the relevance features from latent features as the element-wise product of query vector and candidate vector. First, we obtain query vector  $\mathbf{g}(q)$  and candidate vector  $\mathbf{g}(c)$

<sup>2</sup> Pre-trained with Wikipedia 2014 + Gigaword 5 (<https://nlp.stanford.edu/projects/glove/>)

for each of the document vector compositions described in Subsection 3.2. Then, we compute the relevance vector of query  $q$  and candidate  $c$  by the following element-wise product.

$$\mathbf{h}(q, c) = \mathbf{g}(q) \odot \mathbf{g}(c) \quad (15)$$

The combination of lexical features and latent features is presented in the query-candidate relevance vector as the concatenation of lexical matching features and the element-wise product of latent feature vectors of the query and the candidate.

$$\text{relevance-vector}(q, c) = [\text{lexical-features}(q, c); \mathbf{h}(q, c)] \quad (16)$$

## 4 Experiments

### 4.1 Summarization

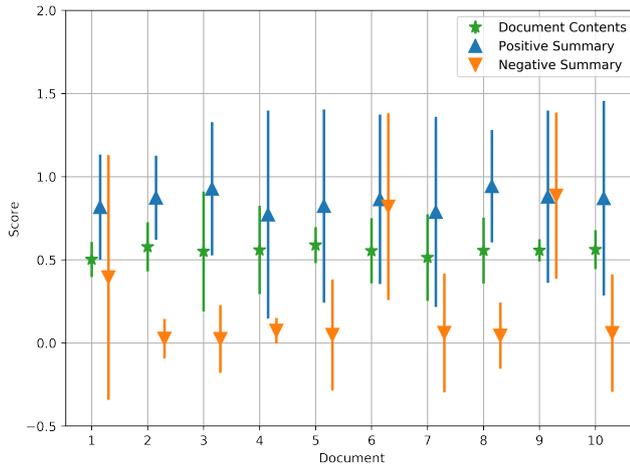
We trained the model with settings shown in Table 3. We use two sets of loss coefficients: (i) the parameters used for COLIEE 2018 submission, (ii) the parameters used for COLIEE 2019 submission. While the parameter set (i) is copied from (Tran et al., 2018), the parameter set (ii) is obtained by random searching around the set (i) for better retrieval performance on COLIEE 2018 dataset.

We report the empirical evaluation of the phrase scoring model applied to case summarization. A predicted summary of a given case is composed according to Section 2.3. We evaluate the predicted summary with length-threshold  $t$  values from 10% to 50% of document length. The evaluation is performed with ROUGE metrics including: ROUGE-1, ROUGE-2, ROUGE-SU. Results of the evaluation are shown in 4.

**Table 3** Phrase scoring model parameters. We use two sets of loss coefficients: (i) the parameters used for COLIEE 2018 submission, (ii) the parameters used for COLIEE 2019 submission. While the parameter set (i) is copied from (Tran et al., 2018), the parameter set (ii) is obtained by random searching around the set (i) for better retrieval performance on COLIEE 2018 dataset.

Parameter	Description
Embeddings (vector size $d$ )	GloVe (Pennington et al., 2014) $d = 300$ <sup>3</sup>
CNN filters $c$	300
CNN window size $l$	5
MLP hidden size	300
Optimizer	Adam(Duchi et al., 2011)
Learning rate	0.0001
Gradient clipping max norm	5.0
Loss coefficients $(a_1, a_2, b_1, b_2, b_3, b_4)$	(i) (1.0, 1.0, 0.5, 0.1, 0.01, 0.02) (ii) (1.0, 1.7, 0.3, 0.7, 0, 0)
Size of negative set $ \{d'\} $	2

<sup>3</sup> Pre-trained with Wikipedia 2014 + Gigaword 5 (<https://nlp.stanford.edu/projects/glove/>)



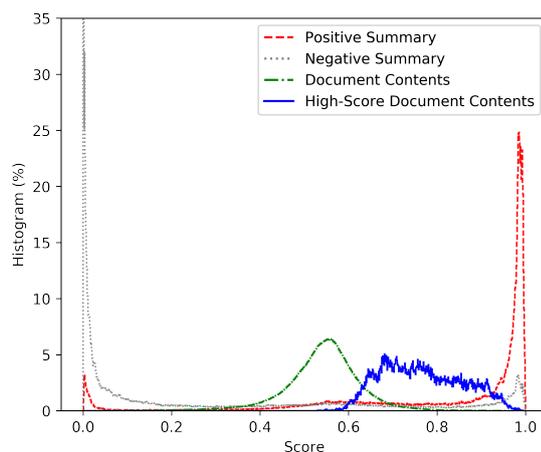
**Fig. 3** Visualization of score distribution (95% confidence) per document showing the comparison among scores of a document’s contents with its summary (positive summary) and other random document’s summary (negative summary). Most of the sample cases, positive summaries have higher mean scores than document contents, and document contents have higher mean scores than negative summaries.

**Table 4** Summarization performance measured in ROUGE scores on dataset from COLIEE 2018 case law retrieval task. The phrase scoring model is trained with loss coefficients (i).

Length Threshold $t$	ROUGE-1			ROUGE-2			ROUGE-SU6		
	Pre	Rec	F1	Pre	Rec	F1	Pre	Rec	F1
10%	0.482	0.409	0.405	0.186	0.152	0.152	0.258	0.199	0.167
20%	0.377	0.592	0.424	0.155	0.244	0.174	0.169	0.388	0.184
30%	0.304	0.687	0.390	0.135	0.311	0.174	0.116	0.511	0.155
40%	0.253	0.745	0.352	0.121	0.364	0.169	0.084	0.592	0.125
50%	0.216	0.784	0.318	0.109	0.407	0.162	0.063	0.651	0.100

The phrase score statistics are shown in Fig. 3 and Fig. 4. Most of the sample cases, positive summaries have higher mean scores than document contents, and document contents have higher mean scores than negative summaries. As shown in Table 5, using our phrase scoring model, we can extract phrases similar to gold summary phrases.

We measure the score distribution but over all data. Similar to per document, positive summaries have higher mean scores than document contents, and document contents have higher mean scores than negative summaries. High-score document contents are selected from top 50 highest phrases for each document. The phrases in high-score document contents affects much to the composition of document vectors, and could also be selected for summarizing documents.



**Fig. 4** Visualization of score distribution over all data. Positive summaries have higher mean scores than document contents, and document contents have higher mean scores than negative summaries. High-score document contents are selected from top 50 highest phrases for each document. The phrases in high-score document contents affects much to the composition of document vectors, and could also be selected for summarizing documents.

<b>Gold Summary</b>	Chan sought judicial review of a visa officer’s decision denying his permanent residence application. He sought an injunction requiring the Minister to allow him to continue working and to take no enforcement action against him until his judicial review application was finally determined. The Federal Court of Canada, Trial Division, dismissed the motion.
<b>High-score Phrases</b>	<ul style="list-style-type: none"> <li>- denying him an immigrant visa</li> <li>- injunction requiring the Minister to</li> <li>- for an immigrant visa moot</li> <li>- judicial review of a visa</li> <li>- existing application for an immigrant</li> <li>- officer ’s decision denying him</li> <li>- application for an immigrant visa</li> <li>- visa officer ’s decision denying</li> <li>- decision denying him an immigrant</li> <li>- immigrant visa has been finally</li> </ul>
<b>Medium-score Phrases</b>	<ul style="list-style-type: none"> <li>- Sciences ( Economics ) from</li> <li>- This will be a matter</li> <li>- might be good reasons that</li> <li>- reasons that could bring s.</li> <li>- of factors listed in column</li> <li>- current employment authorization .</li> <li>- will arise for the applicant</li> <li>- applied for permanent residence in</li> <li>- his application issued on September</li> <li>- its application , even though</li> </ul>
<b>Low-score Phrases</b>	<ul style="list-style-type: none"> <li>- harm will arise as a</li> <li>- turn then to the question</li> <li>- While the applicant has demonstrated</li> <li>- Immigration Regulations are : “</li> <li>- the applicant has demonstrated a</li> <li>- established that irreparable harm will</li> <li>- I turn then to the</li> <li>- of the non-issuance of an</li> <li>- he has not established that</li> <li>- a serious question to be</li> </ul>

**Table 5** Example outputs of phrase scoring model.

For comparison, we evaluate on the dataset provided by Hachey and Grover (2004), which is a collection of 47 judgments of the House of Lord<sup>4</sup> (HOLJ) from 2001 to 2003. We compare the results of sentence selection with the methods of Hachey and Grover (2004) and Kim et al. (2013). Since HOLJ corpus has only 47 documents, the phrase scoring model is trained on COLIEE 2018 dataset. For the task of sentence selection, given a document  $d = \{s\}$ , we select top  $t$  sentences with highest scores computed by sum of sentence (n-gram) phrase scores.

- Hachey and Grover (2004): develop a sentence classification method using models trained on several labor linguistic features: cue phrase, location, entities, sentence length, quotations, and thematic words.
- Kim et al. (2013): develop a graph-based algorithm which selects sentences towards the conclusion/decision of the case. The sentences are connected based on the embedding probability, the probability that a sentence is embedded in another.

**Table 6** Sentence selection results by selection F-score on HOLJ corpus.

Top $t$ Sentences	Pre	Rec	F1
10%	0.197	0.136	0.155
20%	0.182	0.245	0.201
30%	0.168	0.344	0.219
40%	0.168	0.460	0.240
50%	0.171	0.579	0.258
Hachey et al.	0.317	0.307	0.312
Kim et al.	0.313	0.364	0.337

**Table 7** Sentence selection results by ROUGE scores on HOLJ corpus.

Top $t$ Sentences	ROUGE-1			ROUGE-2			ROUGE-SU6		
	Pre	Rec	F1	Pre	Rec	F1	Pre	Rec	F1
10%	0.523	0.715	0.583	0.313	0.424	0.347	0.302	0.530	0.342
20%	0.365	0.846	0.494	0.258	0.592	0.348	0.155	0.739	0.236
30%	0.289	0.896	0.424	0.221	0.685	0.325	0.097	0.824	0.164
40%	0.247	0.931	0.380	0.205	0.770	0.315	0.071	0.882	0.126
50%	0.220	0.957	0.350	0.194	0.838	0.307	0.056	0.928	0.103

Even though, using the phrase scoring model, we can select sentences with high overlap with the gold sentences (Table 7), the accuracy of selecting the labeled sentences is low (Table 6). The results are understood as our phrase scoring model focuses on evaluating the importance of phrases, and is not directly learned to score sentences. Besides, there are two factors our phrase scoring model does not have during inference: (1) any explicit linguistic features other than word embedding, (2) statistical information: term frequency-inverse document frequency. Furthermore, the phrase scoring model is trained on a different corpus. The common of

<sup>4</sup> <https://www.parliament.uk/business/lords/>

the training corpus (COLIEE 2018) with the test corpus (HOLJ) is essentially captured through the use of word embedding.

## 4.2 Retrieval

In the data used in our experiments, the legal cases are sampled from a database of predominantly Federal Court of Canada case laws, provided by Compass Law. The data are provided by COLIEE competition (Kano et al., 2018) held in two years 2018 and 2019. In each of both the datasets, the data contain 285 queries, each query is attached with 200 candidate cases. Each candidate case is presented as a raw text document file which describes the details of the case. While a summary is presented in the query case, the candidate cases may not have summary section.

We formulate the task as bipartite ranking problem and devise the learning to ranking method to solve it. We utilize pair-wise ranking strategy: pairing each noticed case with an irrelevant case from the candidate list. We adopt Linear-SVM as the learning algorithm for solving the optimization problem. The input of the learning-to-rank algorithm is the query-candidate relevance vectors obtained from Equation 16 in Section 3.3. After obtaining the scored candidates as a ranked list, we proceed to select top  $k$  highest scored candidates as the predicted noticed cases.

The phrase scoring model was trained on only COLIEE 2018 dataset, and then adopted to generate encoded summarization vectors for case documents, and text summaries for the candidate cases in COLIEE 2019 dataset. For generating the text summaries, the summary length threshold  $t$  (Section 2.3) is set to  $t = 20\%$  document-length. As shown in Table 1, the average length of summaries is  $\approx 10\%$  document-length for COLIEE 2018 dataset, and  $\approx 9\%$  document-length for COLIEE 2019 dataset. Thus, with a threshold  $t = 20\%$  document-length, we could expect to cover potential information with good recall rate ( $\approx 70\%$ ) while keeping an acceptable summary length.

We evaluated our approach by performing leave-one-out validation where we tested on each and every query from the provided 285 queries and the rest as training data.

We reported our system’s validation results with the following metrics:

- MAP: Mean average precision.
- P, R, F1: Precision, Recall, F-measure whose values are averaged by query. This is straightforward as we average the results of all folds in the leave-one-out validation.

The results in Tables 8, 10, 12, 13, 14, and 15 show that **Lexical+EncSum**, the combination of lexical features with encoded summarization, achieves the best performance.

The validation results of COLIEE 2018 (Table 8) and COLIEE 2019 (Table 10) show that lexical features and latent features complement each other really well. The highest performance with either lexical or latent features is lower than the lowest performance of the combination. The improvement by the combination hints the existence of important information captured by latent features but not captured by lexical features.

*WordEmb*-Hierarchical-pooling performs better than *WordEmb*-Max-pooling and *WordEmb*-Avg-pooling. The hierarchical pooling consists of *WordEmb*-Max-pooling,

**Table 8** Validation results on COLIEE 2018 dataset. We select top 10 highest scored candidates when measuring precision, recall and f-measure. “(summary)” indicates that the corresponding encoding method is applied only on the summary part of the document.

Model	MAP	P	R	F1
Lexical	0.530	0.420	0.520	0.398
<i>WordEmb</i> -Avg-pooling	0.452	0.386	0.440	0.356
<i>WordEmb</i> -Max-pooling	0.325	0.306	0.326	0.275
<i>WordEmb</i> -Hierarchical-pooling	0.528	0.434	0.481	0.400
<i>doc2vec</i>	0.552	0.438	0.533	0.415
<i>WordEmb</i> -Avg-pooling (summary)	0.515	0.444	0.499	0.410
<i>WordEmb</i> -Max-pooling (summary)	0.400	0.370	0.362	0.324
<i>WordEmb</i> -Hierarchical-pooling (summary)	0.619	0.503	0.570	0.469
<i>doc2vec</i> (summary)	0.422	0.367	0.407	0.334
EncSum(i)	0.659	0.510	0.584	0.478
<b>EncSum(ii)</b>	<b>0.690</b>	<b>0.529</b>	<b>0.608</b>	<b>0.494</b>
Lexical+ <i>WordEmb</i> -Avg-pooling	0.686	0.522	0.653	0.502
Lexical+ <i>WordEmb</i> -Max-pooling	0.687	0.515	0.642	0.494
Lexical+ <i>WordEmb</i> -Hierarchical-pooling	0.772	0.565	0.705	0.545
Lexical+ <i>doc2vec</i>	0.684	0.518	0.644	0.496
Lexical+ <i>WordEmb</i> -Avg-pooling (summary)	0.688	0.528	0.646	0.505
Lexical+ <i>WordEmb</i> -Max-pooling (summary)	0.711	0.544	0.677	0.524
Lexical+ <i>WordEmb</i> -Hierarchical-pooling (summary)	0.783	0.579	0.725	0.560
Lexical+ <i>doc2vec</i> (summary)	0.704	0.539	0.675	0.520
Lexical+EncSum(i)	0.849	0.601	0.761	0.583
Lexical+EncSum(ii)	<b>0.888</b>	<b>0.623</b>	<b>0.788</b>	<b>0.607</b>

**Table 9** Lexical feature impact analysis by validation results on COLIEE 2018 dataset. We select top 10 highest scored candidates when measuring precision, recall and f-measure. The coding for lexical features is in the form of q-c, where q is a subset of query components including a summary (s) and paragraphs (p), c is a subset of candidate components including a summary (s) and paragraphs (p). For example, the lexical method sp-plc (q=sp, c=sp) means we perform all 4 matching options, and the lexical method s-p (q=s, c=p) means we only compare the summary of a query with the paragraphs of a candidate.

Lexical Combination	MAP	P	R	F1
s-s	0.372	0.331	0.378	0.302
s-p	0.482	0.386	0.486	0.367
p-s	0.435	0.356	0.434	0.331
p-p	0.469	0.372	0.463	0.355
sp-s	0.458	0.371	0.45	0.346
sp-p	0.510	0.403	0.506	0.384
sp-sp	<b>0.530</b>	<b>0.420</b>	<b>0.520</b>	<b>0.398</b>

*WordEmb*-Avg-pooling features and further sentence-level pooling which regards the sentence information boundary.

As shown in Table 8, when limiting the document to only the summary part than the whole content, most of the models using *WordEmb* or *doc2vec* perform better, except *doc2vec* without lexical features. This suggests the important of summarization in legal case retrieval task.

The suggestion strongly presents in the results of the models using encoded summarization. The models with encoded summarization features outperforms other latent feature generation candidates including *WordEmb*, *doc2vec* on either

**Table 10** Validation results on COLIEE 2019 dataset. We select top 5 highest scored candidates when measuring precision, recall and f-measure.

Model	MAP	P	R	F1
Lexical	0.715	0.495	0.641	0.485
<i>WordEmb</i> -Avg-pooling	0.218	0.177	0.210	0.161
<i>WordEmb</i> -Max-pooling	0.270	0.223	0.260	0.206
<i>WordEmb</i> -Hierarchical-pooling	0.417	0.331	0.405	0.311
<i>doc2vec</i>	0.567	0.404	0.540	0.398
EncSum(i)	0.542	0.430	0.516	0.402
EncSum(ii)	0.576	0.436	0.534	0.410
Lexical+ <i>WordEmb</i> -Avg-pooling	0.733	0.508	0.658	0.496
Lexical+ <i>WordEmb</i> -Max-pooling	0.750	0.526	0.679	0.513
Lexical+ <i>WordEmb</i> -Hierarchical-pooling	0.782	0.549	0.704	0.534
Lexical+ <i>doc2vec</i>	0.725	0.493	0.638	0.482
Lexical+EncSum(i)	0.792	0.552	0.700	0.533
Lexical+EncSum(ii)	<b>0.833</b>	<b>0.579</b>	<b>0.724</b>	<b>0.557</b>

**Table 11** Lexical feature impact analysis by validation results on COLIEE 2019 dataset. We select top 5 highest scored candidates when measuring precision, recall and f-measure. The coding for lexical features is in the form of q-c, where q is a subset of query components including summary (s) and paragraphs (p), c is a subset of candidate components including paragraphs (p), lead sentences (l), and generated summary (e) (described in Section 2.3). For example, the lexical method sp-ple (q=sp, c=ple) means we perform all 6 matching options, and the lexical method s-p (q=s, c=p) means we only compare the summary of a query with the paragraphs of a candidate.

Lexical Combination	MAP	P	R	F1
s-p	0.690	0.484	0.620	0.470
s-l	0.589	0.420	0.528	0.405
s-e	0.561	0.401	0.517	0.390
p-p	0.680	0.476	0.601	0.461
p-l	0.619	0.443	0.563	0.429
p-e	0.588	0.413	0.534	0.402
sp-p	0.712	0.490	0.635	0.480
sp-l	0.634	0.448	0.570	0.435
sp-e	0.602	0.429	0.553	0.416
sp-pl	0.713	0.493	0.639	0.483
sp-pe	0.709	0.485	0.633	0.476
sp-ple	<b>0.715</b>	<b>0.495</b>	<b>0.641</b>	<b>0.485</b>

the summary part or the whole document. Furthermore, the improvement of the encoded summarization suggests that this feature type not only embeds the summary properties of the document but also carries selectively important information from the document content.

The above points also suggest that the summary of a case contains important information but may not contain all relevant information for case retrieval. This is intuitively seen as that the whole case may discuss various legal points besides the main points. Since the encoded summarization weights the case content based on the summary which contains the main points of the case, the other various legal points which are potentially related to the main points may be captured. Hence, the selectively carried information by the encoded summarization could be the related points to the main points of the case.

**Table 12** Results on test data of COLIEE 2018. We select top 10 highest scored candidates when measuring precision, recall and f-measure. “(summary)” indicates that the corresponding encoding method is applied only on the summary part of the document.

Model	P	R	F1
Lexical	0.458	0.429	0.443
<i>WordEmb</i> -Avg-pooling	0.417	0.391	0.404
<i>WordEmb</i> -Max-pooling	0.331	0.310	0.320
<i>WordEmb</i> -Hierarchical-pooling	0.493	0.463	0.477
<i>doc2vec</i>	0.466	0.437	0.451
<i>WordEmb</i> -Avg-pooling (summary)	0.490	0.459	0.474
<i>WordEmb</i> -Max-pooling (summary)	0.432	0.405	0.418
<i>WordEmb</i> -Hierarchical-pooling (summary)	0.585	0.548	0.566
<i>doc2vec</i> (summary)	0.444	0.417	0.430
EncSum(i)	0.598	0.561	0.579
EncSum(ii)	<b>0.608</b>	<b>0.571</b>	<b>0.589</b>
Lexical+ <i>WordEmb</i> -Avg-pooling	0.569	0.534	0.551
Lexical+ <i>WordEmb</i> -Max-pooling	0.566	0.531	0.548
Lexical+ <i>WordEmb</i> -Hierarchical-pooling	0.607	0.569	0.587
Lexical+ <i>doc2vec</i>	0.571	0.536	0.553
Lexical+ <i>WordEmb</i> -Avg-pooling (summary)	0.578	0.542	0.559
Lexical+ <i>WordEmb</i> -Max-pooling (summary)	0.598	0.561	0.579
Lexical+ <i>WordEmb</i> -Hierarchical-pooling (summary)	0.637	0.598	0.617
Lexical+ <i>doc2vec</i> (summary)	0.622	0.583	0.602
Lexical+EncSum(i)	0.676	0.634	0.655
Lexical+EncSum(ii)	<b>0.690</b>	<b>0.647</b>	<b>0.668</b>

**Table 13** Results on test data of COLIEE 2019. We select top 5 highest scored candidates when measuring precision, recall and f-measure.

Model	P	R	F1
Lexical	0.485	0.448	0.466
<i>WordEmb</i> -Avg-pooling	0.157	0.145	0.151
<i>WordEmb</i> -Max-pooling	0.239	0.221	0.230
<i>WordEmb</i> -Hierarchical-pooling	0.334	0.309	0.321
<i>doc2vec</i>	0.403	0.373	0.387
EncSum(i)	0.413	0.382	0.397
EncSum(ii)	0.426	0.394	0.409
Lexical+ <i>WordEmb</i> -Avg-pooling	0.489	0.452	0.469
Lexical+ <i>WordEmb</i> -Max-pooling	0.541	0.500	0.520
Lexical+ <i>WordEmb</i> -Hierarchical-pooling	0.590	0.545	0.567
Lexical+ <i>doc2vec</i>	0.475	0.439	0.457
Lexical+EncSum(i)	0.544	0.503	0.523
Lexical+EncSum(ii)	<b>0.600</b>	<b>0.555</b>	<b>0.576</b>

The validation results (Tables 9, and 11) of lexical features with various combinations (from the 4 matching options for COLIEE 2018 and 6 matching options for COLIEE 2019) described in Section 3.1 show that the combination of lexical matching options does have positive effect to improve the performance on both COLIEE 2018 and COLIEE 2019 datasets. On one hand, it is meaningful to have expert summaries for lexical matching as in COLIEE 2018, and on the other hand, pseudo/generated summaries could also help boost retrieval performance in COLIEE 2019 where candidate summaries are not available.

**Table 14** Participants’ results on test data of COLIEE 2018. We participated in the competition under the name ”JNLP”. ”JNLP-k=10” is our best system utilizing the combination of lexical and encoded summarization using the base parameters.

Model	P	R	F1
HUKB1	0.497	0.308	0.381
HUKB2	0.405	0.304	0.347
JNLP-r=2.5	0.546	0.655	0.596
JNLP-k=10	0.676	0.634	<b>0.655</b>
Smartlaw	0.287	0.431	0.345
UA	0.372	0.323	0.346
UA-postproc	0.348	0.404	0.374
UA-smote	0.354	0.393	0.372
UBIRLED-1	0.133	0.623	0.219
UBIRLED-2	0.196	0.720	0.308
UBIRLED-3	0.561	0.102	0.172
UL	0.564	0.302	0.393

**Table 15** Participants’ results on test data of COLIEE 2019. We participated in the competition under the name ”JNLP”. ”JNLP.task\_1.p” is our best system utilizing the combination of lexical and encoded summarization using the pre-trained phrase scoring model.

Team	Run name	P	R	F1
CACJ	submit_task1.CACJ01	0.212	0.585	0.311
CLArg	CLarg	0.927	0.306	0.460
HUKB	task1.HUKB	0.702	0.400	0.510
IITP	task1.IITPB25	0.626	0.385	0.477
IITP	task1.IITPd2v	0.465	0.346	0.397
IITP	task1.IITPdocBM	0.637	0.388	0.482
ILPS	BERT_Score_0.946	0.681	0.433	0.530
ILPS	BERT_Score_0.96	0.819	0.342	0.483
ILPS	BM25_Rank_6	0.467	0.518	0.491
<i>JNLP</i>	<i>JNLP.task_1.p</i>	0.593	0.549	<i>0.570</i>
<b>JNLP</b>	<b>JNLP.task_1.pl</b>	0.600	0.555	<b>0.576</b>
<b>JNLP</b>	<b>JNLP.task_1.ple</b>	0.600	0.555	<b>0.576</b>
UA	UA_0.52	0.351	0.336	0.344
UA	UA_0.54	0.364	0.324	0.343
UA	UA_0.57	0.356	0.333	0.344

The encoded summarization (EncSum) approach alone achieves MAP of 0.576 and F1 of 0.410 on COLIEE 2019 dataset, lower performance than the best lexical combination. The effect is different from the observation in COLIEE 2018 dataset where the performance of encoded summarization (MAP of 0.690 and F1 of 0.494) is higher than lexical matching approach. Since the encoded summarization model is trained on only COLIEE 2018 dataset, some summary phenomena in COLIEE 2019 dataset may not be well captured.

The combination of encoded summarization and lexical features does improve performance. The improvement by the combination of show that, even though the encoded summarization may not perform well alone, it still provides useful information for identifying relevant cases.

Since our system use similarity as features for predicting supporting relationship, it has the limitation when it comes to non-supporting (unnoticed) but highly similar cases. As shown in Table 16, in the top 10, many retrieved cases are about

**Table 16** An output example of our system from COLIEE 2018 test data. “/\* ... \*/”: omitted.

Query		
Candidates	Noticed	Ranked by our system
AstraZeneca applied for judicial review of a decision by the Minister of Health to disclose certain information related to AstraZeneca’s supplementary new drug submission for LOSEC tablets (omeprozole magnesium) for the treatment of dyspepsia. AstraZeneca argued that: (1) the decision was a nullity because it was made by a person who lacked the authority to make the decision; and (2) the information requested was exempt from disclosure pursuant to ss. 20(1)(b) and (c) of the Access to Information Act or, alternatively, because it was either irrelevant to the request or had previously been severed. /* ... */		
The Minister of Health released certain records related to the applicant’s new drug submission. /* ... */	NO	1
Cyanamid applied under s. 44 of the Access to Information Act to review the Minister’s decision to disclose the product monographs of two drugs and certain severed documents relating to the new drug submission for one of the drugs. /* ... */	YES	2
/* ... */ All the proceedings were brought under the provisions of the Patented Medicines (Notice of Compliance) Regulations and concerned a drug containing the medicine known as omeprazole. /* ... */	NO	3
The Minister of National Health and Welfare refused to permit Apotex Inc. to add information to its New Drug Submissions. /* ... */	NO	4
Allergan Inc. commenced an application under the Patented Medicines (Notice of Compliance) Regulations respecting the ’691 patent. /* ... */	NO	5
The Minister of Transport decided to disclose records which included information about City Express Airline, operated by Air Atonabee Ltd. Air Atonabee applied for judicial review of the decision under s. 44 of the Access to Information Act. /* ... */	YES	7
Brookfield Lepage Johnson Controls Facility Management Services (BLJC) provided professional facility management services to property owners and tenants across Canada. /* ... */	NO	8
Sandoz moved under ss. 6(5)(a) and 6(5)(b) of the Patented Medicines (Notice of Compliance) Regulations, SOR/1993-133, for the dismissal in part of Abbott Laboratories’ prohibition application in respect of clarithromycin patents which included patent 2,387,361. /* ... */	NO	9
The Canadian Tobacco Manufacturers Council et al. applied pursuant to s. 44 of the Access to Information Act (Can.) for an order prohibiting the Minister of National Revenue from disclosing third party information. /* ... */	YES	10
Matol Botanical International Inc. applied to have four decisions authorizing disclosure of information relating to its business reviewed and set aside under s. 44 of the Access to Information Act. /* ... */	YES	12

medicines even though they are not noticed. Aside from that, our system does retrieve the noticed cases (ranked 2, 7, and 10) which are actually diverse in topics.

## 5 Related Work

Legal case retrieval or retrieval of prior cases is an important research topic for decades where approaches to solve the corresponding task involve performing linguistics analysis, logical analysis, common lexical matching, and distributed vector representation with both common and legal expertise knowledge (Bench-Capon et al., 2012). In (Jackson et al., 2003), they build a system called “History Assistant” which extracts rulings from court opinions and retrieves relevant prior cases from a citator database by combining partial parsing techniques with domain knowledge and discourse analysis to extract information from the free text of court opinions. In (Zeng et al., 2005), they develop a knowledge representation model for the intelligent retrieval of legal cases involving decomposing issues into sub-issues, and categorizing factors into pro-claimant, pro-responder and neutral factors. In (Saravanan et al., 2009), they overcome the problem of keyword-based search due to synonymy and ambivalence of words by developing an ontological framework to enhance the user’s query and ensure efficient retrieval by enabling inferences based on domain knowledge. Other works related to building legal ontology are (Wyner, 2008; Wyner and Hoekstra, 2012; Getman and Karasiuk, 2014). Aside of linguistics approaches which are expensive to develop because of the required expertise knowledge, other approaches utilizes the emerging effectiveness of neural networks for natural language processing with the pioneer method of mapping texts to continuous vector space (Mikolov et al., 2013; Le and Mikolov, 2014). In (Mandal et al., 2017a), the authors measure legal document similarity considering structural information of the document including paragraphs, summary and utilizing various representation methods including lexical features: TF-IDF, and topic modeling, and distributed vector representational features: *word2vec*, and *doc2vec*. They, however, do not perform the combination of those features. As lexical features and representational features may potential embed different information since they are extracted by different methodologies, the combination of them is promising.

## 6 Conclusion

We have presented our approach for modeling document summary into continuous vector space. We showed that our approach has positive signs in building an effective legal case retrieval system. The results show the importance of exploiting the summary for solving legal case retrieval task. Furthermore, the improvement by the encoded summarization suggests that this feature type not only embeds the summary properties of the given case but also carries selectively important information from the case content which could be potentially related legal points to the main points of the case. Furthermore, the combination of lexical features and latent features generated with neural networks yields positive results for solving the legal case retrieval task. The experimental results show that lexical features and latent features complement each other. The highest performance with either

lexical or latent features is lower than the lowest performance of the combination. The improvement of the combination hints the existence of latent features not captured by lexical approach. We have also showed that the phrase scoring model trained from COLIEE 2018 dataset can provide useful features for representing documents in COLIEE 2019 dataset. There are several directions for improving the performance of legal case retrieval systems. One is that we can use the documents having a summary in COLIEE 2019 dataset for fine-tuning the phrase scoring model. Besides, the lexical matching has not yet considered the statistical information of terms in the corpus, which can be modeled by term frequency-inverse document frequency for example. Including such information may improve the matching by recognizing the statistically typical words for each document.

**Acknowledgements** This work was supported by JST CREST Grant Number JPMJCR1513, Japan.

## References

- Bench-Capon T, Araszkiwicz M, Ashley K, Atkinson K, Bex F, Borges F, Bourcier D, Bourguine P, Conrad JG, Francesconi E, et al. (2012) A history of ai and law in 50 papers: 25 years of the international conference on ai and law. *Artificial Intelligence and Law* 20(3):215–319
- Chen Q, Zhu X, Ling ZH, Wei S, Jiang H, Inkpen D (2017) Enhanced LSTM for natural language inference. In: *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Association for Computational Linguistics, Vancouver, Canada, pp 1657–1668, DOI 10.18653/v1/P17-1152, URL <https://www.aclweb.org/anthology/P17-1152>
- Duchi J, Hazan E, Singer Y (2011) Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research* 12(Jul):2121–2159
- Galgani F, Compton P, Hoffmann A (2012a) Citation based summarisation of legal texts. In: *Pacific Rim International Conference on Artificial Intelligence*, Springer, pp 40–52
- Galgani F, Compton P, Hoffmann A (2012b) Towards automatic generation of catchphrases for legal case reports. In: *Proceedings of the 13th International Conference on Computational Linguistics and Intelligent Text Processing - Volume Part II*, Springer-Verlag, Berlin, Heidelberg, CILing’12, pp 414–425, DOI 10.1007/978-3-642-28601-8\_35, URL [http://dx.doi.org/10.1007/978-3-642-28601-8\\_35](http://dx.doi.org/10.1007/978-3-642-28601-8_35)
- Getman AP, Karasiuk VV (2014) A crowdsourcing approach to building a legal ontology from text. *Artificial intelligence and law* 22(3):313–335
- Hachey B, Grover C (2004) A rhetorical status classifier for legal text summarisation. *Text Summarization Branches Out*
- Jackson P, Al-Kofahi K, Tyrrell A, Vachher A (2003) Information extraction from case law and retrieval of prior cases. *Artificial Intelligence* 150(1-2):239–290
- Johnson R, Zhang T (2015) Effective use of word order for text categorization with convolutional neural networks. In: *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Hu-*

- man Language Technologies, Association for Computational Linguistics, Denver, Colorado, pp 103–112, URL <http://www.aclweb.org/anthology/N15-1011>
- Kalchbrenner N, Grefenstette E, Blunsom P (2014) A convolutional neural network for modelling sentences. In: Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Baltimore, Maryland, pp 655–665, URL <http://www.aclweb.org/anthology/P14-1062>
- Kano Y, Kim MY, Yoshioka M, Lu Y, Rabelo J, Kiyota N, Goebel R, Satoh K (2018) Coliee-2018: Evaluation of the competition on legal information extraction and entailment. Twelfth International Workshop on Juris-informatics (JURISIN), COLIEE
- Kim MY, Xu Y, Goebel R (2013) Summarization of legal texts with high cohesion and automatic compression rate. In: Motomura Y, Butler A, Bekki D (eds) *New Frontiers in Artificial Intelligence*, Springer Berlin Heidelberg, Berlin, Heidelberg, pp 190–204
- Kim Y (2014) Convolutional neural networks for sentence classification. In: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), Association for Computational Linguistics, Doha, Qatar, pp 1746–1751, URL <http://www.aclweb.org/anthology/D14-1181>
- Le Q, Mikolov T (2014) Distributed representations of sentences and documents. In: Proceedings of the 31st International Conference on Machine Learning (ICML-14), pp 1188–1196
- Levy O, Goldberg Y (2014) Dependency-based word embeddings. In: *ACL* (2), pp 302–308
- Liu Y, Zhang M (2018) Neural network methods for natural language processing. *Computational Linguistics* 44(1):193–195, DOI 10.1162/COLI\_r\_00312, URL [https://doi.org/10.1162/COLI\\_r\\_00312](https://doi.org/10.1162/COLI_r_00312), [https://doi.org/10.1162/COLI\\_r\\_00312](https://doi.org/10.1162/COLI_r_00312)
- Mandal A, Chaki R, Saha S, Ghosh K, Pal A, Ghosh S (2017a) Measuring similarity among legal court case documents. In: Proceedings of the 10th Annual ACM India Compute Conference on ZZZ, ACM, pp 1–9
- Mandal A, Ghosh K, Pal A, Ghosh S (2017b) Automatic catchphrase identification from legal court case documents. In: Proceedings of the 2017 ACM on Conference on Information and Knowledge Management, ACM, New York, NY, USA, CIKM '17, pp 2187–2190, DOI 10.1145/3132847.3133102, URL <http://doi.acm.org/10.1145/3132847.3133102>
- Mikolov T, Sutskever I, Chen K, Corrado GS, Dean J (2013) Distributed representations of words and phrases and their compositionality. In: *Advances in neural information processing systems*, pp 3111–3119
- Olsson L, of Judicial Administration AI (1999) *Guide to Uniform Production of Judgments*. Australian Institute of Judicial Administration, URL <https://books.google.co.jp/books?id=mKnAAQAACAAJ>
- Pennington J, Socher R, Manning C (2014) Glove: Global vectors for word representation. In: Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP), pp 1532–1543
- Saravanan M, Ravindran B, Raman S (2009) Improving legal information retrieval using an ontological framework. *Artificial Intelligence and Law* 17(2):101–124, DOI 10.1007/s10506-009-9075-y, URL <https://doi.org/10.1007/s10506-009-9075-y>

- Severyn A, Moschitti A (2015) Learning to rank short text pairs with convolutional deep neural networks. In: Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval, ACM, pp 373–382
- Tran VD, Nguyen ML, Satoh K (2018) Automatic catchphrase extraction from legal case documents via scoring using deep neural networks. In: Workshop on MIning and REasoning with Legal texts
- Wyner A (2008) An ontology in owl for legal case-based reasoning. *Artificial Intelligence and Law* 16(4):361
- Wyner A, Hoekstra R (2012) A legal case owl ontology with an instantiation of *popov v. hayashi*. *Artificial Intelligence and Law* 20(1):83–107
- Zeng Y, Wang R, Zeleznikow J, Kemp E (2005) Knowledge representation for the intelligent legal case retrieval. In: Khosla R, Howlett RJ, Jain LC (eds) *Knowledge-Based Intelligent Information and Engineering Systems*, Springer Berlin Heidelberg, Berlin, Heidelberg, pp 339–345