# A Summary of the ALQAC 2021 Competition

Nguyen Ha Thanh[1,*], Bui Minh Quan[1], Chau Nguyen[1], Tung Le[1], Nguyen Minh Phuong[1],
Dang Tran Binh[1], Vuong Thi Hai Yen[2], Teeradaj Racharak[1], Nguyen Le Minh[1], Tran Duc Vu[3],
Phan Viet Anh[4], Nguyen Truong Son[5], Huy Tien Nguyen[5], Bhumindr Butr-indr[6],
Peerapon Vateekul[7], Prachya Boonkwan[8],

[1]*Japan Advanced Institute of Science and Technology*
Ishikawa, Japan
[2]*VNU University of Engineering and Technology*
Hanoi, Vietnam
[3]*The Institute of Statistical Mathematics (ISM), Japan*
Tokyo, Japan
[4]*Le Quy Don Technical University (LQDTU), Vietnam*
Hanoi, Vietnam
[5]*Ho Chi Minh University of Science (VNU-HCMUS), Vietnam*
Ho Chi Minh City, Vietnam
[6]*Faculty of Law, Thammasat University (TU), Thailand*
Bangkok, Thailand
[7]*Chulalongkorn University (CU), Thailand*
Bangkok, Thailand
[8]*National Electronics and Computer Technology Center (NECTEC), Thailand*
Pathumthani, Thailand

*Abstract*—We summarize the evaluation of the first Automated Legal Question Answering Competition (ALQAC 2021). The competition this year contains three tasks, which aims at processing the statute law document, which are Legal Text Information Retrieval (Task 1), Legal Text Entailment Prediction (Task 2), and Legal Text Question Answering (Task 3). The final goal of these tasks is to build a system that can automatically determine whether a particular statement is lawful. There is no limit to the approaches of the participating teams. This year, there are 5 teams participating in Task 1, 6 teams participating in Task 2, and 5 teams participating in Task 3. There are in total 36 runs submitted to the organizer. In this paper, we summarize each team's approaches, official results, and some discussion about the competition. Only results of the teams who successfully submit their approach description paper are reported in this paper.

*Index Terms*—ALQAC 2021, summary, legal processing, competition

## I. Introduction

The achievements of natural language processing in recent years are remarkable. Their applications are present in almost every industry. Automated Legal Question Answering Competition (ALQAC 2021) is held with the goal of building a research community in legal text processing and collecting the idea of applying the most advanced techniques to solve law-related problems. ALQAC 2021 is the first time this competition take place, co-located with KSE, International Conference on Knowledge and Systems Engineering Conference. Knowledge and system engineering are the two factors that mostly affect the performance of the system in legal text processing. Inspired by the Competition on Legal Information Extraction/Entailment (COLIEE) [1] for English and Japanese, ALQAC is designed for low-resource language with its own language challenges.

In ALQAC 2021, there are three tasks in statute law processing (information retrieval, entailment, and question answering). The competition's data is prepared in Vietnamese and Thai language. However, this year, all of the participants only choose to work on the Vietnamese dataset. Task 1 is a statute law retrieval task, given a legal query, the systems need to retrieve the most relevant articles. We evaluate the systems using the F2 score metric to balance between precision and recall with precision is weighted higher. In Task 2, the systems can use the retrieval results from Task 1 to make a prediction of whether the articles entail the given query or not, and from that, answer if the statement in the query is lawful or non-lawful. In Task 3, the systems need to directly predict the lawfulness of the query without using the relevant articles. The evaluation metric for Task 2 and Task 3 is accuracy.

The rest of the paper is organized as follows: Section II describes the dataset and evaluation metrics, Sections III, IV, V describe each task, presenting their definitions, list of approaches submitted by the participants, and results attained. Section VI presents some final remarks.

## II. Dataset

Legal data used in ALQAC 2021 are Vietnamese and Thai legal documents. Based on the legal text, we pose questions that can be answered using the relevant articles alone without the need for additional sources such as legal theory or supporting evidence. The questions as well as the relevant

* Corresponding: nguyenhathanh@jaist.ac.jp.

arXiv:2204.10717v2 [cs.CL] 25 Apr 2022

articles are verified by legal experts. Finally, the dataset is formatted to make it most convenient for data processing and result evaluation.

The dataset file formats are shown via examples as follows.

1) **Legal Articles:** Details about each article are in the following format:

```
[
    {
        "id": "45/2019/QH14",
        "articles": [
            {
                "text": "The content of
                    legal article",
                "id": "1"
            }
        ]
    }
]
```

2) **Annotation Samples:** Details about each sample are in the following format:

```
[
    {
        "question_id": "q-1",
        "text": "The content of question
            or statement",
        "label": true,
        "relevant_articles": [
            {
                "law_id": "45/2019/QH14",
                "article_id": "1"
            }
        ]
    }
]
```

$$Precision_i = \frac{\text{the number of correctly retrieved articles of query } i^{th}}{\text{the number of retrieved articles of query } i^{th}} \tag{1}$$

$$Recall_i = \frac{\text{the number of correctly retrieved articles of query } i^{th}}{\text{the number of relevant articles) of query } i^{th}} \tag{2}$$

$$F2_i = \frac{(5 \times Precision_i \times Recall_i)}{(4 \times Precision_i + Recall_i)} \tag{3}$$

$$Marcro - F2 = \text{Average of}(F2_i) \tag{4}$$

$$Accuracy = \frac{\text{(the number of queries which were correctly confirmed as true or false)}}{\text{(the number of all queries)}} \tag{5}$$

The evaluation metrics for Task 1 are precision, recall, and Macro-F2 score as in Equations 1-4. Accuracy is used to evaluate results in Task 2 and Task 3 with respect to whether the yes/no question was correctly confirmed in Equation 5.

In ALQAC 2021, the method used to calculate the final F2 score of all queries is macro-average (evaluation measure is calculated for each query and their average is used as the final evaluation measure) instead of micro-average (evaluation measure is calculated using results of all queries).

## III. TASK 1 - LEGAL INFORMATION RETRIEVAL

### A. Task Description

Task 1's goal is to return articles that are related to a giving statement. An article is considered "relevant" to a statement *iff* the statement rightness can be entailed (as Yes/No) by the article. This task requires the retrieval of all the articles that are relevant to a query.

Specifically, the input samples consist of:

1) **Legal Articles:** whose format is the same as Legal Articles described in Section II.

2) **Question:** whose format is in JSON as follows:

```
[
    {
        "question_id": "q-1",
        "text": "The content of question
            or statement"
    }
]
```

The system should retrieve all the relevant articles as follows:

```
[
    {
        "question_id": "q-193",
        "relevant_articles": [
            {
                "law_id": "100/2015/QH13",
                "article_id": "177"
            }
        ]
    },
    ...
]
```

In which, "relevant_articles" are the list of all relevant articles of the questions/statements.

| Team | Run ID | F2 Score |
|------|--------|----------|
| AimeLaw | run1_bm25 | 0.7842 |
| AimeLaw | run2_bm25_bert_cnn_09_threshold | 0.8061 |
| AimeLaw | run3_bm25_bert_domain_09_threshold | 0.8061 |
| Aleph | result_task_1 | **0.8807** |
| Kodiak | Task_1_run_1 | 0.7855 |
| Kodiak | Task_1_run_2 | 0.7919 |
| Kodiak | Task_1_run_3 | 0.7955 |
| Dline | result_task_1_ver1 | 0.7766 |
| Dline | result_task_1_ver2 | 0.7776 |
| Dline | result_task_1_ver3 | 0.7936 |

### B. Approaches

There are in total 10 validated runs submitted in this task.

1) **AimeLaw (3 runs)** [2] propose approaches of combining scores of BM25 with Domain Invariant Supporting Model and Deep CNN Supporting Model using a weighted sum function.

2) **Aleph (1 run)** [3] build an article ranking model by finetuning their own pretrained model VNLawBERT [4] with a binary classification problem. The authors make negative samples by choosing the closest candidate with the gold samples.

3) **Kodiak (3 runs)** [5] use the lexical matching methods along with the semantic search models with the goal of reaching a balance between the lexical and semantic features.

4) **Dline (3 runs)** [6] preprocess the data to extract lexical features and word embedding features and feeds them into an XGBoost classifier to predict whether a pair of a query and an article is relevant or not.

### C. Results

The final results in this task can be seen in Table I. The performance of the teams is generally not too disparate. All teams surpass the TF-IDF baseline which is 0.6392 F2 Score. With only one run, **Aleph** surprisingly achieve state-of-the-art performance with 0.8807 F2 Score, followed by two runs of **AimeLaw** with 0.8061 F2 Score.

## IV. TASK 2 - LEGAL TEXTUAL ENTAILMENT

### A. Task Description

Task 2's goal is to construct Yes/No question answering systems for legal queries, by entailment from the relevant articles. Based on the content of legal articles, the system should answer whether the statement is true or false.

Specifically, the input samples consist of the pair of question/statement and relevant articles ($>= 1$) as follows:

```json
[
    {
        "question_id": "q-1",
        "text": "The content of question or
            statement",
        "relevant_articles": [
            {
                "law_id": "45/2019/QH14",
                "article_id": "1"
            }
        ]
    }
]
```

The system should answer whether the statement is true or false via "label" in JSON format as follows:

```json
[
    {
        "question_id": "q-193",
        "label": false
    },
    ...
]
```

The evaluation measure is accuracy, with respect to whether the yes/no question was correctly confirmed as in Equation 5.

### B. Approaches

The participants submit 11 valid runs for Task 2 (legal textual entailment).

1) **AimeLaw (3 runs)** [2] propose to fine-tune a BERT model on the sentence pair data created by applying their data augmentation and text matching technique on the annotated data.

2) **Aleph (2 runs)** [3] also utilize a BERT model to be fine-tuned on sentence pair data. However, instead of augmenting data based on the annotated data, they collect external data from law-related websites.

3) **Kodiak (3 runs)** [5] propose to fine-tune BERT models (Multilingual-BERT and PhoBERT [7]) on the annotated data, with/without an asymmetric truncation technique (to address the limitation of BERT models when dealing with long sentences).

4) **Dline (3 runs)** [6] propose two directions: (i) utilize classifiers (Random Forest and SVM) on TF-IDF features, (ii) fine-tune BERT models on the sequence classification task.

### C. Results

We use the proportion of the majority class as a baseline with an accuracy of 0.5455. Table II provides the results in accuracy of participants' models for Task 2. Most of the teams' performance surpasses the baseline with a significant gap. The highest accuracy, 0.6989, is achieved by **AimeLaw** and **Aleph**, and the second-highest accuracy belongs to **Kodiak** with 0.6818. In this task, all participants make use of *BERT* pretrained models. It indicates the robust ability of large-scale pretrained models in language understanding, and their adaptation ability to the legal domain.

| Team | Run ID | Accuracy |
|---|---|---|
| AimeLaw | aimelaw_predictions_1_ | **0.6989** |
| AimeLaw | aimelaw_predictions_2_ | 0.6761 |
| AimeLaw | aimelaw_predictions_3_ | 0.4318 |
| Aleph | aleph_single | 0.5398 |
| Aleph | aleph_pair | **0.6989** |
| Kodiak | Task_2_run_1 | 0.6364 |
| Kodiak | Task_2_run_2 | <u>0.6818</u> |
| Kodiak | Task_2_run_3 | 0.5568 |
| Dline | result_task_2_ver1 | 0.4034 |
| Dline | result_task_2_ver2 | 0.4148 |
| Dline | result_task_2_ver3 | 0.4148 |

| Team | Run ID | Accuracy |
|---|---|---|
| AimeLaw | run1_bm25+bert | <u>0.6477</u> |
| AimeLaw | run2_bm25_bert_domain | 0.6136 |
| AimeLaw | run3_bm25_bert_cnn | 0.6136 |
| Aleph | result_task_3_pair | **0.7102** |
| Aleph | result_task_3_single | 0.5398 |
| Kodiak | Task_3_run_1 | 0.5739 |
| Kodiak | Task_3_run_2 | 0.5682 |
| Kodiak | Task_3_run_3 | 0.6250 |
| Dline | result_task_3_ver1 | 0.5114 |
| Dline | result_task_3_ver2 | 0.5057 |
| Dline | result_task_3_ver3 | 0.5170 |

## V. TASK 3 - LEGAL QUESTION ANSWERING

### A. Task Description

Task 3's goal is to construct Yes/No question answering systems for legal queries.

Given a legal statement or legal question, the task is to answer "Yes" or "No", in other words, to determine whether it is true or false. This question answering could be a concatenation of Task 1 and Task 2, but not necessarily so, e.g. using any knowledge source other than the results of Task 2.

Specifically, the input samples consist of the question/statement as follows:

```
[
    {
        "question_id": "q-1",
        "text": "The content of question or
            statement"
    }
]
```

The system is not provided with the relevant articles. As a result, the teams can use their own results in Task 1 or another source of knowledge to automated answer whether the question/statement is true or false via "label" in JSON format as same as in Task 2:

```
[
    {
        "question_id": "q-194",
        "label": true
    },
    ...
]
```

### B. Approaches

In this task, we receive 9 different valid runs from 4 participants. The summary information about each team's approach is as below:

1) **AimeLaw (3 runs)** [2] leverage their system in Task 1. They use Task 1 pretrained BERT [8] model to classify whether the giving query is Yes or No. By using a Text Matching technique, they find the most relevant clause for a giving query within the articles, then use this clause and the giving query as the input of BERT classifier.

2) **Aleph (2 runs)** [3] use a pretrained model as the classifier with the training data prepared as in the two following methods:
   - Utilizing Task 1 and Task 2 outputs for generating training data.
   - Using provided training data and generating more positive samples by filtering clauses in provided articles.

3) **Kodiak (3 runs)** [5] use output of Task 1 and Task 2 to train a sentence classifier model based on BERT [8], and use this model to make predictions on Task 3.

4) **Dline (3 runs)** [6] suggest an idea of separating an article into meaningful sentences. Based on this segmentation, Dline use a pair of a single sentence and a giving query as the sample to fine-tune on PhoBert.

### C. Results

The baseline for Task 3 is also 0.5455 in accuracy, we can see in Table III, 7 of 11 runs are better than the baseline. The common idea of participated teams for Task 3 is to use outputs of Task 1 and Task 2 for training a classifier. Therefore, **Aleph** achieves first place on Task 3 with *0.7102* accuracy due to impressive results on Task 1 *0.8807* on F2 score, and Task 2 *0.6989*.

## VI. FINAL REMARKS

We have summarized the results of the ALQUAC 2021 competition. Task 1 is about retrieving articles for the purpose of verifying the lawfulness of a given legal question. Task 2 and Task 3 to confirm whether the legal question is lawful or not with and without the relevant articles. We received results from 5 teams for Task 1 (a total of 11 runs), 6 teams for Task 2 (a total of 13 runs), and 5 teams for Task 3 (a total of 12 runs). There are 4 teams successfully submitting technical reports on their methods. This year, we have 1 first prize, 1 second prize for Task 1; 2 first prize, 1 second prize for Task 2; and 1 first prize, 2 second prize for Task 3.

A variety of methods were used for Task 1: retrieval on only the lexical features, taking advantage of pretrained model, combining the lexical and semantic information, and using

boosting technique for classification. For Task 2, transformer-based [9] models are prevalent, but other classical machine learning models are also applied. There is much room for improvement in this task. For Task 3, pretrained models are also utilized in the approaches. Participating teams are not provided with a relevant article list for each query. Interestingly, though, the highest and lowest results in Task 3's rankings are higher than the highest and lowest results in Task 2's rankings.

In COLIEE 2021 on English and Japanese statute law data, pretrained models have an advantage over traditional methods [10]–[12]. This phenomenon is also repeated at ALQAC 2021 on Vietnamese data. This opens up new possibilities in the application of pretrained models in the field of law. In the coming years, we plan to add more data on legal information retrieval, legal question answering as well as introduce to the community other interesting tasks related to automated legal processing.

## REFERENCES

[1] Rabelo Juliano, Goebel Randy, Kano Yoshinobu, Kim Mi-Young, Yoshioka Masaharu, and Satoh Ken. Summary of the competition on legal information extraction/entailment (coliee) 2021. *Proceedings of the COLIEE Workshop in ICAIL*, 2021.

[2] Quang Huy Ngo, Manh Duc Tuan Nguyen, Anh Duong Nguyen, and Quang Nhat Minh Pham. Aimelaw at alqac 2021: Enriching neural network models with legal-domain knowledge. *Proceedings of the 1st Automated Legal Question Answering Competition. ALQAC'2021*, 2021.

[3] Tieu Truong Thinh, Chieu Nguyen Chau, Bui Nguyen-Minh-Hoang, Nguyen Truong Son, and Nguyen Le Minh. Apply bert-based models and domain knowledge for for automated legal question answering tasks at alqac 2021. *Proceedings of the 1st Automated Legal Question Answering Competition. ALQAC'2021*, 2021.

[4] Chieu-Nguyen Chau, Truong-Son Nguyen, and Le-Minh Nguyen. Vn-lawbert: A vietnamese legal answer selection approach using bert language model. In *2020 7th NAFOSTED Conference on Information and Computer Science (NICS)*, pages 298–301. IEEE, 2020.

[5] Dinh Truong Do. Kodiak@alqac2021: Deep learning for vietnamese legal information processing. *Proceedings of the 1st Automated Legal Question Answering Competition. ALQAC'2021*, 2021.

[6] Hoai Linh Luu, Hai Long Nguyen, and Hai Yen Nguyen. Vietnamese legal question answering with combined features and deep learning. *Proceedings of the 1st Automated Legal Question Answering Competition. ALQAC'2021*, 2021.

[7] Dat Quoc Nguyen and Anh Tuan Nguyen. PhoBERT: Pre-trained language models for Vietnamese. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1037–1042. Association for Computational Linguistics, November 2020.

[8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

[9] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.

[10] Sabine Wehnert, Viju Sudhi, Shipra Dureja, Libin Kutty, Saijal Shahania, and Ernesto W De Luca. Legal norm retrieval with variations of the bert model combined with tf-idf vectorization. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Law*, pages 285–294, 2021.

[11] Masaharu Yoshioka, Yasuhiro Aoki, and Youta Suzuki. Bert-based ensemble methods with data augmentation for legal textual entailment in coliee statute law task. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Law*, pages 278–284, 2021.

[12] Ha-Thanh Nguyen, Vu Tran, Phuong Minh Nguyen, Thi-Hai-Yen Vuong, Quan Minh Bui, Chau Minh Nguyen, Binh Tran Dang, Minh Le Nguyen, and Ken Satoh. Paralaw nets–cross-lingual sentence-level pretraining for legal text processing. *Proceedings of the COLIEE Workshop in ICAIL*, 2021.