*Article*

# Contractor's Risk Analysis of Engineering Procurement and Construction (EPC) Contracts Using Ontological Semantic Model and Bi-Long Short-Term Memory (LSTM) Technology

So-Won Choi [1] and Eul-Bum Lee [1,2,*]

1 Graduate Institute of Ferrous and Energy Materials Technology, Pohang University of Science and Technology (POSTECH), Pohang 37673, Korea; smilesowon@postech.ac.kr
2 Department of Industrial and Management Engineering, Pohang University of Science and Technology (POSTECH), Pohang 37673, Korea
* Correspondence: dreblee@postech.ac.kr; Tel.: +82-54-279-0136

**Abstract:** The development of intelligent information technology in the era of the fourth industrial revolution requires the EPC (engineering, procurement, and construction) industry to increase productivity through a digital transformation. This study aims to automatically analyze the critical risk clauses in the invitation to bid (ITB) at the bidding stage to strengthen their competitiveness for the EPC contractors. To this end, we developed an automated analysis technology that effectively analyzes a large amount of ITB documents in a short time by applying natural language processing (NLP) and bi-directional long short-term memory (bi-LSTM) algorithms. This study proposes two models. First, the semantic analysis (SA) model is a rule-based approach that applies NLP to extract key risk clauses. Second, the risk level ranking (RLR) model is a train-based approach that ranks the risk impact for each clause by applying bi-LSTM. After developing and training an artificial intelligent (AI)-based ITB analysis model, its performance was evaluated through the actual project data. As a result of validation, the SA model showed an F1 score of 86.4 percent, and the RLR model showed an accuracy of 46.8 percent. The RLR model displayed relatively low performance because the ITB used in the evaluation test included the contract clauses that did not exist in the training dataset. Therefore, this study illustrated that the rule-based approach performed superior to the training-based method. The authors suggest that EPC contractors should apply both the SA and RLR modes in the ITB analysis, as one supplements the other. The two models were embedded in the Engineering Machine-learning Automation Platform (EMAP), a cloud-based platform developed by the authors. Rapid analysis through applying both the rule-based and AI-based automatic ITB analysis technology can contribute to securing timeliness for risk response and supplement possible human mistakes in the bidding stage.

**Keywords:** AI; EPC contract risk extraction; NLP; ontological semantic model; EPC contract lexicon; deontic logic; bi-LSTM; risk level ranking; digital transformation

## 1. Introduction

The lump sum turn key (LSTK) contract for engineering, procurement, and construction (EPC) projects is a typical contract type used in large-scale and complex plant projects [1]. The EPC plant project combines manufacturing and services such as knowledge service, design, equipment, and construction. In addition, it is a complex industry with various front and back sectors. Furthermore, it includes global supply chains throughout the entire cycle, from bidding to maintenance [2]. In particular, the LSTK contract, in which the EPC contractor bears all liabilities related to design, purchase, construction, and commissioning, is an unbalanced contract as it pays more risks to the EPC contractor due to the increase in complexity when the size of the project expands [1]. Overseas EPC plant projects of Korean companies have been growing in earnest since the mid-2000s, and the

number of orders has continued to increase due to market expansion [3]. However, the need to improve project risk management emerged as the EPC plant industry experienced an earnings shock, such as a decline in cut yields due to a decrease in oil prices. The development of intelligent information technology in the Fourth Industrial Revolution is currently evolving into the digital transformation of all industries. Thus, it is required to increase productivity and strengthen competitiveness using convergence technology to respond to the EPC business environment that is becoming more extensive and complex. Accordingly, the authors' research team considered applying artificial intelligence (AI) technology to manage the risk of bidding documents in the bidding stage of the project.

Invitation to bid (ITB), contract, and claim, mainly used during the EPC project, are text-based unstructured data that describe the client's requirements and significant contractual issues. Furthermore, failure to adequately review the risks of ITB in the bidding stage may result in future disputes. Nevertheless, EPC contractors struggle with ITB analysis and detection of risk clauses due to the large volume of documents, tight schedules, and lack of experienced practitioners during the bidding phase. To this end, research on a system that can analyze bidding documents, especially ITB risk factors, is required at the project bidding stage. It is necessary to convert text data in natural language form into a script that the computer recognizes. In addition, there were not many cases of NLP and AI in the EPC industry compared to other fields, such as medicine; therefore, it is relatively incomplete. Since EPC project documents consist of a large portion of unstructured text data, there is ample space for the NLP application. NLP is a branch of AI that utilizes AI to enable computers to process natural language text [4]. This paper explores a novel approach to automating risk analysis of EPC contracts and computational developments in NLP.

The purpose of the study is to effectively analyze a vast amount of ITB documents in a short period and reduce the uncertainty of decision-making based on human experience and judgment. In addition, it is aimed to support the quick decision-making of EPC contractors and enhance competitiveness by automatically analyzing the critical risks of the ITB in the bidding stage of the EPC project. In this paper, a novel framework of the NLP-based semantic analysis (SA) model and the bi-directional long short-term memory (bi-LSTM) method-based risk level ranking (RLR) model is proposed to analyze the contract risk clause of EPC ITB automatically.

The proposed SA model is an approach that applies the EPC contract lexicon to SVO tuples to develop semantic rules. Then, it extracts risks according to whether the analysis target sentence matches the rules. This study applied the ontology-based semantic information extraction (IE) technique, which maps heterogeneous contract clauses to ontology-based lexicons. Ontology expresses the relationship between objects in a form that a computer can process, and by linking domain knowledge, it becomes the basis for developing semantic rules. In ontology-based semantic IE, the lexicon configuration is significant because it determines the risk clauses by considering the semantic relationship of sentence elements based on the information stored in the lexicon, rather than using a simple keyword search method. Therefore, ontology-based semantic IE performs better than syntactic IE [5]. The EPC contract lexicon was developed along with EPC contract experts through this study. In addition, a PDF structuralization module that recognizes and formalizes text data in documents separately was designed to improve the accuracy of text data analysis. The RLR model was created to address the issue of using the bi-LSTM algorithm, check the risk of each sentence in the ITB, and classify the risk class. Furthermore, the RLR model classifies and extracts each sentence of the EPC contract document into five levels according to the degree of risk. Moreover, a dataset for model training was developed, and hyperparameters were optimized to maximize model performance.

This paper consists of eight sections. Section 1 includes the background and the necessity of the study. The prior research review on knowledge-based risk extraction, contract analysis using AI, and text classification is in Section 2. Section 3 describes the overall architecture and process of the SA and RLR models. Data collection and data conversion for analysis are shown in Section 4. Sections 5 and 6 are the core of this paper, as

they illustrate the development of EPC contract lexicon, introduce the semantic rules for the SA model, and explain the RLR model using bi-LSTM. In addition, the sections discuss the SA and RLR model development, model testing with actual data, and validation. Section 7 analyzes the system application of the developed model, while Section 8 describes the conclusion and implications of this paper. Additionally, Section 8 discusses the limitations and future research directions. In essence, it is expected that rapid risk analysis through AI-based automatic ITB analysis technology will secure timeliness for EPC project bidding risk response.

## 2. Literature Review

Previous studies reviewed three aspects: (1) a knowledge-based risk extraction method in a construction project, (2) automatic extraction of contract risk by applying AI technology to an EPC project, and (3) text classification. Although this study targets the EPC plant project, the prior research also included the construction field.

### 2.1. Knowledge-Based Risk Extraction for EPC Projects

Ebrahimnejada et al. [6] proposed the extended VIKOR method based on the fuzzy set theory as a new risk evaluation approach in large-scale projects. They applied it to the Iranian power plant project to compare the differences with the traditional version. Hung and Wang [7] conducted a study to identify the main risk factors that cause delays in hydropower construction projects in Vietnam and analyze the degree of impact of each risk factor on construction. Jahantigh and Malmir [8] identified, evaluated, and prioritized significant financial risks of EPC projects in terms of national development in developing countries. Furthermore, their work was based on the fuzzy TOPSIS model and they applied the refinery project as a case study. Kim et al. [9] developed the Detail Engineering Completion Rating Index System (DECRIS) that minimizes the rework of EPC contractors and supports schedule optimization for offshore EPC projects. This model improved existing theories, such as the Project Definition Rating Index (PDRI) and front end loading (FEL). Their study verified the effect of schedule and cost through 13 megaprojects. Kabirifar and Mojtahedi [10] studied the most critical factors in EPC project execution by applying the TOPSIS method to a large-scale residential construction project in Iran. In addition, they derived that procurement is the most vital risk factor. Gunduz and Almuajebh [11] ranked 40 critical success factors (CSFs) after reviewing the literature on CSFs considering stakeholder impacts in construction projects. Their collected data were analyzed using the relative importance index (RII) and analytic hierarchy process (AHP) method with Saaty random index. Koulinas et al. [12] proposed a simulation-based approach to estimate the project schedule's delay risk and predict in-time project completion. This approach, implemented through a hotel renovation project, showed better uncertainty expression and superior predictions in comparison to the classic PERT method when estimating budget and time-critical overruns. Okudan et al. [13] developed a knowledge-based risk management tool (namely, CBLisk) using case-based reasoning (CBR). As a web-based tool, this system is characterized by applying the project similarity list in the form of fuzzy linguistic variables for effective case search.

### 2.2. Automatic Extraction of Contract Risks Using AI Technology in EPC Projects

In recent years, research on extracting contract risk from legal documents has been actively conducted by applying AI technology. Surden [14] studied the method of representing specific contractual obligations in computer data for financial contracts, such as stock option contracts. Automated manual comparison has significantly reduced transaction costs associated with contract monitoring compared to traditional written contracts as it applies a technology that transforms specific contract terms into a set of computer-processable rules. In 2018, LawGeex [15] collaborated with 20 experienced lawyers educated in the United States to conduct a study of a contract review platform developed with the AI application. The study, which looked at non-disclosure agreements (NDAs), showed that

AI was 94 percent accurate compared to experienced lawyers, who were 85 percent accurate. Their study improved the quality of legal human resources through faster and more reliable contract management. Cummins and Clack [16] reviewed the concept of "computable contracts", which both humans and computers can understand as the concept exists in text form in natural language. Furthermore, they proposed an integrated framework of various technologies and approaches to model their concepts. Dixon Jr. [17] described the application cases of various AI technologies used in the legal field, such as crime prediction, prevention, detection, and contract drafting and review. Clack [18] studied the problems of converting natural language into computer code that occurred when developing a "smart legal contract", which automates legal contracts using computer technology. His study explained the importance of language design in smart contracts, such as computable language, natural language, and the meaning of the language expression. Salama and El-Gohary [19] studied an automated compliance-checking model that applied deontic logic to the construction domain.

EPC documents consist of a significant portion of text-based unstructured data, while NLP technology is mainly used for text information extraction and retrieval. NLP is an AI-related field of human–computer interaction that enables a computer to interpret human language through machine learning [20]. Zhang and El-Gohary [21] presented a semantic rule-based NLP approach using information extraction (IE) from complex construction regulations. Their study was meaningful as it allowed an advancement in the existing method of selectively extracting only some information from documents. Williams and Gong [22] proposed a risk model to predict cost overruns using data-mining and classification algorithms in bidding documents for construction projects. However, there was a limitation in analyzing only simple keyword-oriented text data, such as project summary information for text analysis. Lee and Yi [23] proposed a bidding risk prediction model using construction project bidding information text mining. However, there was no quantitative explanation of how much the cost should reflect. Zoua et al. [24] proposed an approach that combines two NLP techniques, a vector space model (VSM) and semantic query expansion, to improve search efficiency for accident cases in a construction project. As a result of the study, the problem of semantic similarity remains a significant challenge. Lee et al. [25] proposed a contract risk extraction model for construction projects by applying NLP's automatic text analysis method to the Fédération Internationale Des Ingénieurs-Conseils (FIDIC) Redbook. Their study showed the performance of extracting only about 1.2 percent of the whole sentence as a risk, and their model cannot be applied to other types of contracts other than FIDIC-based, such as offshore plants. Moon et al. [26] proposed an information extraction framework that used Word2Vec and named entity recognition (NER) to develop an automatic review model for construction specifications when bidding for infrastructure projects. Their model targeted only the text data of the construction specification document and it could not analyze the text data shown in the tables or drawings included in the document. Choi et al. [27] developed the Engineering Machine Learning Automation Platform (EMAP). This integrated platform supports decision-making by applying AI and machine learning (ML) algorithms based on data generated throughout the EPC project cycle. Their study is meaningful because it is the first integrated platform for risk extraction of the entire EPC project life cycle. Choi et al. [28] developed a model for checking the presence of a risk clause in an EPC contract using NER and a phrase-matcher. Park et al. studied an ML-based model to extract technical risks from EPC technical specification documents [29]. Choi et al. [27] and Park et al. [29] were interrelated as they created the parts of the sub-element constituting the *EMAP* system. Fantoni et al. [30] utilized state-of-the-art computer language tools with an extensive knowledge base to automatically detect, extract, split, and assign information from technical documents when tendering for a railway project. The implementation of the methodology was utilized during a high-speed train project.

## 2.3. Text Classification

Text classification classifies text data into meaningful categorical classes and is one of the leading research areas of NLP [31]. Traditional text classification methods include dictionary-based and basic machine learning methods [31]. Since the 2000s, it has been replaced by deep learning such as recurrent neural network (RNN), long short-term memory (LSTM), and convolutional neural network (CNN) [32]. Currently, a more powerful text classification technique, such as BERT, has emerged [33]. RNN is one of the neural network architectures used for text mining and classification. Additionally, RNN is a kind of artificial neural network in which directed edges connect hidden nodes to form a directed cycle [34]. Furthermore, it is suitable for processing time-series data that appear sequentially, such as speech and text [35]. However, RNNs have a problem of long-term dependencies in which past learning results disappear. Thus, LSTM was designed to overcome this issue of RNNs [36,37]. The LSTM model proposed by Hochreiter and Schmidhuber [37] is internally controlled by the gating mechanism called input gate, output gate, and forget gate. By improving the long-term dependency problem of RNN, it processes massive data such as time-series data without any problem. However, the unidirectional LSTM has the disadvantage of preserving only past information [38]. Schuster and Paliwal [38] proposed a bi-LSTM model that extends the unidirectional LSTM through introducing a second hidden layer to compensate for this problem of LSTM. Bi-LSTM uses LSTM cells in both directions, therefore past and future information can be exploited [39]. In addition, it is mainly used for text classification due to its excellent performance on sequential modeling problems [33]. Li et al. [40] reviewed text classification methods from 1961 to 2021 and created a taxonomy for text classification tasks from traditional models to deep learning. They also introduced the datasets with a summary table and provided the quantitative results of the leading models. Minaee et al. [41] provided a comprehensive review of deep-learning-based models for text classification developed in recent years and discussed their technical contributions, similarities, and strengths. They also explained a summary of more than 40 popular datasets for text classification.

The research that analyzes the risks of EPC contracts by applying AI technology is relatively insufficient. As a result of the review of previous studies, a majority of the research is focused on construction projects. In addition, most studies have selected and analyzed either rule-based or training-based approaches. However, research that applies both techniques to contract analysis is lacking. This study automatically extracts the contract risk of the EPC contract document by using the rule-based approach of NLP. In addition, it is a study on a training-based method of deep learning that ranks the risk level of contract clauses in the ITB by applying the text classification technique of bi-LSTM.

## 3. Research Framework and Process

### 3.1. Model Framework and Development Process

There are two main types of approaches used in NLP: rule-based and training-based [42]. This study applied both rule-based and training-based methods. The rule-based approach uses manually coded rules for text processing. Although it requires a large amount of human effort, it tends to show higher performance than training-based models [42]. The training-based approach is to learn a text processing model through an ML algorithm so that the machine can understand the meaning of the text expressed in natural language. Recently, the proportion of training-based methods has gradually increased. Furthermore, the SA model was developed by applying a rule-based approach, while the RLR model is training-based and applies the bi-LSTM.

The SA model uses NLP to extract risk clauses that are not found by keyword search through two steps. First, the risk clause is extracted by applying the relationship of the lexicon and subject–verb–object (SVO) tuples. Second, it utilizes deontic logic to formulate the extracted risk clause into obligation, permission, and forbidden. The SA model was developed by a rule-based approach and is expected to have high performance. When applying the rule-based method, ITB composed of a proper sentence shows higher sentence

extraction results than informal sentences such as blogs. The RLR model is classified into five levels, according to the degree of risk, by applying bi-LSTM. This training-based model developed a training dataset and was embedded as a sub-element of the *EMAP* system, an integrated platform. Figure 1 below is the overall architecture and development procedure of the SA model and RLR model.
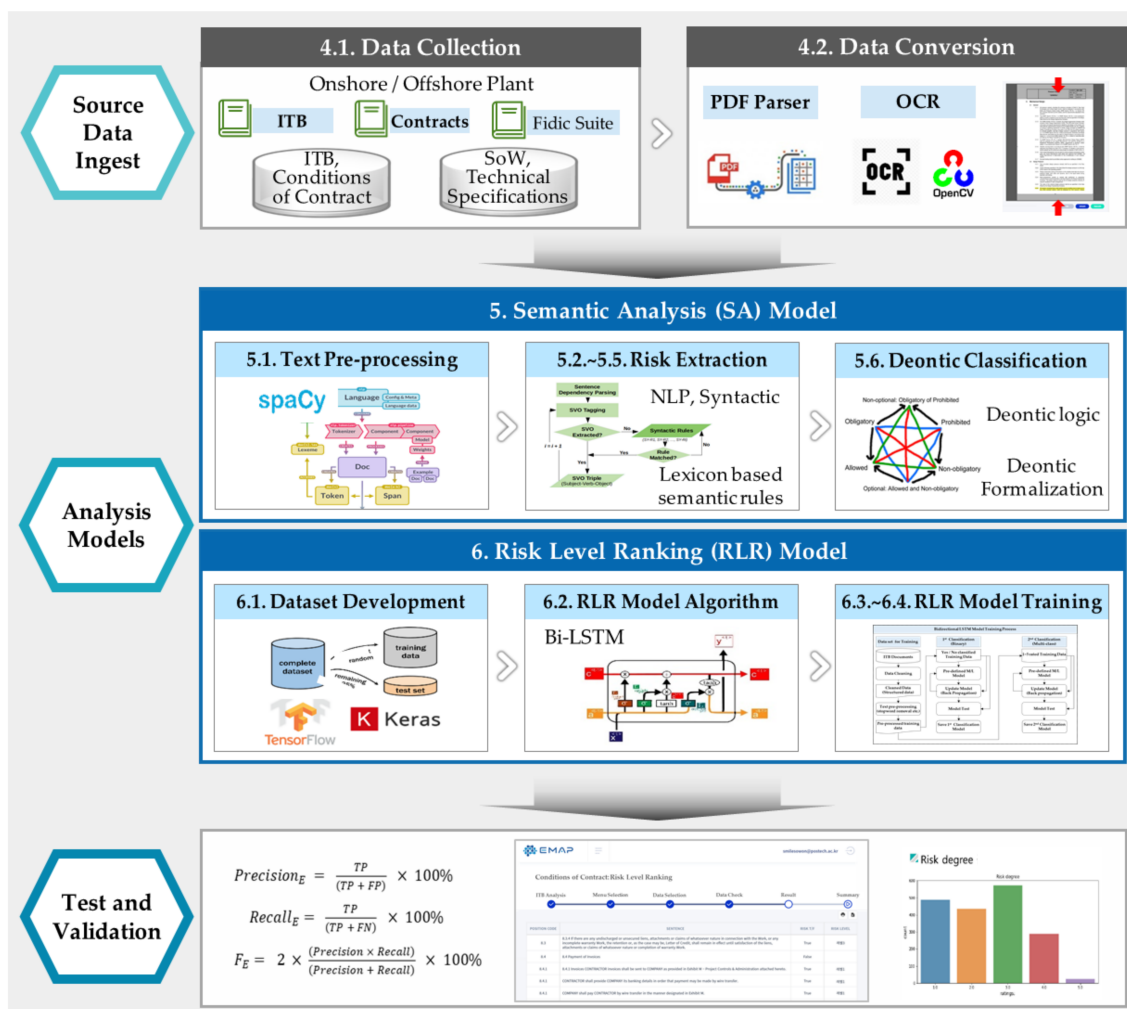


**Figure 1.** The overall architecture of the rule-based SA model and training-based RLR model.

*3.2. Research Scope and Algorithm Development Environments*

The scope of this study was limited as follows. First, this study analyzes the critical risk clauses of the EPC ITB to support the quick and accurate decision-making of the EPC contractor. Second, the subject of this study was restricted to contract documents for onshore and offshore plants of the EPC project. Third, the contract risk defined in this study was set as the key clauses to be confirmed by the EPC contractor in the bidding or contract stage. Fourth, this study applied an AI technique to extract information from contracts to analyze text data. In particular, an ontology-based NLP approach and bi-LSTM of deep-learning were used to propose an automatic risk extraction model for ITB. Fifth, this study does not include the documents generated during the design, procurement, and construction stages of the EPC project. In addition, only contract documents, excluding technical documents, were analyzed among ITB. Sixth, unstructured data such as drawings, tables, images, and videos were excluded. Moreover, only text data among unstructured data were analyzed. These SA and RLR models were applied with NLP and bi-LSTM algorithms, respectively (Table 1). The SA model adapted the NLP technique and spaCy's

2.3.1 library, while the bi-LSTM, Keras 2.6.0, and Tensorflow 2.6.0 libraries were used in the RLR model. The entire procedure was performed in Windows 10 OS and Python 3.7 environments. The system configuration and package information are shown in Table 1.

**Table 1.** Summary of algorithm package information.

| Module | Semantic Analysis | Risk Level Ranking |
|---|---|---|
| AI technology | NLP | Bi-LSTM |
| Libraries | spaCy's 2.3.1 | Keras 2.6.0, Tensorflow 2.6.0 |
| Language | Python 3.7.7 | Python 3.7.11 |
| Input data | EPC Contracts | EPC contracts |
| Operation system | Window 10 | Window 10 |
| Purpose | To extract the risk clauses using the semantic rules based on the lexicon | To classify each sentence of the EPC contracts into five levels by risk degree |

## 4. Data Collection and Conversion

### 4.1. Data Collection

Data were collected for contracts with legal contents to extract the risk of EPC contracts. Among EPC projects ordered over the past 18 years (2003–2020) in North Sea, Australia, Middle East, South and North America, and Africa—eight onshore projects, ten offshore projects, and three Fédération Internationale Des Ingénieurs-Conseils (FIDIC) contracts—a total of 21 contract documents were collected. After converting the data through the PDF structuralization module, the collected contract documents were used for risk analysis. Table 2 summarizes and lists the EPC contracts collected for this study.

**Table 2.** List of the collected ITBs from EPC projects and FIDIC contracts.

| Category | No. | Project Type | Location | Year |
|---|---|---|---|---|
| | 1 | Refinery | Kuwait | 2005 |
| | 2 | Coal-fired Power Plant | Chile | 2007 |
| | 3 | Refinery | Peru | 2008 |
| Onshore | 4 | Combined Cycle Power Plant | Kuwait | 2008 |
| | 5 | Petrochemical | Saudi Arabia | 2011 |
| | 6 | LNG Terminal | USA | 2012 |
| | 7 | Thermal Power Plant | Bangladesh | 2015 |
| | 8 | Combined Cycle Power Plant | Georgia | 2020 |
| | 9 | FPSO [1] | Nigeria | 2003 |
| | 10 | Drillship | For Chartering | 2007 |
| | 11 | FPSO | Angola | 2009 |
| | 12 | FLNG [2] | Brazil | 2010 |
| | 13 | FPSO | Angola | 2011 |
| Offshore | 14 | FPSO | Nigeria | 2012 |
| | 15 | FPSO | Australia | 2012 |
| | 16 | TLP [3] | Congo | 2012 |
| | 17 | Semi-submersible | Gulf of Mexico (US) | 2012 |
| | 18 | Fixed Platform | Norway | 2012 |
| | 19 | FIDIC Red 2017 | | 2017 |
| FIDIC [4] | 20 | FIDIC Silver 2017 | Standard form of Contract | 2017 |
| | 21 | FIDIC Yellow 2017 | | 2017 |

[1] FPSO: floating production storage and offloading. [2] FLNG: floating liquefied natural gas. [3] TLP: tension leg platform. [4] FIDIC: Fédération Internationale Des Ingénieurs-Conseils.

### 4.2. Data Conversion through PDF Structuralization

In order to analyze a document composed of text, such as a contract, it is necessary to delete unnecessary ITB information and structure the data [25]. In this study, a separate

module was developed to extract data for analysis from portable document format (PDF). Furthermore, the documents were composed of text and the module was named PDF structuralization. PDF structuralization is used for removing noise data such as headers, footers, page numbers, and watermarks that are not required for analysis from documents and converting text data into a data frame with sentence units. The data extracted in PDF format were then used for analysis. This study performed data conversion through PDF format before SA and RLR analysis. Figure 2 shows the process of extracting text data from a PDF document.
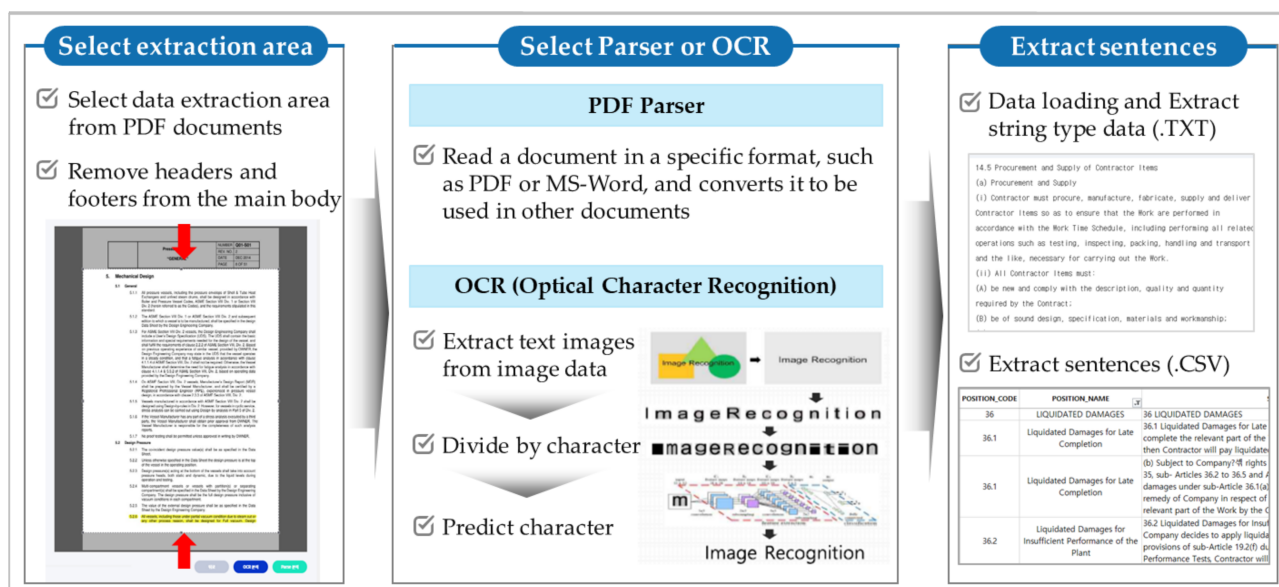


**Figure 2.** Text data extraction process from PDF documents using PDF structuralization module.

The process of extracting text from PDF files was divided into two methods: PDF parser and optical character recognition (OCR). The easy-to-recognize parser is used first, while the OCR is used for text data that the Parser cannot decipher. PDF parser is a compiler that builds in the PDF file structure and allows internalized data to be imported. Through parser, PDF metadata or text data can be read into a computer [43]. OCR is a technology that reads text engraved on a document using light [44]. It extracts text from a scanned image through an OCR character recognition and is utilized for ITB analysis. Furthermore, the structuralized text data were converted into Excel or CSV format for DB. Figure 3 showcases the result of the structuralization of PDF documents in CSV format using PDF parser.



(**a**)



(**b**)

**Figure 3.** Example of text data conversion using PDF structuralization (PDF to CSV). (**a**) Original text contents from the ITB (.PDF); (**b**) automatically extracted text contents (.CSV).

The parser technique was used to define relationships and identify the composition of sentences in text data. The position code shown in Figure 3b is a number for tracking the table of contents. The position name corresponds to the table of contents in the document as well. This structuralization makes it easy to find the affiliation of each sentence. The data conversion result was automatically generated as a CSV file and was used as input data for ITB analysis.

## 5. Semantic Analysis Model

In recent years, research on information extraction of construction contracts using AI has been actively conducted. Various studies applying the NLP technique have been attempted, such as a study to check the presence or absence of a risk clause using a phrase-matcher [28] and an NER model that finds similar texts through training [26]. A limitation of the NER study is that an error in similar phrase tagging occurs in the case of a label with insufficient learning due to the lack of a training dataset. Most of the risk clauses in the contract can be extracted by keyword search through phrase-matcher. However, risk clauses cannot be extracted through a simple keyword search, such as the structuralization Fail-Safe clause. Thus, the SA model is used to extract specific contract clauses that are not extracted with NER or phrase-matcher. This study proposes an ontology-based lexicon and semantic rule mapping approach. The SA model is a method of automatically extracting the key risk clauses based on NLP rules and the knowledge of the EPC contract. Although this model has low efficiency in analysis time compared to ML, it has an advantage in higher information extraction accuracy, as it applies human knowledge to the system [21].

The model also provides a high-level description of the knowledge base (KB), the lexicon, that we built for risk extraction. In this study, the lexicon classifies the contract clauses of the EPC contract and shows the relationship of subsumptions with related terms. In addition, the SA model follows the rule-based NLP pipeline. Text data that have undergone preprocessing is in the order of syntactic analysis, EPC contract lexicon development, semantic rule development, rule matching, risk clause extraction, deontic formalization, and classification. Figure 4 illustrates an analysis procedure of the SA model that extracts contract risk based on this approach.
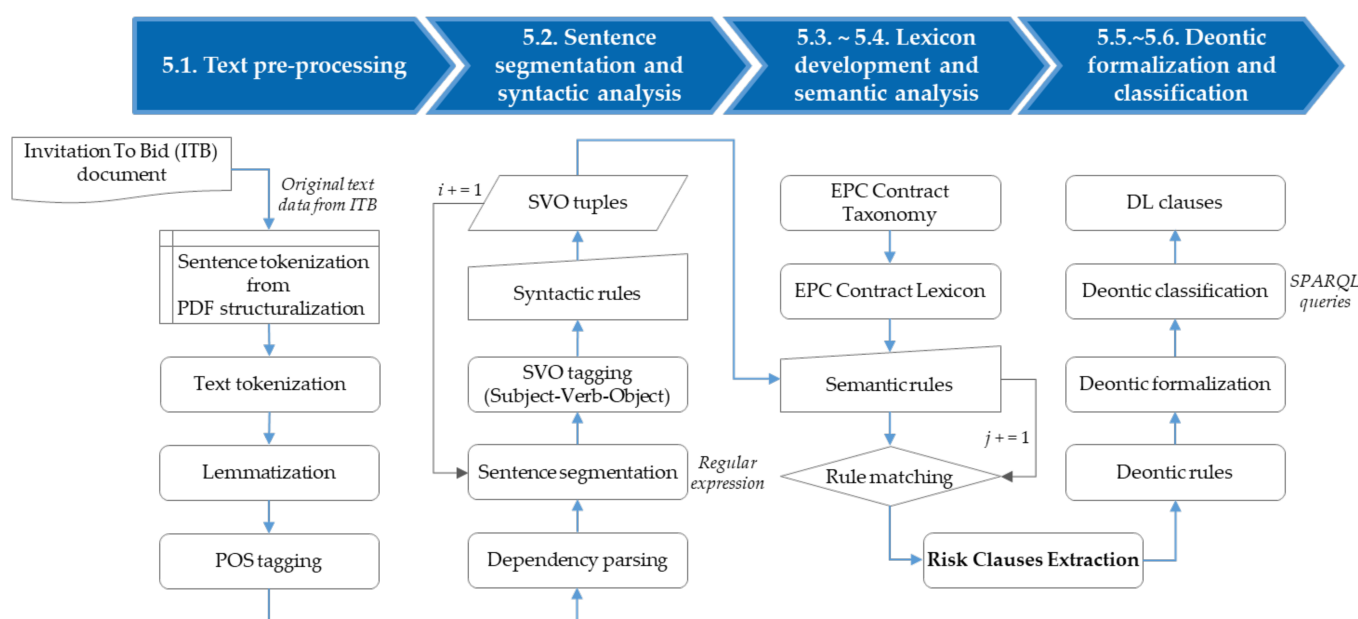


**Figure 4.** The SA model algorithm and implementation process.

As shown in Figure 4, the raw data of the ITB text were converted into data that can be analyzed through the PDF structuralization module. For the structured text data,

preprocessing was performed, such as text tokenization, lemmatization, POS tagging, and dependency parsing. SVO tuples in the sentence were extracted by performing syntactic analysis on the preprocessed data to understand the sentence's grammatical structure. We developed a semantic rule by applying the EPC contract lexicon to the extracted SVO tuples and removed the risk according to whether the rule matched. Then, the procedure of formalizing the extracted risk clause resulted in obligation, permission, and forbidden, which was determined by using deontic logic.

### 5.1. Text Data Preprocessing

Text data that have completed PDF structuralization cannot be used directly for syntactic analysis. A separate preprocessing was required for the basic structure analysis of sentences to perform syntactic analysis [20,25]. Preprocessing techniques for text data consist of various computing techniques such as tokenization and stop word removal [45]. The purpose of preprocessing for syntactic analysis is to refine the text data converted through the PDF structuralization module into the smallest sentence unit. Preprocessing such as text tokenization, lemmatization, POS tagging, and dependency parsing was performed. The above process used the programming language Python and spaCy, an open-source NLP library.

Tokenization is the division of the text into meaningful segments, which are called tokens [46]. Furthermore, it is a type of document segmentation in which a given text is divided into smaller and more specific informational segments. It not only divides a document into paragraphs but also divides a paragraph into sentences, a sentence into phrases (or clauses), and a phrase into tokens (words) and punctuation marks, which are all types of segmentation [5]. In linguistics, a lemma means the original form of a specific word (e.g., the original form be of is and are). Lemmatization means determining and normalizing the lemma of words inflected in various forms in a sentence [5]. Lee et al. [25] applied lemmatization to reduce sentence noise using SyntaxNet. However, unlike previous studies, this study, using spaCy's NLP model, identified each word's part of speech (POS) through an artificial neural network model when performing lemmatization. Then, the lemma was restored based on that information [47].

After tokenization, spaCy predicts the POS of each word through an artificial neural network [48]. The NLP model provided by spaCy uses an artificial neural network trained to predict the correct POS of each word in the context based on numerous sentence data. As a result, one word is assigned one POS information. When analyzing POS in English, spaCy utilizes the Universal Dependencies v2 POS tag set applied, regardless of language, and the Onto-Notes 5 version of the Penn Treebank tag set specialized for English. In this study, the Universal POS tag was applied, and when difficulties arise, OntoNote 5 version Penn treebank tag was used additionally.

Dependency parsing is necessary for preprocessing syntactic analysis. The syntax analysis method uses dependency parsing and phrase structure parsing [49]. Dependency parsing analyzes the relationship and dependency of POS tagged words [5]. In addition, it analyzes the role of each word in the sentence and receives information about it. This study used dependency parsing, in which word order is relatively free, and subject or object can be omitted. In spaCy's basic pipeline, when tokenizer and POS tagging are completed, syntactic dependency, which is syntactic dependence between each word in a sentence, can be identified through another artificial neural network different from POS [50]. Since the dependency relationship between words is not determined by one scheme, clearNLP's CLEAR Style was used to classify English grammatical dependence [51]. Although the above procedure yields syntactic analysis results similar to human interpretation, continuous efforts to improve accuracy are required.

### 5.2. Syntactic Analysis with Sentence Segmentation

Based on the depth and level of analysis, NLP techniques can be classified into lexical, syntactic, and semantic analysis [52]. Syntactic analysis is a sentence-level analysis that

identifies the grammatical relationship of each word in a sentence [52]. The SA model requires the understanding and analysis of the entire context of a sentence rather than extracting only a few risk words from ITB. As a result of dependency parsing, the primary dependency relationship of a sentence can be identified, but the analysis of complex sentences such as ITB has limitations. Therefore, this study applied a syntactic analysis method to define the subject, verb, and object as information extraction factors to understand the semantic relationship between each word in a sentence.

In terms of sentence structure, there are short sentences with only one subject, verb, and object; however, most of the sentences in ITB documents are complex with a parallel structure. Given such a complex sentence, if it is not divided into several simple sentences with one verb and one object, the accuracy of SVO-based grammatical analysis is lowered. To reduce the complexity, sentence segmentation was performed before syntactic analysis to separate and simplify complex sentences. Then the subject, verb, object, and modifiers were separated from the isolated simple sentences. Since it is unrealistic to investigate all grammars in syntactic analysis, this study was limited to finding the main elements of a sentence, such as the subject, verb, and object. Furthermore, this study determines modifiers and clauses that modify each element.

The sentence segmentation rules used various logic based on "if–then". Additionally, regular expressions are utilized to separate with delimiters such as (a), (b), (c), or (i), (ii), (iii): *\({0,1}[a-h]\) or *\((?!\))(?:m{0,4}(?:cm | cd | d?c{0,3})(?:xc | xl | l?x{0,3}) (?:ix | iv | v?i{0,3}))(?<!\()\).

For the data that have completed sentence segmentation, the SVO is separated to identify the related entities in the sentence through dependency parsing. Afterwards, all elements of the sentence are defined to be assigned to at least one of the SVOs. In particular, the principal and conditional clause SVO were extracted for sentences containing conditional clauses such as if and unless, which are often used in contracts. If the main clause and the conditional clause were analyzed without separation, it will result in an error occurring when tagging the verb and object of the if clause.

In this paper, the Fail-Safe clause related to liquidated damages (LD) in the EPC contract is set as a proof of concept (POC) for applying the SA model. LD is one of the most severe risks in EPC contracts because it is a crucial contractual clause in which the EPC contractor promises to compensate the owner for losses if the contractually promised delivery date or performance is not met [1]. Among them, the Fail-Safe clause is the owner's two-tier backstop if the LD fails to work correctly, and it is one of the extremely dangerous clauses for the contractor. What makes the Fail-Safe clause more severe for the contractor is that the word Fail-Safe does not exist in the ITB; therefore, it cannot be found through a simple keyword search. This paper focused on extracting risk clauses that could not be found through keyword search and used the Fail-Safe clause as an example. In addition, it was applied as a case to the lexicon development in Chapter 5.3 and the semantic rule development in Chapter 5.4. The following is an example of the original sentence of Fail-Safe shown in the actual ITB and the rules used to separate the sentences.

[Original sentence]

*If liquidated damages are found not to be payable or the Articles in this Contract in relation to liquidated damages are found to be invalid or unenforceable for any reason, then the Parties agree that Contractor's liability to Company will instead be for general damages at law for Contractor's failure to comply with the relevant obligation.*

[Sentence segmentation rules]

If

If-clause [If <subject> + (<MD>) + <verb> + <object>] + Main-clause [<subject 1> + (i) + <verb 1> + <object 1>]

Then,

<subject> + (<MD>) + <verb> + <object>, <subject 1> + <verb 1> + <object 1>

When the syntactic analysis is completed, each contract sentence is divided into key elements such as subject, verb, object, and modifiers. Figure 5 shows the results of the original PoC sentences and extracted SVO tuples separated by syntactic analysis.
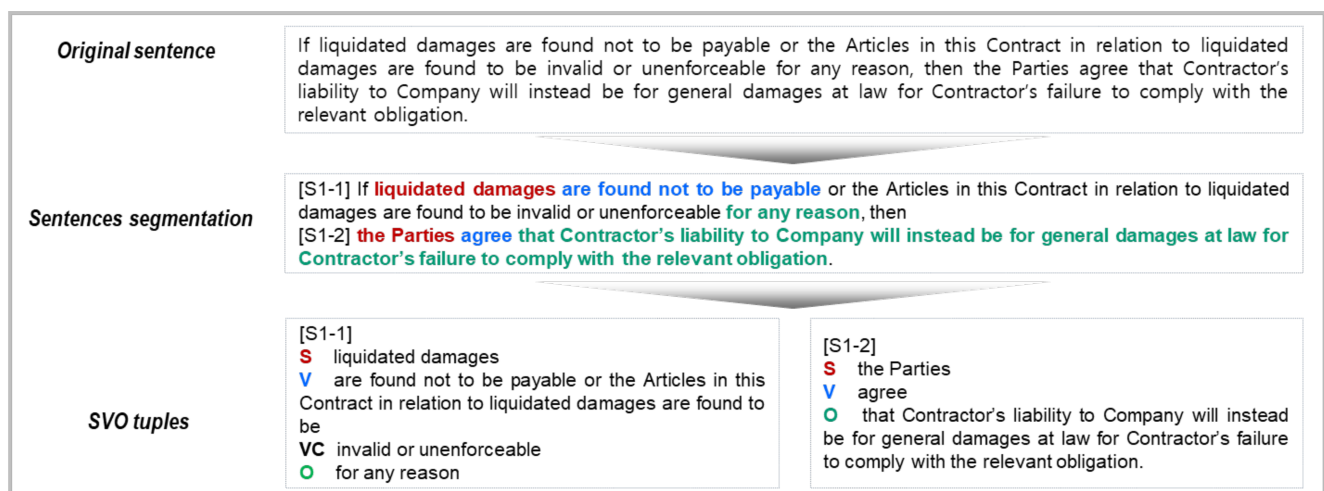
| | |
|---|---|
| ***Original sentence*** | If liquidated damages are found not to be payable or the Articles in this Contract in relation to liquidated damages are found to be invalid or unenforceable for any reason, then the Parties agree that Contractor's liability to Company will instead be for general damages at law for Contractor's failure to comply with the relevant obligation. |
| ***Sentences segmentation*** | [S1-1] If **liquidated damages** **are found not to be payable** or the Articles in this Contract in relation to liquidated damages are found to be invalid or unenforceable **for any reason**, then<br>[S1-2] **the Parties** **agree** **that Contractor's liability to Company will instead be for general damages at law for Contractor's failure to comply with the relevant obligation.** |
| ***SVO tuples*** | **[S1-1]**<br>**S** liquidated damages<br>**V** are found not to be payable or the Articles in this Contract in relation to liquidated damages are found to be<br>**VC** invalid or unenforceable<br>**O** for any reason<br><br>**[S1-2]**<br>**S** the Parties<br>**V** agree<br>**O** that Contractor's liability to Company will instead be for general damages at law for Contractor's failure to comply with the relevant obligation. |

**Figure 5.** The result of sentence segmentation and SVO tuples applied for syntactic analysis.

### 5.3. Ontology-Based EPC Contract Lexicon

5.3.1. EPC Contract Taxonomy

Taxonomy is a classification system that classifies structural relationships between concepts used in a specific domain according to a hierarchical structure [53]. Research on taxonomy and lexicon for ontology-based construction contracts have only been studied until the taxonomy development [53]. Alternatively, the glossary was defined by adding the vocabulary used to the existing lexical dictionary [25]. These prior studies on construction contracts are not taxonomy through professional analysis. Therefore, when constructing a lexicon, the vocabulary is limited due to the lack of lexical analysis by domain. To overcome this limitation, we first organized the EPC contract taxonomy through a workshop targeting subject matter experts (SMEs) with 10 to 20 years of experience in the EPC field. The EPC contract taxonomy in this study was classified into seven categories beneath Class 1 (Figure 6).
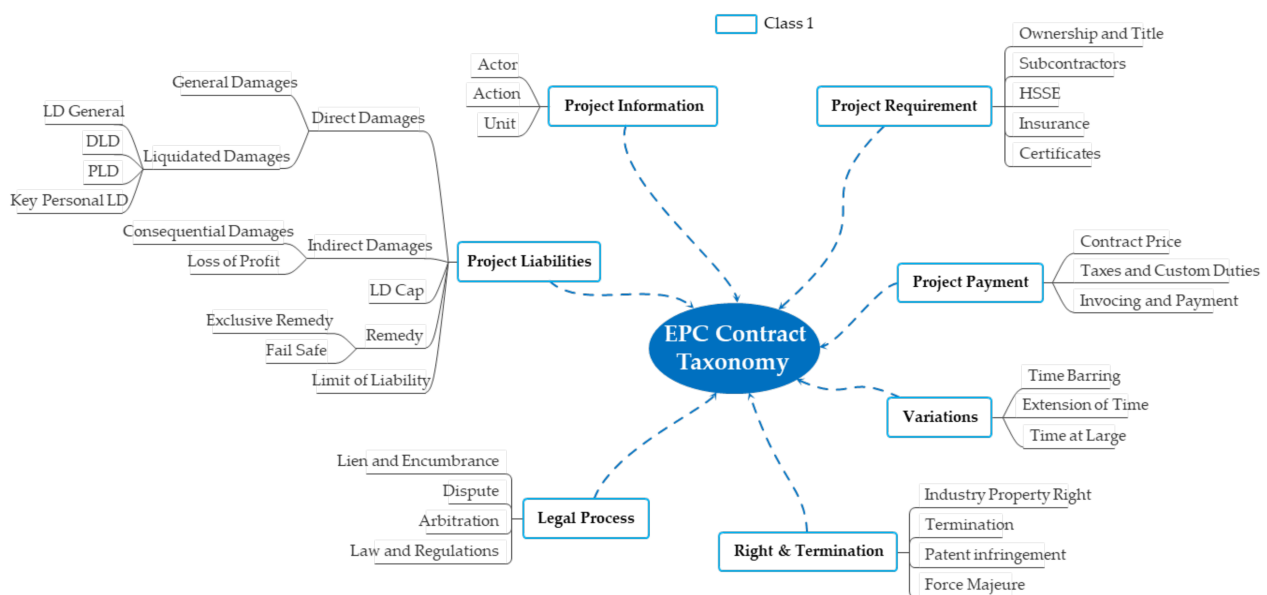
**Figure 6.** Class 1 of the EPC contract taxonomy for the EPC contract lexicon.

Seven SMEs participated in the workshop, consisting of EPC executives, academia, and EPC lawyers. These SMEs also participated in developing a gold standard for verification

of the SA model. Details for SMEs are described in the model test. A total of seven groups were set as Class 1, including:

1.  Project information including the subject and general matters of the contract.
2.  Project requirement for contractual requirements.
3.  Project liabilities for contractor's liability and damage compensation.
4.  Project payment for progress.
5.  Variations including construction changes.
6.  Project rights and termination of contracts.
7.  Legal process, which includes disputes between contractors and owners.

5.3.2. Development of EPC Contract Lexicon

For the EPC contract taxonomy, Class 4 was set as the lowest level. A lexical dictionary was developed among them by finding a synonym of the word corresponding to the lowest level. Then, the lexicon was defined by connecting the lowest level to the lexical dictionary. The EPC contract lexicon consists of 32 items and 79 details in this study. Numeric terms are excluded from this lexicon because their expression is infinite.

The EPC contract lexicon developed from this study is to be optimized for NLP application through the analyzation of the semantic correlation between contract phrases. EPC contract experts verified the developed lexicon. As a result, compared to the existing construction glossary, a highly detailed lexicon was designed and supported a more accurate ITB risk extraction model. In particular, it is significant as it was the first to systematize the entire EPC contract for AI applications. Table 3 is a part of the EPC contract lexicon that was newly developed through this study. It shows the terms related to Fail-Safe set by PoC which belong within the taxonomy.

**Table 3.** An example of EPC contract lexicon based on the EPC contract taxonomy liquidated damages and fail-safe.

| Class 1 | Class 2 | Class 3 | Class 4 | Terms |
|---|---|---|---|---|
| | | | Liquidated Damages General | Liquidated damages, LD, reasonable and genuine pre-estimate of loss, damages for any loss, pre-estimate of loss, not a penalty, not as a penalty, not meet, fail to complete, not ready for delivery |
| | Direct Damages | Liquidated Damages | Delay Liquidated Damages | Delay liquidated damages, DLD, liquidated damages for delay, delay damages, liquidated damages for such delay, liquidated damages for any delay, liquidated damages for late completion |
| Project Liabilities | | | Performance Liquidated Damages | Performance liquidated damages, performance of the plant, PLD, liquidated damages for insufficient failure of the plant to achieve the performance tests, performance liquidated damages, damages for failure to pass tests on completion |
| | Remedy | Exclusive Remedy | | Exclusive remedy, sole and exclusive remedy, sole and exclusive financial remedy, remedy, remedies, obligation to complete the work |
| | | Fail Safe | | Invalid, unenforceable, validity, enforceability, no challenge, remit, refund, reimburse |

*5.4. Semantic Analysis Modeling Based on Rules*

Semantic analysis is a method to understand the meaning of language. It is also used to grasp the contextual meaning and morphological elements of a natural language [42]. The SA model is an approach that maps SVO tuples extracted from syntactic analysis and

the lexicon. Syntactic analysis and EPC contract lexicon alone have limitations in extracting risk sentences; therefore, a semantic rule was developed. The semantic rule sets the subject, verb, and object elements to match the lexicon class. Then, it defines the class using various aspects of domain ontology (i.e., concepts and attributes). Applying these rules makes it possible to grasp the sentence's contextual meaning and identify the correlated risks. The SA model was developed in the following two steps.

First, we developed semantic grammars to match SVO tuples and contract lexicons.

A semantic grammar is required to make a semantic rule by matching each sentence's SVO tuples and lexicons. Additionally, the semantic grammar in this study is a standard for applying lexicon and corresponds to the SVO when developing semantic rules.

Semantic grammar uses the following regular expressions:

- <Class name>: True if there is a term included in the class of lexicon (e.g., <Liquidated Damages>).
- ( ): Give preference to operations in parentheses (e.g., (<employer-side> or <contractor-side>) and "contract").
- or: True if any of the preceding and following elements is True (e.g., <contractor-side> or <both party>).
- and: True only when both preceding and following elements are True (e.g., <Liquidated Damages> and <Exclusive Remedy>).
- <POS: XXX>: True if there is a word corresponding to the part-of-speech of XXX (e.g., <POS: PRON>).

Next, we developed a semantic rule that matches the lexicon's subject, verb, and object by applying the regular expression of the semantic grammar. The general form of a semantic rule is as follows: *IF (Precondition) THEN (Postcondition).*

The precondition (v) indicates the conditions for rule application. The specific class name corresponding to the subject, verb, and object (SVO) is designated to determine the risk. The precondition of the rule is based on the comparison of specific terms specified by SVO and lexicon. As a result, the postcondition extracts only cases where the SVO and the corresponding term match the class to which it belongs. The following are examples of various semantic rules for Fail-Safe sentence extraction.

- Example 1 [S1-1]:
  IF,
  Subject == <Liquidated Damages>.
  Verb == <Liquidated Damages> or ("not" and <Legal-action>).
  Object == <Fail-Safe> or <General Damages>.
  THEN,
  "Fail-Safe" clause is extracted.

In the above Example 1 matching rule, the subject is a term included in the liquidated damages class. While the verb consists of the liquidated damages class and legal-action class. Furthermore, the word "not" is included as well in the sentence. If the Fail-Safe class or *General Damages* class is included in the object; then the rule is to extract the sentence as a Fail-Safe clause. Here, subject, verb, and object do not necessarily all exist in the semantic rule. The rule is established even if only some of the SVOs are present.

- Example 2 [S1-1]:
  IF,
  Subject == <Liquidated Damages>.
  Verb == <Liquidated Damages> or ("not" and <Legal-action>) or <Legal-action>).
  Object == <Fail Safe> or <General Damages>.
  THEN,
  "Fail-Safe" clause is extracted.

- Example 3 [S1-2]:
  IF,
  Subject == <Both-party>.

Verb == <Liquidated Damages> or <Legal-action>.
Object == <Fail-Safe> or <General Damages>.
THEN,
"Fail-Safe" clause is extracted.

Each sentence may have simple or complex semantic matching rules. Based on Examples 2 and 3 above, even for one sentence, one or several semantic matching rules can be defined according to SVO tuples elements' diversity. We developed 79 lexicon classes and 87 corresponding semantic rules through this study.

### 5.5. Risk Clauses Extraction

For each sentence of ITB, the lexicon item corresponding to the SVO tuples of the semantic rule is checked. If the lexicon item in each SVO tuples element is not registered in the lexicon, then the semantic rule is not applied. The result of checking the lexicon for each SVO appears as true/false. When the detailed items of each subject, verb, and object in the semantic rule are all true, the entire rule appears as true; otherwise, it is false. Taking the above Fail-Safe rule as an example, if the subject has a vocabulary belonging to liquidated damages class, then the verb will also contain a term belonging to liquidated damages. If a sentence with words belonging to Fail-Safe class is entered in the object, the result will appear as true. The semantic rules are sequentially applied to all sentences during document analysis. If a sentence matches the practice, it is recognized as a risk clause. Therefore, it requires extraction and needs to be set to the following rule.

### 5.6. Deontic Classification

Deontology is a theory about rights and obligations. In addition, deontic logic (DL) is a branch of modal logic dealing with obligations, permissions, and prohibitions [54]. This section is focused on formalization by applying DL to the risk clauses of the extracted EPC contract, and the extracted contract clauses are divided into O (obligation), P (permission), and F (prohibition/forbidden). An off-the-shelf deontic logic reasoner could not be used; thus, this study utilized DL for the logical formulation of contract risk clauses.

The deontic model helps distinguish between legal contracts and contractual clauses by evaluating whether a particular action or condition is correct, incorrect, permissible, or prohibited [54]. Furthermore, the DL uses deontic operators such as O, P, and F to indicate whether a subject complies with [19].

A deontic rule of this study was classified into four modalities and was matched with the class of SVO tuples and EPC contract lexicon.

- Agent: the accountable agent, corresponding to the actor in the lexicon (e.g., contractor, owner).
- Predicate: represent concepts, relations between objects, corresponding to action in the lexicon (e.g., legal-action, obligated-action, permitted-action, payable-action).
- Topic: the topic it addresses, corresponding to class level 2 in the lexicon (e.g., safety, environment, cost, quality).
- Object: the object it applies to, corresponding to the class level 3 in the lexicon.

A deontic classification (DC) was used to formalized DL statements based on these deontic rules and classify them into three types: O, P, and F.

The basic notation of DL is to express the normative form, and it was used in this study for the logical formulation of contract sentences [55]. Logical formalization means logical expression through deontic logic. Additionally, DL statements consist of predicates or functions combined using two types of operators [19]. The DL formalization of this paper is expressed using two types of operators: deontic operators and first-order logic (FOL) operators [56]. For example, $P\alpha$ in deontic operators means that $\alpha$ is a member of permitted actions. The detailed descriptions of the two types of operators and their representations are given in Table 4.

**Table 4.** Deontic representations with the corresponding descriptions.

| Operator Type | Deontic Representation | Descriptions | Examples |
|---|---|---|---|
| Deontic Operators [1] | $O$ | Obligation | '$O\alpha$' means $\alpha$ is obligated |
| | $P$ | Permission | '$P\alpha$' means $\alpha$ is permitted |
| | $F$ | Forbidden/Prohibition | '$F\alpha$' means $\alpha$ is forbidden |
| | $I$ | Indifferent | '$I\alpha$' means $\alpha$ is indifferent |
| First-Order Logic Operators [2] | $\wedge$ | Conjunction | 'A $\wedge$ B' means A is true and B is true |
| | $\vee$ | Disjunction | 'A $\vee$ B' means A is true or B is true |
| | $\neg$ | Negation | '$\neg$A' means A is not true |
| | $\supset \cap$ | Implication | 'A $\supset \cap$ B' means A implies B (if A is true then B is true) |

[1] McNamara 2022; [2] Salama and El-Gohary 2013.

In the universe of discourse, true (T) and false (F) can be discriminated according to the range of a variable. DL statements specify the variables range by using quantifiers (i.e., $\forall$ and $\exists$) [19]. A universal quantifier ($\forall$ or for all) asserts true (T) only if all instances of a variable are valid, whereas an existential quantifier ($\exists$ or exists) evaluates to true (T) if at least one instance of the variable satisfies true. In addition, logical equivalence is denoted by $\equiv$. The DL formal representation expressed is based on the quantifier, and two types of operators for the S1-1 Fail-Safe clause are as follows.

- Deontic formalization for [S1-1]:
  $\forall x, y, z, h$ (Liquidated Damages $(x)$
  $\wedge$ Both-party $(y)$
  $\wedge$ Permitted-action $(z)$
  $\wedge$ General Damages $(h) \supset \cap$ P(Fail Safe $(z, h)$)

Deontic formalization supports the reasoning of contract sentences; however, a systematic deontic reasoner has not yet been developed [55,57]. As an alternative to a deontic reasoner for classification applying DL, this study proposes a DC approach that converts DL statements into SPARQL queries, then classifies them into O, P, and F. SPARQL queries are suitable for extracting risk clauses in contracts because of their capabilities of semantic understanding and contextual building information in a knowledge base [58]. The DC approach proposed in this study formalizes the previously developed semantic rules into DL statements and then converts them into SPARQL queries. In the future, research that enables direct DL reasoning through an FOL-based reasoner without additional conversion to SPARQL is required. This study classifies by limiting it to only the obligation sentence. The following is the expression of SPARQL queries for the S1-1 Fail-Safe clause.

- SPARQL queries for [S1-1]:
  SELECT ?x ?y ?z ?h
  WHERE {
  ?x a subject:Liquidated Damages. ?y a subject:Both-party. ?z a predicate:Legal-action. ?h a object: General Damages. FILTER (?z= obligated-action).
  }

The SA model proposed uses Python for model implementation, the spaCy library for risk extraction, and SPARQL queries for deontic classification.

*5.7. Implementation and Validation of the SA Model*

This study tested using the actual ITB to verify the performance of the SA model.

5.7.1. Gold Standard and Test Dataset for the SA Model

The extraction accuracy of the SA model was verified by comparing the extraction results of the automatic extraction model with the "Gold Standard" [42]. The Gold Standard is a compilation of target information in text sources by experts such as SMEs. It evaluates a model's performance and extracts information intuitively constructed by humans [25]. The Gold Standard for this study is a standard manually developed by a group of experts, not the semiautomatic gold standard of the study by Zhang and El-Gohary [42]. The Gold Standard was developed through a workshop with a group of SMEs. The SMEs group consists of seven people with more than 10 to 30 years of experience in the EPC plant field and are experts in EPC contracts, including EPC executives, academia, and EPC lawyers. The SMEs analyzed the presence or absence of risks for each contract clause, an affiliation of risk, and the degree of risk impact through workshops. The taxonomy classification and the class of the EPC contract lexicon were also set in workshops. The developed Gold Standard was used to verify the performance of the model by comparing it with the automatic extraction results of the SA model. Table 5 below summarizes information on SMEs that participated in developing the gold standard in this study.

**Table 5.** Information on SMEs that participated in the Gold Standard development.

| Expert Code | Category | Discipline | Year of Experiences | Affiliation |
|---|---|---|---|---|
| A | Offshore | Contract | 32 | EPC company |
| B | Offshore | Planning | 16 | EPC company |
| C | Offshore | Engineering | 18 | EPC company |
| D | Onshore/Power plant | PM | 17 | EPC company |
| E | Offshore/Onshore | Contract | 28 | Law firm |
| F | Onshore/Infra | PM&IT | 22 | Academia |
| G | Onshore/Infra | PM&IT | 24 | Academia |

For the test of the SA model, four offshore plant contracts, a total of 7119 records (sentences), were selected as the test dataset and used for model evaluation. The dataset information for model testing is shown in Table 6.

**Table 6.** Dataset information for testing risk extraction performance of the SA model.

| Dataset No. | Project Name | Domain | Owner | No. of Records |
|---|---|---|---|---|
| 1 | 'I' project | Offshore FPSO | I & T companies consortium | 1864 |
| 2 | 'C' project | Offshore FPSO | T company | 1894 |
| 3 | 'M' project | Offshore TLP | T company | 1371 |
| 4 | 'P' project | Offshore FLNG | P company | 1990 |
| | Total | | No. of Records | 7119 |

5.7.2. Test Results and Validation of the SA Model

The SA model was evaluated using the test data, and the results of the model were compared with the Gold Standard for the model validation. To verify the performance of the risk clause extraction of the SA model, precision, recall, and F-measure indicators were applied [24,27,42]. A confusion matrix is used to calculate precision, recall, and F-measure. Table 7 shows the confusion matrix and the four variables that makeup it: true positive (TP), false positive (FP), false negative (FN), and true negative (TN).

**Table 7.** An evaluation criteria (confusion matrix) for the SA model.

| Criteria | Type | The Results of | The Gold Standard |
|---|---|---|---|
| | | Positive | Negative |
| The actual | Extracted | True positive (TP) | False negative (FN) |
| extraction results | Not extracted | False positive (FP) | True negative (TN) |

True positive (TP) means that the risk results of both the gold standard and the SA model are equally extracted as risks in the confusion matrix. False positive (FP) indicates that a clause that is not an actual risk is extracted as a risk, and false negative (FN) shows a result that the model did not extract, even though it is a risk clause. True negative (TN) is a case in which the machine does not extract non-risk clauses and identifies irrelevant information. TP and TN mean ground truth values, while FP and FN indicate errors due to incorrect extraction. Utilizing the information from the confusion matrix, the equations for precision, recall, and F-measure are as follows. $Precision_E$ is the ratio of the values extracted correctly divided by the risk from the results extracted through the model. $Recall_E$ is the ratio of the values correctly extracted by the model divided by the total positive results. Lastly, $F\text{-}measure_E$ is defined as the harmonic mean between precision and recall. The given equations are equal to (1), (2), and (3), respectively [24,27,42]:

$$Precision_E = \frac{TP}{(TP + FP)} \times 100\% \tag{1}$$

$$Recall_E = \frac{TP}{(TP + FN)} \times 100\% \tag{2}$$

$$F - measure_E = 2 \times \frac{(Precision \times Recall)}{(Precision + Recall)} \times 100\% \tag{3}$$

High precision implies that the extraction results of the SA model are highly correlated with the Gold Standard. High recall illustrates that risk clauses defined by experts in the Gold Standard are relatively extracted the same as from the model. Table 8 summarizes the validation results for the SA model.

**Table 8.** Validation results of the SA model for risk clause extraction.

| Risk Extraction | No. of | Extractions | | | Performance | | |
|---|---|---|---|---|---|---|---|
| | TP | FP | FN | TN | Precision (%) | Recall (%) | F-measure (%) |
| | 765 | 98 | 142 | 6444 | 88.6 | 84.3 | 86.4 |

The F-measure of the SA model was 86.4 percent, which showed higher risk extraction performance than the previous study [25]. Of the 7119 sentences, 765 were extracted as risk sentences (TP). Although 142 is a value extracted as a risk, in actuality it is not a risk (FP). However, a value of 98 is a risk sentence even though the sentence is not extracted (FN). The result of extracting non-risk sentences was 6444 (TN). In comparison to the 84.3 percent recall, the precision of 88.6 percent can be interpreted as extracting most of the risks specified by experts, but not extracting approximately 4 percent.

A high precision value indicates that the results extracted from the model are highly correlated with sentences confirmed as the risk of the gold standard. A high recall value means that most of the risk clauses of the gold standard are extracted the same way in the model. In the results of Table 8, FP and FN show errors in risk clauses extraction. There were several types of extraction and causes of FN errors in the SA model. When adding a completely different new sentence type not defined by the semantic rule into the test, the model will not be able to extract the sentence. Furthermore, when developing a rule for risk

clauses within the defined data, it is analyzed that the extraction accuracy for the existing type of sentence is high. However, there are still difficulties in extracting a completely different new kind of sentence. A case in which risk is incorrectly extracted even though it is not a risk is called an FP error. When reducing the FP error, the FN error increases, and conversely, when lowering the FN error, the FP error increases. Number-related terms were excluded from the lexicon when developing the EPC contract lexicon. This is the main cause of the lower recall value, and as a result, it is analyzed as a factor affecting the F-measure.

The SA model proposed a more accurate risk extraction model based on the lexicon defined by the EPC contract taxonomy to which the ontology concept was applied as the knowledge base (KB). In particular, the contract clause extraction mechanism of the SA model provided accurate and advanced performance compared to the poisonous clause's detection model proposed by Lee et al. [25]. Compared with the previous study of Lee et al., only 1.6% of toxin sentences were extracted out of 708 sentences, resulting in a risk clause extraction rate of about 12% and a relatively high F-measure. It is interpreted that the cause is due to detailed lexicon development and improvement of syntax analysis accuracy. In many NLP studies, the parsing performance deteriorated as the complexity of the sentence structure increased. Specifically, when the sentence structure is complex, such as in an EPC contract, there are more errors in parsing than in typical sentences. Incorrect parsing affects subsequent steps in applying semantic rules, which can directly affect the performance of risk clause extraction. Syntax analysis of this study using the spaCy library shows a higher extraction rate and accuracy than the previous study by Lee et al. due to the improvement of spaCy's performance. However, parsing of complex sentences still shows somewhat poor results.

EPC ITB's automatic risk clause extraction model is a challenging attempt to change the contract review process performed manually by humans to an AI-based automatic extraction process. Therefore, continuous improvement is required.

This study measured the performance of deontic classification in terms of accuracy, and *Accuracy*$_{DC}$ was calculated using Equation (4) [59].

$$Accuracy_{DC} = \frac{TP + TN}{TP + TN + FP + FN} \times 100\% \tag{4}$$

DC's evaluation data were for one FPSO project' ITB, and the total number of records was 1864. Table 9 displays the test results for classifying only DC's obligation statement.

**Table 9.** Validation results of the SA model for deontic classification (Class_O).

| | **The** | **Classified** | **Results** | **(Class_O)** | **Performance** | | | |
|---|---|---|---|---|---|---|---|---|
| Deontic Classification | TP | FP | FN | TN | Precision (%) | Recall (%) | F-measure (%) | Accuracy (%) |
| | 572 | 57 | 198 | 1037 | 90.9 | 74.3 | 81.8 | 86.3 |

Concerning the performance of DC, we achieved an accuracy of 86.3 percent. These results showed the effectiveness of our approach in interpreting and classifying contract clauses. The relatively low recall value (74.3 percent) was analyzed due to errors and omissions in DL formalization. Another cause of the error is the existence of a syntactic or grammatical factor that does not correspond to the four modal components of the deontic rule in the contract sentence. In essence, recall errors may appear due to the omission of extraction.

The SA model improves accuracy by describing linguistic phenomena through rules; however, it is impossible to extract all the numerous linguistic phenomena from the usage of rules. Both the risk clause extraction function and the deontic classification function have the problem of generating a rule for every sentence in the contract. Since language

has many exceptional aspects, it is impossible to cover all linguistic phenomena using a simple rule-based approach. In addition, it is challenging to express the entire sentence of the contract document as a rule. As a result, it is expected that extraction cannot occur when a sentence that is not in the rule is inputted.

## 6. Risk Level Ranking Model

### 6.1. Preprocessing for RLR Model

ITB sentences that were text data preprocessed before using the RLR model were converted through the usage of the PDF structuralization module. Preprocessing was embedded through word embedding after stop-word removal. Word embedding is a vectorization process that maps a word to a specific R-dimensional vector in NLP analysis [20]. Embedding plays a vital role in NLP because high-quality embeddings increase the document classification accuracy and learning speed [60].

In this study, the text data preprocessing was performed by applying the Keras framework in the Python library. In addition, the bi-LSTM model, a deep learning model [61], was implemented and the Keras tokenizer was utilized for word embedding [62]. Keras is a deep learning framework for Python that makes it convenient to develop and train almost any deep learning model [63]. Backend engines such as Theano, Google's Tensorflow, and Microsoft's CNTK can be seamlessly integrated with Keras [63]. In this study, Keras is designed to operate on Tensorflow. An index was assigned to the Keras tokenizer in the order of the word frequencies. For example, "1" was assigned to the out-of-vocabulary (OOV). However, this specific vocabulary is not assigned an index. As a preprocessing for inputting sentences into the artificial neural network, it is converted into an integer sequence that lists the indexes assigned to the words of each sentence. Furthermore, Keras embedding is a function that converts an integer sequence into a dense vector of a fixed size by inputting a vector that has completed one-to-one correspondence with a word through the Keras tokenizer [61].

### 6.2. Risk Level Ranking Modeling with Bi-LSTM

The RLR model is a model developed by applying the bi-LSTM algorithm. The model performs the first classification of ITB sentences according to the presence of risk and true/false (T/F). Then, it divides the sentences classified as true into five levels. The first T/F classification is a binary classification that divides the preprocessed ITB sentences into risky sentences and non-risky sentences. The second classification, risk level, is a multi-class classification that categorizes the true sentences ranked as a risk in the first model into five risk levels. Figure 7 showcases the process of the RLR model that classifies the risk level of sentences by applying the bi-LSTM architecture.

As shown in Figure 7, the trained model for first classification classifies T/F when inputted to an RLR model that has completed PDF structuralization. Sentences classified as true are categorized into five risk levels through a trained secondary classification model. For risk level, the five-point Likert scale was applied to increase statistical reliability [64]. The five-point scale applied to the training dataset in this study is "very high, high, moderate, low, and very low", with very high being level 5 and very low being level 1. In summary, the RLR model is learned through a training dataset and sequentially performs first and second classification.

The bi-LSTM algorithm used in the RLR model includes two LSTM layers (i.e., forward and backward), and the output sequences of the two layers are combined using a concatenating function ($\sigma$) as depicted in Figure 8 [65]. The bi-LSTM uses a forward LSTM layer that sequentially reads the sentence from the word on the left, then uses a backward LSTM layer that reads from the right in reverse order to consider its context. As a result, it shows better performance than RNN or LSTM by minimizing the loss of the output value and learning all parameters simultaneously. Traditional machine learning algorithms, such as neural networks, support vector regression (SVR), and ensembles, cannot handle sequential data or hierarchical representation learning of time series; thus, using bi-LSTM is

more desirable. Moreover, the bi-LSTM approach is suitable for sequential data modeling, based on previous observations, which performs the learning process in both directions [66]. The advantage of the bi-LSTM is that its performance does not deteriorate even if the data length is long. However, the disadvantages of bi-LSTM are high cost, because it has double LSTM cells, and that it is not suitable for speech recognition.
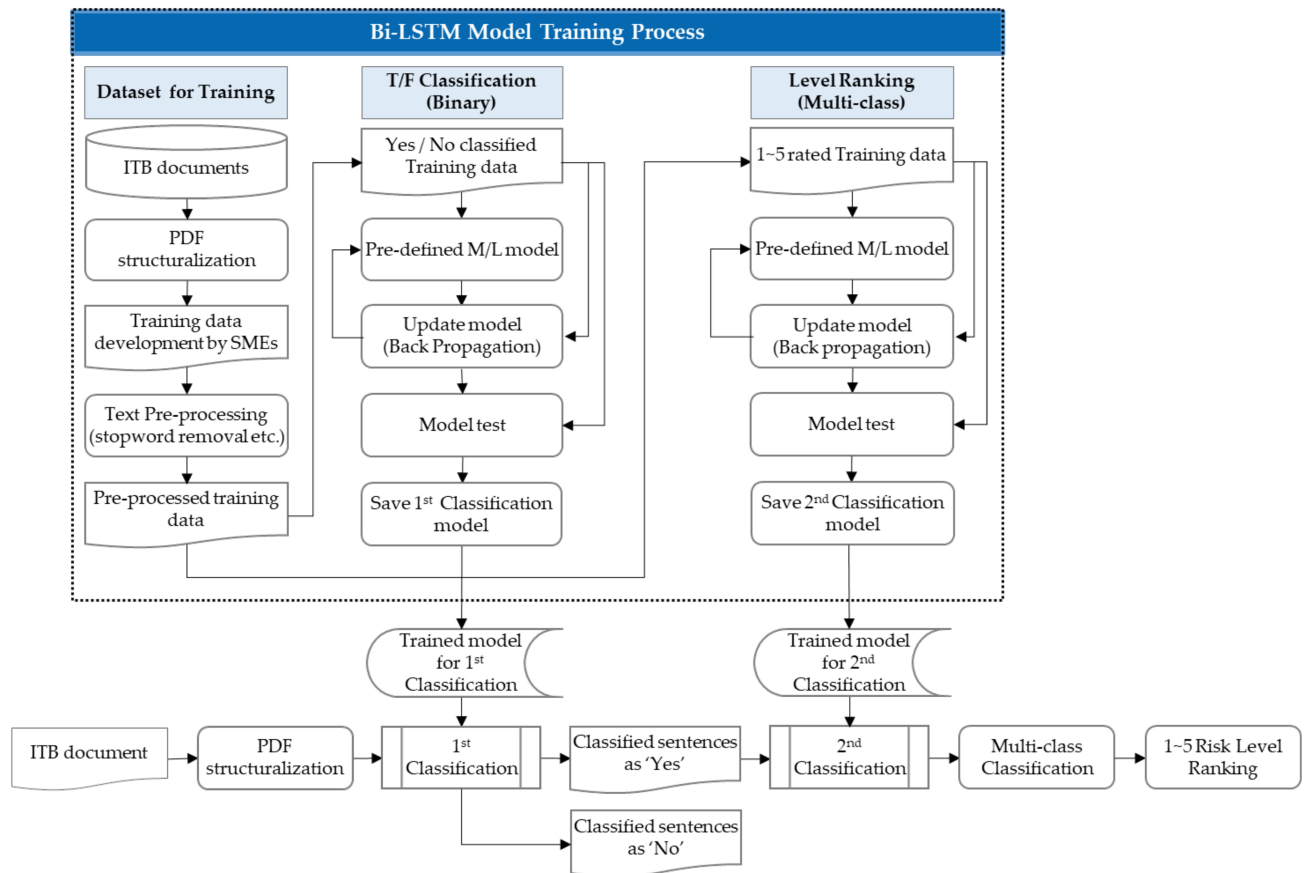


**Figure 7.** The RLR model with bi-LSTM algorithm and implementation process.
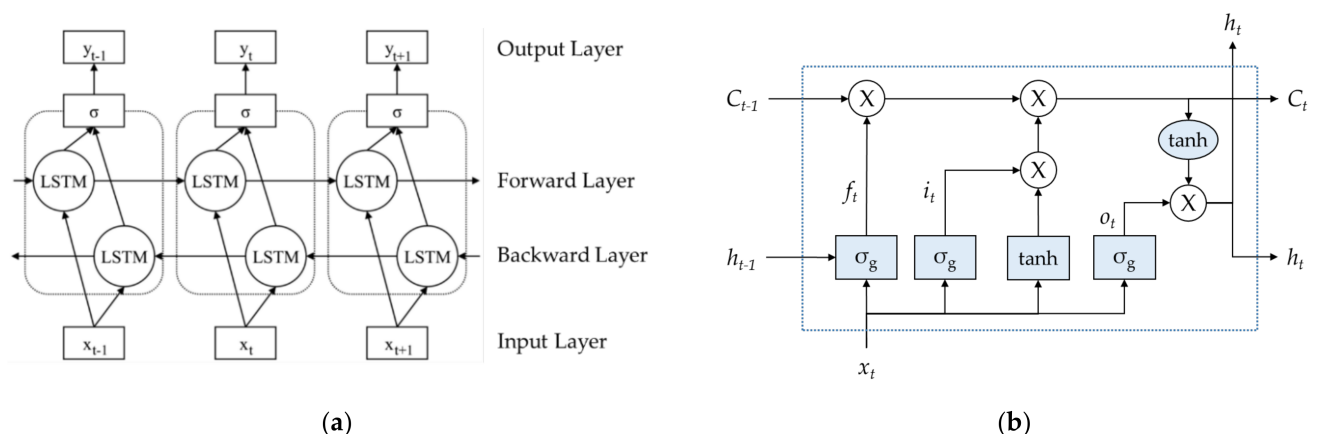


(**a**)                    (**b**)

**Figure 8.** The architecture of bi-LSTM: (**a**) Unfolded architecture of bi-LSTM with three consecutive steps; (**b**) the detailed structure within an LSTM cell. Source: Cui et al. 2018; Moon et al., 2021; modified by the authors.

Looking at the detailed structure of the LSTM cell, it consists of three gates (forget gate, input gate, and output gate) and two memory states (cell state and hidden state).

If it is time $t$, then $x_t$ of the LSTM unit is the input vector and $h_t$ is the layer output. Additionally, $f_t$ denotes a forget gate, $i_t$ represents an input gate, and $o_t$ represents an output gate (Figure 8b). In the cell state, the previous cell output is $C_{t-1}$, the new cell output is $C_t$, and the cell input is $\widetilde{C}_t$, where tanh is a value between $-1$ and 1. Furthermore, $C_t$, a memory cell, stores all necessary information from the past to time t. The $h_{t-1}$ is the hidden state, and $h_t$ is the next hidden state. The $W_s$ and $U_s$ are the weight matrices, and $b_s$ are bias vectors. Lastly, σg stands for the activation function, and X plays the role of opening and closing information, showing the gate mechanism of the LSTM [26,65] (Figure 8).

In the RLR model of this study, the data are inputted to the $t$-th cell, while $x_t$, is a vector in which the $t$-th word of a given sentence of ITB is embedded. The output $y_t$ of the $t$-th cell has the following values (Figure 8a):

1. T/F classification (the first classification): Classify whether the sub-sentence cut to the $t$-th word of the input sentence is a risk.

2. Risk level classification (the second classification): Classify where the risk level of the sub-sentence cut up to the $t$-th word of the input sentence belongs among the five steps. In this study, classification was performed up to $y_t$, the output of the last cell.

The RLR model was developed using Python and embedded as a component technology of the *EMAP* system [27]. The RLR model is meaningful in that it attempted an AI technique to analyze the risk impact of an EPC contract and developed a system for it.

### 6.3. Development of Training Dataset

The training dataset generation of the RLR model was performed in two ways: (1) the development of the risk analysis method for the ITB contract clause, and (2) the training dataset development by the expert group. Considering the complexity of the EPC contract, the adjustment index was added to the evaluation criteria for the risk assessment. As for the risk assessment method, a three-dimensional risk matrix with a new standard of coordination was added to the probability and impact (PI) matrix. In addition, a traditional two-axis evaluation method was applied [67]. The feature can remove the redundancy of risk ranking at the point where $P$ and $I$ meet; however, this is a disadvantage of the *PI* matrix [68]. Jang et al. showed that the main influence variables were concentrated on impact because of the *PI* 2-axis evaluation. However, according to the *PIC* 3-axis evaluation, it was shown that the main influence variables were evenly distributed in probability, impact, and coordination compared to 2-axis [68]. The equation for the *PIC* 3-axis evaluation is as follows (5).

$$\text{Rik Degree} = \sqrt{(P2 + I2 + C2)} \tag{5}$$

where

$P$: Risk probability, likelihood that a risk will occur;

$I$: Risk impact, impact on project objectives;

$C$: Coordination index.

Each contract clause of the ITB converted through the PDF structuralization module was evaluated on a five-point Likert scale according to its effect on $P$, $I$, and C. Figure 9 is a part of the evaluating results of the contract clause of the "I" FPSO by the PIC three-axis method.

An expert group conducted a non-face-to-face workshop to develop the training dataset. The experts consisted of seven SMEs, including EPC lawyers and contract practitioners. Each SME was asked to score risk probability and impact for each sentence in the ITB document. The analysis results were received by e-mail.

Figure 10 below is an example of a training dataset for the RLR model. The training dataset consists of 9520 records with ten integrated EPC ITB documents, and it was converted to DB in CSV format. Risk T/F means either true or false, depending on the presence of risk in the sentence, while risk degree is a classification of the level according to the degree of risk impact among the true sentences where risks exist. After ranking the risk

by the final risk degree score for each sentence through the PIC three-axis evaluation, it was converted into five levels. The risk degree level approaches 5 when the risk is high and approaches 1 when the risk is low (Figure 10).

| Section | | Contract Clauses | Probability | Impact | Coordination Index | Total Risk Degree |
|---|---|---|---|---|---|---|
| 13.1 | Acquaintance with the Contract Documents | (b) Contractor warrants that it will perform the Work in accordance with the Work Time Schedule, failing which Company has the remedies specified under the Contract, including those provided for under Article 51 and sub-Article 36.1 concerning defective performance by Contractor and liquidated damages for late completion. | 5 | 5 | 5 | 8.66 |
| 36.2 | Liquidated Damages for Insufficient Performance of the Plant | If Company decides to apply liquidated damages, in accordance with the provisions of sub-Article 19.2(f) due to a failure of the Plant to achieve the Performance Tests, Contractor will pay to Company liquidated damages in accordance with Exhibit B. | 4 | 5 | 5 | 8.12 |
| 54.3 | Termination of the Contract for Force Majeure | (a) If the event of Force Majeure continues for more than sixty (60) consecutive days or ninety (90) days in aggregate:<br>(i) Company may at any time terminate the Contract by giving notice to Contractor; | 1 | 4 | 4 | 5.74 |

**Figure 9.** An example of risk degree analysis of EPC contract applying *PIC* 3-axis method.

| Doc. code | Sentence in ITB | Risk T/F | Risk degree |
|---|---|---|---|
| 4 | 15 5. | F | - |
| 4 | Effect of Rescission: It is expressly understood and agreed by the parties that in any case, if the BUYER rescinds this Contract under this Article, the BUYER shall not be entitled to any liquidated damages, or any other recourse unless by means of the provisions of Article X hereof. | T | 4 |
| 4 | (End of Article) 16 ARTICLE IV - APPROVAL OF PLANS AND DRAWINGS AND INSPECTION DURING CONSTRUCTION 1. | F | - |
| 4 | Approval of Plans and Drawings: Approved plans and drawings of the BUILDER's HN. 1674 (including all amendments, additions, deletions and variations incorporated into the Specification up to the date of this Contract signing) shall be deemed to be approved by the BUYER and shall be applied to the Drillship. | T | 1 |
| 4 | The BUILDER shall be exempted from the approval of the BUYER for the plans and drawings in accordance with the Specifications. | T | 2 |

**Figure 10.** An example of the training dataset for the RLR model.

### 6.4. Fine Tuning for the Risk Level Ranking Model

Hyperparameters refer to the values set by the model user directly in the deep learning model and are used to control the training process [69]. Furthermore, hyperparameters are optimized to maximize the performance of deep learning models. Examples include values such as epoch, the number of nodes (neurons) inside the cell, learning amount, and learning rate that mainly determine how far to proceed. In addition, the performance of the deep learning model varies depending on the combination of these values. In this way, exploring the combination of hyperparameters to maximize model performance is called hyperparameter optimization [70]. Epoch indicates the number of passes in the entire training dataset that the algorithm has completed. When epoch = 1, it indicates that training has been completed once for the entire dataset [71]. It is possible to prevent underfitting and overfitting only by setting an appropriate epoch value when training the model. Overfitting is when a deep learning model learns the training data in too much detail. It occurs when the training data are insufficient, or the model is too complex for the characteristics of the data. Moreover, the result of this problem is that the general model's performance is reduced due to excessive adaptation to the training data. This results in excellent learning performance, but poor responsiveness to untrained data [72]. Keras supports the early stopping of training via a callback called EarlyStopping, even if the specified epoch is not filled [73].

This study performed early stopping after 14 epochs were categorized to secondary classification. A loss function refers to a function that calculates the error between the expected and current output of the algorithm [74]. The closer to the actual value, the smaller the value appears. Cross entropy is one of the methods to measure the difference between two probability distributions [75]. Binary cross entropy is used when only two label classes are in cross-entropy and performs best in equal data distribution among class scenarios, while categorical cross-entropy is used when two or more label classes exist

among cross-entropy. Binary cross-entropy is derived from the Bernoulli distribution, and categorical cross-entropy is derived from the multinoulli distribution [75]. In deep learning, optimization refers to the process of finding the extremum of a specific objective function [76]. When optimizing hyperparameters, the parameter that changes the model performance most dramatically and easily is the optimizer. Adaptive moment estimation (Adam) is an optimization algorithm that improves the accuracy of deep learning and, thus, was used as the optimizer in this study. The Adam function is a function that finds the minimum value for the objective function by applying optimization according to size [76]. The bi-LSTM model of this study uses the text sentence of the contract as input, the primary output is returned as T/F classification, and the second is returned as a five-level multi-class. Of the total 9520 records of 10 ITBs, 80 percent were used as training data, and the remaining 20 percent were used for the test. Table 10 shows the hyperparameters used for the RLR model.

**Table 10.** Hyperparameters of the RLR model using bi-LSTM.

| Type of Model | Hyperparameters | Value Determined |
|---|---|---|
| T/F Classification (Binary) | Epoch | 10 |
| | Early stopping | - |
| | Loss function | Binary cross entropy |
| | Optimizer | Adam |
| | Train data: Test data | 8:2 |
| Degree Ranking (Multi-class) | Epoch | 100 |
| | Early stopping | 14 epoch |
| | Loss function | Categorical cross entropy |
| | Optimizer | Adam |
| | Train data: Test data | 8:2 |

*6.5. Implementation and Validation of Risk Level Ranking Model*

We used 2380 test data records to evaluate the RLR model, and we verified the test results. Based on the testing of the bi-LSTM model, 1806 true values and 574 false values were classified in the first classification. As a result of the second classification, risk level 1 was classified into 489, 434 for level 2, 572 for level 3, 288 for level 4, and 23 for level 5. Table 11 illustrates the test results for the first and second classification of the RLR model.

**Table 11.** Test results for the first and second classification of the RLR model.

| Category | Type of Model | | | Test | Result | | |
|---|---|---|---|---|---|---|---|
| 1st | T/F Classification | T/F | | True | | False | |
| | (Binary) | No. of sentences | | 1806 | | 574 | |
| 2nd | Risk Level Classification | Risk level | 1 | 2 | 3 | 4 | 5 |
| | (Multi-class) | No. of sentences | 489 | 434 | 572 | 288 | 23 |

The performance test was implemented by inputting the actual EPC ITB. When a user inputs a new ITB, the level of risk is extracted for each sentence. Figure 11 showcases the result of inputting the actual ITB into the RLR model after data conversion in the PDF structuralization module. Furthermore, the figure displays how each sentence is classified according to the level of risk.
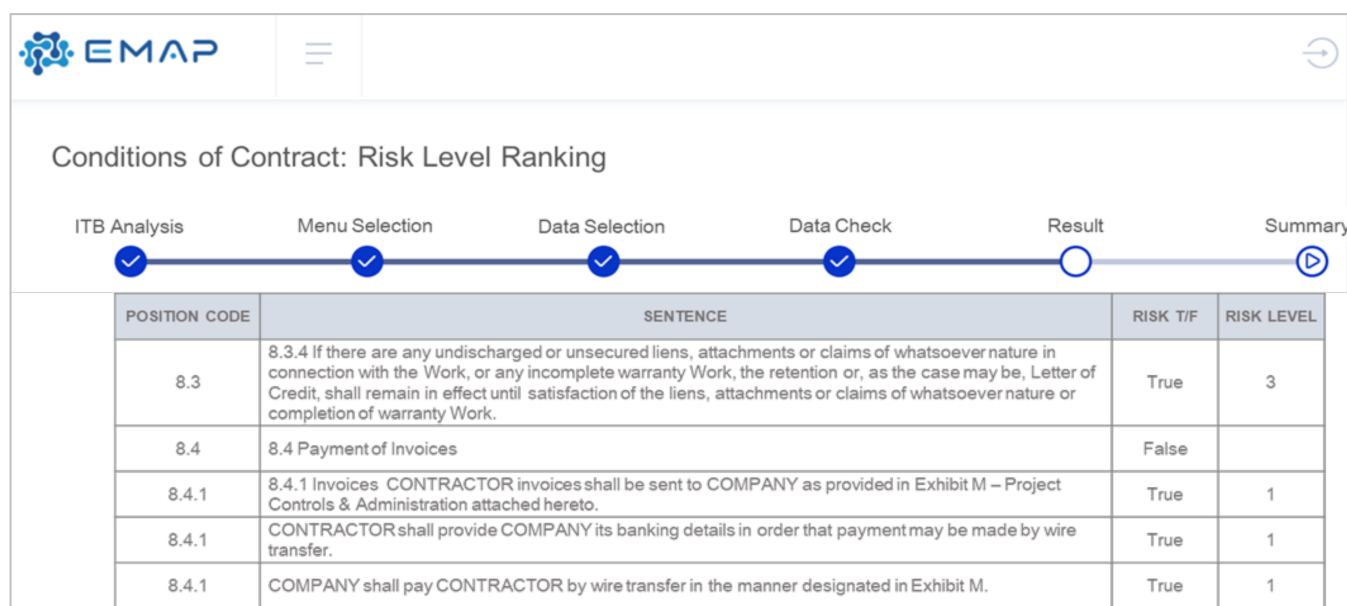
**Figure 11.** An example of analysis results of the RLR model in the *EMAP*.

The position code makes it easy to find the affiliation of a sentence, and it matches the table of contents number in the document. RISK T/F indicates the risk of each sentence separated based on the position code. If there is no risk, it is displayed as false, and the RISK LEVEL does not exist. The output result on the UI of *EMAP* can be downloaded as a CSV file, and it is set to show the summary of the risk level classification results in a bar graph. Accuracy was used as a performance evaluation index since the class ratio of the dataset used in the RLR model test is the same at 1:1. In addition, the accuracy formula is the same as the above Equation (4).

As a result of measuring the accuracy of the trained bi-LSTM model, the performance of the second classification model was lower than that of the first classification model. In particular, the second classification, multi-class classification, showed an accuracy rate of approximately 47 percent. Table 12 displays the validation results for the first and second classification of the RLR model.

**Table 12.** Validation results for the first and second classification of the RLR model.

| Category | Type of Model | Performance | |
|---|---|---|---|
| | | **Train Set** | **Test Set** |
| 1st | T/F Classification (Binary) | Loss: 0.141 Accuracy: 0.955 | Loss: 0.356 Accuracy: 0.882 |
| 2nd | Risk Level Classification (Multi-class) | Loss: 0.547 Accuracy: 0.888 | Loss: 2.522 Accuracy: 0.468 |

The rule-based approach generally tends to show better text processing performance than the train-based approach [42]. The low performance of multi-class classification, the second classification of the RLR model, can be attributed to the lack of a training dataset. Specifically, when a sentence that did not exist in the model training process occurred during the test, accuracy was lowered due to the untrained sentence. Moreover, the lack of table recognition of the ITB is also interpreted as the cause of the low performance. Documents of a specific domain, especially contract documents such as EPC ITB, have limitations from data collection. In addition, developing a training dataset can be difficult. For example, a task such as assigning a risk rating to a contract sentence by a legal expert may pose a challenge in terms of both accessibility and economic feasibility. Although this

study developed the training dataset together with SMEs specialized in the contract field, the performance evaluation result recognized the limitations of the lack of the training dataset. Nevertheless, when first- and second-year junior engineers performed ITB analysis, the RLR model derived analysis results similar to those of senior engineers with 8 to 9 years of experience. In essence, the RLR model is significant in that it was developed by simultaneously applying traditional risk assessment methods and AI techniques to the risk analysis of EPC contracts.

## 7. System Application on Cloud Platform

In this study, the EMAP system, a cloud-based engineering machine learning integrated analysis platform, was developed to support the decision-making of the EPC project [27]. Software development, system integration, and cloud computing are the core of EPMA implementation. The software development was integrated with various software such as text analysis tools, web application server (WAS), engineering machine learning platform, decision support system, and data open system. EMAP was developed using HTML, Cascading Style Sheet (CSS), JavaScript, and Tomcat applied to WAS [77]. The data entering interface was created for a local server through the open-source MySQL DataBase Management System (DBMS) to access the database [78]. The system-to-system linkage used the application programming interface (API) [79], and the data were linked in the form of JavaScript Object Notation (JSON) [80]. The detailed analysis module of the decision support system was developed using the Python programming language. In addition, user convenience was improved by applying cloud services to allow users to analyze on the web.

The *EMAP* system consists of five main modules.: (M1) ITB Analysis, (M2) Design Cost Estimation, (M3) Design Error Check, (M4) Change Order Forecast, and (M5) Equipment Predictive Maintenance [27]. Among them, the SA model (Section 5) and RLR model (Section 6) described in this paper are included in the (M1) ITB Analysis module. The risk clause analysis model and the risk frequency analysis model are discussed in another paper [28]. The user interface (UI) of the submenu composing the ITB analysis module is shown in Figure 12.



**Figure 12.** UI of the ITB analysis module from the *EMAP* system.

Each model is analyzed according to its purpose, and the results are the output. The results shown on the screen can be downloaded in Excel or CSV format. When selecting the UI menu shown in Figure 12, the CSV format document that has completed PDF structuralization requires uploading. Errors and bugs that occurred during system development were improved through corrections.

## 8. Conclusion and Future Works

*8.1. Summary and Contributions*

This study supports EPC contractor decision-making by analyzing ITB in the bidding stage of the EPC project and applying NLP and machine learning technology. For this purpose, two models were proposed. The SA model that extracts the critical risk clauses of ITB is a rule-based approach to which the ontology concept is applied. The RLR model classifies each ITB sentence into five levels, according to the degree of risk, through the classification model to which bi-LSTM is applied. These two models were developed in the following steps.

First, the contract documents of the EPC plant project were collected for this study. Twenty-one contract documents, including the completed onshore and offshore plant projects and FIDIC, were collected and used as primary data for model development. Second, a PDF structuralization module was developed and used for analysis to convert the text data from the original ITB into the digital data frame for NLP analysis, which was then used as input for SA and RLR models by applying PDF parser and OCR technology. Third, the converted data after the PDF structuralization were used for model development and training through data preprocessing. The SA model automatically extracted key contract clauses using NLP's rule-based approach. The RLR model was developed and trained based on deep learning technology by applying bi-LSTM with word embedding. Fourth, collaboration with SMEs in the EPC field was carried out from the early stage of model development to enhance the reliability of the model's performance. To develop an ontology-based EPC contract lexicon, the authors established a taxonomy of EPC contracts with experts such as EPC lawyers and developed a very detailed lexicon. For the RLR model, a training dataset was developed by seven through reviewing the risk impact rating of each sentence of the contract.

Each model was validated through a pilot test using an actual EPC contract that was not used in the model development. As a result of the model validation, the SA model achieved 86.4 percent F-measure, which showed a higher performance of risk extraction than the previous study [25]. Furthermore, the accuracy of DC reached 86.3 percent. In particular, the even distribution of precision of 88.6 and 84.3 percent of recall can be interpreted that the ground truth value and the automatic extraction result of the SA model match more than 80 percent. The RLR model showed an accuracy of 88 percent due to the first classification but 46.8 percent accuracy in the second classification. It is expected that the accuracy of the RLR model will improve through additional data collection and training. The SA and RLR models are embedded in the *EMAP* system, an engineering integration platform. *EMAP* was provided on a cloud basis considering user convenience.

It is believed that this research has theoretical contributions. It has been confirmed in previous studies that most studies apply either a rule-based or a training-based approach when analyzing contracts and extracting information using AI techniques. This study has theoretical significance in that both rule-based and training-based approaches were used for contract analysis. It was applied to a cloud-based engineering machine learning integrated analysis platform and provided a solution. In addition, this study contributes to the literature on the contract lexicon development for the first time in the EPC field. The EPC contract lexicon was based on the contract taxonomy systematized by contract experts, and it can be said that it is the first detailed and systematic contract lexicon in the EPC plant field. The lexicon developed through this study consisted of 32 items, 79 details, and was created with 87 semantic rules. The SA model, which shows significantly improved risk clause extraction accuracy compared to the previous study proposed by Lee et al. [25], is analyzed due to the accuracy of the detailed lexicon and semantic rule associated with it. In the study of information recognition and extraction of documents, the development of the PDF structuralization module suggests a practical alternative to effective information recognition before text preprocessing. In most NLP studies, text data recognition of documents is read as a text file. However, the data read into the text file contain noise data, such as headers, footers, and page numbers, that are not required

for analysis. It is also challenging to convert text data into data frames of a chapter and sentence units. The PDF structuralization module developed in this study suggests a method to resolve this problem. Additionally, it was explained to assist other researchers who aim to build a similar tool in the future. From a risk management perspective, this study contributes to an attempt to manage contract risk using AI. In particular, the RLR model is significant as it developed a risk management model that analyzes the risk impact of EPC contracts by applying AI techniques to the traditional risk management method. The RLR model can improve the accuracy of the analysis result through additional data and training of the model.

The practical implications of this study are as follows. First, the technical system that can effectively support the EPC ITB document analysis work is developed. The senior engineers tested the ITB analysis, and there was a high deviation in analysis accuracy, depending on the individual engineer's competency and experience. This study established a technical system that prevents errors due to human mistakes and contributes to improving the accuracy of risk analysis tasks. As a result, it is expected to improve the work efficiency (time, quality) of inexperienced junior engineers in the risk analysis of EPC contracts. Second, it was confirmed that the owner's requirement clarification time could be significantly reduced through the onsite voice of customer (VOC) of the engineer in charge of ITB analysis. The automated management of the ITB analysis tasks that previously depended only on engineers' experience can shorten ITB analysis time and reduce engineers' workload in charge. Third, in general, when senior engineers change their roles or reposition to another department, their knowledge and skills will not transfer over. In comparison, the knowledge learned by AI can be converted into organizational learning for a sustainable EPC industry. Lastly, this study contributes to the intelligence of the sustainable EPC industry and the establishment of a digital workforce based on AI technology.

*8.2. Limitations and Further Study*

Limitations exist in this study and will be discussed in the following. First, this study's scope was to support the EPC contractor's decision-making in the bidding stage of the EPC project. The subject of analysis was limited to the conditions of the contract among the ITB documents. In addition, this study does not include technical specification documents. Second, the SA model has some errors even though a rule-based approach improves the automatic extraction accuracy. It was challenging to generate all the contract sentences as a rule, and it was also impossible to extract numerous linguistic phenomena through the rule. Continuous research is required to improve the performance of rule-based NLP. Third, this study determined the risk through consultation with SMEs in the EPC field. However, there was a difference in risk judgment for each expert, and accordingly, the training dataset could not be considered perfect. In addition to adding data, research to improve the reliability of training data will continue in the future. Fourth, there was a limitation in data collection. Due to the availability of the collected data, the size of the training dataset used for the RLR model was limited. Data collection is a big challenge in the case of documents in a specific domain, especially with contracts. Despite the collection of more than 20 EPC contracts, there were limitations to model training. The RLR model is expected to increase performance accuracy through additional data collection and training. Overall, research that can improve machine learning performance based on small data is required in the future.

The discussion points for further research are as follows. First, the EMAP system was developed as part of a three-year-long engineering research project. The ITB analysis module was coded by a research team rather than a professional software company, and there are some shortcomings related to response time. In particular, the phenomenon was conspicuous in the training-based RLR model. Further research is required to solve the response time delay of the RLR model. Second, the EPC contract lexicon was developed based on an extensive literature review and brainstorming of experts. The developer can

modify this lexicon if necessary, and there is still room for improvement. Third, this study only targeted unstructured data, such as text format. Tables and drawings in ITB were excluded from the analysis. In the future, active research is needed for further analysis of the automation of preprocessing of tables and drawings, and improvement of recognition accuracy. The integration of unstructured data preprocessing automation technology and various unstructured data analyses is expected to enable the realization of a trustworthy big-data-based engineering machine learning platform.

In recent years, the accuracy of text classification models such as the RLR model has improved. However, it cannot indicate whether the model "understands" the text from the semantic level in the same way as human beings. Therefore, the model's semantic representation ability and decision confidence robustness require further improvement. Finally, despite the difficulty of collecting EPC contracts, the RLR model is expected to increase the analysis accuracy through the additional collection and learning of EPC contract data. Additionally, the RLR model's performance can be further enhanced because it is trained with an extensive training dataset. To conclude, research that can improve machine learning performance based on small data is required in the future.

## Abbreviations

The following abbreviations and parameters are used in this paper:

| | |
|---|---|
| Adam | Adaptive moment estimation |
| AI | Artificial intelligence |
| ANN | Artificial neural network |
| API | Application programming interface |
| CNN | Convolutional neural network |
| CRC | Critical risk check |
| CSV | Comma-separated values |
| DC | Deontic classification |
| DL | Deontic logic |
| *EMAP* | *Engineering Machine-learning Automation Platform* |
| EPC | Engineering, procurement, construction |
| FIDIC | Fédération Internationale Des Ingénieurs-Conseils |
| FOL | First-order logic |
| IE | Information extraction |
| ITB | Invitation to bid |
| JSON | JavaScript object notation |
| LD | Liquidated damages |
| LSTM | Long short-term memory |
| ML | Machine learning |

| NER | Named entity recognition |
| NLP | Natural language processing |
| OCR | Optical character recognition |
| OOV | Out-of-vocabulary |
| OPF | Obligation, permission, and prohibition/forbidden |
| PDF | Portable document format |
| PI | Probability and impact |
| PM | Project management |
| POC | Proof of concept |
| POS tagging | Part of speech tagging |
| RNN | Recurrent neural network |
| RLR | Risk level ranking |
| SA | Semantic analysis |
| SMEs | Subject matter experts |
| SVO | Subject–verb–object |
| SVR | Support vector regression |
| WAS | Web application server |

## References

1. DLA Piper. EPC Contracts in the Process Plant Sector. Available online: www.dlapiper.com (accessed on 5 February 2022).
2. Ritsche, F.-P.; Wagner, R.; Schlemmer, P.; Steinkamp, M.; Valnion, B.D. *Innovation Project EPC 4.0 'Unleashing the Hidden Potential'*; ProjectTeam: Hamburg, Germany, 2019.
3. International Trade Administration. South Korea-Construction Services. Available online: https://www.trade.gov/country-commercial-guides/south-korea-construction-services (accessed on 7 February 2022).
4. Vogl, R. *The Coming of Age of Legal Technology*; Stanford University: Stanford, CA, USA, 2016.
5. Lane, H.; Hapke, H.; Howard, C. *Natural Language Processing in Action: Understanding, Analyzing, and Generating Text with Python*; Simon and Schuster: New York, NY, USA, 2019.
6. Ebrahimnejad, S.; Mousavi, S.M.; Tavakkoli-Moghaddam, R.; Heydar, M. Evaluating high risks in large-scale projects using an extended VIKOR method under a fuzzy environment. *Int. J. Ind. Eng. Comput.* **2012**, *3*, 463–476. [CrossRef]
7. Hung, M.S.; Wang, J. Research on Delay Risks of EPC Hydropower Construction Projects in Vietnam. *Int. J. Power Energy Eng.* **2016**, *4*, 8. [CrossRef]
8. Jahantigh, F.F.; Malmir, B.; Avilaq, B.A. Economic risk assessment of EPC projects using fuzzy TOPSIS approach. *Int. J. Ind. Syst. Eng.* **2017**, *27*, 161–179. [CrossRef]
9. Kim, M.-H.; Lee, E.-B.; Choi, H.-S. Detail Engineering Completion Rating Index System (DECRIS) for Optimal Initiation of Construction Works to Improve Contractors' Schedule-Cost Performance for Offshore Oil and Gas EPC Projects. *Sustainability* **2018**, *10*, 2469. [CrossRef]
10. Kabirifar, K.; Mojtahedi, M. The impact of Engineering, Procurement and Construction (EPC) Phases on Project Performance: A Case of Large-scale Residential Construction Project. *Buildings* **2019**, *9*, 15. [CrossRef]
11. Gunduz, M.; Almuajebh, M. Critical Success Factors for Sustainable Construction Project Management. *Sustainability* **2020**, *12*, 1990. [CrossRef]
12. Koulinas, G.K.; Xanthopoulos, A.S.; Tsilipiras, T.T.; Koulouriotis, D.E. Schedule delay risk analysis in construction projects with a simulation-based expert system. *Buildings* **2020**, *10*, 134. [CrossRef]
13. Okudan, O.; Budayan, C.; Dikmen, I. A knowledge-based risk management tool for construction projects using case-based reasoning. *Expert. Syst. Appl.* **2021**, *173*, 114776. [CrossRef]
14. Surden, H. Computable contracts. *UC Davis Law Rev.* **2012**, *46*, 629.
15. LawGeex. Comparing the Performance of AI to Human Lawyers in the Review of Standard Business Contracts. Available online: https://ai.lawgeex.com/rs/345-WGV-842/images/LawGeex%20eBook%20Al%20vs%20Lawyers%202018.pdf (accessed on 10 January 2022).
16. Cummins, J.; Clack, C. Transforming Commercial Contracts through Computable Contracting. *arXiv* **2020**, arXiv:2003.10400. [CrossRef]
17. Dixon, H.B., Jr. What judges and lawyers should understand about artificial intelligence technology. *ABA J.* **2020**, *59*, 36–38.
18. Clack, C.D. Languages for Smart and Computable Contracts. *arXiv* **2021**, arXiv:2104.03764.
19. Salama, D.A.; El-Gohary, N.M. Automated compliance checking of construction operation plans using a deontology for the construction domain. *J. Comput. Civ. Eng.* **2013**, *27*, 681–698. [CrossRef]
20. Chopra, D.; Joshi, N.; Mathur, I. *Mastering Natural Language Processing with Python*; Packt Publishing Ltd.: Birmingham, UK, 2016.
21. Zhang, J.; El-Gohary, N. Automated reasoning for regulatory compliance checking in the construction domain. In *Construction Research Congress 2014: Construction in a Global Network*; ASCE: Reston, VA, USA, 2014; pp. 907–916.
22. Williams, T.P.; Gong, J. Predicting construction cost overruns using text mining, numerical data and ensemble classifiers. *Autom. Constr.* **2014**, *43*, 23–29. [CrossRef]

23. Lee, J.; Yi, J.-S. Predicting project's uncertainty risk in the bidding process by integrating unstructured text data and structured numerical data using text mining. *Appl. Sci.* **2017**, *7*, 1141. [CrossRef]

24. Zou, Y.; Kiviniemi, A.; Jones, S.W. Retrieving similar cases for construction project risk management using Natural Language Processing techniques. *Autom. Constr.* **2017**, *80*, 66–76. [CrossRef]

25. Lee, J.H.; Yi, J.-S.; Son, J.W. Development of automatic-extraction model of poisonous clauses in international construction contracts using rule-based NLP. *J. Comput. Civ. Eng.* **2019**, *33*, 04019003. [CrossRef]

26. Moon, S.; Lee, G.; Chi, S.; Oh, H. Automated construction specification review with named entity recognition using natural language processing. *J. Constr. Eng. Manag.* **2021**, *147*, 04020147. [CrossRef]

27. Choi, S.-W.; Lee, E.-B.; Kim, J.-H. The Engineering Machine-Learning Automation Platform (*EMAP*): A Big-Data-Driven AI Tool for Contractors' Sustainable Management Solutions for Plant Projects. *Sustainability* **2021**, *13*, 10384. [CrossRef]

28. Choi, S.J.; Choi, S.W.; Kim, J.H.; Lee, E.-B. AI and Text-Mining Applications for Analyzing Contractor's Risk in Invitation to Bid (ITB) and Contracts for Engineering Procurement and Construction (EPC) Projects. *Energies* **2021**, *14*, 4632. [CrossRef]

29. Park, M.-J.; Lee, E.-B.; Lee, S.-Y.; Kim, J.-H. A Digitalized Design Risk Analysis Tool with Machine-Learning Algorithm for EPC Contractor's Technical Specifications Assessment on Bidding. *Energies* **2021**, *14*, 5901. [CrossRef]

30. Fantoni, G.; Coli, E.; Chiarello, F.; Apreda, R.; Dell'Orletta, F.; Pratelli, G. Text mining tool for translating terms of contract into technical specifications: Development and application in the railway sector. *Comput. Ind.* **2021**, *124*, 103357. [CrossRef]

31. Jang, B.; Kim, M.; Harerimana, G.; Kang, S.-U.; Kim, J.W. Bi-LSTM model to increase accuracy in text classification: Combining Word2vec CNN and attention mechanism. *Appl. Sci.* **2020**, *10*, 5841. [CrossRef]

32. Pouyanfar, S.; Sadiq, S.; Yan, Y.; Tian, H.; Tao, Y.; Reyes, M.P.; Shyu, M.-L.; Chen, S.-C.; Iyengar, S. A survey on deep learning: Algorithms, techniques, and applications. *ACM Comput. Surv. (CSUR)* **2018**, *51*, 1–36. [CrossRef]

33. Devlin, J.; Chang, M.-W.; Lee, K.; Toutanova, K. Bert: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv* **2019**, arXiv:1810.04805.

34. Zhou, S.K.; Rueckert, D.; Fichtinger, G. *Handbook of Medical Image Computing and Computer Assisted Intervention*; Academic Press: Cambridge, MA, USA, 2019.

35. Kowsari, K.; Jafari Meimandi, K.; Heidarysafa, M.; Mendu, S.; Barnes, L.; Brown, D. Text Classification Algorithms: A Survey. *Information* **2019**, *10*, 150. [CrossRef]

36. Pascanu, R.; Mikolov, T.; Bengio, Y. On the difficulty of training recurrent neural networks. In Proceedings of the 30th International Conference on Machine Learning, Atlanta, GA, USA, 17–21 June 2013; pp. 1310–1318.

37. Hochreiter, S.; Schmidhuber, J. Long short-term memory. *Neural. Comput.* **1997**, *9*, 1735–1780. [CrossRef] [PubMed]

38. Schuster, M.; Paliwal, K.K. Bidirectional recurrent neural networks. *IEEE Trans. Signal Process.* **1997**, *45*, 2673–2681. [CrossRef]

39. Zhou, P.; Qi, Z.; Zheng, S.; Xu, J.; Bao, H.; Xu, B. Text classification improved by integrating bidirectional LSTM with two-dimensional max pooling. *arXiv* **2016**, arXiv:1611.06639.

40. Li, Q.; Peng, H.; Li, J.; Xia, C.; Yang, R.; Sun, L.; Yu, P.S.; He, L. A Survey on Text Classification: From Traditional to Deep Learning. *ACM Trans. Intell. Syst. Technol.* **2022**, *13*, 364. [CrossRef]

41. Minaee, S.; Kalchbrenner, N.; Cambria, E.; Nikzad, N.; Chenaghlu, M.; Gao, J. Deep Learning—Based Text Classification. *ACM Comput. Surv.* **2022**, *54*, 1–40. [CrossRef]

42. Zhang, J.; El-Gohary, N.M. Semantic NLP-based information extraction from construction regulatory documents for automated compliance checking. *J. Comput. Civ. Eng.* **2016**, *30*, 04015014. [CrossRef]

43. PDF Parser. Available online: https://py-pdf-parser.readthedocs.io/en/latest/overview.html (accessed on 13 January 2022).

44. Fernández-Caballero, A.; López, M.T.; Castillo, J.C. Display text segmentation after learning best-fitted OCR binarization parameters. *Expert. Syst. Appl.* **2012**, *39*, 4032–4043. [CrossRef]

45. Vijayarani, S.; Ilamathi, M.J.; Nithya, M. Preprocessing techniques for text mining-an overview. *Int. J. Comput.* **2015**, *5*, 7–16.

46. spaCy. Tokenization. Available online: https://spacy.io/usage/linguistic-features#tokenization (accessed on 15 January 2022).

47. spaCy. Lemmatization. Available online: https://spacy.io/usage/linguistic-features#lemmatization (accessed on 15 January 2022).

48. spaCy. Part-of-Speech Tagging. Available online: https://spacy.io/usage/linguistic-features#pos-tagging (accessed on 15 January 2022).

49. Wu, Y.; Zhang, Q.; Huang, X.-J.; Wu, L. Phrase dependency parsing for opinion mining. In Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing, Singapore, 6–7 August 2009; Association for Computational Linguistics: Singapore, 2009; pp. 1533–1541.

50. spaCy. Dependency-Parsing. Available online: https://spacy.io/usage/linguistic-features#dependency-parse (accessed on 17 January 2022).

51. Google. ClearNLP. Available online: https://github.com/clir/clearnlp-guidelines (accessed on 17 January 2022).

52. Tiwary, U.; Siddiqui, T. *Natural Language Processing and Information Retrieval*; Oxford University Press, Inc.: Oxford, UK, 2008.

53. Niu, J.; Issa, R.R. Developing taxonomy for the domain ontology of construction contractual semantics: A case study on the AIA A201 document. *Adv. Eng. Inform.* **2015**, *29*, 472–482. [CrossRef]

54. Prisacariu, C.; Schneider, G. A formal language for electronic contracts. In *International Conference on Formal Methods for Open Object-Based Distributed Systems*; Springer: Berlin/Heidelberg, Germany, 2007; pp. 174–189.

55. Xu, X.; Cai, H. Ontology and rule-based natural language processing approach for interpreting textual regulations on underground utility infrastructure. *Adv. Eng. Inform.* **2021**, *48*, 101288. [CrossRef]

56. McNamara, P.; Van De Putte, F. The Stanford Encyclopedia of Philosophy. Available online: https://plato.stanford.edu/entries/logic-deontic/ (accessed on 14 February 2022).

57. Cheng, J. Deontic relevant logic as the logical basis for representing and reasoning about legal knowledge in legal information systems. In *International Conference on Knowledge-Based and Intelligent Information and Engineering Systems*; Springer: Berlin/Heidelberg, Germany, 2008; pp. 517–525.

58. Hu, S.; Zou, L.; Yu, J.X.; Wang, H.; Zhao, D. Answering natural language questions by subgraph matching over knowledge graphs. *IEEE Trans. Knowl. Data Eng.* **2017**, *30*, 824–837. [CrossRef]

59. Chen, Y.-H.; Lu, E.J.-L.; Ou, T.-A. Intelligent SPARQL Query Generation for Natural Language Processing Systems. *IEEE Access* **2021**, *9*, 158638–158650. [CrossRef]

60. Lai, S.; Liu, K.; He, S.; Zhao, J. How to generate a good word embedding. *IEEE Intell. Syst.* **2016**, *31*, 5–14. [CrossRef]

61. Keras. Text Data Preprocessing. Available online: https://keras.io/api/preprocessing/text/ (accessed on 18 January 2022).

62. TensorFlow. Tokenizer. Available online: https://www.tensorflow.org/api_docs/python/tf/keras/preprocessing/text/Tokenizer (accessed on 2 March 2022).

63. Manaswi, N.K. Understanding and working with Keras. In *Deep Learning with Applications Using Python*; Springer: Berlin/Heidelberg, Germany, 2018; pp. 31–43.

64. Dawes, J. Do data characteristics change according to the number of scale points used? An experiment using 5-point, 7-point and 10-point scales. *Int. J. Mark. Res.* **2008**, *50*, 61–104. [CrossRef]

65. Cui, Z.; Ke, R.; Pu, Z.; Wang, Y. Deep bidirectional and unidirectional LSTM recurrent neural network for network-wide traffic speed prediction. *arXiv* **2018**, arXiv:1801.02143.

66. Shen, S.-L.; Atangana Njock, P.G.; Zhou, A.; Lyu, H.-M. Dynamic prediction of jet grouted column diameter in soft soil using Bi-LSTM deep learning. *Acta Geotech.* **2021**, *16*, 303–315. [CrossRef]

67. Renn, O. Three decades of risk research: Accomplishments and new challenges. *J. Risk Res.* **1998**, *1*, 49–71. [CrossRef]

68. Jang, W.-S.; Hong, H.-U.; Han, S.-H. Risk Identification and Priority method for Overseas LNG Plant Projects-Focusing on Design Phase. *Korean J. Constr. Eng. Manag.* **2011**, *12*, 146–154. [CrossRef]

69. Bergstra, J.; Bengio, Y. Random search for hyper-parameter optimization. *J. Mach. Learn. Res.* **2012**, *13*, 281–305.

70. Amir-Ahmadi, P.; Matthes, C.; Wang, M.-C. Choosing prior hyperparameters: With applications to time-varying parameter models. *J. Bus. Econ. Stat.* **2020**, *38*, 124–136. [CrossRef]

71. Afaq, S.; Rao, S. Significance of Epochs On Training A Neural Network. *Int. J. Sci. Technol. Res.* **2020**, *19*, 485–488.

72. TensorFlow. Overfit and Underfit. Available online: https://www.tensorflow.org/tutorials/keras/overfit_and_underfit (accessed on 7 March 2022).

73. TensorFlow. EarlyStopping. Available online: https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/EarlyStopping (accessed on 7 March 2022).

74. Aretz, K.; Bartram, S.M.; Pope, P.F. Asymmetric loss functions and the rationality of expected stock returns. *Int. J. Forecast.* **2011**, *27*, 413–437. [CrossRef]

75. Jadon, S. A survey of loss functions for semantic segmentation. In Proceedings of the 2020 IEEE Conference on Computational Intelligence in Bioinformatics and Computational Biology (CIBCB), Vina del Mar, Chile, 27–29 October 2020; pp. 1–7.

76. Kingma, D.P.; Ba, J. Adam: A method for stochastic optimization. In Proceedings of the 3rd International Conference on Learning Representations (ICLR 2015), San Diego, CA, USA, 7–9 May 2015.

77. Apache. Tomcat Software. Available online: http://tomcat.apache.org (accessed on 22 March 2022).

78. Oracle. MySQL. Available online: https://www.oracle.com/mysql/ (accessed on 22 March 2022).

79. Santoro, M.; Vaccari, L.; Mavridis, D.; Smith, R.; Posada, M.; Gattwinkel, D. *Web Application Programming Interfaces (APIs): General-Purpose Standards, Terms and European Commission Initiatives*; European Union: Luxembourg, 2019. [CrossRef]

80. Gunnulfsen, M. Scalable and Efficient Web Application Architectures: Thin-Clients and Sql vs. Thick-Clients and Nosql. Master's Thesis, The University of Oslo, Oslo, Norway, 2013.