*Article*

# Design and Development of Knowledge Graph for Industrial Chain Based on Deep Learning

**Yue Li** [ID]**, Yutian Lei \*, Yiting Yan, Chang Yin and Jiale Zhang**

School of Computer Science and Technology, Donghua University, Shanghai 201620, China; frankyueli@dhu.edu.cn (Y.L.); 2222843@mail.dhu.edu.cn (Y.Y.); 2222779@mail.dhu.edu.cn (C.Y.); 2222748@mail.dhu.edu.cn (J.Z.)
\* Correspondence: 2222849@mail.dhu.edu.cn

**Abstract:** This paper aims to structure and semantically describe the information within the industrial chain by constructing an Industry Chain Knowledge Graph (ICKG), enabling more efficient and intelligent information management and analysis. In more detail, this paper constructs a multi-domain industrial chain dataset and proposes a method that combines the top-down establishment of a semantic expression framework with the bottom-up establishment of a data layer to build an ICKG. In the data layer, a deep learning algorithm based on BERT-BiLSTM-CRF is used to extract industry chain entities from relevant literature and reports. The results indicate that the model can effectively identify industry chain entities. These entities and relationships populate a Neo4j graph database, creating a large-scale ICKG for visual display and aiding cross-domain applications.

**Keywords:** industrial chain; knowledge graph; information extraction; entity recognition; NLP

## 1. Introduction

The industrial chain is one of the cornerstones of the modern economic system, describing the entire process of products and services from the starting point of raw materials to the final consumer. Currently, many scholars are dedicated to constructing and refining industry chains in various sectors. Among them, Lou et al. [1] proposed a framework for the traditional Chinese medicine industry chain, which includes upstream, midstream, and downstream segments. Lu et al. [2] conducted an in-depth analysis of the technical efficiency and influencing factors of the energy industry chain in China. Yang et al. [3] explained how studying the carbon footprint in the industry chain would promote hydrogen development in specific sectors. Song et al. [4] highlighted the importance and challenges of the rice industry chain. The authors of [5] emphasized the comprehensive industry chain perspective on mineral resource security, highlighting the importance of industry chains in safeguarding the supply of strategic mineral resources and sustainable development.

A knowledge graph is a technological means capable of mining, organizing, and efficiently managing knowledge from massive data, providing qualitative improvement for information services, and offering users more intelligent and personalized services [6]. In various fields, knowledge graphs have shown a wide range of application prospects, including search engine optimization in medical care [7–12] intelligent recommendation systems [13–15], intelligent transportation and urban planning, financial risk management [16–21], etc. Knowledge graphs can integrate scattered information within the industrial chain, establish comprehensive entity relationships, achieve efficient sharing and management of information, and thereby break information silos. Utilizing knowledge graphs can intelligently analyze the relationships between various links in the industrial chain, helping to deepen the understanding of business processes and optimize decisions. Knowledge graphs can also integrate data from multiple sources, enabling real-time monitoring and prediction of potential risks in the industrial chain, thereby improving the accuracy and efficiency of risk management. Based on knowledge graph-based data

analysis, more comprehensive and accurate information support can be provided for decision-making in the industrial chain, helping managers make wiser decisions.

The current industrial chain is becoming increasingly complex, covering multiple stages including raw material procurement, manufacturing, logistics, inventory management, and sales channels, with participants including manufacturers, suppliers, and logistics companies [22]. However, these participants typically use different information systems and data formats, leading to information asymmetry and inefficient processes [23–25]. Existing industrial chain datasets have limited coverage [26], are outdated, and struggle to handle unstructured data. Traditional industrial chains also face issues such as information asymmetry, supply chain risks, and inefficiencies [27]. Particularly when dealing with complex texts, traditional entity recognition algorithms and relationship extraction rules may encounter difficulties [28]. Texts often contain noise and ambiguity, making entity recognition and relationship extraction even more complex.

To address these challenges, this research aims to construct a cross-disciplinary, comprehensive industrial chain dataset by integrating data from multiple sources with strong real-time capabilities and various dimensions. This integration aims to enhance the comprehensive understanding of industrial chains and provide decision support. Simultaneously, the introduction of knowledge graphs into industrial chains aims to promote digital transformation, improve overall operational efficiency and management levels, and achieve efficient sharing and management by eliminating information silos. In order to extract knowledge more accurately, we plan to optimize the entity recognition algorithm and relationship extraction rules and use deep learning and natural language processing technology to process unstructured data to improve the quality of data extraction. This comprehensive strategy aims to address current challenges and promote the intelligent and informatized development of industrial chains.

This paper collected a large amount of structured data and unstructured text related to the industry chain and preprocessed these sample data to construct a comprehensive industry chain dataset covering multiple domains. Based on this comprehensive dataset, the study investigated a deep learning-based method for industry chain entity recognition, comparing the effectiveness of different deep learning algorithms on the industry chain dataset to obtain the best method for industry chain entity recognition. The experiments showed that the BERT-BiLSTM-CRF model significantly improved the efficiency of industry chain entity recognition. Finally, the extracted standard data were stored in the Neo4j graph database, successfully constructing a large-scale ICKG. This knowledge graph provides strong support for cross-domain applications, including knowledge retrieval, intelligent question answering, smart decision-making, intelligent marketing, and intelligent recommendation systems. The main contributions of our research are as follows:

- This paper integrated structured data (such as Shenwan Industry, Shenzhen Stock Exchange, Shanghai Stock Exchange, etc.) and unstructured data (such as encyclopedias, news, annual reports, etc.) to build a comprehensive industry chain dataset. This comprehensive dataset helps provide more comprehensive and multidimensional industry information and knowledge;
- The entity recognition algorithm based on the BERT-BiLSTM-CRF model proposed in this paper performed excellently on the industry chain dataset. The macro-F1, macro-P, and macro-R of the BERT-BiLSTM-CRF model were 97.10%, 96.80%, and 96.95%, respectively, showing the best performance among the three models. The macro-F1 of this model is 0.44% higher than that of the BERT model and 24.78% higher than that of the BERT-CRF model;
- This paper proposes a method of combining top-down and bottom-up approaches to construct an ICKG. This ICKG serves as a core resource for cross-domain applications, providing extensive and robust support for various fields, including knowledge retrieval, intelligent question answering, smart decision-making, intelligent marketing, and intelligent recommendation systems.

The remaining sections of the paper are arranged as follows: Section 2 outlines the relevant research on entity recognition and relation extraction. Section 3 introduces the concept and overall architecture of the ICKG. Sections 4 and 5, respectively, discuss the construction and improvement of the model layer and data layer of the ICKG. Section 6 employs various deep learning models to perform entity recognition on the industrial chain dataset and thoroughly analyze the experimental results. Finally, we summarize the findings of the study and propose directions for future research.

## 2. Related Work

Named entity recognition (NER) is a technique used to identify entities from a large amount of text [29]. NER can be mainly categorized into three types [28]. The first type is rule-based or knowledge-based methods [30], which have the advantages of strong interpretability and good adaptability in specific domains. Due to their reliance on manually designed rules and domain expert knowledge, these methods do not require a large amount of annotated data. However, their generalization ability is limited, making it difficult to handle different languages and complex grammars, and the maintenance cost is high when facing new domains or changes. Rule-based methods may perform poorly in processing long texts or complex contexts due to their relatively weak capability of handling contextual dependencies and complex grammars. The second type is learning-based methods [31], which utilize machine learning algorithms, such as deep learning models, to learn patterns for identifying named entities from text through a large amount of annotated data. Their advantages lie in their strong adaptability, the ability to handle diverse language structures, and the absence of manual rule design. The third type is neural network methods based on feature inference [32], which infer named entities by learning features from the text. Their advantages include efficient utilization of contextual information, adaptation to diverse language structures, and no need for manual feature design. However, they depend heavily on a large amount of annotated data, have higher model complexity, are relatively difficult to interpret, and require more computational resources.

The goal of entity relation extraction models is to identify relationships between entities from text [33]. Commonly used methods for entity relation identification include rule-based methods [34,35], pattern-based methods [36], machine learning-based methods [37,38], and deep learning-based methods [39,40].

After a comprehensive analysis of relevant research work, this paper addresses the following issues when constructing an ICKG: how to deal with data scarcity and labeling difficulties; how to improve the model's domain adaptability and generalization ability; and explore how to more effectively utilize contextual information. This paper constructs a comprehensive industrial chain dataset, and the proposed BERT-BiLSTM-CRF-based model can significantly improve the efficiency of industrial chain entity recognition.

## 3. The Overall Architecture of the ICKG

The industrial chain refers to the entire production and delivery process of a product, from the production of raw materials, processing and manufacturing, and distribution to the final consumer. The industrial chain is typically divided into three parts: upstream, midstream, and downstream. Upstream refers to the part of the industrial chain responsible for providing raw materials, components, and primary processing. Midstream is the part of the industrial chain responsible for processing, manufacturing, and assembly. Downstream is the part of the industrial chain responsible for product distribution, sales, and final consumption. Upstream provides raw materials and basic substances; midstream processes and manufactures products; and downstream interacts with end-users, driving product flow to the market. This division helps to understand and analyze the entire industrial chain, thereby optimizing decisions related to supply chain management, product design, and market positioning.

This paper employs data processing and structured extraction techniques to successfully construct a comprehensive dataset of the industrial chain. Entity recognition methods

based on BERT-BiLSTM-CRF models and rule-based methods are utilized for knowledge extraction, and the data are stored in the Neo4j graph database, forming a large-scale ICKG. This knowledge graph provides powerful support for cross-disciplinary applications, covering various aspects such as knowledge retrieval, intelligent question-answering, intelligent decision-making, intelligent marketing, and intelligent recommendation. The overall architecture design is illustrated in Figure 1.
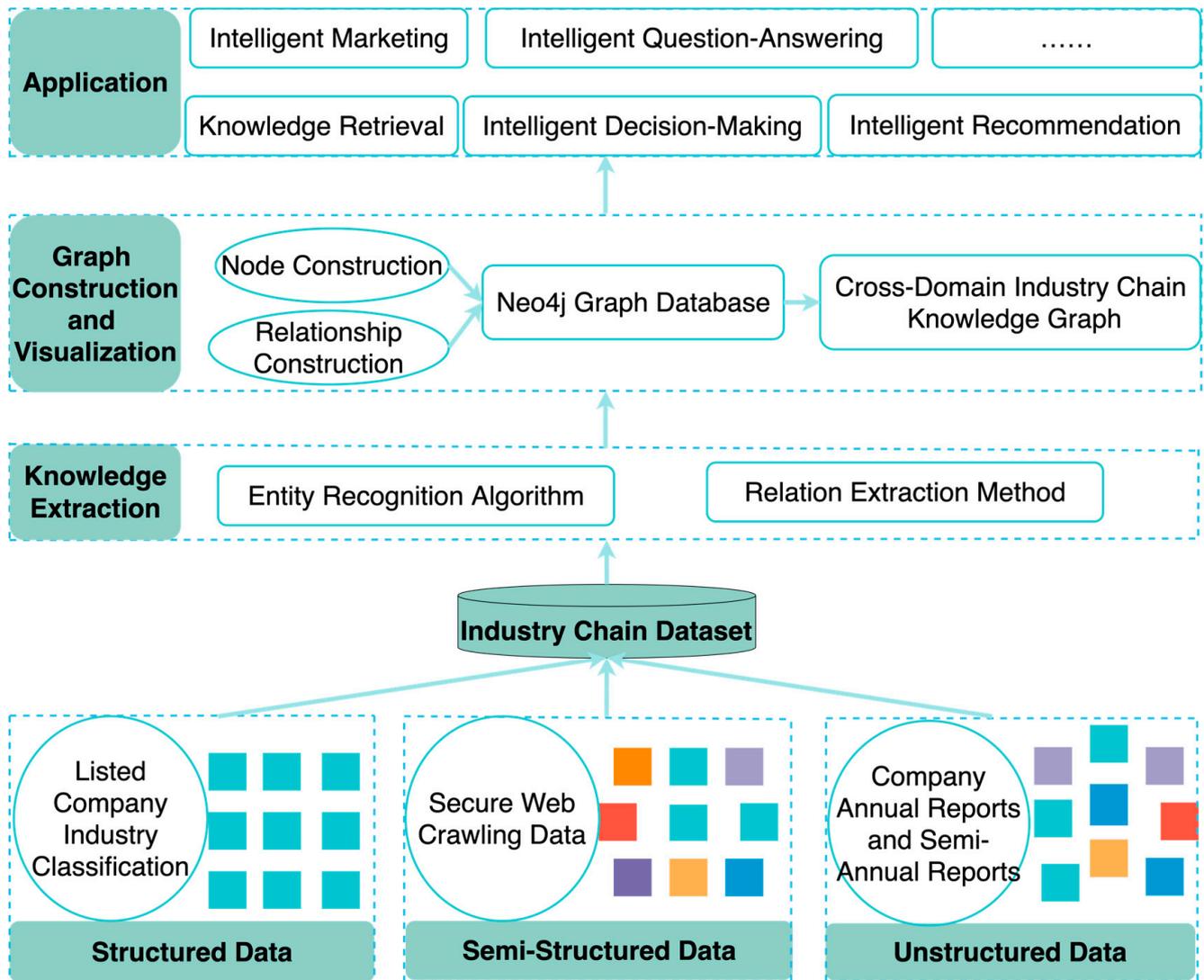


**Figure 1.** Framework for constructing the ICKG.

The construction of the ICKG adopts a combined approach of top-down and bottom-up methodologies, as illustrated in Figure 2, depicting the basic process of constructing the pattern layer and data layer. The top-down construction of the pattern layer provides a theoretical framework and abstract model, guiding the construction of the knowledge graph. The bottom-up construction of the data layer ensures the authenticity and practicality of the knowledge graph through the extraction and integration of actual data. The integration of top-down and bottom-up approaches, through modeling and validation from multiple perspectives, enhances the accuracy and credibility of the knowledge graph. This enables it to meet the requirements of different domains and scenarios, flexibly applying it to various practical situations.
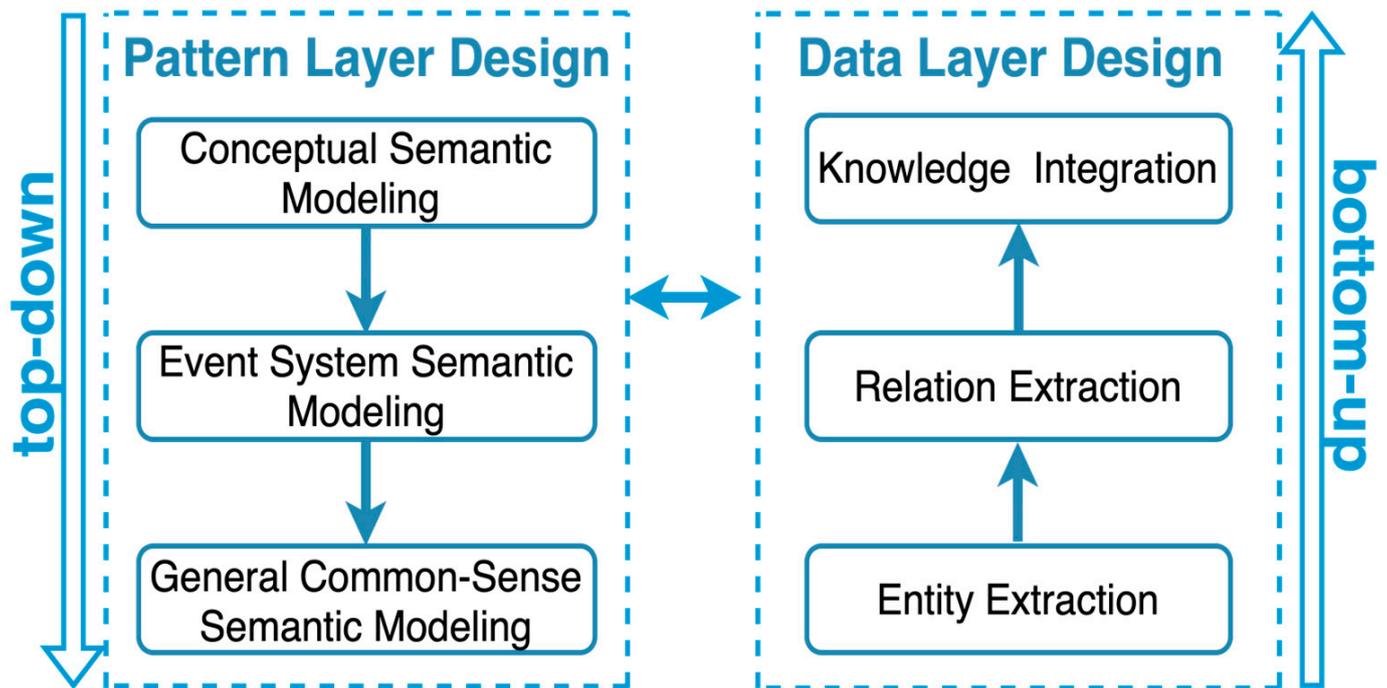
**Figure 2.** Basic process of constructing the pattern layer and data layer of the ICKG.

## 4. Construction of the Pattern Layer of the ICKG

This paper constructs the pattern layer of the ICKG through three steps: conceptual semantic modeling, event semantic modeling, and general common-sense semantic modeling. Conceptual semantic modeling aims to define and establish key concepts, entities, relationships, and their logic within the business domain. Event semantic modeling follows conceptual semantic modeling, relying on the defined conceptual model to derive semantic relationships between events, actions, or processes. General common-sense semantic modeling builds upon the first two steps, relying on the models of conceptual and event semantics to add more universal concepts and relationships to the system's knowledge base. By top-down construction of the pattern layer of the ICKG through conceptual, event, and general common-sense semantic modeling, comprehensive modeling and logical derivation of key concepts, events, and universal relationships in the industrial chain are achieved, enhancing the richness and universality of the knowledge graph.

### 4.1. Conceptual Semantic Modeling

The first layer of conceptual semantic modeling in the pattern layer aims to associate entities with relevant concepts, with the goal of describing and representing concepts, entities, and their relationships in the real world to better understand, analyze, and process complex information and data. It helps organize complex concepts and information into a structured format for knowledge sharing, automated reasoning, and data integration. In the context of the industrial chain, conceptual semantic modeling involves modeling industry standards, domain-specific terminology, and so on. As depicted in Figure 3, the conceptual model adopts a hierarchical tree-like classification system and forms a concept dictionary by hierarchically encoding the concept tree. The benefit of this design is that when renaming concepts, only the concept dictionary information needs to be updated without the need to update the data of indexes or relationships. Since concepts are associated with many entities, any change in a concept will involve changes throughout the tree. Using a concept dictionary effectively addresses this issue.
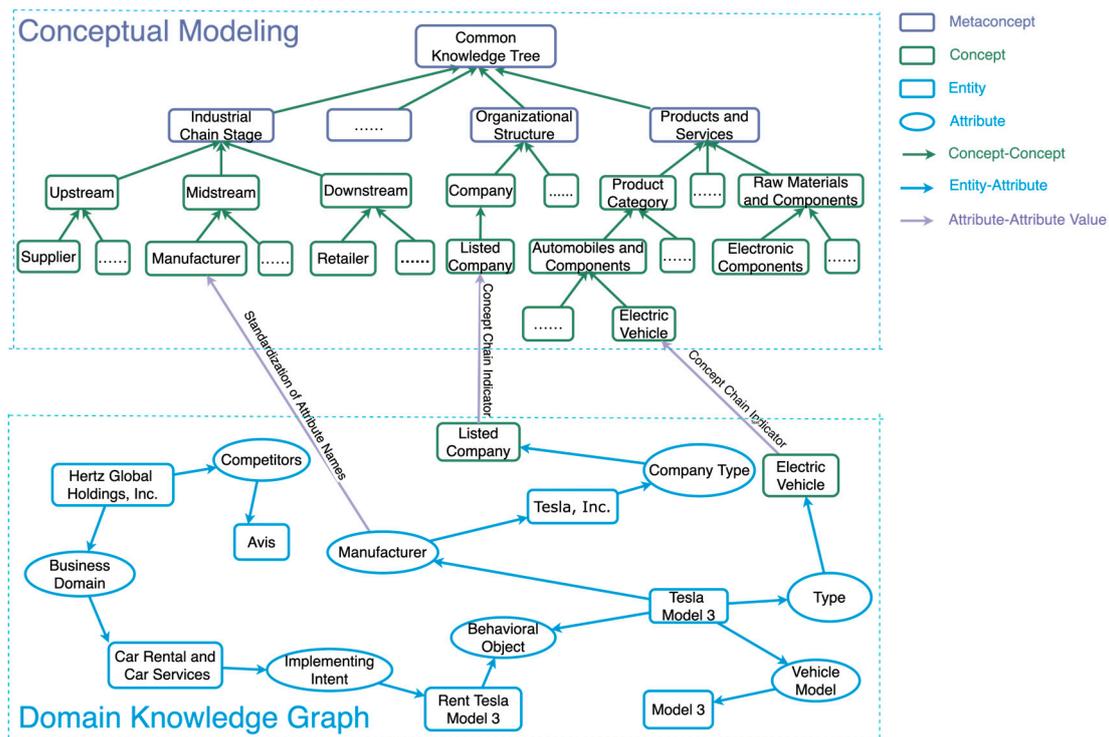
**Figure 3.** Conceptual semantic modeling.

### 4.2. Event Semantic Modeling

The second layer of semantic modeling in the pattern layer focuses on defining and modeling events, determining relationships between events, and associating events with entities to capture dynamic information in the knowledge graph. In the ICKG, this involves defining and modeling events within the industrial chain, such as market changes and company mergers.

The industrial chain encompasses numerous types of events, including supply chain event modeling, market event modeling, environmental event modeling, marketing and sales modeling, as well as service and after-sales modeling. Introducing each new event type necessitates designing a new schema, which adds complexity to knowledge extraction and management. To streamline this process, this paper establishes a set of standardized event modeling specifications, as shown in Table 1. This standardized event modeling framework defines uniform modeling patterns for different event types, categorizing event attributes into different categories such as basic elements, temporal elements, spatial elements, subject elements, and object elements.

**Table 1.** Industrial chain event schema model.

| Element Type | Attribute Name | Attribute Type | Value Examples |
|---|---|---|---|
| Basic Elements | Type | String | Product Launch, Mergers and Acquisitions, Supply Chain Disruption |
| | Description | Text | New Model Launch, Company Acquisition, Supplier Bankruptcy |
| | Level | String | High, Medium, Low |
| Temporal Elements | Occurrence_Time | Datetime | 15 January 2023 10:00:00 |
| | Duration | Integer | 2, 5, 10 (units in days) |
| | Deadline | Datetime | 28 February 2023 23:59:59 |

**Table 1.** *Cont.*

| Element Type | Attribute Name | Attribute Type | Value Examples |
|---|---|---|---|
| Spatial Elements | Location | String | Shanghai, Beijing, Shenzhen |
| | Impact_Area | String | Global, Regional, Local |
| | Coordinates | Geographical Coordinates | (31.2304° N, 121.4737° E) |
| Subject Elements | Subject | String | Company, Individual, Government |
| | Participant | String | Investor, Supplier, Customer |
| | Responsible_Party | String | Manufacturer, Distributor, Government Agency |
| Object Elements | Product_Model | String | Tissue Paper |
| | Raw_Material | String | Steel, Semiconductor |
| | Transaction_Amount | Currency | 500,000 |

*4.3. General Common-Sense Semantic Modeling*

The third layer of the pattern layer focuses on general common-sense semantic modeling, which involves modeling the association, attributes, and events of general common sense with specific industrial chains. This allows the system to better understand and infer various relationships within the industrial chain. Such models can cover multiple aspects, including industry interrelations, supply-demand relationships, market trends, and more, ensuring that the knowledge graph not only relies on domain-specific information but also comprehends a broader context.

This paper first determines the scope and objectives of modeling, clarifying the common-sense knowledge content that the system needs to acquire and understand, including the involved industrial chain links, related entities, events, etc. Next, it collects and organizes common-sense knowledge within the industrial chain, including industry standards, expert knowledge, document materials, case analyses, historical data, and other sources. Knowledge is extracted and integrated from various sources and organized into a form that the system can comprehend and utilize. Finally, based on the extracted and integrated common-sense knowledge, a general common-sense model is established.

**5. Construction of the Data Layer of the ICKG**

This paper constructs the data layer of the ICKG from the bottom up through three steps: entity relationship modeling, entity recognition and relation extraction, and knowledge fusion. Entity relationship modeling is the process of defining and modeling entities and their relationships within the industrial chain. This step clarifies the connections between various entities in the industrial chain, establishing a clear framework for subsequent entity recognition and relation extraction, making them more targeted and accurate. The second step involves identifying specific entities and relationships from text. Through this step, textual information is transformed into structured data, enriching the content of the ICKG. Finally, knowledge fusion integrates and combines knowledge obtained from different sources to make the knowledge graph more complete and comprehensive. Knowledge fusion helps address the limitations of a single data source, improving the accuracy and credibility of the knowledge graph.

*5.1. Entity Relationship Modeling*

Entity relationship modeling defines and models entities in the knowledge graph, determines their relationships, and specifies entity attributes. Its aim is to create a structured representation of things and their relationships in the real world to better organize and understand complex data structures. The main goal is to represent complex real-world data structures in a clear and maintainable manner for data storage, querying, and analysis.

The schema architecture in the knowledge graph plays a role in defining and delineating entities, attributes, and relationships. It abstracts and standardizes the data model within the domain, providing the foundation for constructing and applying the knowledge graph. By using types and properties, the schema specifies feasible categories and con-

straints, enabling the knowledge graph to organize, query, and utilize domain knowledge more effectively. The enterprise schema model is defined in Table 2.

**Table 2.** Enterprise schema definition.

| Attribute Name | Attribute Type | Attribute Value |
|:---:|:---:|:---:|
| Id | String | 4873 |
| Code | String | 300636 |
| City | String | Jiangxi Province |
| Fullname | String | China Jiangxi Tonghe Pharmaceutical Co., Ltd. |
| Location | String | Shenzhen Stock Exchange |
| Comp_Name | String | Tonghe Pharmaceutical |
| Reg_Capital | Numeric | 350 million |
| Setup_Date | Datetime | 31 March 2017 |

As shown in Figure 4, the ICKG encompasses five types of entities: products, companies, industries, supply chains, and customers. It defines ten types of entity relationships, including upstream material, higher-level industry, downstream product, product subcategory, core product, business sector, purchase relationship, service relationship, procurement relationship, and partnership relationships.
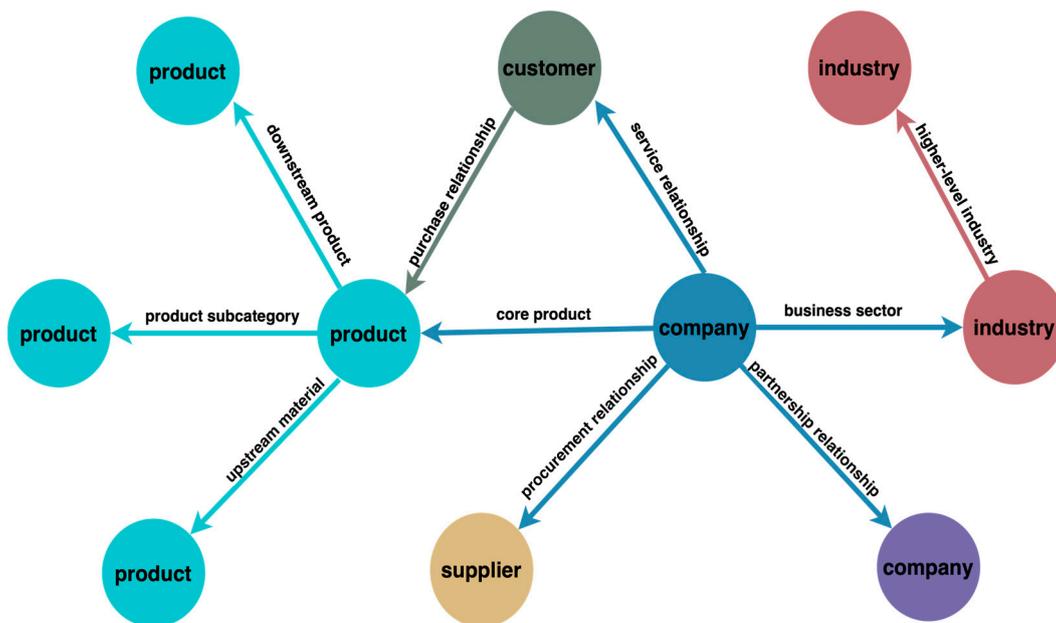


**Figure 4.** Entity and relationship definitions.

## 5.2. Entity Recognition Methods

NER is the process of identifying named entities with specific meanings from text and classifying them into predefined categories. This section analyzes and compares various deep learning algorithms' applications in entity recognition in the industrial chain context from three aspects: semantic representation, context awareness, and sequence labeling.

### 5.2.1. Semantic Representations

The role of semantic representations is to convert textual data into vector representations with semantic information. Its main function is to map natural language text into a high-dimensional vector space, where texts with similar semantics are represented closer to each other in the vector space. In these models, each word or subunit is mapped to a vector, and the entire text is represented by the combination of these vectors.

This paper utilizes the Bidirectional Encoder Representations from Transformers (BERT) model to learn semantic representations of textual data. The BERT model employs bidirectional transformer encoders, which can consider context information simultaneously during the pre-training phase, thereby better capturing the context and semantics of words. This architecture can effectively capture long-distance dependencies, which is advantageous for processing longer text sequences. During pre-training, BERT performs two tasks simultaneously: the Masked Language Model (MLM) and Next Sentence Prediction (NSP). This multi-task learning approach helps the model learn richer semantic representations.

### 5.2.2. Context-Aware Models

Context-aware models refer to models or components used to capture contextual information from textual data. Their main function is to consider the context of textual data before and after, enabling a more comprehensive understanding of the text content. The role of context-aware models is to not only consider the current word or sentence when processing textual data but also take into account the context before and after, leading to more accurate understanding and processing.

This paper utilizes the Bidirectional Long Short-Term Memory (BiLSTM) model to study its application effectiveness in entity recognition in the industrial chain context. Long Short-Term Memory (LSTM) is a variant of Recurrent Neural Network (RNN) that effectively addresses the vanishing gradient and exploding gradient problems in traditional RNNs by introducing internal gate mechanisms, thus better capturing long-term dependencies. BiLSTM extends LSTM by considering both past and future context information at each time step. BiLSTM consists of two LSTMs, one responsible for extracting information from the forward sequence and the other for extracting information from the backward sequence. Their outputs are then concatenated to obtain richer contextual information.

### 5.2.3. Sequence Tagger

The main function of a sequence tagger is to label input sequence data by assigning a label or category to each element in the input sequence (such as words, characters, etc.). Conditional Random Field (CRF) is a probabilistic graphical model commonly used for sequence tagging tasks. CRF models the label sequences in the sequence, considering the dependency relationships between labels, thereby improving the accuracy of sequence tagging. In industrial entity recognition, CRF models can utilize previously extracted features (such as word embeddings, part-of-speech tagging, character-level features, etc.) to predict the label for each word while considering the transition probabilities between labels, thus ensuring the consistency and rationality of the recognition results.

### *5.3. Relation Extraction*

The ICKG constructed in this paper includes five main entity types, including product, company, industry, supplier, and customer. It defines ten types of entity relationships. Firstly, focusing on the industry relationships of listed companies, based on publicly available results of industry classification for listed companies, we construct the business sector of listed companies. These relationships demonstrate clear associations with the industries in which companies operate, providing an important foundation for the establishment of the ICKG. Secondly, according to the latest industry classification table, we establish higher-level industry relationships between industries. These relationships describe the hierarchical structure between different industries, revealing the hierarchical relationships between industries and aiding in better understanding and analyzing the structure and hierarchy of the entire industrial chain. Additionally, by formulating corresponding rules and applying the method of rule pattern matching, we successfully extract rich relationship information from the semi-annual and annual reports published by companies each year. This information provides detailed descriptions of companies' main businesses, the association between products and raw materials, the connections between products and downstream products, and the subdivision of product categories, among other aspects,

thereby providing important information and perspectives for a deeper understanding of the interactions between various links in the industrial chain.

*5.4. Knowledge Fusion*

The ICKG integrates information from different data sources to eliminate duplicates, consolidate identical entities, and ensure the completeness and accuracy of the information in the knowledge graph. Different descriptions of the same entity exist across different data sources, and knowledge fusion eliminates this heterogeneity, thus unifying entity information in the knowledge graph. Various links, entities, and relationships within the industrial chain need to be comprehensively considered, and knowledge fusion helps integrate information from multiple domains or stages to achieve a holistic view of the entire industrial chain. This paper primarily conducts knowledge fusion on industrial chain data from two aspects: semantic knowledge fusion and entity fusion.

5.4.1. Semantic Knowledge Fusion

Semantic knowledge fusion in the industrial chain mainly involves semantic mapping and ontology fusion, which can enhance the consistency of data across different stages of the industrial chain. For instance, when constructing a knowledge graph covering the electronic product supply chain, it includes data such as different suppliers, part information, production stages, and market demands from multiple sources. These data employ different terminologies and standards to describe the same concepts. Supplier A may describe a component as an "aluminum bracket", while Supplier B may describe the same component as an "aluminum support frame". Through semantic mapping, these two descriptions are mapped to a unified semantic model, ensuring consistent descriptions of the components in the knowledge graph and establishing their associations. In the electronic product supply chain, different stages may use different ontologies or knowledge representation models to describe product specifications, production processes, quality standards, etc. Production stages in the supply chain use a specific set of terms and concepts to describe the manufacturing process and technical requirements, while sales stages use different terms. Through ontology fusion, these different ontologies are integrated to establish a common ontology representation model, ensuring consistent and interrelated descriptions of product manufacturing and sales stages in the knowledge graph.

5.4.2. Entity Fusion

Entity fusion refers to identifying the correspondence between descriptions of the same entity in different data sources and aligning them with the same entity. This includes resolving differences in entity names, eliminating multiple descriptions of the same entity, and associating dispersed entity information to ensure consistency of entities in the knowledge graph. For example, in the aviation industry, there are multiple similar entities, such as airlines, aerospace engineering organizations, aviation logistics service providers, and aviation maintenance service providers. These entities describe different aspects of the same industry, but they may differ in naming and description. In the process of constructing the ICKG, it is necessary to standardize entity names, unify naming formats, or use consistent naming conventions, such as unifying "aviation logistics service" and "aviation cargo transportation" into "aviation logistics". Next, text similarity algorithms are used to analyze the descriptions of these entities to determine their degree of similarity. Entities with high similarity in descriptions are identified and associated with the same entity, indicating that they refer to the same domain or service. Finally, integrate other attribute information about these entities, such as service scope, technical characteristics, customer base, etc., to ensure they are associated with the same entity. Such entity alignment and linking processes ensure the consistency of entities in relevant domains of the knowledge graph, providing a consistent and reliable entity information foundation for comprehensive analysis of the industrial chain. Table 3 below shows entity fusion between some similar industries when constructing the ICKG.

**Table 3.** Schematic representation of closely related industries.

| Industry | Related Industry |
|---|---|
| Solar Cells | Wind Power Generation, Hydroelectric Power Generation |
| Medical Devices | Pharmaceutical Manufacturing, Medical Equipment Components |
| Artificial Intelligence (AI) | Machine Learning, Deep Learning |
| Catering Industry | Fast Food Chain, Bar, and Entertainment |
| E-commerce | Internet Finance, Online Payment |
| Environmental Technology | Waste Management, Renewable Energy |
| Aviation Industry | Aerospace Engineering, Air Logistics |
| Human Resources (HR) | Training and Development, Recruitment Services |

## 6. Experiment

This section discusses the application of deep learning algorithms to entity recognition within industrial chains. Through deep learning models, the automatic identification and extraction of entity information relevant to specific industrial chains from text data are achieved, laying the foundation for constructing an ICKG.

### 6.1. Experimental Data and Experimental Environment

We utilized web crawling technology to collect various data sources related to industrial chains, including structured data from industries such as the Shenwan Industry Index, the Shenzhen Stock Exchange, and the Shanghai Stock Exchange, as well as unstructured data extracted from sources like encyclopedias, news, and annual reports. Our dataset contains structured information from 9300 companies, including fields such as comp_name, ts_code, chairman, manager, secretary, reg_capital, setup_date, province, city, website, email, office, employees, etc. Additionally, we obtained unstructured text data from these companies, including introductions, business scopes, main businesses, and company essentials, totaling 37,200 texts. We conducted experiments using the Python language and employed the BIO annotation method to label the positions of entity names such as companies and products in the text. In the BIO labeling method, 'B' is used to mark the beginning word of an entity, 'I' is used to mark words inside the entity, and 'O' is used to mark words that do not belong to any entity. We manually annotated 108,462 sentences from the crawled unstructured text, randomly selecting 80% as the training set, 10% as the validation set, and 10% as the test set.

The experimental environment utilized a computational server equipped with NVIDIA GeForce RTX 4090 and NVIDIA GeForce RTX 4060 GPUs sourced from NVIDIA Corporation, located in Santa Clara, CA, USA. TensorFlow 2.5.0, a deep learning framework, and the Python 3.8.17 programming language were used for entity extraction experiments.

In industrial chain entity recognition experiments, adjusting experimental parameters is aimed at optimizing model performance and training process efficiency. These parameters include the learning rate, number of iterations, maximum sequence length, and batch size. The learning rate controls the speed of parameter updates, and an appropriate learning rate can accelerate model convergence and improve performance. The choice of iteration number depends on factors such as dataset size, model complexity, and training time. Typically, increasing the number of epochs can enhance model performance but may also increase the risk of overfitting. The maximum sequence length limits the length of input sequences, helping to save computational resources and avoid inefficiencies in handling excessively long sequences. Batch size determines the number of training samples in each iteration. With the initial learning rate set to $2 \times 10^{-5}$, the maximum sequence length set to 10, and the batch size set to 16, by iteratively adjusting these parameters, the goal is to enhance the accuracy and training speed of the model in entity recognition tasks to achieve optimal results while avoiding resource waste and overfitting risk.

### 6.2. Deep Learning Method Selection and Model Evaluation Criteria

BERT, as a base model, possesses strong language representation capabilities and excels in natural language processing tasks, effectively addressing named entity recognition

tasks. CRF and BiLSTM-CRF, as commonly used sequence labeling models, demonstrate excellent performance in named entity recognition tasks, effectively considering contextual information and entity relationships. By combining BERT with CRF or BiLSTM-CRF, one can fully leverage BERT's understanding of contextual nuances and the modeling capabilities of sequence labeling models to enhance the accuracy and robustness of entity extraction. This paper selects three models, namely BERT, BERT-CRF, and BERT-BiLSTM-CRF, to investigate the application effectiveness of deep learning algorithms in industrial chain named entity recognition tasks. Through multiple rounds of comparative experiments, a thorough evaluation of different models' performance in industrial chain named entity recognition is conducted. The selection and framework of entity extraction models are shown in Figure 5.
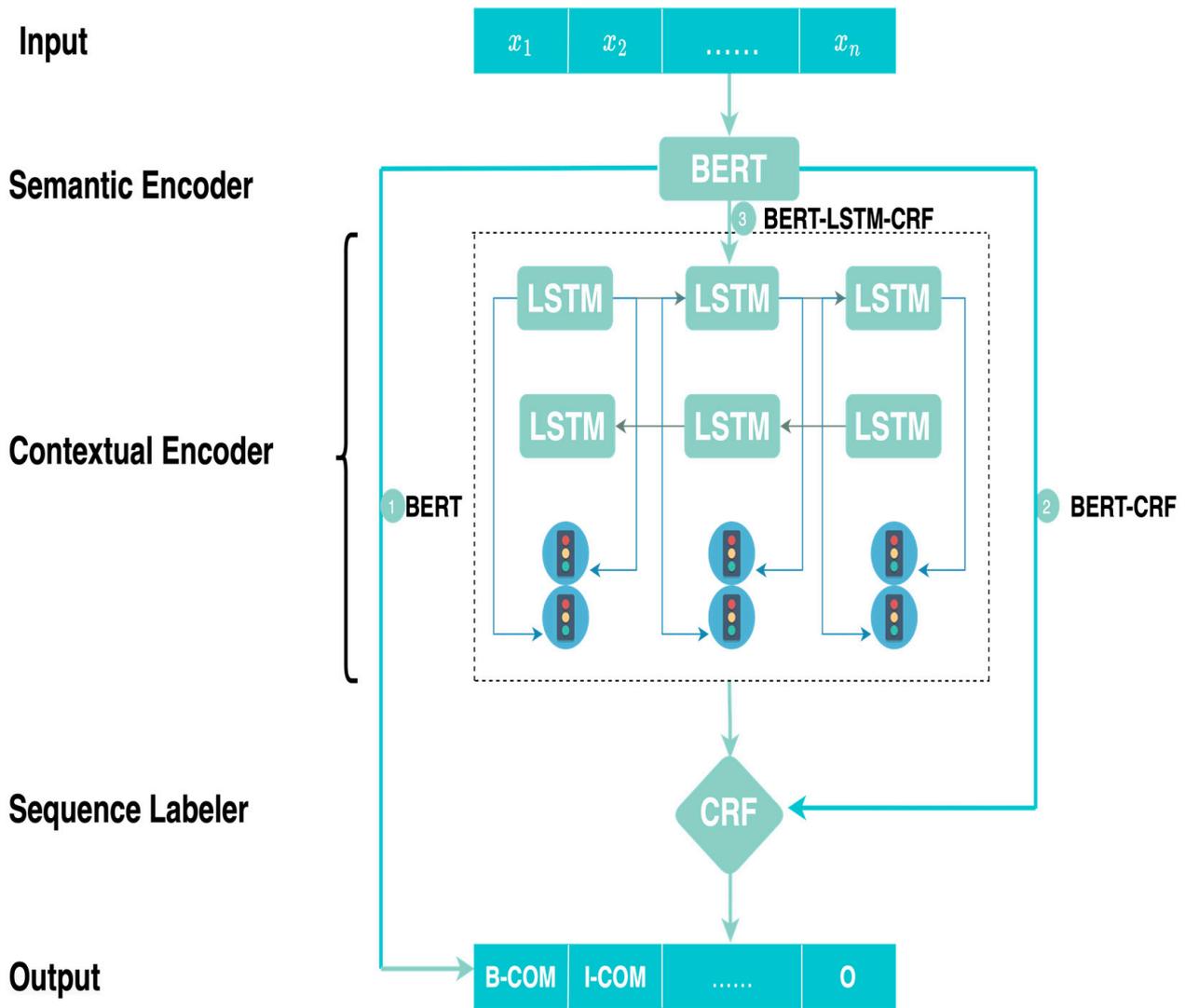


**Figure 5.** Model selection and framework.

This study selects macro-averaging as the evaluation metric for the industrial chain entity recognition experiments. Macro-averaging balances the influence of different categories by considering the performance of each category and averaging them. In this multi-classification problem, macro-averaging can independently evaluate the performance of each category, providing a better understanding of the performance of each category. The evaluation metrics include macro-average F1 score (Macro-F1), macro-average precision (Macro-P), and macro-average recall (Macro-R). Where Macro-P, Macro-R, and Macro-F1

are, respectively, the arithmetic averages of Precision (*P*), Recall (*R*), and F1-score (*F1*). The formulas for calculating *P*, *R*, and *F*1 are as follows (Equations (1)–(3)):

$$P = \frac{TP}{TP + FP} \tag{1}$$

$$R = \frac{TP}{TP + FN} \tag{2}$$

$$F1 = \frac{2 * R * P}{R + P} \tag{3}$$

Among them, *TP* refers to the number of samples that the model correctly identifies as positive examples, that is, the number of positive samples for which the entity recognition is correct. *FP* refers to the number of samples in which the model incorrectly identifies negative examples as positive examples, that is, the number of positive samples that the model misjudges. *FN* refers to the number of samples in which the model incorrectly identifies positive examples as negative examples, that is, the number of positive samples that the model misses.

*6.3. Experimental Results and Discussion*

This paper conducted comparative experiments among BERT, BERT-CRF, and BERT-BiLSTM-CRF and obtained the experimental results shown in Table 4 after multiple parameter adjustments. Based on the experimental results in the table, we can draw the following conclusions:

Using different learning rates, the performance of the BERT model does not vary significantly, but overall performance is good, with macro-F1 scores ranging from 0.9666 to 0.9695. macro-P and macro-R scores of the BERT model both remain at a high level. The BERT-CRF model, which combines BERT with CRF, performs significantly lower in the F1 score compared to the pure BERT model. The BERT-BiLSTM-CRF model further extends the BERT-CRF model by adding BiLSTM and fine-tuning the learning rate of CRF. This model achieves a macro-F1 score of 0.9710, an increase of 0.44% compared to the BERT model and 24.78% compared to the BERT-CRF model. However, there is a slight decrease in macro-P and macro-R compared to the pure BERT model, although they still remain at a high level.

**Table 4.** Experimental results of entity recognition for different models.

| Model | Macro-F1 | Macro-P | Macro-R | Learning Rate | Epoch | Maxlen | Batch_Size |
|---|---|---|---|---|---|---|---|
| | 0.9674 | 0.9735 | 0.9709 | $3 \times 10^{-5}$ | 3 | 10 | 16 |
| BERT | 0.9695 | 0.9735 | 0.9715 | $4 \times 10^{-5}$ | 3 | 10 | 16 |
| | 0.9666 | 0.9693 | 0.9679 | $5 \times 10^{-5}$ | 3 | 10 | 16 |
| BERT-CRF | 0.7232 | 0.7658 | 0.8611 | $5 \times 10^{-5}$ (bert), $3 \times 10^{-3}$ (crf) | 3 | 10 | 16 |
| BERT-BiLSTM-CRF | 0.9710 | 0.9680 | 0.9695 | $2 \times 10^{-5}$ (bert), $3 \times 10^{-3}$ (crf) | 3 | 10 | 16 |
| | 0.9710 | 0.9671 | 0.9690 | $3 \times 10^{-5}$ (bert), $3 \times 10^{-3}$ (crf) | 3 | 10 | 16 |

Overall, the BERT-BiLSTM-CRF model achieves the best performance, with an F1 score higher than the other two models, and it also performs well in terms of precision and recall, making it suitable for the entity recognition task of industrial chain naming addressed in this paper. In this model, BERT learns bidirectional contextual relationships between words to generate word-level contextual representations. BiLSTM has two directions of hidden states, one for reading sequences from left to right and the other for reading sequences from right to left, which enables capturing semantic information over longer distances.

CRF is a sequence labeling model that, based on the output of the BiLSTM layer, labels the entire sequence and considers dependencies between labels, thereby improving the accuracy of labeling.

The advantages of the BERT-BiLSTM-CRF model in the task of entity name recognition in the industrial chain mainly include the following points: First, because the industrial chain involves multiple fields and professional terms, the semantic representation learned by BERT during pre-training can better understand the terms and relationships in specific fields, thereby improving the accuracy of entity name recognition, especially when dealing with proprietary terms and abbreviations in the industrial chain. Secondly, BiLSTM can capture the influence of context on entity names, while BERT provides a more comprehensive understanding of context, further enhancing the precision of entity recognition. In addition, the CRF model globally models the labeled sequences, considering the dependency relationships between entity labels, thus improving consistency and accuracy, especially when dealing with complex entity relationships in the industrial chain. Finally, the BERT-BiLSTM-CRF model has strong versatility and generalization ability, adapting to entity name recognition tasks in different industrial chain scenarios and providing robust performance.

## 7. ICKG Instance Display

Through complex data processing and extraction work, a large-scale knowledge graph has been successfully established, with data stored in the graph database, Neo4j. This knowledge graph encompasses hundreds of thousands of entities and their relationships within the industrial chain and has been visualized. A partial display of the ICKG is shown in Figure 6.
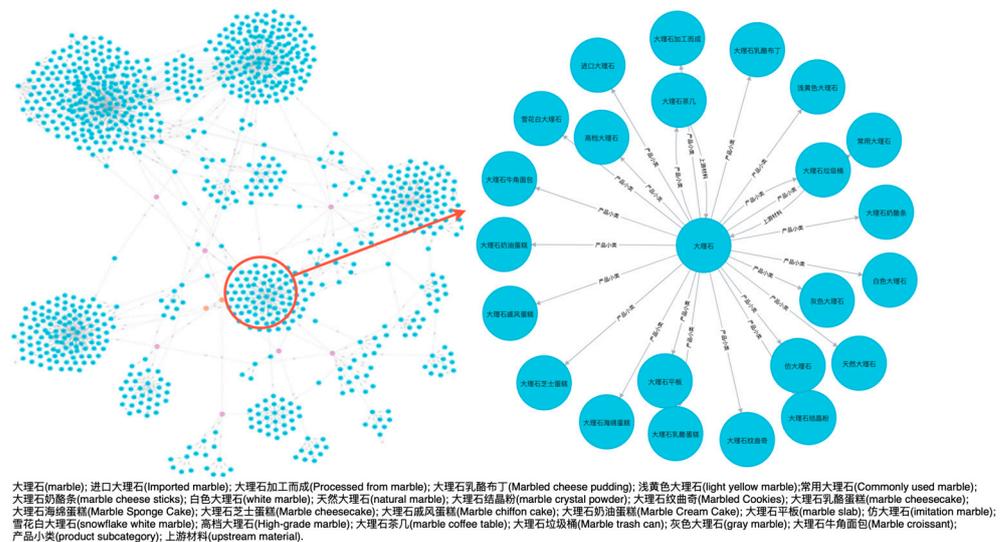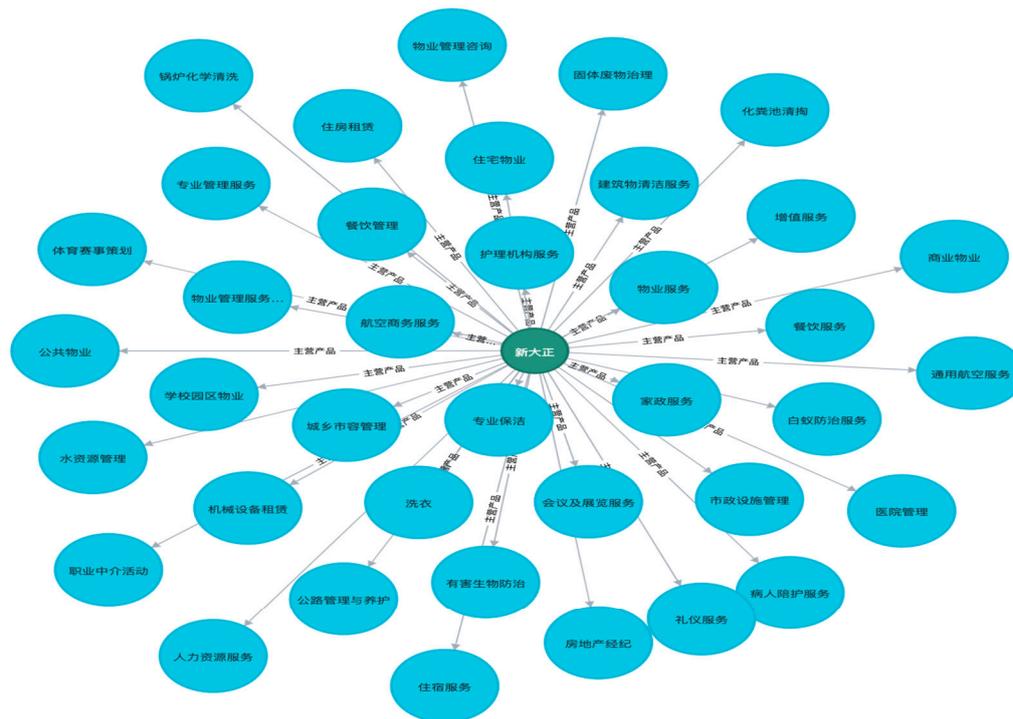


大理石(marble); 进口大理石(Imported marble); 大理石加工而成(Processed from marble); 大理石乳酪布丁(Marbled cheese pudding); 浅黄色大理石(light yellow marble);常用大理石(Commonly used marble);
大理石奶酪条(marble cheese sticks); 白色大理石(white marble); 天然大理石(natural marble); 大理石结晶粉(marble crystal powder); 大理石纹曲奇(Marbled Cookies); 大理石乳酪蛋糕(marble cheesecake);
大理石海绵蛋糕(Marble Sponge Cake); 大理石芝士蛋糕(Marble cheesecake); 大理石戚风蛋糕(Marble chiffon cake); 大理石奶油蛋糕(Marble Cream Cake); 大理石平板(marble slab); 仿大理石(imitation marble);
雪花白大理石(snowflake white marble); 高档大理石(High-grade marble); 大理石茶几(marble coffee table); 大理石垃圾桶(Marble trash can); 灰色大理石(gray marble); 大理石牛角面包(Marble croissant);
产品小类(product subcategory); 上游材料(upstream material).

**Figure 6.** Cross-domain knowledge graph (partial).

The query statement is entered in the graph database Neo4j, and the ICKG system queries the entity information through the semantic network built in the knowledge graph and visually displays the knowledge structure of this type. Figure 7 displays all the main products of China Chongqing Xinda Zheng Property Group Co., Ltd. Information about the main products can be used in various internal applications of the enterprise, such as marketing, product development, and supply chain management. Understanding a company's main products helps analyze its market position, competitive advantages, and product positioning, providing support for business decisions.

新大正(Xinda Zheng); 锅炉化学清洗(Boiler chemical cleaning); 物业管理咨询(Property management consulting); 固体废物治理(solid waste management); 化粪池清掏(Septic tank cleaning); 专业管理服务(Professional management services); 住房租贷(housing rental loan); 住宅物业(residential property); 建筑物清洁服务(building cleaning services); 增值服务(Value-added services); 体育赛事策划(Sports event planning); 餐饮管理(Catering management); 护理机构服务(nursing facility services); 物业服务(commercial service); 商业物业(commercial property); 公共物业(public property); 物业管理服务(Property management services); 航空商务服务(Aviation business services); 餐饮服务(Catering Services); 水资源管理(water resources management); 学校园区物业(School campus property); 城乡市容管理(Urban and rural city appearance management); 专业保洁(Professional cleaning); 家政服务(Housekeeping); 白蚁防治服务(Termite control services); 通用航空(General aviation services); 职业中介活动(Employment agency activities); 机械设备租赁(Machinery and equipment rental); 洗衣(Laundry); 会议及展览服务(Conference and exhibition services); 市政设施管理(Municipal facilities management); 医院管理(hospital management); 人力资源服务(Human resources services); 公路管理与养护(Highway management and maintenance); 有害生物防治(pest control); 住宿服务(Accommodation services); 房地产经纪(Real estate brokerage); 礼仪服务(ceremonial service); 病人陪护服务(Patient escort service); 主营产品(core product)

**Figure 7.** ICKG query display.

## 8. Conclusions

This paper utilized both structured data from sources like the Shenwan Industry Index, Shenzhen Stock Exchange, Shanghai Stock Exchange, etc., and unstructured data from sources like encyclopedias, news articles, annual reports, etc., to construct a comprehensive dataset for the industrial chain. Through data processing and structured extraction techniques, we built a comprehensive industrial chain dataset. After analyzing the performance of various deep learning algorithms in industrial chain NER, we selected an entity recognition algorithm based on BERT-BiLSTM-CRF and rule-based relationship extraction for industrial chain knowledge extraction. Regarding graph generation and visualization, we constructed the framework as shown in the figure through node and relationship construction and stored the data in the Neo4j graph database, successfully creating a comprehensive ICKG. This knowledge graph provides strong support for cross-disciplinary applications, including knowledge retrieval, intelligent question answering, smart decision-making, intelligent marketing, and recommendations.

Despite constructing a full-scale ICKG and providing support for cross-disciplinary applications, there are still some shortcomings and areas for further improvement. The following aspects are considered for improvement in future work:

- Enhancing Precision: Further optimize the entity recognition algorithm and relationship extraction rules to improve the recall and precision of knowledge extraction. Enhance its generalization ability to adapt to entity and relationship extraction tasks in different contexts.
- Data Completeness and Updates: Ensure the completeness of knowledge graph data and regularly update it. Continuously collect the latest data and promptly incorporate it into the knowledge graph to maintain its timeliness and comprehensiveness.

- Optimization of Intelligent Application Functions: Further develop and optimize intelligent application functions based on the knowledge graph to enhance the intelligence level in areas such as knowledge retrieval, intelligent question answering, decision support, marketing, and recommendations.
- User Experience and Interface Optimization: Design user-friendly interfaces and interactive experiences to make it easier for users to access and utilize the information and functionalities provided by the knowledge graph.

## References

1. Lou, Q.; Xin, T.Y.; Song, J.Y. Application of DNA barcoding technology in the whole industrial chain of traditional Chinese medicine. *Acta Pharm. Sin.* **2020**, *12*, 1784–1791.
2. Lu, H.; Peng, J.; Lu, X. Do Factor Market Distortions and Carbon Dioxide Emissions Distort Energy Industry Chain Technical Efficiency? A Heterogeneous Stochastic Frontier Analysis. *Energies* **2022**, *15*, 6154. [CrossRef]
3. Yang, Y.; Tong, L.; Yin, S.; Liu, Y.; Wang, L.; Qiu, Y.; Ding, Y. Status and Challenges of Applications and Industry Chain Technologies of Hydrogen in the Context of Carbon Neutrality. *J. Clean. Prod.* **2022**, *376*, 134347. [CrossRef]
4. Song, H.; Lu, B.; Ye, C.; Li, J.; Zhu, Z.; Zheng, L. Fraud Vulnerability Quantitative Assessment of Wuchang Rice Industrial Chain in China Based on AHP-EWM and ANN Methods. *Food Res. Int.* **2021**, *140*, 109805. [CrossRef] [PubMed]
5. Haizhong, A.N.; Huajiao, L.I. Theory and Research Advances in Whole Industrial Chain of Strategic Mineral Resources. *Resour. Ind.* **2022**, *24*, 8.
6. Chen, X.; Jia, S.; Xiang, Y. A Review: Knowledge Reasoning over Knowledge Graph. *Expert Syst. Appl.* **2020**, *141*, 112948. [CrossRef]
7. Fu, D.; Zhou, S.; Shen, B.; Chen, Y. Enhancing Semantic Search of Crowdsourcing IT Services Using Knowledge Graph. In Proceedings of the International Conference on Software Engineering and Knowledge Engineering, Lisbon, Portugal, 10 July 2019.
8. Sarrafzadeh, B.; Vechtomova, O.; Jokic, V. Exploring Knowledge Graphs for Exploratory Search. In Proceedings of the 5th Information Interaction in Context Symposium, Regensburg, Germany, 26 August 2014.
9. Fernández, M.; Cantador, I.; López, V.; Vallet, D.; Castells, P.; Motta, E. Semantically Enhanced Information Retrieval: An Ontology-Based Approach. *J. Web. Semant.* **2011**, *9*, 434–452. [CrossRef]
10. Wang, Q.; Li, M.; Wang, X.; Parulian, N.; Han, G.; Ma, J.; Tu, J.; Lin, Y.; Zhang, R.H.; Liu, W.; et al. COVID-19 Literature Knowledge Graph Construction and Drug Repurposing Report Generation. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Demonstrations, Online, 6–11 June 2021.
11. Yu, T.; Li, J.; Yu, Q.; Tian, Y.; Shun, X.; Xu, L.; Zhu, L.; Gao, H. Knowledge Graph for TCM Health Preservation: Design, Construction, and Applications. *Artif. Intell. Med.* **2017**, *77*, 48–52. [CrossRef] [PubMed]
12. Shi, L.; Li, S.; Yang, X.; Qi, J.; Pan, G.; Zhou, B. Semantic Health Knowledge Graph: Semantic Integration of Heterogeneous Medical Knowledge and Services. *BioMed Res. Int.* **2017**, *2017*, 1–12. [CrossRef] [PubMed]
13. Guo, Q.; Zhuang, F.; Qin, C.; Zhu, H.; Xie, X.; Xiong, H.; He, Q. A Survey on Knowledge Graph-Based Recommender Systems. *IEEE Trans. Knowl. Data Eng.* **2022**, *34*, 3549–3568. [CrossRef]
14. Wang, H.; Zhang, F.; Wang, J.; Zhao, M.; Li, W.; Xie, X.; Guo, M. RippleNet: Propagating User Preferences on the Knowledge Graph for Recommender Systems. In Proceedings of the 27th ACM International Conference on Information and Knowledge Management, Torino, Italy, 22–26 October 2018.
15. Wang, H.; Zhao, M.; Xie, X.; Li, W.; Guo, M. Knowledge Graph Convolutional Networks for Recommender Systems. In Proceedings of the World Wide Web Conference, San Francisco, CA, USA, 13–17 May 2019.
16. Tan, J.; Qiu, Q.; Guo, W.; Li, T. Research on the Construction of a Knowledge Graph and Knowledge Reasoning Model in the Field of Urban Traffic. *Sustainability* **2021**, *13*, 3191. [CrossRef]

17. Ahmed, U.; Srivastava, G.; Djenouri, Y.; Lin, J.C.W. Knowledge Graph Based Trajectory Outlier Detection in Sustainable Smart Cities. *Sustain. Cities Soc.* **2022**, *78*, 103580. [CrossRef]

18. Luettin, J.; Monka, S.; Henson, C.; Halilaj, L. A Survey on Knowledge Graph-Based Methods for Automated Driving. In Proceedings of the Knowledge Graphs and Semantic Web, Madrid, Spain, 21–23 November 2022.

19. Miao, R.; Zhang, X.; Yan, H.; Chen, C. A Dynamic Financial Knowledge Graph Based on Reinforcement Learning and Transfer Learning. In Proceedings of the 2019 IEEE International Conference on Big Data (Big Data), Los Angeles, CA, USA, 9–12 December 2019.

20. Zehra, S.; Mohsin, S.F.M.; Wasi, S.; Jami, S.I.; Siddiqui, M.S.; Syed, M.K.U.R.R. Financial Knowledge Graph Based Financial Report Query System. *IEEE Access* **2021**, *9*, 69766–69782. [CrossRef]

21. Yang, B.; Liao, Y. Research on Enterprise Risk Knowledge Graph Based on Multi-Source Data Fusion. *Neural Comput. Appl.* **2022**, *34*, 2569–2582. [CrossRef]

22. Gundlach, G.T.; Bolumole, Y.A.; Eltantawy, R.A.; Frankel, R. The Changing Landscape of Supply Chain Management, Marketing Channels of Distribution. *Logist. Purchasing. J. Bus. Ind. Mark.* **2006**, *21*, 428–438. [CrossRef]

23. Mao, D.; Wang, F.; Hao, Z.; Li, H. Credit Evaluation System Based on Blockchain for Multiple Stakeholders in the Food Supply Chain. *Int. J. Environ. Res. Public Health* **2018**, *15*, 1627. [CrossRef] [PubMed]

24. Agrawal, T.K.; Kumar, V.; Pal, R.; Wang, L.; Chen, Y. Blockchain-Based Framework for Supply Chain Traceability: A Case Example of Textile and Clothing Industry. *Comput. Ind. Eng.* **2021**, *154*, 107130. [CrossRef]

25. Nasar, Z.; Jaffry, S.W.; Malik, M.K. Named Entity Recognition and Relation Extraction: State-of-the-Art. *ACM Comput. Surv.* **2021**, *54*, 20:1–20:39. [CrossRef]

26. Avinadav, T.; Shamir, N. The Effect of Information Asymmetry on Ordering and Capacity Decisions in Supply Chains. *Eur. J. Oper. Res.* **2021**, *292*, 562–578. [CrossRef]

27. Peck, H. Reconciling Supply Chain Vulnerability, Risk and Supply Chain Management. *Int. J. Logist. Res. Appl.* **2006**, *9*, 127–142. [CrossRef]

28. Al-Moslmi, T.; Gallofré Ocaña, M.L.; Opdahl, A.; Veres, C. Named Entity Extraction for Knowledge Graphs: A Literature Overview. *IEEE Access* **2020**, *8*, 32862–32881. [CrossRef]

29. Sun, J.; Liu, Y.; Cui, J.; He, H. Deep Learning-Based Methods for Natural Hazard Named Entity Recognition. *Sci. Rep.* **2022**, *12*, 4598. [CrossRef] [PubMed]

30. Gallofré Ocaña, M.; Opdahl, A. Supporting Newsrooms with Journalistic Knowledge Graph Platforms: Current State and Future Directions. *Technologies* **2022**, *10*, 68. [CrossRef]

31. Evtimova-Gardar, M.; Mellouli, N. An Overview of Methods and Tools for Extraction of Knowledge for COVID-19 from Knowledge Graphs. In Proceedings of the Pattern Recognition and Artificial Intelligence, Paris, France, 1–3 June 2022.

32. Braşoveanu, A.M.P.; Andonie, R. Integrating Machine Learning Techniques in Semantic Fake News Detection. *Neural Process. Lett.* **2021**, *53*, 3055–3072. [CrossRef]

33. Shi, P.; Lin, J. Simple BERT Models for Relation Extraction and Semantic Role Labeling. *arXiv* **2019**, arXiv:1904.05255.

34. Stapley, B.J.; Benoit, G. Biobibliometrics: Information retrieval and visualization from co-occurrences of gene names in Medline abstracts. In Proceedings of the Biocomputing 2000, Honolulu, HI, USA, 4–9 January 2000.

35. Jenssen, T.K.; Lægreid, A.; Komorowski, J.; Hovig, E. A Literature Network of Human Genes for High-Throughput Analysis of Gene Expression. *Nat. Genet* **2001**, *28*, 21–28. [CrossRef] [PubMed]

36. Bunescu, R.; Ge, R.; Kate, R.J.; Marcotte, E.M.; Mooney, R.J.; Ramani, A.K.; Wong, Y.W. Comparative Experiments on Learning Information Extractors for Proteins and Their Interactions. *Artif. Intell. Med.* **2005**, *33*, 139–155. [CrossRef] [PubMed]

37. Raychaudhuri, S.; Chang, J.T.; Sutphin, P.D.; Altman, R.B. Associating Genes with Gene Ontology Codes Using a Maximum Entropy Analysis of Biomedical Literature. *Genome Res.* **2002**, *12*, 203–214. [CrossRef] [PubMed]

38. Leonard, J.E.; Colombe, J.B.; Levy, J.L. Finding Relevant References to Genes and Proteins in Medline Using a Bayesian Approach. *Bioinformatics* **2002**, *18*, 1515–1522. [CrossRef] [PubMed]

39. Choi, S.P. Extraction of Protein–Protein Interactions (PPIs) from the Literature by Deep Convolutional Neural Networks with Various Feature Embeddings. *J. Inf. Sci.* **2018**, *44*, 60–73. [CrossRef]

40. Zhao, Z.; Yang, Z.; Lin, H.; Wang, J.; Gao, S. A Protein-Protein Interaction Extraction Approach Based on Deep Neural Network. *IJDMB* **2016**, *15*, 145. [CrossRef]