

Ontology-Based Knowledge Graph Framework for Industrial Standard Documents via Hierarchical and Propositional Structuring

Jiin Park*
jiinpark@hanyang.ac.kr
Hanyang University
Seoul, Republic of Korea

Hyuna Jeon*
hyunajeon@hanyang.ac.kr
Hanyang University
Seoul, Republic of Korea

Yoonseo Lee
yoonseolee@hanyang.ac.kr
Hanyang University
Seoul, Republic of Korea

Jisu Hong
jisu.hong@themiraclesoft.com
The Miraclesoft
Seoul, Republic of Korea

Misuk Kim†
misukkim@hanyang.ac.kr
Hanyang University
Seoul, Republic of Korea

Abstract

Ontology-based knowledge graph (KG) construction is a core technology that enables multidimensional understanding and advanced reasoning over domain knowledge. Industrial standards, in particular, contain extensive technical information and complex rules presented in highly structured formats that combine tables, scopes of application, constraints, exceptions, and numerical calculations, making KG construction especially challenging. In this study, we propose a method that organizes such documents into a hierarchical semantic structure, decomposes sentences and tables into atomic propositions derived from conditional and numerical rules, and integrates them into an ontology-knowledge graph through LLM-based triple extraction. Our approach captures both the hierarchical and logical structures of documents, effectively representing domain-specific semantics that conventional methods fail to reflect. To verify its effectiveness, we constructed rule, table, and multi-hop QA datasets, as well as a toxic clause detection dataset, from industrial standards, and implemented an ontology-aware KG-RAG framework for comparative evaluation. Experimental results show that our method achieves significant performance improvements across all QA types compared to existing KG-RAG approaches. This study demonstrates that reliable and scalable knowledge representation is feasible even for industrial documents with intertwined conditions, constraints, and scopes, contributing to future domain-specific RAG development and intelligent document management. Code is available at: https://anonymous.4open.science/r/ontology_based_kg_paper

*Both authors contributed equally to this research.

†Corresponding Author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
Conference acronym 'XX, Woodstock, NY

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-XXXX-X/2018/06
<https://doi.org/XXXXXXX.XXXXXXX>

CCS Concepts

• Computing methodologies → Knowledge representation and reasoning.

Keywords

Ontology, Knowledge Graph, Industrial Standard, Retrieval Augmented Generation, Large Language Model

ACM Reference Format:

Jiin Park, Hyuna Jeon, Yoonseo Lee, Jisu Hong, and Misuk Kim. 2018. Ontology-Based Knowledge Graph Framework for Industrial Standard Documents via Hierarchical and Propositional Structuring. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 13 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 Introduction

Industrial standards are essential specifications used across industries, including marine and shipbuilding, to ensure consistency and quality in materials and manufacturing processes. These documents encompass thousands of interrelated clauses that include tables, scopes of application, constraints, exceptions, and numerical specifications, making interpretation challenging. In particular, toxic clauses or differences in interpretation can lead to disputes, which therefore require careful examination. Currently, domain experts manually compare and review these documents, resulting in high costs and inefficiency. Therefore, there is an increasing need for systematic and automated methods to structure and utilize industrial standards. A promising solution is the use of ontology and knowledge graphs (KGs). Ontology formally represents domain concepts and relationships to organize and reason over complex knowledge [14], while a knowledge graph (KG) extends these ontological concepts and relations into concrete entities and facts represented in a graph structure [21]. In other words, ontology provides the conceptual framework of knowledge, and KG connects it with real-world facts, enabling exploration, reasoning, and application. By integrating the two, industrial standards can be processed more efficiently and understood more intuitively through a knowledge-based system. Thanks to these properties, ontology and KGs have become powerful tools in applications such as semantic search, recommendation systems, and legal document analysis [28].

Traditional ontology learning studies have combined subtasks such as concept discovery and relation extraction to structure domain knowledge [2]. With the recent advancement of large language models (LLMs), there have been growing attempts to automate these processes; however, such approaches remain limited in capturing the complex structures inherent in industrial standards. [3]. Moreover, conventional KG construction methods often suffer from redundancy, sparsity, and inconsistency caused by synonymy and expression variability which ultimately degrade KG embedding and retrieval-augmented generation (RAG) performance [24, 33]. While general-purpose KGs such as Wikidata, DBpedia, and YAGO have been widely used for information retrieval and reasoning, they are limited in representing domain-specific rules and table-based knowledge. Therefore, a reliable and consistent ontology-based KG construction methodology is required for highly structured technical documents such as industrial standards.

Against this background, this study aims to answer the following question: **Can the complex rules and tabular structures of industrial standards be effectively represented in an ontology-based KG framework?** To this end, we propose a novel method that first organizes documents into a hierarchical semantic structure, subsequently constructs a KG by extracting triples from sentence-level ontologies that reflect rule conditions, and ultimately establishes explicit links between the ontology and the KG. Unlike conventional approaches, our framework captures domain-specific semantics, enhances table and conditional reasoning, and supports practical applications such as automated toxic clause detection. However, building such a framework entails several technical challenges:

- **Challenge 1:** Industrial standards often exhibit highly complex and non-standardized formats, making it difficult to accurately parse their hierarchical structures and consistently normalize both text and tables.
- **Challenge 2:** A reliable method is required to extract meaningful triples from conditional rules (if–then), scopes of application and exceptions, numerical conversions, and complex tables inherent in industrial documents.
- **Challenge 3:** To align the triples extracted by LLMs with the ontology structure, a refinement process is needed, including synonym normalization and relation pruning to remove redundant or irrelevant links.

To address these challenges, we propose a three-stage pipeline: (1) **Hierarchical Ontology Modeling**, which maps the hierarchical structure of industrial standards onto an ontology schema; (2) **Atomic Proposition Ontology Modeling**, which decomposes conditional, scope, and numerical clauses into atomic propositions and represents them as triples using an LLM; and (3) **Ontology-based Knowledge Graph Construction**, which links the extracted triples to the ontology schema and refines entities and relations through synonym normalization and pruning to ensure consistency and efficiency.

To verify the effectiveness of the proposed pipeline, we constructed an ontology-based KG using the industrial standards ASTM

A578/A578M¹, API_2W² and ASTM A6/A6M³. Based on these resources, we constructed a QA dataset including rule, table, and multi-hop types, as well as a toxic clause dataset, and implemented an ontology-aware KG-RAG framework. Finally, we conducted comparative experiments with LLM only, LLM with full document, BM25/Ada-002 based RAG, and existing KG-RAG [6, 31] to evaluate its performance.

In summary, the main contributions of this study are as follows:

- We developed an ontology construction method tailored for industrial standards and devised a way to extract triples from textual elements such as constraints, exceptions, scopes of application, and tables. By integrating these components, we propose a comprehensive ontology-based KG construction methodology.
- We built new QA datasets—including rule, table, and multi-hop types—as well as a toxic clause dataset derived from industrial standards, which can serve as benchmarks for future research.
- We implemented an ontology-aware KG-RAG framework and demonstrated its effectiveness through comparative experiments with various retrieval and reasoning baselines.

2 Related Works

Ontology ranges from simple concept collections to complex reasoning systems that encompass relations, constraints, and axioms. Representative examples include WordNet [20], which contains 117,659 concepts and 89,089 hierarchical relations, and Gene Ontology [1], which contains 42,255 concepts and 66,810 relations. Ontology Learning aims to automatically construct ontology structures. Early studies addressed the limitations of manual ontology construction by proposing semi-automatic and fully automatic approaches that integrate statistical analysis, rule-based parsing, and machine learning techniques [5, 7, 15]. Nevertheless, challenges in scalability and quality remain, and recent research has increasingly explored the use of LLMs. [12] leveraged GPT-3.5 to iteratively expand an “is-a” hierarchy without predefined schemas, demonstrating the potential of LLMs in ontology construction, although limited to a single relation type. Similarly, [18] proposed an end-to-end framework that generates the overall ontology structure, but it focuses solely on taxonomic relations and overlooks non-taxonomic relations.

Early research on KG construction primarily relied on rule-based or statistical methods to extract triples from web text. Systems such as KnowItAll [11] and TextRunner [32] demonstrated scalability in large-scale data processing but exhibited notable limitations in generalization. Subsequently, neural network-based models such as CasRel [30] and GPNN [25] emerged, which refined relation extraction and syntactic parsing. However, these approaches still required extensive annotation and domain-specific retraining. More recently, large language models (LLMs) have been actively explored as a means to alleviate these constraints. Leveraging the vast linguistic knowledge acquired during pretraining, LLMs can perform

¹<https://webstore.ansi.org/standards/astm/astma578a578m17>

²https://store.accuristech.com/standards/api-spec-2w?product_id=2032271&srsId=AfmBOooF1rn33A4SZWUXIamtWFgBlmp41eyVnRjGZVAZdH7V61Bjts76

³<https://webstore.ansi.org/standards/astm/astma6a6m24b>

KG construction without requiring large-scale domain-specific annotations. For example, REBEL [17] proposed a T5-based model that directly generates triples from text, while EDC [34], LLM-TIKG [16], and ACKG-LLM [29] extended this idea through prompt engineering and adaptive learning across various domains. Collectively, these studies demonstrate that LLMs enable an efficient and generalizable paradigm for KG construction.

3 Ontology Construction

The ontology construction pipeline comprises two stages: (1) Hierarchical Ontology Modeling, which defines a schema reflecting the hierarchical structure of industrial standards, and (2) Atomic Proposition Ontology Modeling, which decomposes conditional statements, exceptions, and numerical rules into atomic sentence-level expressions structured for LLM-based triple extraction.

3.1 Hierarchical Ontology Modeling

3.1.1 Ontology Design. We define a schema that reflects the hierarchical organization of standard documents. To this end, based on PDF data parsed through OCR, we extract both text and structure while preserving positional and layout information within the document. This schema reveals that a document is not merely a collection of text but a semantically organized structure governed by parent-child relationships, which can be explicitly represented at the ontology level. Each document is divided into various units such as Section, Subsection, Subsubsection, text, table, and footnote, all of which form a hierarchical structure. This can be represented as a set, as shown in Eq. 1:

$$\mathcal{H} = \{\text{Section} \rightarrow \text{Subsection} \rightarrow \text{Subsubsection} \rightarrow (\text{Text} \cup \text{Table}) \rightarrow \text{Footnote}\}. \quad (1)$$

3.1.2 Hierarchical Structure Recognition. The parsed document data are normalized and mapped to the previously defined schema. To achieve this, regular expressions are used to recognize various section numbering formats such as ‘1. SCOPE’, ‘5.6 SCANNING’, and ‘ANNEX A’, which are then converted into a consistent hierarchical structure. Tables and footnotes are also linked to their corresponding parent sections using surrounding context, forming a hierarchical path of Section \rightarrow Subsection \rightarrow Subsubsection \rightarrow Text/Table \rightarrow Footnote. As a result, each textual unit preserves its position and context within the document, allowing hierarchical semantics to be reflected in subsequent knowledge extraction.

3.1.3 Knowledge Unit Structuring. The text with an assigned hierarchical structure is segmented into minimal sentence units suitable for knowledge extraction. Formally, the resulting knowledge units can be represented as **title-text pairs** as follows:

$$\mathcal{D} = \{U_1, U_2, \dots, U_N\}, \quad U_i = (t_{\text{title}}^{(i)}, t_{\text{text}}^{(i)}), \quad (2)$$

where U_i denotes an individual knowledge unit extracted from a document, and \mathcal{D} represents the set of such units that collectively constitute the entire document. $t_{\text{title}}^{(i)}$ indicates the hierarchical context (e.g., Section, Subsection), while $t_{\text{text}}^{(i)}$ refers to the corresponding text—such as paragraphs, tables, or notes—associated

with that context. Since industrial standards often include domain-specific patterns such as ‘3-in. [75-mm]’ and ‘ASTM A6/A6M’, we applied tokenization rules that preserve units, abbreviations, and table identifiers. Ultimately, the title serves as an anchor defining the contextual scope, and the text provides concrete facts within that scope. The combination of these two elements transforms the document into a structured input suitable for knowledge graph construction and establishes a foundation for effective LLM-based triple extraction. An actual example is shown below.

Ontology-Aligned Title-Text Representation

Title ($t_{\text{title}}^{(i)}$): Procedure \rightarrow Scanning

Text ($t_{\text{text}}^{(i)}$): the major plate axis, on 3-in. [75-mm] or smaller centers. Measure the lines from the center or one corner of the plate with an additional path within 2 in. [50 mm] of all edges of the plate on the searching surface.

3.2 Atomic Proposition Ontology Modeling

3.2.1 Propositional Decomposition of Sentences. After hierarchically structuring the paragraphs, a sentence-level ontology was constructed for each sentence $t_{\text{text}}^{(i)}$. Industrial standards frequently describe manufacturing environments and technical specifications, and thus often adopt the form of conditional statements, such as “if A then B ”. This characteristic is particularly evident in tabular formats, where both row and column conditions jointly determine the applicable rule within each cell. Accordingly, a sentence-level ontology schema was defined to represent such conditional logic explicitly. In propositional logic, an **atomic proposition** denotes an indivisible unit of meaning, while a **compound proposition** is composed of multiple atomic propositions connected through logical operators such as \wedge , \vee , and \rightarrow [8]. Formally, a conditional statement “if A then B ” can be represented as $A \rightarrow B$, where A and B are independent sub-propositions that can be recursively decomposed into atomic propositions. Hence, any sentence can ultimately be represented as a hierarchical composition of atomic propositions. By recursively decomposing a sentence based on its logical connectors, a tree-structured representation is obtained, where the leaf nodes correspond to atomic propositions. Let \mathcal{P} be a set of propositions, and $p_j^{(\ell)} \in \mathcal{P}$ denote the j -th proposition at decomposition depth ℓ . Then, the sentence structure can be formally expressed as follows:

$$\begin{aligned} t_{\text{text}}^{(i)} &= p_1^{(1)} ::= f^{(1)}(p_1^{(2)}, p_2^{(2)}) = p_1^{(2)} \rightarrow p_2^{(2)}, \\ p_k^{(2)} &::= f^{(2)}(p_{k_1}^{(3)}, p_{k_2}^{(3)}, \dots, p_{k_{m_2}}^{(3)}), \\ &\vdots \\ p_s^{(n-1)} &::= f^{(n-1)}(p_{s_1}^{(n)}, p_{s_2}^{(n)}, \dots, p_{s_{m_{n-1}}}^{(n)}), \end{aligned} \quad (3)$$

where $f^{(\ell)}(\cdot)$ denotes the logical operator that combines propositions at layer ℓ , with $f^{(\ell)} \in \wedge, \vee, \rightarrow$, and m_ℓ represents the number of sub-propositions forming each compound proposition

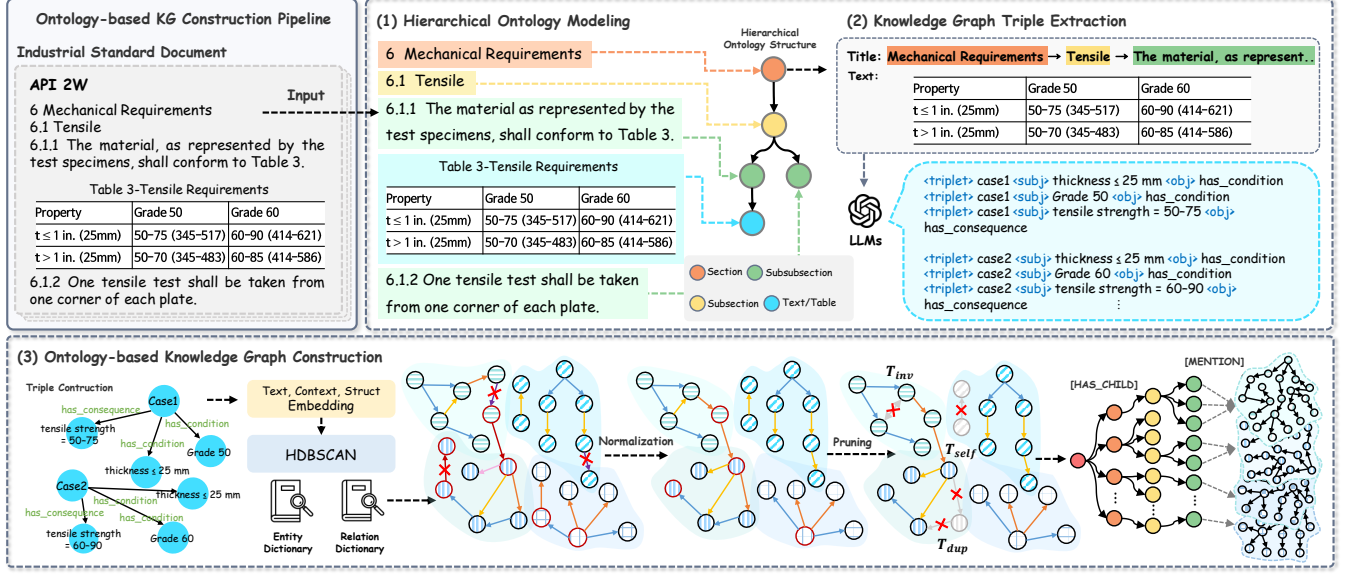


Figure 1: Ontology-based Knowledge Graph Construction Pipeline. (1) Hierarchical Ontology Modeling: The document is organized into a hierarchical ontology schema (Section, Subsection, Text/Table) to represent its structure explicitly. (2) Knowledge Graph Triple Extraction: An LLM extracts condition-consequence triples from both textual and tabular content. (3) Ontology-based Knowledge Graph Construction: The extracted triples are refined through embedding-based synonym normalization, HDBSCAN clustering, and pruning, then stored in Neo4j for use in the Ontology-aware KG-RAG framework.

$p_j^{(\ell)}$ at level ℓ . Each compound proposition $p_j^{(\ell)}$ is defined as a composition of its unique set of sub-propositions $p_{j_1}^{(\ell+1)}, \dots, p_{j_{m_\ell}}^{(\ell+1)}$. Each $p_j^{(\ell)}$ may represent either an atomic or a compound proposition, while $p_j^{(\ell)}$ always denotes an atomic proposition.

3.2.2 Triple Extraction Function. Following the decomposition, knowledge triples were extracted according to this hierarchical logical structure. We define the triple extraction function as

$$T : X \mapsto \{\text{triples}\}.$$

This function converts the hierarchical propositional structure of the text into a graph representation by connecting parent and child nodes across layers. In the structure $A \rightarrow B$, elements on the left-hand side (A), and those on the right-hand side (B) via `has_condition`. For each layer $T(\ell \geq 2)$, the triple extraction function $T(\ell)$ is defined as follows.

(i) Atomic proposition:

$$T(p_j^{(\ell)}) = \{(n_{group}^X, \text{has}_{\alpha,\beta}, p_{j_t}^{(\ell+1)})\}, \forall t \in \{1, \dots, m_{\ell-1}\},$$

$$\text{where } T(p_1^{(1)}) = \{(n_{case}, \text{has}_{\alpha,\beta}, p_j^{(2)})\}, \forall j \in \{1, 2\}.$$
(4)

(ii) Compound proposition:

$$T(p_j^{(\ell)}) = \{(n_{group}^X, \text{has}_{\alpha,\beta}, n_{group}^Y)\},$$

$$\text{where } T(p_1^{(1)}) = \{(n_{case}, \text{has}_{\alpha,\beta}, n_{group}^Y)\}.$$
(5)

Here, n_{case} denotes a conditional unit in the standard document (e.g., a table row or a conditional clause in a paragraph) and serves

as the root node of the decomposition tree. n_{group} represents an intermediate node generated during decomposition, grouping atomic or lower-level propositions connected through logical operators such as AND, OR. At each depth, an atomic proposition is represented directly as a node, while a compound proposition forms a group node. Here, $\alpha \in \text{condition, consequence}$ and $\beta \in \wedge, \vee$. The complete set of triples for a text is given by:

$$T(t_{\text{text}}^{(i)}) = \bigcup_{\ell=1}^{n-1} \bigcup_{j=1}^{m_{\ell-1}} T(p_j^{(\ell)}), \quad (6)$$

where $T(p^{(i)})$ denotes the set of triples extracted from each hierarchical level. To operationalize the proposed framework on real-world data, a rule-based triple extraction process was formulated as an LLM prompt template. In addition, few-shot examples from actual standard documents were incorporated to guide the model toward consistent logical interpretation. Conversion rules for SI unit normalization were also applied to ensure numerical consistency. As a result, the system automatically generated triple sets that faithfully reflect the hierarchical logical structure of each sentence. An example of the extracted triples is illustrated in Figure 2.

4 Ontology-Based Knowledge Graph Construction

The ontology-based KG construction pipeline proposed in this study comprises three stages: (1) Synonym Dictionary Construction, (2) Knowledge Graph Refinement and Pruning, and (3) Graph Structure Design.

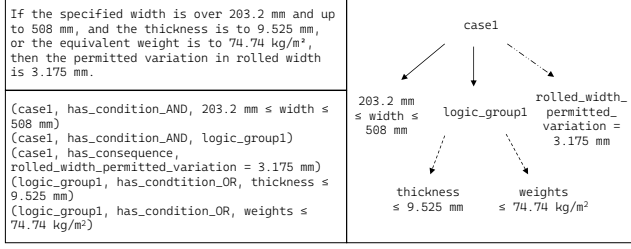


Figure 2: Example using the Atomic Proposition Ontology. (top-left) An example of table data from the A6/A6M document represented in sentence form, showing its conditional structure. (bottom-left) The set of triples extracted using the Atomic Proposition Ontology. (right) An example of a knowledge graph generated from the Atomic Proposition Ontology. Solid lines represent [has_condition_AND], dashed lines represent [has_condition_OR], and dash-dot lines represent [has_consequence] relations.

4.1 Synonym Dictionary Construction

In this stage, a synonym dictionary for entities and relations is constructed to ensure representational consistency and eliminate redundancy in the KG. To achieve this, embeddings that integrate textual, contextual, and structural information were clustered.

For the entity dictionary, the final representation of each entity e is defined by concatenating three types of embeddings, as shown in Eq. 7:

$$E_{\text{final}}(e) = [E_{\text{text}}(e); E_{\text{ctx}}(e); E_{\text{pos}}(e)], \quad (7)$$

where the textual semantic embedding $E_{\text{text}}(e)$ captures the semantic similarity of each entity using Sentence-BERT [26]. The local contextual embedding $E_{\text{ctx}}(e)$ is computed through Personalized PageRank [13] to represent the contextual importance around a specific node. Finally, the global structural embedding $E_{\text{pos}}(e)$ reflects the overall structural position within the graph by employing Spectral Embedding [4], which is based on the eigenvectors of the graph Laplacian. The resulting embeddings $E_{\text{final}}(e)$ are then clustered using HDBSCAN [19] to form synonym entity groups. For the relation dictionary, we adopt the same embedding approach for each relation r , with an additional Selectional Preference embedding $E_{\text{prefer}}(r)$ that captures the type distribution of subjects and objects predominantly connected by relation r (see Eq. 8).

$$E_{\text{final}}(r) = [E_{\text{text}}(r); E_{\text{ctx}}(r); E_{\text{pos}}(r); E_{\text{prefer}}(r)]. \quad (8)$$

In other words, we represent the likelihood of a relation occurring between specific types of entity pairs (u, v) as a probability distribution $p(t|u, v)$ to capture their functional similarity. Finally, we cluster the relations using HDBSCAN and select the member closest to the cluster center as the canonical predicate, thereby standardizing relation representations.

4.2 Knowledge Graph Refinement and Pruning

4.2.1 Conceptual Normalization. Based on the previously constructed synonym dictionary, dispersed relation labels are replaced with canonical predicates, and diverse entity notations are unified into standardized forms. This process eliminates conceptual redundancy and prevents identical terms from being unnecessarily separated into multiple nodes.

4.2.2 Structural Pruning. To enhance the structural efficiency and density of the normalized graph, a multi-step pruning procedure is applied to remove redundant or low-information elements. The graph is represented by a node set V and a triple set $T \subseteq U \times R \times V$, and the pruning process is defined as follows. Through this structural pruning, the knowledge graph achieves both computational efficiency and analytical clarity.

(i) **Self-loop Pruning:** Relations referring to the same entity, $T_{\text{self}} = \{(u, r, v) \in T \mid u = v\}$, are semantically irrelevant and therefore removed.

(ii) **Inverse-edge Pruning:** When bidirectional edges coexist, the one with the lower weight according to Eq. 9 is removed to ensure directionality. If the weights are identical, $w(u, v) = w(v, u)$, one edge is removed according to a predefined node ordering rule to maintain consistency.

$$T_{\text{inv}} = \left\{ \arg \min_{t \in \{(u, v), (v, u)\}} w(t) \mid (u, v), (v, u) \in T \right\} \quad (9)$$

(iii) **Duplicate Triple Removal:** When identical triples are extracted multiple times from the same document, as shown in Eq. 10, only one instance is retained.

$$T_{\text{dup}} = \{(u, r, v)_1, \dots, (u, r, v)_k \mid (u, r, v)_i \in T, k > 1\} \quad (10)$$

After this triple removal process, isolated nodes with a degree of 0 are eliminated, resulting in the final triple set defined in Eq. 11.

$$T_{\text{final}} = T \setminus (T_{\text{self}} \cup T_{\text{inv}} \cup T_{\text{dup}}) \quad (11)$$

4.3 Graph Structure Design

After refining the KG representations through Sections 4.1 and 4.2, we designed a unified graph schema that integrates the document's hierarchical structure (ontology) with the KG. Finally, the unified graph model was implemented in Neo4j, enabling both hierarchical navigation and ontology-aware retrieval within the ontology-KG integrated framework. In addition, a [MENTION] relation was established to link each section node to the entities it contains. This dual structure enables effective exploration of knowledge within specific document regions and supports query answering that considers hierarchical context. Finally, the

5 Ontology-Aware KG-RAG

While conventional KGs can be evaluated using predefined answer sets, our ontology-based KG lacks explicit ground truths, making direct evaluation difficult. Therefore, we adopted the ontology-aware KG-RAG framework (Figure 3) to verify whether the proposed KG contributes to QA performance improvement. To this end, we selected two baselines with different retrieval and reasoning strategies: MindMap [31] and KG-Retriever [6], and modified them to incorporate the ontology-aware KG for comparative evaluation.

Ontology-aware MindMap is a method designed to enhance cooperative reasoning between the KG and the LLM. It extracts key entities from the question, matches them with KG embeddings to establish anchors, and collects condition-consequence triples along the document's hierarchical structure. It then identifies relevant evidence by exploring shortest paths between entities and comparing triple-query embedding similarities. Finally, the selected

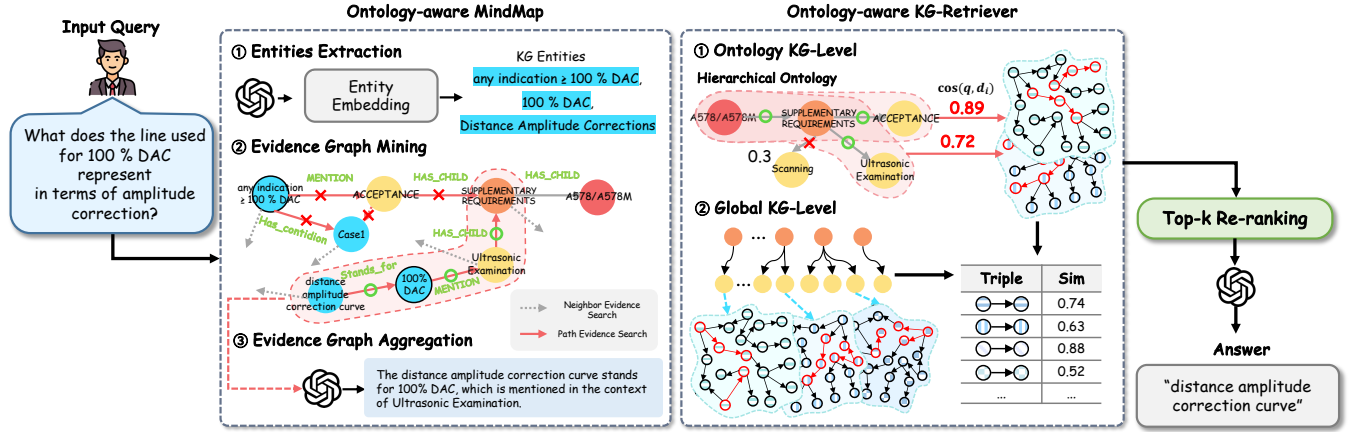


Figure 3: Ontology-aware KG-RAG framework. (Left) The ontology-aware MindMap selects evidence by reflecting the document’s hierarchical ontology structure. (Right) The Ontology-aware KG-Retriever combines ontology-based retrieval with global KG exploration. Finally, triples are extracted and inserted into the prompt through a re-ranking process.

evidence is integrated into a reasoning graph and inserted into the LLM prompt to enable structurally coherent answer generation.

In contrast, ontology-aware KG-Retriever integrates the hierarchical structure and global context of documents to extract precise evidence. At the ontology KG-Level, each section pair $(t_{\text{title}}^{(i)}, t_{\text{text}}^{(i)})$ is embedded, and the cosine similarity $\cos(\mathbf{q}, \mathbf{d}_i)$ between the query \mathbf{q} and section embedding \mathbf{d}_i is computed. Triples are then retrieved from the top- K sections. At the global KG-Level, the guided search module calculates the similarity between the query and triples to extract supplementary evidence. Finally, the retrieved results are integrated, reranked, and inserted into the LLM prompt for answer generation. Detailed descriptions of the framework are provided in Appendix A.

6 Experiments

In this section, we empirically analyze the performance of the proposed ontology-aware KG-RAG from multiple perspectives. We conduct quantitative experiments to answer the following research questions (RQs):

- **RQ1:** To what extent does the proposed ontology-aware KG-RAG improve QA performance compared to baseline RAG models?
- **RQ2:** Can the proposed method effectively generalize to real-world domain problems such as toxic clause detection?
- **RQ3:** How does QA performance vary depending on the combination of input information types (Text, KG, Ontology)?
- **RQ4:** What are the individual and combined effects of the synonym dictionary and pruning modules on performance?
- **RQ5:** How does the knowledge graph exploration range (number of hops) affect QA performance, and what is the optimal depth?

6.1 Experimental Setup

IndusSpec-QA Dataset. In this study, we constructed a dataset that reflects the structural characteristics of industrial documents. Considering the hierarchical nature of industrial standards, both

Table 1: Summary statistics of the IndusSpec-QA dataset constructed from three industrial standards. ‘Page’ indicates the number of pages in the original documents. The dataset consists of three task types—rule, multi-hop, and table—and is built upon custom knowledge graphs. The table summarizes the number of QA pairs per task and the scale of the graph (nodes, relations, triples).

Dataset	A578/A578M	API_2W	A6/A6M
Page	5	27	63
Question	501	384	663
Rule	414	179	200
Multi-hop	79	86	63
Table	8	119	400
Node	606	1320	10976
Relation	52	121	216
Triplet	1986	5128	53508

section titles and subordinate texts were utilized to enable natural question generation. The dataset was built based on the hierarchical ontology model defined in Section 3.1, where within U_i , (1) $t_{\text{title}}^{(i)}$ was used to define the scope of the query, and (2) $t_{\text{text}}^{(i)}$ (including text and tables) was referenced for automatic question generation. The generated questions were then manually reviewed by the research team to correct inconsistencies, ambiguity, and duplication, followed by final validation by domain experts. Dataset statistics are presented in Table 1, and see Appendix Table 8.

Baselines. In this study, we established four categories of baseline models: (1) cases without external knowledge (LLM-only), (2) cases where the full document is directly provided to the model (LLM w/full doc), (3) cases utilizing unstructured text retrieval (Text-only RAG), and (4) KG-based approaches that do not incorporate ontology structures (KG-RAG). Comparative experiments were conducted across these categories using various models. See Appendix B.

Evaluation Metrics. In this study, model performance was evaluated using F1-Score, BERTScore, BLEU, and Rouge. F1 measures

Table 2: Performance comparison across datasets. Four categories of QA models were compared using three industrial standard documents. The results represent weighted averages of rule, table, and multi-hop task performances, with the best scores shown in bold and the second-best underlined.

Type	Method	A578/A578M				API_2W				A6/A6M			
		F1	BERTScore	BLEU	Rouge	F1	BERTScore	BLEU	Rouge	F1	BERTScore	BLEU	Rouge
LLM-only	GPT-4o-mini	0.013	0.806	0.004	0.020	0.010	0.800	0.003	0.014	0.002	0.790	0.002	0.003
	gemini-2.0-flash	0.000	0.802	0.000	0.004	0.000	0.797	0.000	0.002	0.003	0.790	0.002	0.003
	DeepSeek-v3.1	0.043	0.817	0.016	0.049	0.040	0.809	0.018	0.046	0.036	0.800	0.016	0.039
LLM w/full doc	GPT-4o-mini	0.308	0.867	0.157	0.357	0.298	0.864	0.120	0.361	0.152	0.826	0.066	0.174
	gemini-2.0-flash	0.282	0.861	0.159	0.315	0.276	0.853	0.142	0.314	0.206	0.834	0.083	0.273
	DeepSeek-v3.1	0.077	0.874	0.165	0.384	0.382	0.884	0.156	0.463	0.130	0.824	0.069	0.152
Text-only RAG	BM25+Qwen1.5-14B	0.306	0.876	0.129	0.381	0.306	0.881	0.124	0.431	0.154	0.837	0.054	0.208
	BM25+Mistral-7B	0.279	0.865	0.118	0.346	0.242	0.859	0.112	0.300	0.128	0.837	0.050	0.183
	BM25+GPT-4o-mini	0.366	0.880	0.193	0.433	0.308	0.868	0.140	0.357	0.160	0.831	0.070	0.195
	BM25+gemini-2.0	0.340	0.876	0.184	0.405	0.291	0.859	0.137	0.350	0.159	0.825	0.076	0.184
	Dense+Qwen1.5-14B	0.264	0.870	0.104	0.322	0.372	0.891	0.134	0.475	0.179	0.852	0.067	0.223
	Dense+Mistral-7B	0.255	0.859	0.100	0.318	0.277	0.864	0.112	0.348	0.152	0.840	0.059	0.223
	Dense+GPT-4o-mini	0.314	0.871	0.499	0.370	0.413	0.888	0.160	0.476	0.175	0.828	0.079	0.205
	Dense+gemini-2.0	0.309	0.870	0.167	0.359	0.348	0.875	0.164	0.471	0.167	0.825	0.077	0.203
KG-RAG (non-finetuned)	MindMap (Wen et al., 2024)	0.370	0.879	0.224	0.426	0.334	0.887	0.145	0.412	0.223	0.851	0.113	0.255
	KG Retriever (Chen et al., 2024)	0.319	0.874	0.170	0.476	0.223	0.862	0.090	0.301	0.190	0.853	0.067	0.281
Ontology-aware KG-RAG (non-finetuned)	Ontology-aware MindMap	0.511	0.911	0.313	0.585	0.494	0.915	0.265	0.569	0.441	0.892	0.186	0.518
	Ontology-aware KG-Retriever	0.639	0.933	0.424	0.653	0.448	0.909	0.268	0.527	0.191	0.854	0.084	0.194

the balance between precision and recall, BERTScore evaluates semantic similarity between sentences, and BLEU and Rouge assess n-gram and phrase overlap. In the toxic clause detection task, Recall was additionally analyzed, as accurately identifying actual toxic clauses is of primary importance.

Implementation Details. All ontology-aware KG-RAG experiments used GPT-4o-mini as the inference model, while GPT-5-mini was employed for KG triple extraction. ontology-aware MindMap was configured with 3-hop exploration and KG-Retriever with 2-hop, with both models inserting up to 30 triples into the LLM prompt. Sentence-BERT (all-MiniLM-L6-v2) was used for embedding-based similarity computation, and experiments with open-source LLMs were conducted on eight NVIDIA RTX 3090 GPUs 24GB.

6.2 Experimental Results

(RQ1) Performance Comparison. We compared the performance of the four categories of baseline models across the three industrial standard datasets, and the results are summarized in Table 2. From these experiments, the following key observations can be drawn:

- **Ontology-aware KG-RAG achieved the highest performance across all datasets.** The average F1 score reached 0.454, representing a significant improvement over existing methods. This demonstrates that ontology-based knowledge graphs are a decisive factor in enhancing QA performance for specialized domains such as industrial standards.
- **The LLM-only approach shows strong limitations in specialized domains.** With an average F1 score of only 0.016, it indicates that LLMs alone, without access to external knowledge, cannot generate meaningful answers in technical document domains.
- **LLM w/full doc suffers from information overload.** Although the average F1 score of 0.235 outperformed the LLM-only approach, it still fell far short of the proposed ontology-aware KG-RAG. This suggests that simply injecting the entire

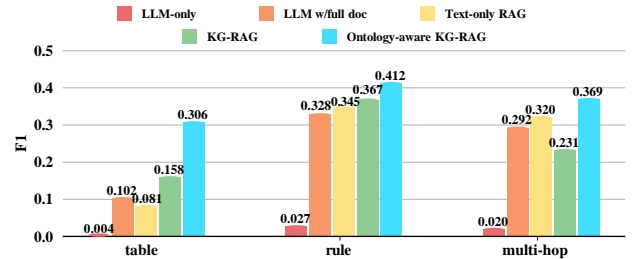


Figure 4: F1 performance comparison by data type. In the QA experiments conducted on three standard documents, the average F1-scores of five model categories were compared across three data types: table, rule, and multi-hop.

document leads to information overload, hindering effective reasoning.

- **KG-RAG shows only marginal improvement over Text-only RAG.** The performance gap between Text-only RAG (average F1: 0.277) and conventional KG-RAG (average F1: 0.304) was minor, implying that previous KG-RAG frameworks failed to fully exploit the potential of structural knowledge. In contrast, the ontology-aware KG-RAG achieved a 64% improvement (0.277 → 0.454), underscoring the crucial role of domain ontology design in industrial QA tasks.
- **Ontology-aware KG-RAG remains stable regardless of document length.** Overall, performance was slightly lower for the longer document A6/A6M (63p) compared to the shorter A578/A578M (5p), suggesting that evidence extraction becomes more complex as document length increases. Nevertheless, ontology-aware KG-RAG maintained relatively stable performance even under such conditions, showing particular strength in long-document QA.

Furthermore, Figure 4 presents the average F1-scores across models within each methodological category, showing that the performance gap between methods varies depending on the task type. The ontology-aware KG-RAG achieved the highest average

Table 3: Performance on API_2W toxic clause detection task (F1, Accuracy, Recall). Ontology-aware KG-Retriever achieved the best results across all metrics, and ontology-aware MindMap also showed improved performance compared to the baseline without ontology integration.

Type	Method	F1	Acc	Recall
LLM-only	GPT-4o-mini	0.558	0.667	0.438
	gemini-2.0-flash	0.000	0.520	0.000
	DeepSeek-v3.1	0.591	0.685	0.471
LLM w/full doc	GPT-4o-mini	0.812	0.833	0.752
	gemini-2.0-flash	0.527	0.667	0.364
	DeepSeek-v3.1	0.850	0.857	0.890
Text-only RAG	BM25+Qwen1.5-14B	0.870	0.869	0.895
	BM25+Mistral-7B	0.837	0.826	0.917
	BM25+GPT-4o-mini	0.878	0.879	0.886
	BM25+gemini-2.0-flash	0.887	0.888	0.895
	Dense+Qwen1.5-14B	0.871	0.872	0.893
	Dense+Mistral-7B	0.834	0.825	0.809
	Dense+GPT-4o-mini	0.851	0.861	0.826
	Dense+gemini-2.0-flash	0.851	0.864	0.818
KG-RAG (non-finetuned)	MindMap (Wen et al., 2024)	0.819	0.812	0.880
	KG Retriever (Chen et al., 2024)	0.900	0.905	0.893
Ontology-aware KG-RAG (non-finetuned)	Ontology-aware MindMap	0.828	0.832	0.883
	Ontology-aware KG-Retriever	0.910	0.913	0.926

F1 in all task types. Specifically, it outperformed the conventional KG-RAG by 93.7% in the table task (0.158 \rightarrow 0.306), by 12.3% in the rule task (0.367 \rightarrow 0.412), and by 59.7% in the multi-hop task (0.231 \rightarrow 0.369). The strong improvement in the table task suggests that tabular, structured data aligns particularly well with ontology-based schema representations, leading to more effective reasoning. These results indicate that the complex rules and constraints of industrial standards can only be effectively handled through the structural knowledge representation provided by ontologies. Overall, this finding strongly underscores that ontology-based, knowledge graph-driven approaches are essential for highly specialized domains such as industrial standards. The substantial improvement over LLM-only models clearly demonstrates the importance of integrating external knowledge and highlights that systematic domain ontology construction will be a key direction for future development.

(RQ2) Performance on Toxic Clause Detection Task. We conducted comparative experiments to evaluate the performance of the proposed method in the toxic clause detection task. This task involves distinguishing whether a given clause is toxic by presenting both the original sentences from the document and their modified versions with altered numerical values or words. The dataset was constructed using the API_2W document and consisted of 252 classification problems, with examples provided in Appendix Table 9. As shown in Table 3, the ontology-aware KG-Retriever achieved the best performance across all evaluation metrics. Furthermore, the ontology-aware MindMap also demonstrated improved results compared to the original MindMap without ontology integration. Since correctly identifying actual toxic clauses is the key objective of this task, the recall metric is particularly important. The results confirm that the ontology-based KG-RAG approach achieved notable performance gains in this aspect as well.

(RQ3) Effect of KG and Ontology Integration. To evaluate the individual contributions of the ontology and KG components, we

Table 4: Ablation study on API_2W rule data. Performance comparison of LLM w/doc, with Ontology, with KG, and with both. Adding both Ontology and KG yields the highest gains across all metrics.

	F1	BERTScore	BLEU	Rouge
LLM w/doc	0.298	0.864	0.120	0.361
+ KG	0.334	0.887	0.145	0.412
+ Ontology	<u>0.413</u>	<u>0.888</u>	<u>0.160</u>	<u>0.476</u>
+ Ontology + KG	0.494	0.915	0.265	0.569

conducted an ablation study. In the experiment using only the ontology module, triple extraction was excluded, and only the hierarchical ontology structure was applied. As shown in Table 4, when both ontology and KG were utilized together, the F1 increased from 0.298 to 0.494, representing an improvement of approximately 65%. Even when compared to using the KG module alone (F1: 0.334), the combined model achieved about a 47% performance gain. Interestingly, the performance with the ontology module alone surpassed that of the KG-only setting, suggesting that the hierarchical understanding of document structure provided by the ontology contributed more substantially to overall performance improvement.

(RQ4) Effect of Synonym Dictionary and Pruning. In this experiment, we evaluated the contributions of the synonym dictionary and pruning modules in the KG construction process using F1 and Rouge metrics. As shown in Figure 5, the performance of the ontology-aware MindMap gradually improved as the two modules were applied step by step. In contrast, for the ontology-aware KG-Retriever, applying either the synonym dictionary or pruning individually resulted in slight decreases in F1 or Rouge, suggesting that applying each module in isolation may lead to partial information loss. However, when both modules were applied together, the performance improved the most, demonstrating their complementary effects. These results confirm that integrating the synonym dictionary and pruning processes is essential for effective KG refinement.

(RQ5) Performance under Different k -Hop Neighborhoods.

Figure 6 presents the experimental results of varying hop sizes for the two proposed ontology-based methods. The ontology-aware MindMap achieved the highest performance at 3-hop, while the ontology-aware KG-Retriever performed best at 2-hop. This difference appears to stem from the distinct structural characteristics of their retrieval mechanisms. Specifically, the ontology-aware KG-Retriever explores multidimensionally connected nodes across ontology and KG levels, whereas the ontology-aware MindMap follows a more linear exploration pattern. Consequently, these structural differences lead to different optimal hop values for the two methods. In subsequent experiments, we applied the respective optimal hop settings for each approach.

7 Conclusion

This study proposed an ontology-based KG construction methodology to effectively structure industrial standards that contain conditional rules, numerical relations, and complex tables. The proposed pipeline hierarchically organizes documents into ontologies, extracts triples from logical statements, and refines the extracted triples through synonym normalization and pruning. To evaluate its

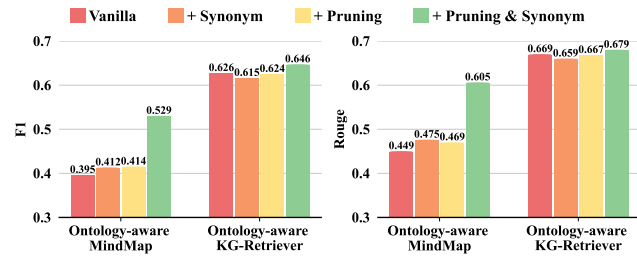


Figure 5: Effect of synonym dictionary and pruning modules on QA performance (F1, ROUGE) using API_2W rule data. The ontology-aware MindMap showed consistent gains with both modules, whereas the ontology-aware KG-Retriever slightly dropped with each module alone but reached peak performance when combined.

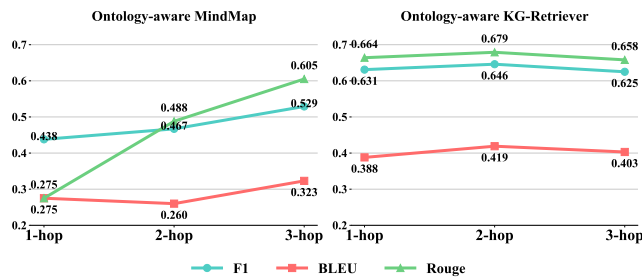


Figure 6: Performance across different k-hop settings using API_2W rule data. For all evaluation metrics, the ontology-aware MindMap achieved the best performance at 3-hop, while the ontology-aware KG-Retriever showed the highest performance at 2-hop.

effectiveness, we constructed QA and toxic clause detection datasets and implemented an ontology-aware KG-RAG framework. The proposed method yielded a 64% average F1 improvement compared to conventional KG-RAG models, with particularly notable gains in table reasoning tasks. These results demonstrate that atomic proposition ontology modeling, hierarchical ontology structuring, and synonym-pruning refinement are complementary components that collectively enhance reasoning performance. Overall, this work establishes that complex industrial standards can be systematically transformed into reliable and scalable knowledge graphs, thereby contributing to domain-specific RAG development and intelligent document management.

References

- [1] M.M. Ashburner, C.A.C. Ball, Judith Blake, David Botstein, Heather Butler, J.M.J. Cherry, Allan Peter Davis, Kara Dolinski, Selina Dwight, Janan Eppig, Midori Harris, D.P. Hill, Laurie Issel-Tarver, A Kasarskis, Suzanna Lewis, John Matese, J.E. Richardson, M Ringwald, and G.M. Rubin. 2000. Gene Ontology: tool for the unification of biology. *The Gene Ontology Consortium. Nat Genet* 25 (05 2000), 25–29.
- [2] Muhammad Asim, Muhammad Wasim, Muhammad Usman Khan, and Waqar Mahmood. 2018. A survey of ontology learning techniques and applications. *Database The Journal of Biological Databases and Curation* 2018 (10 2018), 1–24. doi:10.1093/database/bay101
- [3] Hamed Babaei Giglou, Jennifer D’Souza, and Sören Auer. 2023. LLMs4OL: Large Language Models for Ontology Learning. In *The Semantic Web – ISWC 2023: 22nd International Semantic Web Conference, Athens, Greece, November 6–10, 2023, Proceedings, Part I* (Athens, Greece). Springer-Verlag, Berlin, Heidelberg, 408–427. doi:10.1007/978-3-031-47240-4_22
- [4] Mikhail Belkin and Partha Niyogi. 2001. Laplacian eigenmaps and spectral techniques for embedding and clustering. In *Proceedings of the 15th International Conference on Neural Information Processing Systems: Natural and Synthetic* (Vancouver, British Columbia, Canada) (*NIPS’01*). MIT Press, Cambridge, MA, USA, 585–591.
- [5] Christopher Brewster. 2006. Ontology Learning from Text: Methods, Evaluation and Applications Paul Buitelaar, Philipp Cimiano, and Bernado Magnini (editors) (DFKI Saarbrücken, University of Karlsruhe, and ITC-irst), Amsterdam: IOS Press (Frontiers in artificial intelligence and applications, edited by J. Breuker et al., volume 123), 2005, v+171 pp; hardbound, ISBN 1-58603-523-1, \$102.00, €85.00, £59.00. *Computational Linguistics* 32 (12 2006), 569–572. doi:10.1162/col.2006.32.4.569
- [6] Weijie Chen, Ting Bai, Jinbo Su, Jian Luan, Wei Liu, and Chuan Shi. 2025. KG-Retriever: Efficient Knowledge Indexing for Retrieval-Augmented Large Language Models. arXiv:2412.05547 [cs.IR] <https://arxiv.org/abs/2412.05547>
- [7] Philipp Cimiano and Johanna Völker. 2005. Text2Onto: a framework for ontology learning and data-driven change discovery. In *Proceedings of the 10th International Conference on Natural Language Processing and Information Systems* (Alicante, Spain) (*NLDB’05*). Springer-Verlag, Berlin, Heidelberg, 227–238. doi:10.1007/11428817_21
- [8] Stephen A. Cook. 1971. The complexity of theorem-proving procedures. In *Proceedings of the Third Annual ACM Symposium on Theory of Computing* (Shaker Heights, Ohio, USA) (*STOC ’71*). Association for Computing Machinery, New York, NY, USA, 151–158. doi:10.1145/800157.805047
- [9] DeepMind / Google. 2025. Gemini 2.0 Flash. <https://deepmind.google/models/gemini/flash/>.
- [10] DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2025. DeepSeek-V3 Technical Report. arXiv:2412.19437 [cs.CL] <https://arxiv.org/abs/2412.19437>
- [11] Oren Etzioni, Michael Cafarella, Doug Downey, Stanley Kok, Ana-Maria Popescu, Tal Shaked, Stephen Soderland, Daniel S. Weld, and Alexander Yates. 2004. Web-scale information extraction in knowitall: (preliminary results). In *Proceedings of the 13th International Conference on World Wide Web* (New York, NY, USA) (*WWW ’04*). Association for Computing Machinery, New York, NY, USA, 100–110. doi:10.1145/988672.988687
- [12] Maurice Funk, Simon Hosemann, Jean Christoph Jung, and Carsten Lutz. 2023. Towards Ontology Construction with Language Models. arXiv:2309.09898 [cs.AI] <https://arxiv.org/abs/2309.09898>
- [13] Johannes Gasteiger, Aleksandar Bojchevski, and Stephan Günnemann. 2022. Predict then Propagate: Graph Neural Networks meet Personalized PageRank. arXiv:1810.05997 [cs.LG] <https://arxiv.org/abs/1810.05997>
- [14] Thomas Gruber. 1994. Toward Principles for the Design of Ontologies Used for Knowledge Sharing. *International Journal of Human-Computer Studies* 43 (08 1994). doi:10.1006/ijhc.1995.1081
- [15] Maryam Hazman, Samhaa El-Beltagy, and Ahmed Rafea. 2013. A Survey of Ontology Learning Approaches. *International Journal of Computer Applications* 22 (09 2013), 36–43. doi:10.5120/2610-3642
- [16] Yuelin Hu, Futai Zou, Jiajia Han, Xin Sun, and Yilei Wang. 2024. LLM-TIKG: Threat intelligence knowledge graph construction utilizing large language model. *Computers & Security* 145 (2024), 103999. doi:10.1016/j.cose.2024.103999
- [17] Pere-Lluís Huguet Cabot and Roberto Navigli. 2021. REBEL: Relation Extraction By End-to-end Language generation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (Eds.). Association for Computational Linguistics, Punta Cana, Dominican Republic, 2370–2381. doi:10.18653/v1/2021.findings-emnlp.204
- [18] Andy Lo, Albert Q. Jiang, Wenda Li, and Mateja Jamnik. 2025. End-to-end ontology learning with large language models. In *Proceedings of the 38th International Conference on Neural Information Processing Systems* (Vancouver, BC, Canada) (*NIPS ’24*). Curran Associates Inc., Red Hook, NY, USA, Article 2767, 42 pages.

- [19] Leland McInnes, John Healy, and Steve Astels. 2017. hdbscan: Hierarchical density based clustering. *Journal of Open Source Software* 2, 11 (2017), 205. doi:10.21105/joss.00205
- [20] George A. Miller. 1995. WordNet: a lexical database for English. *Commun. ACM* 38, 11 (Nov. 1995), 39–41. doi:10.1145/219717.219748
- [21] Belinda Mo, Kysen Yu, Joshua Kazdan, Proud Mpala, Lisa Yu, Chris Cundy, Charilaos Kanatsoulis, and Sanmi Koyejo. 2025. KGGGen: Extracting Knowledge Graphs from Plain Text with Language Models. arXiv:2502.09956 [cs.CL] <https://arxiv.org/abs/2502.09956>
- [22] OpenAI. 2022. New and improved embedding model. <https://openai.com/index/new-and-improved-embedding-model/>.
- [23] OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, and Lama Ahmad. 2024. GPT-4 Technical Report. arXiv:2303.08774 [cs.CL] <https://arxiv.org/abs/2303.08774>
- [24] Marinela Parović, Ze Li, and Jinhua Du. 2025. Generating Domain-Specific Knowledge Graphs from Large Language Models. In *Findings of the Association for Computational Linguistics: ACL 2025*, Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (Eds.). Association for Computational Linguistics, Vienna, Austria, 11558–11574. doi:10.18653/v1/2025.findings-acl.602
- [25] Siyuan Qi, Wenguan Wang, Baoxiong Jia, Jianbing Shen, and Song-Chun Zhu. 2018. Learning Human-Object Interactions by Graph Parsing Neural Networks. arXiv:1808.07962 [cs.CV] <https://arxiv.org/abs/1808.07962>
- [26] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. arXiv:1908.10084 [cs.CL] <https://arxiv.org/abs/1908.10084>
- [27] Stephen Robertson and Hugo Zaragoza. 2009. The Probabilistic Relevance Framework: BM25 and Beyond. *Foundations and Trends® in Information Retrieval* 3, 4 (2009), 333–389. doi:10.1561/15000000019
- [28] Philipp Sorg and Philipp Cimiano. 2012. Exploiting Wikipedia for cross-lingual and multilingual information retrieval. *Data & Knowledge Engineering* 74 (2012), 26–45. doi:10.1016/j.datak.2012.02.003
- [29] Qingwang Wang, Chaohui Li, Yi Liu, Qiubai Zhu, Jian Song, and Tao Shen. 2025. An Adaptive Framework Embedded With LLM for Knowledge Graph Construction. *IEEE Transactions on Multimedia* 27 (2025), 2912–2923. doi:10.1109/TMM.2025.3557717
- [30] Zhepei Wei, Jianlin Su, Yue Wang, Yuan Tian, and Yi Chang. 2020. A Novel Cascade Binary Tagging Framework for Relational Triple Extraction. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (Eds.). Association for Computational Linguistics, Online, 1476–1488. doi:10.18653/v1/2020.acl-main.136
- [31] Yilin Wen, Zifeng Wang, and Jimeng Sun. 2024. MindMap: Knowledge Graph Prompting Sparks Graph of Thoughts in Large Language Models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, Bangkok, Thailand, 10370–10388. doi:10.18653/v1/2024.acl-long.558
- [32] Alexander Yates, Michele Banko, Matthew Broadhead, Michael Cafarella, Oren Etzioni, and Stephen Soderland. 2007. TextRunner: Open Information Extraction on the Web. In *Proceedings of Human Language Technologies: The Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL-HLT)*, Bob Carpenter, Amanda Stent, and Jason D. Williams (Eds.). Association for Computational Linguistics, Rochester, New York, USA, 25–26. <https://aclanthology.org/N07-4013/>
- [33] Bowen Zhang and Harold Soh. 2024. Extract, Define, Canonicalize: An LLM-based Framework for Knowledge Graph Construction. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, Miami, Florida, USA, 9820–9836. doi:10.18653/v1/2024.emnlp-main.548
- [34] Bowen Zhang and Harold Soh. 2024. Extract, Define, Canonicalize: An LLM-based Framework for Knowledge Graph Construction. arXiv:2404.03868 [cs.CL] <https://arxiv.org/abs/2404.03868>

A Ontology-Aware KG-RAG Details

A.1 Ontology-Aware MindMap

This methodology consists of four main stages: **Entity Extraction**, **Evidence Graph Mining**, **Evidence Graph Aggregation**, and **LLM-integrated Reasoning**.

In the **Entity Extraction** stage, core entity candidates are identified from the question using an LLM-based module. The extracted entities are then refined into a set of the most relevant entities by matching them against the preconstructed KG embeddings using

cosine similarity. This refined entity set serves as the starting point (anchor) for subsequent graph exploration.

In the **Evidence Graph Mining** stage, the graph is expanded around the matched entities to collect evidence. First, *Neighbor Search* explores relationships with adjacent nodes to gather local triples and strengthen contextual understanding. Then, *Path Search* identifies the shortest paths between entity pairs to obtain path evidence. Finally, *Semantic Verification* applies triple embedding-based semantic ranking to filter out irrelevant triples and retain only those that are semantically aligned with the question’s intent.

The **Evidence Graph Aggregation** stage integrates the triples obtained from the exploration phase into a unified reasoning graph after deduplication and structural normalization. The resulting reasoning graph is then transformed into a natural language evidence description suitable for LLM input.

In the **LLM-integrated Reasoning** stage, the reasoning graph is inserted into the LLM prompt and utilized for reasoning. During this process, the LLM integrates semantic diversity among KG entities, generates coherent answers following the reasoning paths provided by the KG, and compensates for missing or incorrect information in the KG using its internal knowledge. This enables cooperative reasoning between the KG and the LLM.

A.2 Ontology-Aware KG-Retriever

The proposed method consists of two stages of evidence retrieval at the ontology KG level and the global KG level, followed by a final stage that integrates the retrieved evidence into the LLM prompt.

At the **Ontology KG-Level**, each title within the document is treated as a semantic unit, and its associated text is embedded to compute cosine similarity with the query. Based on this, the top three titles are selected to form a coarse-grained candidate set. Then, for the KG triples contained within each title, 1-hop triples are first retrieved, and if insufficient, expansion continues to 2-hop triples to sequentially collect evidence. Specifically, 20, 10, and 5 triples are gathered from each title, respectively, constructing a fine-grained evidence set of 35 triples in total. This approach allows the model to capture both local semantic contexts and the hierarchical structure of the document.

The **Global KG-Level** excludes title segmentation and directly searches the entire KG for triples semantically similar to the query. Both 1-hop and 2-hop searches are performed to collect 30 triples, thereby providing global structural cues beyond local constraints. Additionally, a **Guided Search module** computes direct similarity between the query embeddings and the textual representations of triples (including table-based ones) to retrieve an additional 20 triples. This guided search functions as an independent complementary module that reinforces semantic relevance alongside structural cues.

In the final stage, the evidence obtained from the ontology KG and global KG explorations (a total of 85 triples) is integrated into a single candidate pool. The top 30 triples are then selected via similarity-based re-ranking with the query and inserted into the LLM prompt for answer generation. Through this multi-stage evidence extraction and re-ranking procedure, the ontology-aware KG-Retriever effectively incorporates both the semantic context

Table 5: A578/A578M performance comparison across different question types. Results of comparing various QA models on table, rule, and multi-hop question types.

Type	Method	table				rule				multi-hop			
		F1	BERTScore	BLEU	Rouge	F1	BERTScore	BLEU	Rouge	F1	BERTScore	BLEU	Rouge
LLM-only	GPT-4o-mini	0.000	0.773	0.000	0.000	0.013	0.806	0.004	0.021	0.017	0.806	0.006	0.016
	gemini-2.0-flash	0.000	0.773	0.000	0.000	0.000	0.802	0.000	0.005	0.000	0.806	0.000	0.000
	DeepSeek-v3.1	0.000	0.773	0.000	0.000	0.048	0.819	0.018	0.055	0.021	0.814	0.005	0.025
LLM w/full doc	GPT-4o-mini	0.000	0.776	0.000	0.031	0.297	0.866	0.154	0.339	0.388	0.881	0.167	0.454
	gemini-2.0-flash	0.000	0.773	0.000	0.000	0.276	0.861	0.163	0.307	0.344	0.872	0.157	0.388
	DeepSeek-v3.1	0.000	0.766	0.000	0.195	0.033	0.876	0.174	0.383	0.317	0.890	0.139	0.407
Text-only RAG	BM25+Qwen1.5-14B	0.000	0.821	0.000	0.219	0.312	0.876	0.137	0.38	0.307	0.880	0.133	0.402
	BM25+Mistral-7B	0.000	0.776	0.000	0.179	0.276	0.865	0.121	0.331	0.320	0.874	0.114	0.441
	BM25+GPT-4o-mini	0.064	0.812	0.009	0.343	0.360	0.880	0.194	0.419	0.430	0.890	0.209	0.516
	BM25+gemini-2.0-flash	0.000	0.800	0.000	0.25	0.332	0.875	0.184	0.391	0.416	0.887	0.201	0.493
	Dense+Qwen1.5-14B	0.000	0.828	0.000	0.324	0.262	0.869	0.109	0.310	0.302	0.882	0.104	0.389
	Dense+Mistral-7B	0.011	0.775	0.001	0.185	0.251	0.859	0.101	0.305	0.290	0.867	0.105	0.390
	Dense+GPT-4o-mini	0.110	0.827	0.018	0.411	0.295	0.868	0.560	0.342	0.432	0.862	0.225	0.516
	Dense+gemini-2.0	0.000	0.803	0.000	0.281	0.296	0.868	0.161	0.337	0.407	0.887	0.216	0.480
KG-RAG (non-finetuned)	MindMap (Wen et al., 2024)	0.500	0.833	0.972	0.633	1.000	0.878	0.238	0.417	0.345	0.883	0.148	0.455
	KG Retriever (Chen et al., 2024)	0.118	0.822	0.730	0.750	0.338	0.878	0.170	0.403	0.253	0.859	0.119	0.529
Ontology-aware KG-RAG (non-finetuned)	Ontology-aware MindMap	0.875	0.978	0.725	1.000	0.511	0.911	0.329	0.573	0.476	0.907	0.226	0.607
	Ontology-aware KG-Retriever	0.351	0.876	0.030	0.375	0.639	0.934	0.440	0.654	0.670	0.935	0.380	0.676

Table 6: API_2W performance comparison across different question types. Results of comparing various QA models on table, rule, and multi-hop question types.

Type	Method	table				rule				multi-hop			
		F1	BERTScore	BLEU	Rouge	F1	BERTScore	BLEU	Rouge	F1	BERTScore	BLEU	Rouge
LLM-only	GPT-4o-mini	0.000	0.800	0.000	0.000	0.010	0.797	0.002	0.013	0.025	0.806	0.007	0.035
	gemini-2.0-flash	0.000	0.800	0.000	0.000	0.000	0.795	0.000	0.001	0.000	0.799	0.000	0.006
	DeepSeek-v3.1	0.005	0.803	0.001	0.005	0.005	0.809	0.003	0.058	0.069	0.818	0.017	0.079
LLM w/full doc	GPT-4o-mini	0.222	0.862	0.055	0.312	0.325	0.862	0.155	0.372	0.346	0.872	0.139	0.408
	gemini-2.0-flash	0.121	0.828	0.037	0.182	0.369	0.867	0.208	0.387	0.295	0.859	0.151	0.326
	DeepSeek-v3.1	0.352	0.888	0.078	0.472	0.412	0.884	0.208	0.471	0.360	0.879	0.155	0.435
Text-only RAG	BM25+Qwen1.5-14B	0.098	0.882	0.017	0.388	0.411	0.879	0.185	0.464	0.375	0.886	0.147	0.424
	BM25+Mistral-7B	0.053	0.836	0.015	0.125	0.352	0.870	0.174	0.402	0.273	0.866	0.119	0.330
	BM25+GPT-4o-mini	0.156	0.847	0.027	0.234	0.404	0.881	0.205	0.439	0.316	0.872	0.162	0.357
	BM25+gemini-2.0-flash	0.110	0.830	0.025	0.218	0.396	0.873	0.191	0.428	0.322	0.870	0.178	0.369
	Dense+Qwen1.5-14B	0.334	0.906	0.061	0.531	0.392	0.883	0.171	0.458	0.383	0.889	0.159	0.434
	Dense+Mistral-7B	0.141	0.854	0.024	0.249	0.359	0.869	0.163	0.406	0.294	0.869	0.127	0.364
	Dense+GPT-4o-mini	0.405	0.897	0.072	0.502	0.439	0.886	0.214	0.468	0.372	0.879	0.171	0.419
	Dense+gemini-2.0-flash	0.222	0.863	0.045	0.496	0.449	0.884	0.234	0.491	0.313	0.872	0.185	0.393
KG-RAG (non-finetuned)	MindMap (Wen et al., 2024)	0.468	0.929	0.173	0.610	0.303	0.873	0.157	0.358	0.214	0.860	0.081	0.256
	KG Retriever (Chen et al., 2024)	0.135	0.850	0.027	0.232	0.344	0.873	0.162	0.420	0.081	0.854	0.027	0.151
Ontology-aware KG-RAG (non-finetuned)	Ontology-aware MindMap	0.531	0.933	0.211	0.621	0.529	0.917	0.323	0.605	0.371	0.888	0.217	0.421
	Ontology-aware KG-Retriever	0.135	0.879	0.036	0.317	0.646	0.933	0.419	0.679	0.470	0.902	0.274	0.600

of the query and the structural knowledge of the KG, leading to enhanced QA performance.

B Baselines

- **LLM-only** generates answers directly from the input query without using any external knowledge sources. For this setting, we employ GPT-4o-mini, Gemini-2.0-Flash, and DeepSeek-V3.1 (DeepSeek-Chat) as representative LLMs [9, 10, 23].
- **LLM w/full doc** divides the entire document into multiple segments according to the prompt length limit. Each segment is processed to generate partial answers, which are then aggregated to produce the final response. The same LLMs are used as in the LLM-only setting, and depending on document length, the document is split into 2 to 20 sub-prompts.

- **BM25 retriever + LLMs** is a traditional information retrieval method that evaluates keyword matching between queries and documents based on term frequency (TF) and inverse document frequency (IDF). The top- k documents are retrieved for each query, where $k = 5$ in this study [27].
- **Ada-002 embedding retriever + LLMs** uses the OpenAI embedding model to vectorize text and retrieve the top- k documents based on semantic similarity, capturing contextual meaning more effectively than BM25. In this experiment, $k = 5$ [22].
- **MindMap** is a reasoning framework that integrates LLMs with KGs. It first explores subgraphs relevant to the query, converts them into natural language representations suitable for LLM input, and uses them for reasoning. A key feature of this method

Table 7: A6/A6M performance comparison across different question types. Results of comparing various QA models on table, rule, and multi-hop question types.

Type	Method	table				rule				multi-hop			
		F1	BERTScore	BLEU	Rouge	F1	BERTScore	BLEU	Rouge	F1	BERTScore	BLEU	Rouge
LLM-only	GPT-4o-mini	0.004	0.784	0.003	0.005	0.000	0.799	0.000	0.000	0.000	0.803	0.000	0.000
	gemini-2.0-flash	0.004	0.784	0.003	0.004	0.000	0.799	0.000	0.000	0.003	0.804	0.001	0.003
	DeepSeek-v3.1	0.007	0.786	0.003	0.008	0.091	0.824	0.044	0.099	0.047	0.817	0.011	0.053
LLM w/full doc	GPT-4o-mini	0.083	0.805	0.016	0.100	0.244	0.857	0.136	0.269	0.294	0.864	0.163	0.347
	gemini-2.0-flash	0.110	0.806	0.008	0.132	0.370	0.879	0.216	0.410	0.290	0.865	0.139	0.350
	DeepSeek-v3.1	0.006	0.792	0.003	0.014	0.343	0.876	0.184	0.386	0.243	0.862	0.121	0.294
Text-only RAG	BM25+Qwen1.5-14B	0.053	0.817	0.009	0.112	0.311	0.867	0.121	0.356	0.295	0.869	0.124	0.344
	BM25+Mistral-7B	0.036	0.819	0.007	0.085	0.266	0.862	0.107	0.321	0.274	0.871	0.138	0.37
	BM25+GPT-4o-mini	0.043	0.803	0.010	0.072	0.333	0.871	0.161	0.376	0.349	0.877	0.164	0.396
	BM25+gemini-2.0-flash	0.016	0.788	0.003	0.030	0.367	0.882	0.195	0.423	0.346	0.876	0.165	0.401
	Dense+Qwen1.5-14B	0.078	0.836	0.013	0.116	0.339	0.876	0.152	0.387	0.314	0.876	0.141	0.378
	Dense+Mistral-7B	0.064	0.827	0.012	0.145	0.292	0.858	0.131	0.336	0.264	0.864	0.127	0.362
	Dense+GPT-4o-mini	0.063	0.799	0.012	0.085	0.340	0.871	0.178	0.380	0.362	0.877	0.191	0.414
	Dense+gemini-2.0-flash	0.044	0.792	0.006	0.075	0.364	0.876	0.192	0.399	0.324	0.874	0.162	0.390
KG-RAG (non-finetuned)	MindMap (Wen et al., 2024)	0.117	0.834	0.034	0.130	0.404	0.877	0.255	0.458	0.322	0.880	0.167	0.405
	KG Retriever (Chen et al., 2024)	0.107	0.843	0.013	0.224	0.356	0.874	0.178	0.413	0.185	0.853	0.057	0.224
Ontology-aware KG-RAG (non-finetuned)	Ontology-aware MindMap	0.468	0.898	0.135	0.562	0.438	0.885	0.285	0.472	0.344	0.878	0.192	0.398
	Ontology-aware KG-Retriever	0.124	0.842	0.029	0.106	0.474	0.897	0.163	0.303	0.355	0.881	0.185	0.405

is its entity-centric graph expansion and reasoning path organization, where the evidence graph is structured into a “mind map” [31].

- **KG-Retriever** consists of two stages: a document-level stage that selects documents semantically related to the query, and a KG-level stage that extracts relevant triples from within the selected documents. The retrieved triples are then inserted into the LLM prompt for answer generation [6].

C Dataset Details

Table 8: Representative question examples from IndusSpec-QA. This table presents representative sample questions extracted from the IndusSpec-QA dataset. Each question type (table, rule, multi-hop) reflects the structural characteristics and query patterns of industrial standard documents, derived directly from the specification texts.

Type	Question	Answer
table	What is the maximum percentage of Nickel allowed in the chemical requirements for Grade 60?	1
rule	What does CE stand for in the context of heat analysis?	carbon equivalent
multi-hop	What processes did TMCP evolve from, as discussed in the bibliographical reference regarding fine-grained steel?	controlled rolling processes

Table 9: Representative toxic clause question examples. This table presents sample questions from the toxic clause dataset. The questions reflect potential risk clauses related to industrial standards, and the answers are derived from the corresponding regulatory provisions. In particular, answers marked as No indicate toxic clauses that conflict with the actual standards.

Question	Answer
Does the regulation allow specimens to be taken without being adjacent to tensile test coupons?	No
If a specimen fails during the drop-weight test, must two additional specimens from each plate be retested?	Yes
Are blisters considered acceptable surface imperfections even if they do not reduce the thickness below minimum?	No

D Prompt Template

Prompt Example for Triplet Extraction

Instruction: You are an expert in constructing knowledge graphs. Your task is to extract knowledge triplets (subject–relation–object) from the given tabular or textual content.

Output Format:

<triplet> {subject} <subj> {object} <obj> {relation_name}

Guidelines (excerpt):

1. Extract only semantically meaningful and domain-relevant relations.
2. Normalize all values into SI units.
3. For plain text, use relations such as covers, applies_to, produced_by, used_for, stands_for, and defined_as.
4. For conditional statements and tabular data, extract triplets using only the relations has_condition[_AND/_OR] and has_consequence[_AND/_OR].
5. For tabular input:

- each non-empty cell becomes a case node (case1, case2, ...).
- Row and column labels are treated as conditions, while cell values are treated as consequences.
- Introduce logic_groupX when both AND/OR appear.

Examples:

Example1 (Input):

If the plate thickness exceeds 2 in. [50 mm], heat treatment shall be required.

Example1 (Output):

<triplet> case1 <subj> plate thickness ≥ 50 mm <obj> has_condition

$$\vdots$$

Now extract triplets from the following section:

Input: {text}

Output:

E Synonym Dictionary Visualization

F Ontology-Based Knowledge Graph Visualization

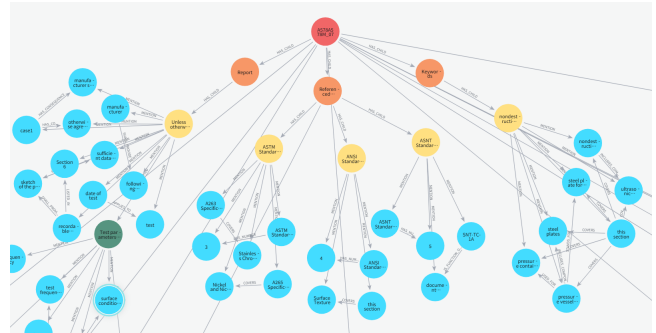


Figure 8: Visualization of an ontology-based knowledge graph using Neo4j. Node colors represent the hierarchical structure of the document: red nodes indicate **document names**, orange nodes represent **sections**, yellow nodes denote **subsections**, and green nodes correspond to **subsubsections**. Blue nodes illustrate specific **entities** mentioned within the document. The edges explicitly express relationships between document structure and knowledge, such as [HAS_CHILD], [MENTION], and [HAS_CONSEQUENCE_AND].

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009



Figure 7: Relation synonym dictionary visualization. This figure presents a partial Sunburst diagram illustrating the hierarchical relationships between canonical predicates (central concepts) and their corresponding paraphrases (synonymous expressions). Each central predicate encompasses a set of outer synonyms, visualized as clusters that integrate semantically similar relation expressions.