
OG-RAG: ONTOLOGY-GROUNDED RETRIEVAL-AUGMENTED GENERATION FOR LARGE LANGUAGE MODELS

A PREPRINT

Kartik Sharma, Peeyush Kumar, Yunqing Li
Microsoft Research
Seattle

December 23, 2024

ABSTRACT

This paper presents OG-RAG, an Ontology-Grounded Retrieval Augmented Generation method designed to enhance LLM-generated responses by anchoring retrieval processes in domain-specific ontologies. While LLMs are widely used for tasks like question answering and search, they struggle to adapt to specialized knowledge, such as industrial workflows or knowledge work, without expensive fine-tuning or sub-optimal retrieval methods. Existing retrieval-augmented models, such as RAG, offer improvements but fail to account for structured domain knowledge, leading to suboptimal context generation. Ontologies, which conceptually organize domain knowledge by defining entities and their interrelationships, offer a structured representation to address this gap. OG-RAG constructs a hypergraph representation of domain documents, where each hyperedge encapsulates clusters of factual knowledge grounded using domain-specific ontology. An optimization algorithm then retrieves the minimal set of hyperedges that constructs a precise, conceptually grounded context for the LLM. This method enables efficient retrieval while preserving the complex relationships between entities. OG-RAG applies to domains where fact-based reasoning is essential, particularly in tasks that require workflows or decision-making steps to follow predefined rules and procedures. These include industrial workflows in healthcare, legal, and agricultural sectors, as well as knowledge-driven tasks such as news journalism, investigative research, consulting and more. Our evaluations demonstrate that OG-RAG increases the recall of accurate facts by 55% and improves response correctness by 40% across four different LLMs. Additionally, OG-RAG enables 30% faster attribution of responses to context and boosts fact-based reasoning accuracy by 27% compared to baseline methods.

Keywords Retrieval Augmented Generation, Knowledge Graphs, Factual Deduction, Context Attribution, Ontology Grounded Retrieval

1 Introduction

Large language models (LLMs) have advanced the capabilities of question-answering systems, search engines, and task-oriented chatbots [Perplexity, 2024, ChatGPT, 2024, Achiam et al., 2023]. However, they face significant challenges with fact-based adaptation, particularly in domains that rely on precise, domain-specific data [Cascella et al., 2023, Thirunavukarasu et al., 2023, Singhal et al., 2023, Guha et al., 2024, Wang et al., 2024, Balaguer et al., 2024]. Consider a precision agriculture system where real-time changes in soil moisture and weather data must influence irrigation decisions. A general-purpose LLM might suggest irrigation plans based on broad knowledge but fail to account for specific soil conditions or plant requirements in that region. This lack of adaptability means the LLM's recommendation could be inaccurate, potentially leading to overwatering or under-irrigation, which can harm crops. Such scenarios highlight a core limitation: the inability of LLMs to reliably adapt to domain-specific decision-making, where accuracy and specialized knowledge are paramount [Rudin, 2019, Sharma et al., 2020].

To overcome these limitations, off-the-shelf LLMs can be either fine-tuned for specific domains [Bommasani et al., 2021] or paired with external tools or documents [Lewis et al., 2020, Zhuang et al., 2023, Schick et al., 2024].

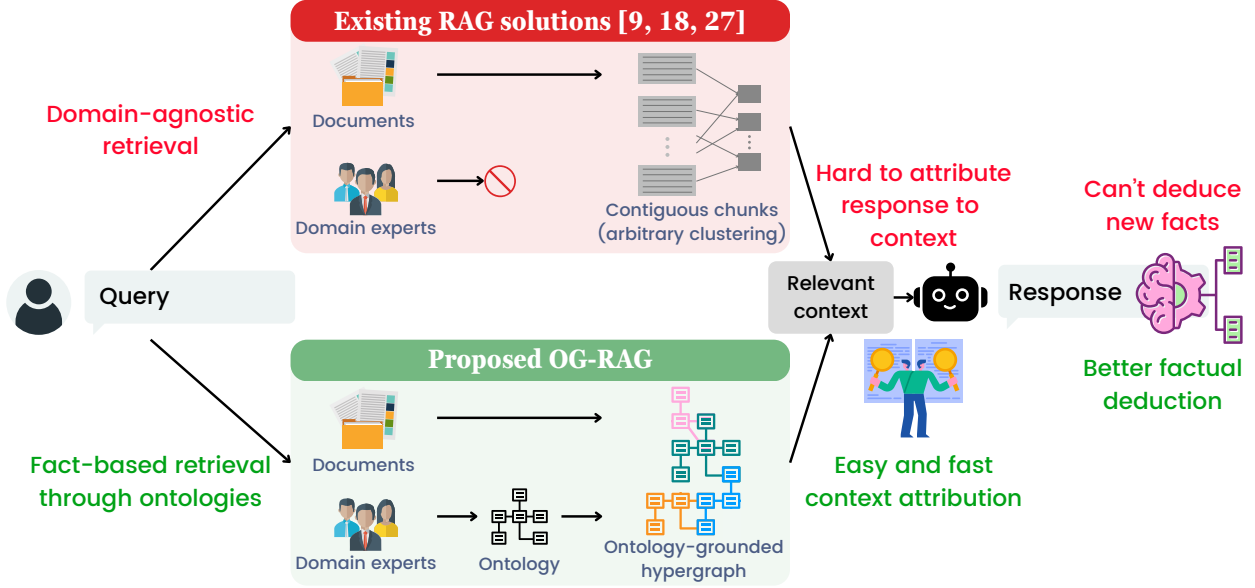


Figure 1: Comparison of the proposed OG-RAG with existing retrieval-augmented generation (RAG) solutions.

However, fine-tuning is computationally expensive and requires extensive data curation, making it a less practical solution [Balaguer et al., 2024, Ovadia et al., 2023]. On the other hand, retrieval-based approaches, such as RAG Lewis et al. [2020], Sarthi et al. [2024], Zhang et al. [2024], Borgeaud et al. [2022], Karpukhin et al. [2020], Edge et al. [2024], use domain-agnostic embeddings to retrieve query-relevant information from domain-specific documents and use the retrieved information for answering. Although promising, these methods fail to capture the deep conceptual relationships and nuanced facts required for accurate domain-specific retrieval.

Each domain organizes its knowledge and terminology in distinct ways, which cannot be generalized across different fields [Mernik et al., 2005]. For example, in industrial workflows, facts and relationships are carefully curated and structured into domain-specific frameworks, while in knowledge work and investigative research, ontologies serve as templates for organizing and analyzing facts and concepts [Jackson, 1990, Guarino et al., 2009]. Current LLMs struggle to adapt to these diverse structures, limiting their accuracy and effectiveness in specialized domains. Another major issue is that users often struggle to trace generated responses back to the relevant context. Furthermore, many specialized domains follow strict procedural rules, and the current techniques fail to reliably deduce accurate conclusions based on this established domain knowledge. This gap presents a major challenge to the wider applicability of LLMs in specialized workflows.

In this paper, we present OG-RAG (Ontology-Grounded Retrieval Augmented Generation) that bridges this gap of existing LLMs by integrating domain-specific ontologies for fact-based adaptation. Ontologies, which define key entities and their relationships within a domain, provide structured representation that is essential for adapting to complex and evolving information landscapes. As shown in Figure 1, OG-RAG leverages these ontologies to enhance LLM responses by grounding retrieval within structured domain knowledge, leading to *improved response accuracy*, supporting *flexible fact-based adaptation*, and enabling *verifiable context attribution*. A key feature of OG-RAG is its use of hypergraph representations of domain documents, which is a more sophisticated and multi-faceted way to model relationships than traditional retrieval approaches. Each hyperedge in the hypergraph represents a collection of related factual knowledge based on the corresponding ontology. Using a greedy algorithm, the engine retrieves a minimal set of hyperedges for a given query that forms a compact context for the LLM. Unlike traditional retrieval approaches, this method effectively distills complex relationships and domain-specific knowledge into a structured context, adapting LLMs to generate context-aware responses without adding significant computational overhead.

OG-RAG applies to a wide set of domains where fact-based adaptation is essential. These include industrial workflows in healthcare, legal, and agricultural sectors, as well as knowledge work such as news journalism, web based investigative research, consulting, and more. Our evaluations of OG-RAG within the agriculture and news domains demonstrate that OG-RAG increases the recall of accurate facts by 55%, and improves the overall correctness of generated responses by 40% across four different LLMs. Our user study shows that attributing LLM responses to the context retrieved by OG-RAG is 30% faster and better. Finally, in a fact-based reasoning task, we find that LLM responses are 27% more

correct when applying pre-defined rules over OG-RAG’s context compared to other methods. These results highlight OG-RAG’s effectiveness in providing more reliable, fact-based answers in specialized workflows.

2 Related Work

Fine-tuning. One approach to overcome the limitations of LLMs is fine-tuning on domain-specific data [Bommasani et al., 2021]. Fine-tuning allows models to adapt to the nuances of a specific domain by retraining the model on specialized datasets. However, this method is computationally expensive, requiring significant resources and extensive data curation, which makes it impractical for many real-world applications [Kumar et al., 2022]. OG-RAG addresses this shortcoming by eliminating the need for costly fine-tuning through retrieval-based solutions.

Hallucination mitigation. LLMs are prone to generating hallucinations, *i.e.*, outputs that are factually incorrect or irrelevant to the input [Ji et al., 2023a]. These hallucinations are especially problematic in domains that require precision, such as scientific research or industrial workflows [Thirunavukarasu et al., 2023]. Existing systems have attempted to mitigate hallucinations through post-generation correction methods and factuality checks, but these often require additional layers of computation and are not foolproof [Ji et al., 2023b, Madaan et al., 2024, Welleck et al., 2024]. OG-RAG reduces hallucinations by transforming data-mapped ontologies into hypergraphs and uses optimized retrieval of relevant fact clusters, ensuring LLM responses are grounded in domain-specific facts.

Retrieval methods In addition to the traditional retrieval augmented generation (RAG) [Lewis et al., 2020], graph-based approaches have also been proposed. These include GraphRAG [Edge et al., 2024], RAPTOR [Sarhi et al., 2024], and other knowledge graph-based frameworks such as Langchain [lan] and Neo4J [neo]. They have advanced LLM performance by leveraging structured knowledge graphs to organize and retrieve contextually relevant information. GraphRAG excels in semantic clustering by organizing entities and relationships, allowing for more efficient handling of complex queries, while RAPTOR uses a hierarchical structure for multi-level abstraction, improving contextual understanding across large documents. However, these approaches rely on ad-hoc extraction of entities and domain-specific information, often without grounding in domain expertise. This results in overly complex workflows for generating the correct structured representation, while still leaving significant gaps in precision. It also leads to weaker context attribution, making it more difficult to trace conclusions back to relevant facts. Improving on these approaches, OG-RAG’s hyperedge construction offers a compact fact representation that enhances transparency through better context attribution, while its hypergraph retrieval mechanism selects optimal fact clusters precisely tailored to the query.

Attribution. To enhance the interpretability and reliability of the LLM responses, it is important to attribute their generation to trustworthy sources. One way is to generate text with citations but prior work has shown limitations of existing zero-shot approaches [Gao et al., 2023] and specially-trained models [Khalifa et al., 2024]. Furthermore, other forms of attribution are also explored since citations require users to search over a full page to verify the claims in the generated response, which is undesirable. Thus, locally-attributable methods [Slobodkin et al., 2024] and human-in-the-loop Kamaloo et al. [2023] strategies have also been proposed. While these approaches provide sentence-level attribution, complementary benefits can be achieved through interpretable RAG contexts. OG-RAG provides easy-to-attribute contexts that require only a little effort from the users to trace the generation of the response.

Deductive reasoning. Traditional rule-based reasoning systems provide interpretable and easily controllable ways to deduce novel conclusions from a given input [Jackson, 1990, Saparov et al., 2023]. However, they lack the flexibility and generalization capabilities of neural models like LLMs. On the other hand, LLMs are prone to arbitrary hallucinations in deductive reasoning, which can be problematic in structured workflows [Wang et al., 2024, Saparov et al., 2023]. OG-RAG combines the structured precision of fact-based reasoning with neural flexibility by anchoring unstructured text to domain-specific vocabulary, enabling LLMs to more effectively apply domain-specific rules while maintaining scalability across multiple domains.

3 Summary of Key Contributions

This paper presents several key contributions that address the challenges of fact-based reasoning and hallucination reduction in large language models (LLMs):

- **Fact-based Context Retrieval Enabling Domain Adaptation:** OG-RAG is a novel framework that integrates domain-specific ontologies into the LLM retrieval process, facilitating precise, fact-based adaptation across domains. OG-RAG achieves this through two core mechanisms: (a) formalization of *facts* by transforming data-mapped ontologies into a hypergraph, and (b) using an optimization-based hypergraph retrieval to extract compact clusters of facts that precisely align with domain-specific queries, resulting in more contextually accurate responses.

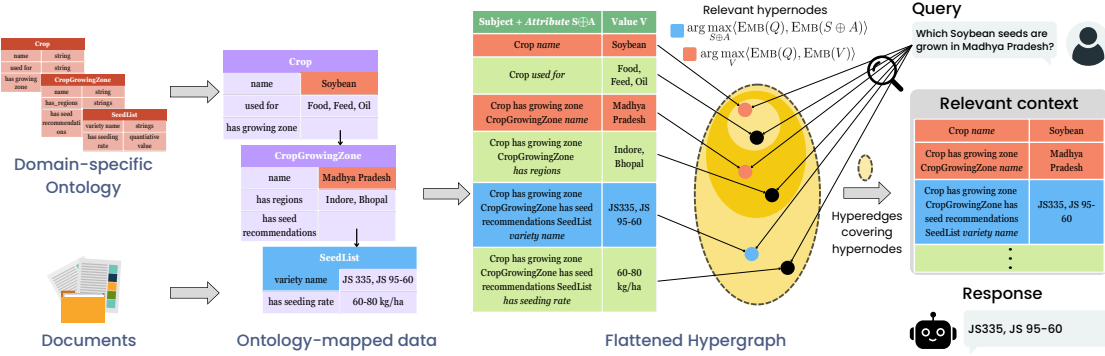


Figure 2: OG-RAG: Ontology-Grounded Retrieval-Augmented Generation

- **Factual Deduction:** OG-RAG enables factual deduction by leveraging domain-adapted facts to infer conclusions based on established knowledge and relationships. This enables the system to generate accurate, domain-specific conclusions dynamically, improving reasoning over complex factual contexts and enhancing the precision of LLM responses.
- **Improved Context Attribution:** OG-RAG enhances response verification by providing precise, fact-based context attribution. This enables both humans and machines to accurately trace conclusions back to their underlying facts, ensuring reliable and transparent fact-checking in domain-specific queries.

Section 6 presents extensive evaluations to support each of the enumerated contributions.

4 OG-RAG

Ontology Grounded Retrieval Augmented Generation (OG-RAG) is a novel framework which integrates ontologies—formal representations of domain-specific concepts and their relationships—into the retrieval process. Unlike existing retrieval-augmented generation (RAG) systems or other ontology based approaches, which rely on general-purpose embeddings or ad-hoc context generation without grounding in domain expertise, OG-RAG leverages ontology-driven hypergraph retrieval to dynamically adapt LLMs to structured knowledge bases and complex domain-specific queries. Figure 2 shows the high-level pipeline of the proposed method while we describe each component in more detail below.

4.1 Hypergraph Construction

The first part includes mapping the general domain-specific documents \mathcal{D} onto a given ontology \mathcal{O} and converting the available information into a usable format for retrieval.

4.1.1 Ontology

Different domains organize decision-making by following specific rules and procedures tailored to their unique workflows. In agriculture, for instance, factual information is organized differently depending on the task. Crop cultivation relies on facts like soil quality, weather patterns, and pest management strategies, where decisions follow a series of steps based on phenotype and environmental data. On the other hand, agricultural budgeting uses the same foundational facts—such as farm size or crop type—but applies them through financial models and cost projections, which require a different decision-making framework.

An **ontology** is a formal representation of key entities and their relationships within a domain. For example, in the agriculture domain, entities like crops, soil, and weather conditions are defined, along with relationships such as "crop is grown in a region" or "soil has moisture level." By defining these entities and relationships, the ontology provides a consistent and clear framework for organizing domain knowledge [Guarino et al., 2009, Jackson, 1990]. It differs from *taxonomy* or *classifications* as it allows for richer relationships between entities that need not be hierarchical. More formally,

Definition 1 (Ontology) An *ontology* $\mathcal{O} \subseteq \mathcal{S} \times \mathcal{A} \times (\mathcal{S} \cup \{\phi\})$ consists of a set of triples that relate a set of entities \mathcal{S} using a set of attributes \mathcal{A} , where $(s, a, v) \in \mathcal{O}$ denotes that the subject entity s has an attribute a , and the value $v := v_{\mathcal{O}}(s, a)$ is either:

- Another entity $s' \in \mathcal{S}$, or
- An unspecified domain value, denoted by ϕ .

Here, $v := v_{\mathcal{O}}(s, a)$ represents the value of the attribute a for entity s , which is either another entity within the ontology or an undefined (unspecified) text or data.

For example, consider a subject entity $s = \text{“Crop”}$, that can have the attribute $a_1 = \text{“is grown in”}$, which maps it to another object entity $v_{\mathcal{O}}(s, a_1) = s' = \text{“Crop Region”}$. Additionally, the same entity s can have another attribute $a_2 = \text{“has name”}$, which maps it to an arbitrary text, denoted as $v_{\mathcal{O}}(s, a_2) = \phi$, indicating that this value is unspecified and can be any relevant text or name in the domain.

Extracting factual information from domain-specific documents \mathcal{D} is challenging due to their specialized language and often underspecified structure. Moreover, relevant facts are frequently scattered across various unrelated documents. An ontology \mathcal{O} provides a structured way to organize key terms and their relationships within a domain.

To address this, we propose leveraging the explicit relationships defined in the ontology to extract factual information from these documents. We claim that since domain-specific facts are often grounded in the underlying ontology, enforcing these relationships can help enrich and disambiguate the information contained in the documents. In particular, we use the documents to find values for attributes by extracting relevant domain-specific text or values from the documents themselves (*i.e.*, when $v_{\mathcal{O}}(s, a) = \phi$). Since domain-specific documents may contain a variety of facts, this value assignment does not have to be unique across all documents. Instead, different parts of the documents may provide distinct yet valid text/data value related to the same ontology entity.

Therefore, we model the extracted information $\mathcal{I} := \mathcal{D}(\mathcal{O})$ using a set of self-contained *factual-blocks* $F \in \mathcal{D}(\mathcal{O})$, where each *factual-block* F consists of relationships that map ontology entities to either a unspecified domain text $\in \mathcal{V}$ or another entity within the same factual-block F . More formally,

Definition 2 (Ontology-mapped data) *Ontology-mapped data* $\mathcal{I} := \mathcal{D}(\mathcal{O})$ is information derived from the documents \mathcal{D} using the ontology \mathcal{O} . It consists of a set of *factual-blocks*, where each *factual-block* F represents a set of ontology relationships. For any relationship $(s, a, v) \in F$, the value v is derived as follows: If value $v_{\mathcal{O}}(s, a) = \phi$ then $v \in \mathcal{V}$ is extracted from the document text; otherwise $v = v_{\mathcal{O}}(s, a)$ is the value provided by the ontology.

Thus, ontology-mapped data represents *self-contained* and *ontology-grounded* information extracted from domain-specific documents. For example, a *factual-block* F might represent that: a term $s = \text{“Seed”}$ is $a_1 = \text{“of crop”}$ $v(s, a_1) = \text{“Soybean”}$ is $a_2 = \text{“is grown in”}$ $v(s, a_2) = (s' = \text{“Crop Region”},$ which $a_3 = \text{“has a name”}$ of $v(s', a_3) = \text{“Northwest Region”}$).

We can apply various pattern-matching heuristics, rule-based strategies, or embedding similarity to generate this ontology-mapped data [Otero-Cerdeira et al., 2015, Jackson, 1990]. However, with the powerful capabilities of LLMs, we leverage their natural language understanding capabilities to map these ontology entities to corresponding document text more effectively [Babaei Giglou et al., 2023]. We perform this task by prompting the LLM to generate the mapped information in JSON-LD format. The complete prompt is provided in Appendix B.1.

One limitation of this method is that domain-specific ontologies may not always be available or sufficiently comprehensive. To address this, we are developing an *ontology learning* method that can automatically generate a robust baseline ontology. This provides domain experts with a starting point, making it easier for them to edit and refine an existing ontology rather than building one from scratch. Additionally, in many fields, rich pre-existing ontologies are already available due to decades of research in data modeling and ontology development, which can be directly leveraged by this method. The details of this work are beyond the scope of the current paper, so we omit them.

4.1.2 Hypergraph Transformation

Due to the nested structure of the definitions in the *factual-block* $F \in \mathcal{D}(\mathcal{O})$, directly processing the information in these blocks is challenging. The combinatorial nature of the multi-layered relationships and dependencies make it difficult to efficiently extract or attribute information, which hampers our goal of providing compact and accurate context attribution. To address this, we flatten the structure so that each *factual-block* F in the ontology-mapped data \mathcal{I} is converted into a set of flattened *factual-blocks* \overline{F} , making the information easier to handle without significant loss of detail. Algorithm 1 outlines the flattening process, which is also illustrated in Figure 2.

Algorithm 1 Flattening a factual block**Require:** Factual block F , Concatenation operator \oplus .**Ensure:** A set of flattened factual-blocks $\overline{F} \leftarrow \text{FLATTEN}(F)$ flattens any nested information present in F .

```

1: procedure FLATTEN( $F$ )
2:    $\overline{F} \leftarrow \{\}$ 
3:    $\overline{F}_0 \leftarrow \{(s \oplus a, v) : (s, a, v) \in F, v \in \mathcal{V}, (s', a', s) \notin F\}$ .       $\triangleright$  no dependencies, can be directly flattened
4:    $\overline{F} \leftarrow \overline{F} \cup \{\overline{F}_0\}$ .
5:   for  $(s, a, s') \in F \setminus \overline{F}_0$  do.
6:     if  $s' \in \mathcal{S}$  then
7:        $F_{s'} \leftarrow \overline{F}_0 \cup \{(s \oplus a \oplus s' \oplus a', v') : (s', a', v') \in F\}$ .
8:        $\overline{F} \leftarrow \overline{F} \cup \text{FLATTEN}(F_{s'})$ .       $\triangleright$  flatten nesting of  $s'$ 
9:     end if
10:  end for
11:  return  $\overline{F}$ 
12: end procedure

```

We define each flattened factual-block as a hyperedge $e \in \mathcal{E}$, where a hyperedge e connects multiple hypernodes $\{n_i \in \mathcal{N}\}$, where each hypernode $n_i \in \overline{F}$ is a primitive set in the flattened-block that can be represented as a key-value pair. Importantly, this flattening process maintains the integrity of the information without introducing data loss. This allows OG-RAG to capture multi-dimensional relationships between facts, unlike simpler graph-based models that only handle pairwise connections. We can now convert the extracted information from the ontology into a more structured hypergraph, defined as follows:

Definition 3 (Hypergraph) A hypergraph $\mathcal{H} := (\mathcal{N}, \mathcal{E})$ consists of hypernodes \mathcal{N} and hyperedges \mathcal{E} , such that each hyperedge $e \in \mathcal{E}$ is a set of nodes with arbitrary length. Defining $\mathcal{P}(X)$ as the power set of X and $\bigoplus X$ as the set that is formed by concatenating the strings within each element of the set X , we have the hyperedges $\mathcal{E} \subseteq \mathcal{P}(\mathcal{N})$ and the hypernodes $\mathcal{N} \subseteq [\bigoplus \mathcal{P}(\mathcal{S} \times \mathcal{A})] \times \mathcal{V}$, where \times is the cartesian product.

Thus, the set of all flattened factual-blocks extracted from the ontology-mapped documents $\overline{\mathcal{I}}$ thus can be seen as a hypergraph. We call this simply $\mathcal{H}(\mathcal{I})$. With this definition, a hypernode is essentially a key-value pair and we declare a hyperedge to be a true *fact* grounded in domain-specific data. Mathematically,

Definition 4 (Fact) A fact is a logical assertion between two entities - subject and object, through a functional attribute, which can be evidentially verified to be either true or false. Formally, it can be expressed as a logical assertion that can be verified to have a value of True or False. For example, consider the assertion: $\text{hasCropYield}(\text{Farm } A) = 500$ tons, where hasCropYield is the functional attribute mapping a farm (subject) to a crop yield (value), and which can be evidentially verified to be either True or False.

Therefore, in OG-RAG a hyperedge can be viewed as a representation of a complex fact. Without loss of generality, consider two hypernodes, $n_1(s_1 \oplus a_1, v_1) = (\text{Crop has name, Soybean})$ and $n_2(p_2 \in \bigoplus \mathcal{P}(\mathcal{S} \times \mathcal{A}), v_2) = (\text{Crop has growing zone CropGrowingZone with name, Northwest})$ forming an hyperedge $e = ((\text{Crop has name, Soybean}), (\text{Crop has growing zone CropGrowingZone with name, Northwest}))$ can be represented as a simplified fact:

$$\text{hasGrowingZone}(\text{Crop has name Soybean}) = \text{Northwest},$$

which can be evidentially verified to be True or False.

In this way, the OG-RAG hypergraph construction enables a compact and accurate representation of *facts* that are adapted to the specific domain. This structure facilitates fact verification by allowing users to inspect the hyperedges, which encapsulate the relationships and dependencies between entities.

4.2 Hypergraph-based retrieval

With the hypergraph constructed on domain-specific information, *i.e.*, $\mathcal{H}(\mathcal{I}(\mathcal{D}, \mathcal{O}))$, OG-RAG is now ready to retrieve relevant context based on user query Q that can support the LLM in generating accurate, domain-specific responses.

Algorithm 2 Ontology-grounded Retrieval Augmented Generation

Require: Query Q , Domain-specific Ontology \mathcal{O} , Documents \mathcal{D} , Sentence embedding function \mathbf{Z} , LLM \mathcal{M}_0 , Maximum length L

Ensure: Retrieved context $\mathcal{C}_{\mathcal{H}}(Q)$ is grounded in the ontology and relevant to the query

```

1: procedure OG-PREPROCESS( $\mathcal{O}, \mathcal{D}, \mathcal{M}_0$ )
2:    $\mathcal{I} \leftarrow$  LLM  $\mathcal{M}_0$ (Ontology Map ( $\mathcal{D}, \mathcal{O}$ )) ▷ See Definition 2
3:    $\mathcal{H}(\mathcal{I}) \leftarrow$  Hypergraph with edges  $\bigcup_{F \in \mathcal{I}} \text{FLATTEN}(F)$ .
4: end procedure
5: procedure OG-RETRIEVE( $Q, \mathcal{H}(\mathcal{I}), \mathbf{Z}, k, L$ )
6:    $\mathcal{N}, \mathcal{E} \leftarrow$  nodes and edges of the hypergraph  $\mathcal{H}(\mathcal{I})$ .
7:    $\mathcal{N}_S(Q) \leftarrow$  top- $k$  arg max $_{(s,a,v) \in \mathcal{N}} \langle \mathbf{Z}(s \oplus a), \mathbf{Z}(Q) \rangle$ .
8:    $\mathcal{N}_V(Q) \leftarrow$  top- $k$  arg max $_{(s,a,v) \in \mathcal{N}} \langle \mathbf{Z}(v), \mathbf{Z}(Q) \rangle$ .
9:    $\mathcal{N}(Q) \leftarrow \mathcal{N}_S(Q) \cup \mathcal{N}_V(Q)$ .
10:   $\mathcal{C}_{\mathcal{H}}(Q) \leftarrow \{\}$ 
11:  while  $(|\mathcal{N}(Q)| > 0) \vee (|\mathcal{C}_{\mathcal{H}}(Q)| < L)$  do
12:     $\mathcal{C}_{\mathcal{H}}(Q) \leftarrow \mathcal{C}_{\mathcal{H}}(Q) \cup \arg \max_{e \in \mathcal{E}} |\{n \in \mathcal{N}(Q) : n \in e\}|$ 
13:  end while
14:  return  $\mathcal{C}_{\mathcal{H}}(Q)$ 
15: end procedure

```

4.2.1 Relevant Nodes

We first identify the set of hypernodes relevant to a given query. Using Definition 3, a hypernode $n \in \mathcal{N}$ can be represented as a key-value pair that comes from the elements in the sets $\mathcal{S}, \mathcal{A}, \mathcal{V}$. A hypernode can then be considered relevant to a query if: (1) the query pertains to an attribute a of the term s , or (2) the query focuses on an object with specific values v . In other words, a hypernode is relevant if either the similarity between the key (representing concatenated entities and attributes) and the query Q is high, or the similarity between v (the value) and the query Q is high. OG-RAG finds two sets of query-relevant hypernodes: $\mathcal{N}_S(Q)$ and $\mathcal{N}_V(Q)$ to represent the two sets respectively. In particular, $\mathcal{N}_S(Q)$ denotes the top k hypernodes with the highest similarity between their attributed term, i.e., $s \oplus a$ and the query Q in the vector space \mathbf{Z} . Similarly, $\mathcal{N}_V(Q)$ represents the top k hypernodes with the highest similarity between their value v and the query Q . Thus, for each query, we extract $2 \cdot k$ relevant hypernodes.

4.2.2 Relevant Hyperedges as Context

We form the relevant context as the set of hyperedges $\mathcal{C}_{\mathcal{H}}(Q) \subset \mathcal{E}$ that minimally cover the relevant hypernodes, $\mathcal{N}(Q) = \mathcal{N}_S(Q) \cup \mathcal{N}_V(Q)$. This is formulated as an optimization problem and solved in a greedy manner. Since the objective of minimizing the number of hyperedges is linear under a matroid constraint, the solution can be shown to be optimal [Korte et al., 2011]. Specifically, we maintain a dictionary that maps each hypernode $n \in \mathcal{N}$ to the set of hyperedges that it is a part of, i.e., $\mathcal{E}(n)$, where $e \in \mathcal{E}(n) \implies n \in e$. In each iteration, we add the hyperedge that covers the largest number of uncovered nodes to the context and remove those nodes from further consideration. This process is repeated until either we have L hyperedges or all the relevant nodes are covered. In this way, the context is constructed as a collection of up to L hyperedges representing *facts* relevant to the given query. By organizing information into hyperedges, OG-RAG is able to group related facts together, ensuring that the retrieved context is both compact and comprehensive, capturing all necessary facts to support accurate LLM responses, while optimizing for efficiency.

4.2.3 Retrieval-Augmented Generation

Finally, given a user query Q and the relevant context as found above, we prompt the LLM \mathcal{M} to use this context to answer the query as $\mathcal{M}(\mathcal{P}(Q, \mathcal{C}_{\mathcal{H}}(Q)))$, where \mathcal{P} denotes the corresponding textual prompt:

Given the context below, generate the answer to the given query. Note that the context is provided as a list of valid facts in a dictionary format.

Context: \langle Line-separated retrieved context $\mathcal{C}_{\mathcal{H}}(Q)$ \rangle

Query: \langle User-defined query Q \rangle

Answer:

4.3 Complexity Analysis

Algorithm 2 outlines the full procedure of the proposed method which consists of two main components: (1) OG-PREPROCESS, applied to the set of documents once, and (2) OG-RETRIEVE, used to retrieve the relevant context for each query.

4.3.1 Query Complexity

Assume the context size for the LLM \mathcal{M}_0 is N_C . The ontology \mathcal{O} , which can be written in a JSON-LD or textual format, has a length $|\mathcal{O}|$, where the attributes are mapped to their corresponding ranges in the natural language vocabulary. OG-PREPROCESS phase may involve several LLM calls depending on the number of document chunks, specifically, $(|\mathcal{D}| + |\mathcal{O}|)/N_C$ number of calls. We do not make any additional LLM calls during the querying time in the OG-RETRIEVE procedure.

4.3.2 Time Complexity

We ignore the time taken by LLM calls while calculating the time complexity, as this is accounted for under query complexity. Thus, the time complexity of the OG-PREPROCESS step only involves the hypergraph transformation by flattening the mapped data. Let us assume we have $|\mathcal{I}|$ factual-blocks derived from the documents, and each factual-block has a maximum length of $|F|_{max} = O(|\mathcal{O}|)$. We consider two cases: **(1) Minimal or No Nesting:** In this case, the time complexity is determined by step 4 in the algorithm, leading to a complexity of $O(|\mathcal{O}||\mathcal{I}|)$, **(2) Maximum Nesting:** In this scenario, step 4 may result in an empty set. Thus, each factual-block F can be recursively flattened $\log |\mathcal{O}|$ times while searching through the entire set, leading to a time complexity of $O(|\mathcal{I}||\mathcal{O}| \log |\mathcal{O}|)$.

4.3.3 Space Complexity

The only storage required is for the hypergraph structure $\mathcal{H}(\mathcal{I})$, which is directly proportional to the number of hyperedges $|\mathcal{E}| = |\overline{\mathcal{I}}|$.

5 Experimental Setup

Datasets. We evaluate OG-RAG across two distinct domain categories that involve specialized workflows: (a) Industrial workflows, with a focus on the agriculture domain, where precise, data-driven decisions are critical for crop management and resource allocation, and (b) Knowledge work, where we evaluate it on research and analysis tasks in the news domain. We avoid general domains like Wikipedia to mitigate potential data contamination in LLM training. For the agriculture domain, we utilize two proprietary high-quality datasets comprising of 85 documents prepared by agriculture experts, focusing on the crop cultivation of Soybean and Wheat in India. For the news domain, we use the publicly available dataset from Multi-hop RAG [Tang and Yang, 2024], filtered for 149 long-form articles (each over 2,000 words) focused on multi-faceted, complex news stories requiring detailed, contextually rich analysis. Please refer to Appendix A for exemplary excerpts from the datasets.

Ontology. We use a semi-automated approach to construct the ontology for both domains, which reflects the broader applicability of OG-RAG in specialized workflows. For the agriculture domain, the ontology was generated using a proprietary ontology learning module, which was then reviewed and verified by multiple experts specializing in crop cultivation. For the news domain, we modify the existing Simple News and Press (SNaP) ontology¹. Specifically, we simplify its structure by excluding certain classes, such as those related to images, videos, and the "stuff" hierarchy. Instead, we allow an asset to be linked to multiple events, and each event can be associated with multiple organizations and persons. The complete ontologies for both domains are provided in Appendix B.2.

Large Language Models. We consider 4 large language models for zero-shot query answering while adding the retrieved context from different methods: 2 closed-box models² (GPT-4o-mini and GPT-4o) and 2 open-source models³ (Llama-3.1-8B and Llama-3.1-70B). These models have been chosen for their remarkable understanding and ability to reason in natural language. We consider 4096 completion tokens and a temperature of 0.

Baselines. We compare OG-RAG against three leading retrieval-based methods, representing state-of-the-art approaches to context retrieval and generation, to demonstrate its effectiveness.

¹<https://iptc.org/thirdparty/snap-ontology/>

²<https://openai.com/index>

³<https://ai.meta.com/blog/meta-llama-3-1/>

Table 1: Quality of contexts retrieved by different methods for domain-specific query-answering. We found the 95% confidence interval to be ≤ 0.05 for all metrics, representing small margin of error. It is not reported here. The symbol ‘-’ denotes that the computation did not complete within 1 day.

Method	Soybean					Wheat					News				
	C-Rec	C-ERec	A-Corr	A-Sim	A-Rel	C-Rec	C-ERec	A-Corr	A-Sim	A-Rel	C-Rec	C-ERec	A-Corr	A-Sim	A-Rel
RAG	0.22	0.08	0.31	0.62	0.29	0.14	0.04	0.29	0.69	0.28	0.01	0.01	0.27	0.67	0.20
RAPTOR	0.54	0.19	0.34	0.68	0.68	0.85	0.29	0.59	0.79	0.89	0.82	0.46	0.58	0.84	0.76
GraphRAG	0.41	0.14	0.26	0.63	0.63	0.78	0.05	-	-	-	-	-	-	-	-
OG-RAG	0.84	0.41	0.48	0.72	0.79	0.95	0.34	0.62	0.79	0.79	0.82	0.52	0.66	0.86	0.73

1. **RAG [Lewis et al., 2020]:** RAG (Retrieval-Augmented Generation) retrieves query-relevant document chunks by embedding them into a vector space and then finding the context based on the maximum chunk-query similarity.
2. **RAPTOR [Sarathi et al., 2024]:** RAPTOR clusters document chunks into hierarchical structures and uses an LLM to summarize the clusters as additional context. For this experiment, we set the tree depth to 3 and use the collapsed-tree retrieval strategy.
3. **GraphRAG [Edge et al., 2024]:** GraphRAG retrieves from a knowledge graph constructed using an LLM by extracting entities and relationships and clustering them into semantic communities. We use default graph construction prompts and local search with community level as 2 for retrieval.

We use the text-embedding-3-small² as the sentence embedding function across all retrieval methods and GPT-4o as the LLM (*i.e.*, \mathcal{M}_0) for pre-processing. For each method, we find $\{2, 5\}$ similar contexts and select the one with the highest performance.

Metrics. Building on the RAGAS framework [RAGAS, 2024], we use the following metrics to assess the quality of the retrieved context and the generated responses while using text-embedding-3-small as the embedding model and GPT-4o as the LLM.

1. **Context Recall (C-Rec):** Proportion of claims in the ground-truth answer that can be attributed to the information present in the retrieved context.
2. **Context Entity Recall (C-ERec):** Proportion of entities in the ground-truth answer that are present in the retrieved context.
3. **Answer Similarity (A-Sim):** Similarity between the generated response and the ground-truth answer in the embedding space.
4. **Answer Correctness (A-Corr):** A combination of answer similarity (defined above) and factual similarity, which is the F1-score between the claims in the ground-truth answer and those in the generated response.
5. **Answer Relevance (A-Rel):** Measures how easily the original question can be inferred from the generated response.

Ethics Consideration. The survey was conducted in compliance with ACM’s Publications Policy on Human Research. No personal data was collected, and all responses were anonymized. For the agriculture and news domain experiments, only publicly available or proprietary datasets were used, with no sensitive data involved, adhering to ethical research standards.

6 Experiments

6.1 Query answering

6.1.1 Question Generation

We generate a set of question/answer pairs using the RAGAS framework [RAGAS, 2024] to validate the factual accuracy of our proposed method. RAGAS prompts off-the-shelf LLM to generate questions of varying difficulty, each with the corresponding ground-truth answers and contexts. Specifically, we generate up to 100 unique questions from RAGAS focused on multi-hop reasoning abilities, which is commonly required in specialized domain tasks. Examples of these generated questions, along with their ground-truth answers, are provided in Appendix C.

6.1.2 Does OG-RAG help in retrieving useful contexts?

A context is deemed useful for a query if it provides sufficient information to derive the ground-truth response. We evaluate this using Context Recall and Context Entity Recall. Table 1 compares the performance of different retrieval

Table 2: Quality of the answers generated by different LLMs using different retrieval methods. We found the 95% confidence interval to be ≤ 0.05 for all metrics, so it is not reported here. The symbol ‘-’ denotes that the computation did not complete within 1 day.

Method	Soybean					Wheat					News				
	C-Rec	C-ERec	A-Corr	A-Sim	A-Rel	C-Rec	C-ER	A-Corr	A-Sim	A-Rel	C-Rec	C-ERec	A-Corr	A-Sim	A-Rel
<i>Llama-3-8B</i>															
RAG	0.22	0.07	0.26	0.59	0.22	0.14	0.05	0.26	0.65	0.23	0.01	0.01	0.15	0.52	0.08
RAPTOR	0.56	0.20	0.34	0.66	0.59	0.84	0.35	0.54	0.76	0.67	0.82	0.47	0.53	0.74	0.68
GraphRAG	0.46	0.11	0.26	0.63	0.52	-	-	0.43	0.35	0.27	-	-	-	-	-
OG-RAG	0.82	0.40	0.40	0.65	0.60	0.95	0.33	0.54	0.73	0.72	0.81	0.51	0.52	0.76	0.69
<i>Llama-3-70B</i>															
RAG	0.24	0.06	0.27	0.59	0.19	0.14	0.03	0.26	0.65	0.14	0.01	0.01	0.17	0.58	0.09
RAPTOR	0.55	0.23	0.41	0.70	0.64	0.85	0.39	0.58	0.77	0.75	0.82	0.47	0.39	0.72	0.64
GraphRAG	-	-	0.30	0.65	0.55	-	-	0.47	0.37	0.29	-	-	-	-	-
OG-RAG	0.84	0.41	0.54	0.75	0.56	0.95	0.31	0.63	0.77	0.73	0.70	0.69	0.51	0.77	0.67
<i>GPT-4o-mini</i>															
RAG	0.24	0.07	0.29	0.66	0.59	0.14	0.05	0.33	0.73	0.66	0.01	0.01	0.34	0.73	0.64
RAPTOR	0.59	0.23	0.34	0.68	0.85	0.84	0.36	0.51	0.77	0.88	0.842995	0.364377	0.51	0.77	0.88
GraphRAG	0.42	0.13	0.25	0.63	0.65	0.78	0.05	0.35	0.70	0.85	-	-	-	-	-
OG-RAG	0.83	0.41	0.48	0.72	0.77	0.95	0.33	0.62	0.78	0.85	0.81	0.51	0.62	0.78	0.85
<i>GPT-4o</i>															
RAG	0.22	0.08	0.31	0.62	0.29	0.14	0.04	0.29	0.69	0.28	0.01	0.01	0.27	0.67	0.20
RAPTOR	0.54	0.19	0.34	0.68	0.68	0.85	0.29	0.59	0.79	0.89	0.83	0.46	0.58	0.84	0.76
GraphRAG	0.41	0.14	0.26	0.63	0.63	0.82	0.52	0.35	0.70	0.86	-	-	-	-	-
OG-RAG	0.84	0.41	0.48	0.72	0.79	0.95	0.34	0.62	0.79	0.79	0.82	0.52	0.66	0.86	0.73

methods across three datasets. OG-RAG outperforms the baselines in almost all cases, boosting the recall of correct claims by 55% and recall of correct entities by 110%. The only exception is the News dataset where OG-RAG matches the context recall performance of RAPTOR but still delivers better performance.

6.1.3 Does OG-RAG help generate factually accurate responses?

A useful context should lead to more factual and precise response when incorporated into the query for various LLMs. We evaluate this by comparing how closely the generated responses/answers align with the ground-truth answer when added as context across different LLMs. Table 2 presents the results of response correctness, similarity, and relevance for the 3 datasets. OG-RAG consistently outperforms the baselines, significantly improving answer correctness by 40%, and answer relevance by 16%. The only notable exception where OG-RAG slightly underperforms is in the Answer Relevance for Wheat and Soybean datasets in GPT-4o and Llama-3-70B. This is likely due to the broad scope of the retrieved context, which can sometimes introduce extraneous information. This can be possibly mitigated through further fine-tuning of the hypergraph retrieval mechanism, adjusting the level of detail to suit the complexity of the queries expected. We leave domain-specific optimization for future work, as the current approach already delivers good responses across all datasets.

6.1.4 Is OG-RAG efficient?

Finally, we demonstrate that OG-RAG is computationally efficient by comparing its pre-processing and per-query retrieval times with other methods across different datasets. Table 3 shows that OG-RAG performs nearly as efficiently as a simple RAG method, with only a minimal increase of at most 2 seconds during querying time despite being at least 100% better in factual accuracy. OG-RAG is also shown to have significantly lower computational time than more competitive baselines such as RAPTOR and GraphRAG at both the pre-processing and query stages, particularly highlighted by a 50% drop in the pre-processing times. This efficiency is particularly critical for real-time applications, such as agricultural monitoring systems, legal research, and automated news fact-checking, where quick retrieval and processing of domain-specific knowledge is essential.

6.2 Context attribution

6.2.1 Survey design

To assess how effectively the proposed method aids humans in verifying facts within LLM-generated responses, we conduct a human study measuring the time taken to verify whether the given context supports the generated response. We randomly select 10 queries from the agriculture dataset and present the responses generated by GPT-4o using both RAG and OG-RAG, each paired with their respective contexts. We exclude RAPTOR due to its content similarity with RAG, and GraphRAG due to its prohibitive context length. Participants are asked to evaluate the level of factual support

Table 3: Efficiency of different retrieval methods on domain-specific query-answering. T_{pre} and T_{query} denote the average pre-processing and query time in seconds. We found the variance to be within 5 seconds, so it is not reported here.

Method	Soybean			Wheat			News		
	$T_{pre} \downarrow$	$T_{query} \downarrow$	$ \mathcal{C} $	$T_{pre} \downarrow$	$T_{query} \downarrow$	$ \mathcal{C} $	$T_{pre} \downarrow$	$T_{query} \downarrow$	$ \mathcal{C} $
RAG	11.41	2.49		10.55	2.36		449.21	3.56	
RAPTOR	71.66	4.81		61.56	4.38		1513.57	5.45	
GraphRAG	157.04	5.95		307.37	5.65		>1 day	-	
OG-RAG	29.61	3.75		47.76	4.09		655.15	4.12	

Table 4: Comparison of the ease with which humans can attribute generated responses to the contexts produced by RAG and OG-RAG, presented with 95% confidence intervals.

Method	Time taken \downarrow	Support [1-5] \uparrow
RAG	61.15 \pm 28.48	2.67 \pm 0.30
OG-RAG	43.50 \pm 18.08	3.46 \pm 0.19

the context provides for the response on a scale of 1-5. We also track the time each participant takes to complete this task. Each participant is shown 10 questions, consisting of 5 random queries, each paired with both RAG and OG-RAG responses and contexts in a randomized order. To ensure fairness, each query is presented an equal number of times across all participants. Examples of the survey design can be found in Appendix D.

6.2.2 Results

A total of 16 participants, aged 18-34 and familiar with LLMs, took part in the survey. Table 4 presents the average time taken and the level of support participants attributed to the contexts. We observed that OG-RAG significantly reduced the time required by 28.8% and increased the human-attributed support by 29.6% on average. These results demonstrate that OG-RAG not only enables faster fact verification but also provides more robust and clear contexts, making the system more user-friendly and reliable for context fact attribution.

6.3 Factual Deduction

6.3.1 Deductive Facts

We assess OG-RAG’s ability to enhance deductive reasoning in LLMs by evaluating how well it can generate new conclusions based on a set of predefined facts. These facts, grounded in domain-specific ontologies, provide the framework for reasoning tasks that require multi-step logic. Specifically, for this experiment we use six agricultural facts to deduce CO2 emissions, as this information is not directly available in the documents. These facts are partially derived from industry sources on the relationship between fossil fuels, pesticides, and greenhouse gases.⁴:

1. Farm area in the North Eastern Hill zone is 1 hectare or ha.
2. Farm area in North Plain Hill zone is 2 hectares or ha.
3. Herbicide production is calculated by multiplying the farm area by the recommended herbicide quantity.
4. 1 kg of herbicide production results in 18.22–26.63 kg of CO2e emissions.
5. 1 kg of insecticide production results in 14.79–18.91 kg of CO2e emissions.
6. 1 kg of fungicide production results in 11.94–29.19 kg of CO2e emissions.

6.3.2 Question Generation

To create the evaluation test set, we prompt GPT-4o following the RAGAS guidelines RAGAS [2024] to generate questions that require the application of deductive facts and a randomly sampled chunk from the ontology-mapped data to generate the responses. Specifically, we use the following prompt:

⁴Adapted in part from <https://www.panna.org/news/linking-fossil-fuels-and-pesticides-to-greenhouse-gases>

Table 5: Comparison of different retrieval methods in their ability to support deductive reasoning from different LLMs.

Method	Soybean			Wheat		
	A-Corr	A-Sim	A-Rel	A-Corr	A-Sim	A-Rel
<i>GPT-4o-mini</i>						
RAG	0.46	0.89	0.66	0.41	0.92	0.64
RAPTOR	0.42	0.89	0.81	0.50	0.92	0.74
GraphRAG	0.44	0.91	0.83	0.49	0.93	0.82
OG-RAG	0.50	0.92	0.75	0.53	0.94	0.83
<i>GPT-4o</i>						
RAG	0.44	0.90	0.56	0.42	0.92	0.54
RAPTOR	0.01	0.11	0.03	0.41	0.91	0.74
GraphRAG	0.48	0.92	0.84	0.44	0.90	0.73
OG-RAG	0.56	0.92	0.75	0.47	0.94	0.83

Given the following data and a set of deductive rules, generate a hard question that requires the application of the rules on the data to generate the answer.

Data: < Domain-specific data >

Rules: < Fixed set of rules >

Question:

Next, we make two additional LLM calls to generate the corresponding answer and assign a rating from 1 to 10, evaluating how well the question tests the application of the rules on the data to derive the answer. We select 10 questions that receive a rating of at least 7. A full list of generated questions is provided in Appendix C.

6.3.3 Results

Table 5 presents the results of factual deductions across two agriculture datasets, using GPT-4o and GPT-4o-mini as the underlying LLMs. In all cases, except two, the OG-RAG context substantially improves the correctness, similarity, and relevance of the generated answers compared to baseline methods. This demonstrates that OG-RAG is more effective at supporting deductive reasoning from a fixed set of facts. One exception is in the Soybean dataset for answer relevance which again points to a slightly less pertinent answer due to a broader retrieved context by OG-RAG. Overall, these results confirm that OG-RAG provides a more robust context for deducing new facts than alternative retrieval methods.

7 Conclusion

In this work, we study the problem of domain adaptation of LLMs using ontology-grounded retrieval-augmented generation. We introduce OG-RAG, a novel hypergraph-based retrieval method that retrieves query-relevant context from documents by structuring their facts as a hypergraph using a domain-specific ontology. OG-RAG has wide applicability in domains which include industrial workflows in healthcare, legal, and agricultural sectors, among others as well as knowledge-driven tasks like news journalism, investigative research, consulting, and more. Through extensive experiments on two agriculture datasets and a news dataset, we demonstrate that OG-RAG significantly improves the factual accuracy of LLM-generated responses, while also enabling faster attribution of answers to their supporting context and more effectively deducing conclusions from domain facts. We recommend that LLMs have better ways to incorporate controlled vocabulary and structured evidence retrieval through fixed ontologies, as this not only enhances user comprehension of generated responses but also facilitates smoother integration of LLMs into industrial workflows and knowledge work. By offering greater flexibility and control over how context is retrieved and utilized, OG-RAG paves the way for more adaptable and reliable language systems. For future work we encourage to explore automated or

semi-automated ontology construction techniques to build these frameworks in an end-to-end fashion, ensuring broader applicability of retrieval-augmented models across diverse domains.

References

- Enhancing RAG-based application accuracy by constructing and leveraging knowledge graphs — [blog.langchain.dev](https://blog.langchain.dev/enhancing-rag-based-applications-accuracy-by-constructing-and-leveraging-knowledge-graphs/). <https://blog.langchain.dev/enhancing-rag-based-applications-accuracy-by-constructing-and-leveraging-knowledge-graphs/>. [Accessed 14-10-2024].
- Using a Knowledge Graph to Implement a RAG Application — [neo4j.com](https://neo4j.com/developer-blog/knowledge-graph-rag-application/). <https://neo4j.com/developer-blog/knowledge-graph-rag-application/>. [Accessed 14-10-2024].
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Hamed Babaei Giglou, Jennifer D’Souza, and Sören Auer. Llms4ol: Large language models for ontology learning. In *International Semantic Web Conference*, pages 408–427. Springer, 2023.
- Angels Balaguer, Vinamra Benara, Renato Luiz de Freitas Cunha, Roberto de M Estevão Filho, Todd Hendry, Daniel Holstein, Jennifer Marsman, Nick Mecklenburg, Sara Malvar, Leonardo O Nunes, et al. Rag vs fine-tuning: Pipelines, tradeoffs, and a case study on agriculture. *arXiv e-prints*, pages arXiv–2401, 2024.
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*, pages 2206–2240. PMLR, 2022.
- Marco Cascella, Jonathan Montomoli, Valentina Bellini, and Elena Bignami. Evaluating the feasibility of chatgpt in healthcare: an analysis of multiple clinical and research scenarios. *Journal of medical systems*, 47(1):33, 2023.
- ChatGPT. <https://chatgpt.com/>, 2024. [Accessed 08-10-2024].
- Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. From local to global: A graph rag approach to query-focused summarization. *arXiv preprint arXiv:2404.16130*, 2024.
- Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. Enabling large language models to generate text with citations. *arXiv preprint arXiv:2305.14627*, 2023.
- Nicola Guarino, Daniel Oberle, and Steffen Staab. What is an ontology? *Handbook on ontologies*, pages 1–17, 2009.
- Neel Guha, Julian Nyarko, Daniel Ho, Christopher Ré, Adam Chilton, Alex Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel Rockmore, Diego Zambrano, et al. Legalbench: A collaboratively built benchmark for measuring legal reasoning in large language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Peter Jackson. *Introduction to expert systems*. Addison-Wesley Longman Publishing Co., Inc., 1990.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38, 2023a.
- Ziwei Ji, Tiezheng Yu, Yan Xu, Nayeon Lee, Etsuko Ishii, and Pascale Fung. Towards mitigating llm hallucination via self reflection. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1827–1843, 2023b.
- Ehsan Kamaloo, Aref Jafari, Xinyu Zhang, Nandan Thakur, and Jimmy Lin. Hagrid: A human-llm collaborative dataset for generative information-seeking with attribution. *arXiv preprint arXiv:2307.16883*, 2023.
- Vladimir Karpukhin, Barlas Öğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906*, 2020.
- Muhammad Khalifa, David Wadden, Emma Strubell, Honglak Lee, Lu Wang, Iz Beltagy, and Hao Peng. Source-aware training enables knowledge attribution in language models. *arXiv preprint arXiv:2404.01019*, 2024.
- Bernhard H Korte, Jens Vygen, B Korte, and J Vygen. *Combinatorial optimization*, volume 1. Springer, 2011.
- Ananya Kumar, Aditi Raghunathan, Robbie Jones, Tengyu Ma, and Percy Liang. Fine-tuning can distort pretrained features and underperform out-of-distribution. *arXiv preprint arXiv:2202.10054*, 2022.

- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474, 2020.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhunoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
- Marjan Mernik, Jan Heering, and Anthony M Sloane. When and how to develop domain-specific languages. *ACM computing surveys (CSUR)*, 37(4):316–344, 2005.
- Lorena Otero-Cerdeira, Francisco J Rodríguez-Martínez, and Alma Gómez-Rodríguez. Ontology matching: A literature review. *Expert Systems with Applications*, 42(2):949–971, 2015.
- Oded Ovadia, Menachem Brief, Moshik Mishaeli, and Oren Elisha. Fine-tuning or retrieval? comparing knowledge injection in llms. *arXiv preprint arXiv:2312.05934*, 2023.
- Perplexity. <https://www.perplexity.ai/>, 2024. [Accessed 08-10-2024].
- RAGAS. <https://docs.ragas.io/en/stable/>, 2024. [Accessed 08-10-2024].
- Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature machine intelligence*, 1(5):206–215, 2019.
- Abulhair Saparov, Richard Yuanzhe Pang, Vishakh Padmakumar, Nitish Joshi, Mehran Kazemi, Najoung Kim, and He He. Testing the general deductive reasoning capacity of large language models using ood examples. *Advances in Neural Information Processing Systems*, 36:3083–3105, 2023.
- Parth Sarthi, Salman Abdullah, Aditi Tuli, Shubh Khanna, Anna Goldie, and Christopher D Manning. Raptor: Recursive abstractive processing for tree-organized retrieval. *arXiv preprint arXiv:2401.18059*, 2024.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36, 2024.
- Abhinav Sharma, Arpit Jain, Prateek Gupta, and Vinay Chowdary. Machine learning applications for precision agriculture: A comprehensive review. *IEEE Access*, 9:4843–4873, 2020.
- Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. Large language models encode clinical knowledge. *Nature*, 620(7972): 172–180, 2023.
- Aviv Slobodkin, Eran Hirsch, Arie Cattan, Tal Schuster, and Ido Dagan. Attribute first, then generate: Locally-attributable grounded text generation. *arXiv preprint arXiv:2403.17104*, 2024.
- Yixuan Tang and Yi Yang. Multihop-rag: Benchmarking retrieval-augmented generation for multi-hop queries. *arXiv preprint arXiv:2401.15391*, 2024.
- Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. Large language models in medicine. *Nature medicine*, 29(8):1930–1940, 2023.
- Siyuan Wang, Zhongyu Wei, Yejin Choi, and Xiang Ren. Can llms reason with rules? logic scaffolding for stress-testing and improving llms. *arXiv preprint arXiv:2402.11442*, 2024.
- Sean Welleck, Amanda Bertsch, Matthew Finlayson, Hailey Schoelkopf, Alex Xie, Graham Neubig, Ilia Kulikov, and Zaid Harchaoui. From decoding to meta-generation: Inference-time algorithms for large language models. *arXiv preprint arXiv:2406.16838*, 2024.
- Tianjun Zhang, Shishir G Patil, Naman Jain, Sheng Shen, Matei Zaharia, Ion Stoica, and Joseph E Gonzalez. Raft: Adapting language model to domain specific rag. *arXiv preprint arXiv:2403.10131*, 2024.
- Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and Chao Zhang. Toolqa: A dataset for llm question answering with external tools. *Advances in Neural Information Processing Systems*, 36:50117–50143, 2023.

Appendix

A Dataset Examples

A.1 Soybean

— title: SOYBEAN: AN INTRODUCTION Classification of States in six major Soybean growing zones: Soybean is majorly grown in the following areas —
 # SOYBEAN: AN INTRODUCTION Classification of States in six major Soybean growing zones: Soybean is majorly grown in the following areas
 ## TILLAGE

- Deep ploughing is essential during summer, after harvesting the Rabi crop. This facilitates exposing the hibernating insects to extreme heat and predatory birds as well as movement of nutrients and infiltration of soil water. Therefore, one deep ploughing once in 3-4 years, otherwise one normal ploughing in summer followed by 2 criss-cross harrowing or cultivation for breaking of soil clods will make ideal seed bed for a good crop of soybean cultivation is recommended. Also, sub-soiling operation once in 4-5 years at an interval of 10 meter, break the compactness of the sub-soil and also facilitate infiltration of rainwater which is useful for un-interrupted crop growth even during drought period also.

A.2 Wheat

CLIMATE

Wheat crop has wide adaptability. It is primarily a temperate crop but is widely cultivated in subtropical regions and is grown even in some tropical countries.

- Ideal temperature for germination of wheat seeds is around 20-23 o C though these can germinate in temperature range of 3.5 to 35 o C.
- During the heading and flowering stages, excessively high or low temperatures and drought are harmful to wheat.
- The temperature conditions at the time of grain filling and development are very crucial for yield.
- Temperatures above 25 o C during grain filling and development period tend to depress grain weight. When temperatures are high, too much energy is lost through the process of transpiration by the plants and the reduced residual energy results in poorer grain formation and lower yields.

A.3 News

Raiders vs. Lions live score, updates, highlights from NFL 'Monday Night Football' game
 author: Dan Treacy
 source: Sporting News
 published_at: 2023-10-30T22:20:03+00:00
 category: sports
 url: <https://www.sportingnews.com/us/nfl/news/raiders-lions-live-score-highlights-monday-night-football/d022b1d62b18af8a70c516f4>

The Lions just needed to get themselves back in the win column after a blowout loss in Baltimore, and they did just that in front of their home fans on Monday night. Detroit rolled to a 26-14 victory over the Raiders in a game that felt much more one-sided than the score indicates.

The Lions thoroughly outplayed the Raiders, out-gaining Las Vegas by 329 yards, but critical mistakes by Detroit left points on the board. The offense struggled to turn red zone opportunities into touchdowns in the first half, and two turnovers deep in Raiders territory – including a pick-six by Jared Goff – kept Las Vegas in the game in the second half. MORE: What to know about Lions' retro 'Honolulu blue' helmets

B Ontology

B.1 Ontology Prompt

Here is a context definition for wheat crop cultivation ontology.

Context Definition:
{context_definition}

Generate a JSON-LD using the following data and the above context definition for crop cultivation ontology.

Use '@graph' object namespace for the data in JSON-LD.

Be comprehensive and make sure to fill all of the data.

Keep nesting to the minimum and still be able to disambiguate.

If there are multiple subfields enumerated in a 'List' namespace then do not combine them in a single subfield, keep them as separate subfields to disambiguate.

Ensure that you populate all items in the 'List' namespace, do not leave any item.

Do not include any explanations or apologies in your response.

Do not add any other text other than the generated JSON-LD in your response.

Generate in Json format.

Data:

{data}

JSON-LD json:

B.2 Examples

Figures 3 and 4 show agriculture and news ontologies used in the work.

C Evaluation Questions

C.1 Soybean

Question: Which pest can be controlled with Imidacloprid 48 FS pesticide in soybean crops?

OG-RAG Answer: The pest that can be controlled with Imidacloprid 48 FS pesticide in soybean crops is the White Grub.

Question: What type of soil should be avoided for growing soybeans due to its medium to high nutrient holding capacity and medium water holding capacity?

RAPTOR Answer: The context does not explicitly mention a specific type of soil that should be avoided for growing soybeans due to its medium to high nutrient holding capacity and medium water holding capacity. However, it does provide recommendations for managing different soil types:

- In saline soils (pH >7.5), gypsum at 150-200 kg/ha should be applied along with organic manure. - In acidic soils (pH <6.5), lime at 600 kg/ha is recommended.

Based on this information, it can be inferred that while soybeans can be grown in various soil types, specific amendments are needed for saline and acidic soils to optimize growth. Therefore, no particular soil type is explicitly advised to be avoided, but appropriate soil management practices are necessary for different soil conditions.

OG-RAG Answer: Sandy soil should be avoided for growing soybeans due to its medium to high nutrient holding capacity and medium water holding capacity.

Question: What is the name of the crop that requires seeds with a germination rate of at least 80%, should be harvested when 85-90% pods have turned brown and lost their green color, and should be stored in a cool, dry place using hermetic storage bags for long-term storage?

RAPTOR Answer: The context provided does not mention a crop that requires seeds with a germination rate of at least 80%, should be harvested when 85-90% pods have turned brown and lost their green color, and should be stored in a cool, dry place using hermetic storage bags for long-term storage. The context specifically discusses soybean harvesting, threshing, and storage practices, including a germination rate of at least 70%, harvesting when 95% of pods change color, and storage recommendations. Therefore, based on the given context, there is no crop that matches the criteria described in the question.

OG-RAG Answer: The name of the crop is Soybean.

Question: Which soybean crop is recommended for the Central Zone and has an early sowing time on January 1st, 2016?

RAPTOR Answer: The context provided does not mention any soybean crop recommended for the Central Zone with an early sowing time on January 1st, 2016. The recommended sowing times for the Central Zone are from 20th June to 5th July. Therefore, there is no soybean crop recommended for the Central Zone with a sowing time of January 1st, 2016.

OG-RAG Answer: The soybean crop recommended for the Central Zone with an early sowing time starting on January 1st, 2016, is the variety "JS 20-69."

Question: What is the temperature range for soybean germination and how does it compare to the temperature requirements for flowering and grain filling?

RAPTOR Answer: The provided context does not specify the exact temperature range required for soybean germination, flowering, or grain filling. Therefore, based on general agronomic knowledge:

1. ****Soybean Germination**:** Soybeans typically germinate best at soil temperatures between 15°C to 30°C (59°F to 86°F). Optimal germination usually occurs around 25°C (77°F).

2. ****Flowering and Grain Filling**:** For flowering and grain filling, soybeans generally require slightly higher temperatures. The optimal temperature range for these stages is typically between 20°C to 30°C (68°F to 86°F). Temperatures above 30°C (86°F) can negatively impact these processes, especially if they are sustained over long periods.

In summary, while soybean germination can occur at a range of temperatures starting from around 15°C, the optimal conditions for flowering and grain filling are generally within a similar but slightly higher range, emphasizing the importance of maintaining moderate temperatures throughout the growing season for optimal yield.

OG-RAG Answer: The temperature range for soybean germination is 15-30°C. For flowering, the temperature requirement is 20-30°C, and for grain filling, it is also 20-30°C. This means that the temperature range for germination is slightly broader on the lower end compared to flowering and grain filling, which have the same temperature requirements.

C.2 Wheat

Question: What is the recommended moisture level for harvesting wheat and how should the grain be dried before storage?

OG-RAG Answer: The recommended moisture level for harvesting wheat is 25%. Before storage, the grain should be dried by spreading it on tarpaulin plastic sheets in bright sunlight to a moisture level below 12

Question: Why conduct a seed germination test for wheat in relation to food production?

OG-RAG Answer: Conducting a seed germination test for wheat is essential to ensure a minimum germination rate of 85%. This is crucial for achieving an optimum plant population, which directly impacts the yield. By verifying the germination status of seeds, farmers can ensure that they plant seeds that are likely to grow successfully, leading to better food production and higher yields.

Question: What are the recommended practices for harvesting and drying wheat grains in the NEP Zone before storage?

OG-RAG Answer: The recommended practices for harvesting and drying wheat grains in the North Western Plains (NEP) Zone before storage are as follows:

1. **Harvesting:** - Harvest when the moisture level of the grain is 25- Wheat should be harvested 4-5 days before it is dead ripe. - Morning is the best time for harvesting. - For quick harvesting, use a combine harvester to avoid losses due to shattering and lodging. - In case of manual harvesting, make bundles and dry them for 3-4 days before threshing with a thresher.
2. **Drying:** - Before storage, dry the grain by spreading it on tarpaulin plastic sheets in bright sunlight. - Ensure the moisture level is reduced to below 12% before storage.

Question: What are the recommended practices for harvesting and drying wheat grains in the Northern Hills Zone before storage?

OG-RAG Answer: In the Northern Hills Zone, the recommended practices for harvesting and drying wheat grains before storage are as follows:

1. **Harvesting:** - Harvest when the moisture level of the grain is 25- Wheat is generally manually harvested, but a combine harvester can be used for quick harvesting to avoid losses due to shattering and lodging. - Harvest 4-5 days before the wheat is dead ripe. - Morning is the best time for harvesting. - In case of manual harvesting, make bundles and dry them for 3-4 days before threshing with a thresher.
2. **Drying:** - Before storage, dry the grain by spreading it on tarpaulin plastic sheets in bright sunlight to reduce the moisture level below 12%.

C.3 News

Question: Which earbuds are stylish and see-through?

OG-RAG Answer: The earbuds that are stylish and see-through are the "Nothing Ear Stick."

Question: What else debuts with 'Hurricane Season' on Netflix Nov. 1?

OG-RAG Answer: On Netflix Nov. 1, along with "Hurricane Season," the following titles debut:

- Locked In (Netflix Film) - Mysteries of the Faith (Netflix Documentary) - Nuovo Olimpo (Netflix Film) - Till Murder Do Us Part: Soering vs. Haysom (Netflix Documentary) - Wingwomen (Netflix Film) - 13 Going on 30 - 13 Hours: The Secret Soldiers of Benghazi - 60 Days In: Season 5 - The Addams Family - The Big Lebowski - The Change-Up - Desperado - Downsizing - Drag Me to Hell - Love in the Wild: Season 1 - Madea's Family Reunion - The Mummy (2017) - Paul Blart: Mall Cop - Paul Blart: Mall Cop 2 - Pitch Perfect - The Scorpion King - Scott Pilgrim vs. the World - Sherlock Gnomes - Six Feet Under: Seasons 1-5 - Sixteen Candles - The Social Network - Ted 2 - Victorious: Season 3 - Whiplash

Question: How do price-match policies curb customer loss in big sales?

OG-RAG Answer: Price-match policies help curb customer loss during big sales by ensuring that customers can get the best price without needing to shop around at multiple retailers. By offering to match lower prices found elsewhere, retailers like Best Buy, Walmart, Dell, eBay, and Lowe's aim to keep customers satisfied and loyal, preventing them from purchasing from competitors. This strategy is particularly effective during major sales events like Black Friday and Cyber Monday, when price competition is intense. By guaranteeing competitive pricing, retailers can maintain their customer base and potentially attract new customers who are looking for the best deals.

Question: Who wrote about AI risks and the creative class?

OG-RAG Answer: Daniel Tencer wrote about AI risks and the creative class.

Question: Who hit the Hawks' game-winning buzzer-beater vs. the Taipans?

OG-RAG Answer: Tyler Harvey hit the Hawks' game-winning buzzer-beater vs. the Taipans.

Question: Which Chrome extension by Steven Tey aids in sharing AI replies?

OG-RAG Answer: The Chrome extension by Steven Tey that aids in sharing AI replies is called "ShareGPT."

Question: Which artists thrived under Michelle Jubelirer at Capitol Music Group?

OG-RAG Answer: The artists that thrived under Michelle Jubelirer at Capitol Music Group include Ice Spice and The Beatles.

Question: Why are JLab earphones seen as a budget-friendly backup with EQ presets?

OG-RAG Answer: JLab earphones are considered a budget-friendly backup because they offer great value by including three different EQ presets directly into the earphones. This feature enhances the listening experience without a high cost, making them an attractive option for those seeking an affordable alternative to their main wireless earphones.

C.4 Factual Deduction

Question: Calculate the total CO₂e emissions from herbicide production for a farm located in the North Eastern Plains Zone, using the recommended herbicide for Phalaris minor. Assume the farm area is 1 hectare. Provide the range of CO₂e emissions based on the given data and rules.

OG-RAG Answer: To calculate the total CO₂e emissions from herbicide production for a farm located in the North Eastern Plains Zone, we need to follow these steps:

1. ****Identify the Recommended Herbicide and Quantity:**** - The recommended herbicide for Phalaris minor is Clodinafop. - The recommended quantity is 60 gm a.i. (active ingredient) per hectare.
2. ****Calculate the Total Herbicide Required:**** - Since the farm area is 1 hectare, the total herbicide required is 60 grams.
3. ****Convert Grams to Kilograms:**** - 60 grams is equivalent to 0.06 kilograms.
4. ****Calculate CO₂e Emissions:**** - The CO₂e emissions for herbicide production range from 18.22 to 26.63 kg of CO₂e per kilogram of herbicide. - For 0.06 kg of herbicide, the CO₂e emissions are calculated as follows: - Lower bound: $0.06 \text{ kg} * 18.22 \text{ kg CO}_2\text{e/kg} = 1.0932 \text{ kg CO}_2\text{e}$ - Upper bound: $0.06 \text{ kg} * 26.63 \text{ kg CO}_2\text{e/kg} = 1.5978 \text{ kg CO}_2\text{e}$

Therefore, the total CO₂e emissions from herbicide production for the farm are in the range of 1.0932 to 1.5978 kg CO₂e.

Question: Calculate the total potential CO₂e emissions from herbicide, insecticide, and fungicide production for a farm located in the North Plain zone, assuming the recommended herbicide quantity is 2 kg per hectare, insecticide quantity is 1.5 kg per hectare, and fungicide quantity is 1 kg per hectare. Use the average CO₂e values for your calculations.

OG-RAG Answer: To calculate the total potential CO₂e emissions from herbicide, insecticide, and fungicide production for a farm located in the North Plain zone, we need to follow these steps:

1. ****Determine the farm area**:** The farm area in the North Plain Hill zone is 2 hectares.
2. ****Calculate the total quantity of each type of chemical used**:** - Herbicide: $2 \text{ kg/ha} * 2 \text{ ha} = 4 \text{ kg}$ - Insecticide: $1.5 \text{ kg/ha} * 2 \text{ ha} = 3 \text{ kg}$ - Fungicide: $1 \text{ kg/ha} * 2 \text{ ha} = 2 \text{ kg}$
3. ****Use the average CO₂e values for each type of chemical**:** - Herbicide: Average CO₂e = $(18.22 + 26.63) / 2 = 22.425 \text{ kg CO}_2\text{e/kg}$ - Insecticide: Average CO₂e = $(14.79 + 18.91) / 2 = 16.85 \text{ kg CO}_2\text{e/kg}$ - Fungicide: Average CO₂e = $(11.94 + 29.19) / 2 = 20.565 \text{ kg CO}_2\text{e/kg}$

4. ****Calculate the total CO₂e emissions for each type of chemical****: - Herbicide: $4 \text{ kg} \times 22.425 \text{ kg CO}_2\text{e/kg} = 89.7 \text{ kg CO}_2\text{e}$ - Insecticide: $3 \text{ kg} \times 16.85 \text{ kg CO}_2\text{e/kg} = 50.55 \text{ kg CO}_2\text{e}$ - Fungicide: $2 \text{ kg} \times 20.565 \text{ kg CO}_2\text{e/kg} = 41.13 \text{ kg CO}_2\text{e}$

5. ****Calculate the total potential CO₂e emissions****: - Total CO₂e = $89.7 \text{ kg} + 50.55 \text{ kg} + 41.13 \text{ kg} = 181.38 \text{ kg CO}_2\text{e}$

Therefore, the total potential CO₂e emissions from herbicide, insecticide, and fungicide production for the farm in the North Plain zone is 181.38 kg CO₂e.

Question: Calculate the total CO₂e emissions from herbicide production for a farm in the North Eastern Hill zone, assuming the recommended herbicide quantity is 3 kg per hectare. Use the range provided for CO₂e emissions per kg of herbicide to determine the minimum and maximum possible emissions.

OG-RAG Answer: To calculate the total CO₂e emissions from herbicide production for a farm in the North Eastern Hill zone, we need to follow these steps:

1. ****Determine the farm area****: According to the context, the farm area in the North Eastern Hill zone is 1 hectare.

2. ****Calculate the total herbicide quantity****: The recommended herbicide quantity is 3 kg per hectare. Therefore, for a 1-hectare farm, the total herbicide quantity is:

$$3 \text{ kg/ha} \times 1 \text{ ha} = 3 \text{ kg}$$

3. ****Calculate the CO₂e emissions****: The CO₂e emissions per kg of herbicide range from 18.22 to 26.63 kg of CO₂e. Therefore, the total CO₂e emissions for 3 kg of herbicide are:

- ****Minimum emissions****:

$$3 \text{ kg} \times 18.22 \text{ kg CO}_2\text{e/kg} = 54.66 \text{ kg CO}_2\text{e}$$

- ****Maximum emissions****:

$$3 \text{ kg} \times 26.63 \text{ kg CO}_2\text{e/kg} = 79.89 \text{ kg CO}_2\text{e}$$

Therefore, the total CO₂e emissions from herbicide production for the farm in the North Eastern Hill zone range from 54.66 kg to 79.89 kg of CO₂e.

D Survey design

Figure 5 shows the design of our survey with the introduction and one exemplary question.

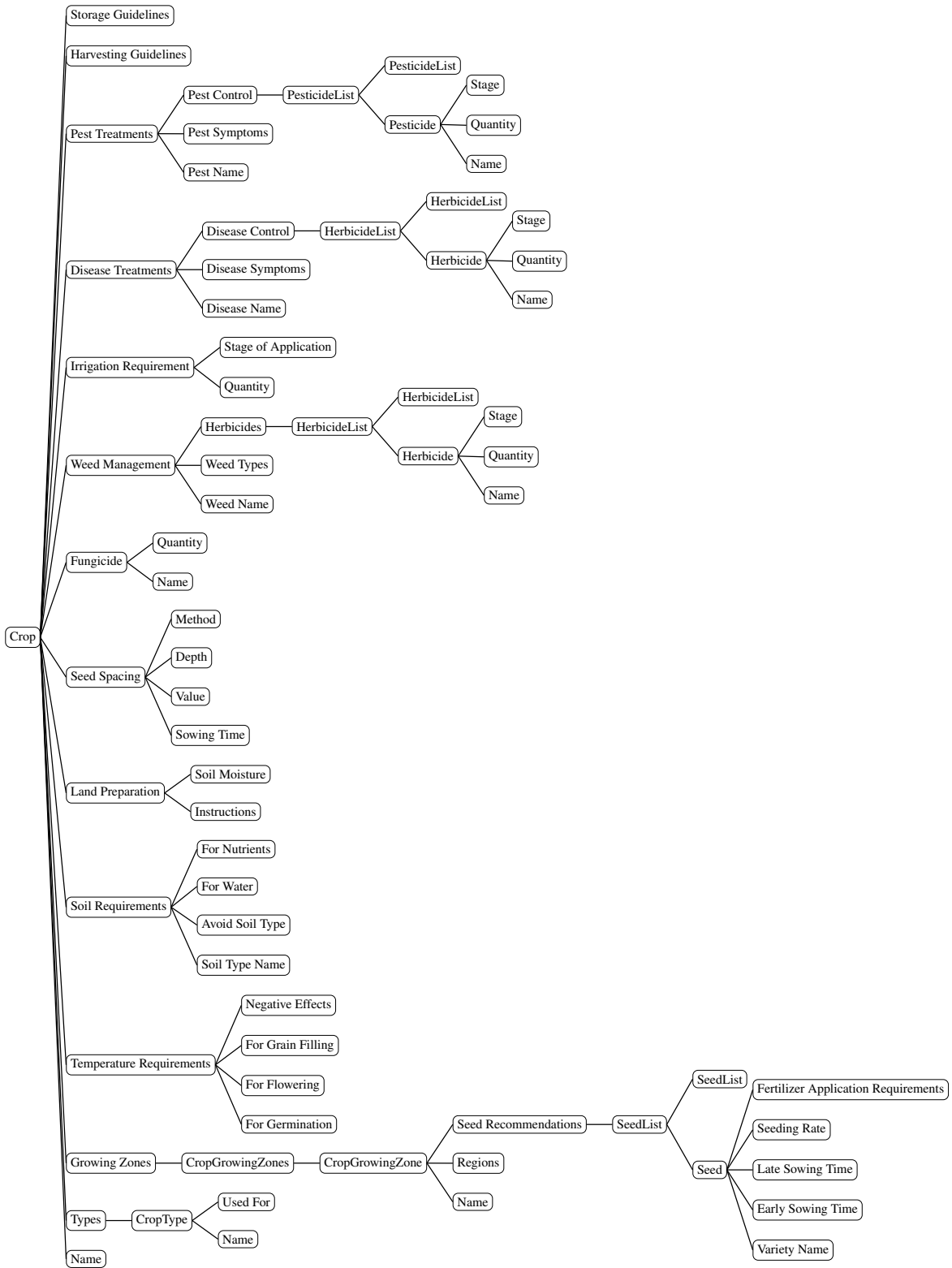


Figure 3: Agriculture ontology

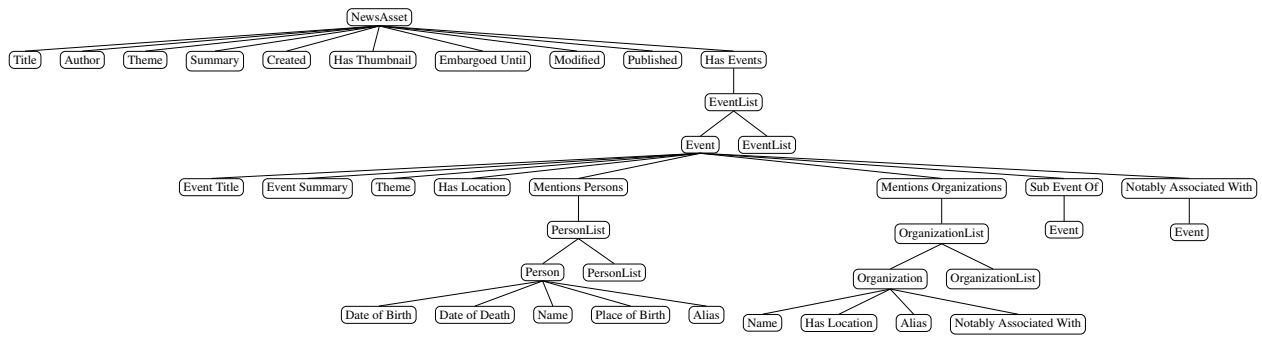


Figure 4: News ontology

Timing

These page timer metrics will not be displayed to the recipient.

First Click: 4.763 seconds
 Last Click: 4.763 seconds
 Page Submit: 0 seconds
 Click Count: 1 clicks

Given the query below, you are supposed to verify the claims in the response using the adjoining context.

Query: What is the name of the crop that requires seeds with a germination rate of at least 80%, should be harvested when 85-90% pods have turned brown and lost their green color, and should be stored in a cool, dry place using hermetic storage bags for long-term storage?

Response

Soybean

Context

HARVESTING & THRESHING

• Optimum time of harvesting is very important for soybean as it causes yield loss due to shattering and seed viability loss due to field weathering. • Initiation of change of pod color (physiological maturity) of soybean is the right indication to go for harvesting. At this time, the moisture percentage of soybean seed is around 14-16%.

Hence, the farmers are advised not to wait till the foliage/plant turns yellow. Harvesting can be done when 90% pod change color from green. • The harvested soybean crop, after sun drying for 2-3 days is ready for threshing. Because of sowing of subsequent crop if threshing is to be performed later, the harvested soybean should be collected preferably on threshing floor and covered with tarpaulin to avoid damage from rain/shattering. To maintain viability of the seed and to avoid low mechanical damage, threshing should be done at 306-400 rpm. • The threshed soybean should again be allowed to sun dry for 3-4 days in order to bring down the moisture up to 10% which is essential to avoid fungal infection during storage. • The storage place should be cool with aeration and insect free. The soybean bags should be kept upright as far as possible. If racking is to be done, it should be only up to 4-5 bags of not more than 5 feet height using platforms in order to maintain the viability/germination of soybean seed. • While moving the seed bags to storage house, it should be carefully placed at the appropriate place/platform. The seed bags should not be in direct contact with floor/wall. The moisture seepage in the walls/floor may be a source of infection of diseases, hence can be avoided for storage.

SEED GERMINATION TEST

Farmers are advised to check germination status of seed purchased available with them before sowing. To ensure optimum plant population and thereby good yield, minimum 70% germination is essential. This can be done through sowing of 100 seeds in 14x21cm plot and it is kept moist. From 5-8 days emergence is counted everyday till the count is stabilized. The germination test can also be done by placing 100 seeds in between two newspaper sheets and rolling them with a moist cloth.

Timing

These page timer metrics will not be displayed to the recipient.

First Click: 0 seconds
 Last Click: 0 seconds
 Page Submit: 0 seconds
 Click Count: 0 clicks

Given the query below, you are supposed to verify the claims in the response using the adjoining context.

Query: What is the name of the crop that requires seeds with a germination rate of at least 80%, should be harvested when 85-90% pods have turned brown and lost their green color, and should be stored in a cool, dry place using hermetic storage bags for long-term storage?

Response (scrollable)

Soybean

Context (both sides are true)

Crop Name: Soybean Harvesting Guidelines: Harvest when 85-90% pods have turned brown and lost their green color. Seed Fungicide Requirements: The moisture content of the seeds should be around 10-12%. Seed Fungicide Requirements: - Fungicide Name: Thiam or Carbendazim - Fungicide Quantity: 2 g per kg of seeds Seed Germination Test Requirements: Seeds should have more than 80% germination. Storage Guidelines: Store in a cool and dry place. Storage Guidelines: Store in a cool, dry place. Use hermetic storage bags for long-term storage.	Crop Name: Soybean Harvesting Guidelines: Harvest when 85-90% pods turn brown and hard. Seed Fungicide Requirements: - Fungicide Name: Thiam or Carbendazim - Fungicide Quantity: 2 g per kg of seeds Seed Germination Test Requirements: Seeds should have more than 80% germination. Storage Guidelines: Store in a cool and dry place.
---	---

How old are you?

Under 18

18-24 years old

25-34 years old

35-44 years old

45-54 years old

55-64 years old

65+ years old

How familiar are you with language models?

Not at all familiar

Somewhat familiar but don't use them often

Quite familiar but don't use them often

Very familiar and use them often

How useful is the context for verifying the claims of the response?

Not at all useful

Can be useful but not specific to the query

Provides support for some claims

Provides support for all claims

How useful is the context for verifying the claims of the response?

Not at all useful

Can be useful but not specific to the query

Provides support for some claims

Provides support for all claims

(a) Introduction page

(b) RAG context attribution

(c) OG-RAG context attribution

Figure 5: Survey design