

---

# Supplementary Materials for “OntoLearner: A Modular Python Library for Ontology Learning with Large Language Models”

---

1	<b>Contents</b>	
2	<b>1 LLMs4OL Paradigm Tasks Definition</b>	<b>2</b>
3	1.1 Term Typing . . . . .	2
4	1.2 Taxonomy Discovery . . . . .	2
5	1.3 Non-Taxonomic Relationship Extraction . . . . .	3
6	<b>2 Complexity Score Analysis</b>	<b>3</b>
7	<b>3 LLMs Contamination Analysis</b>	<b>4</b>
8	<b>4 Experimental Resources and Setup</b>	<b>5</b>
9	4.1 Selected Datasets Statistics . . . . .	5
10	4.2 How to Run OntoLearner for LLMs4OL Tasks . . . . .	6
11	4.3 Standardized Prompts for LLMs . . . . .	7
12	4.4 Computational Details . . . . .	7
13	<b>5 Detailed Results</b>	<b>10</b>
14	5.1 Retrievers . . . . .	10
15	5.2 Reranking for the RAG Pipeline . . . . .	11

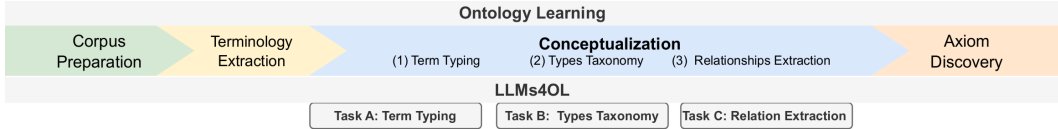


Figure 1: Overview of the OL pipeline and the role of the LLMs4OL paradigm tasks. The workflow spans from corpus preparation and terminology extraction to conceptualization and axiom discovery. Within the conceptualization phase, LLMs4OL focuses on three core tasks: (A) term typing, (B) taxonomy discovery, and (C) non-taxonomic relation extraction. These tasks collectively support the construction and enrichment of ontologies using LLMs.

## 1 LLMs4OL Paradigm Tasks Definition

In this section, we formalize the core ontology learning (OL) tasks addressed in the LLMs4OL paradigm (see Figure 1) and illustrate each task with concrete examples [6, 18, 19]. These tasks capture fundamental ontology engineering operations and serve as standardized benchmarks for evaluating Large Language Models (LLMs) within OntoLearner, which serves as a real-world use case.

The LLMs4OL paradigm focuses on three primary tasks: (i) term typing, which involves assigning semantic categories or types to domain-specific terms; (ii) taxonomy discovery, which aims to organize these terms into hierarchical structures based on “is-a” relationships; and (iii) non-taxonomic relationship extraction, which identifies meaningful semantic relations beyond hierarchies, such as functional or associative links between concepts. Together, these tasks capture the essential components of ontology learning—concept classification, hierarchical organization, and relational enrichment. They form the foundation of the LLMs4OL evaluation framework, enabling systematic benchmarking of both retrieval-based and generative approaches across diverse domains and datasets.

### 1.1 Term Typing

**Definition.** Term typing is defined as the task of assigning one or more semantic types (classes) to a given lexical term. Formally, let  $\mathcal{L}$  denote a lexical term and  $\mathcal{T} = \{t_1, t_2, \dots, t_N\}$  the set of candidate ontology types. The goal is to predict a subset  $\mathcal{T}^*(\mathcal{L}) \subseteq \mathcal{T}$  such that each selected type correctly characterizes the meaning of  $\mathcal{L}$ .

**Example.**

- *Input term:* “aspirin”
- *Candidate types:* {Drug, ChemicalCompound, Disease, Procedure}
- *Output:* {Drug, ChemicalCompound}

This is a multi-class classification problem where a term may belong to multiple semantic categories. The task is essential for populating ontologies and enabling downstream reasoning.

### 1.2 Taxonomy Discovery

**Definition.** Taxonomy discovery aims to identify hierarchical (“is-a”) relationships between concepts. Given two terms or classes  $(c_i, c_j)$ , the task is to determine whether a subsumption relation holds, i.e., whether  $c_i \sqsubseteq c_j$  (“ $c_i$  is a subclass of  $c_j$ ”).

**Example 1.**

- *Input pair:* (“red wine”, “wine”)
- *Output:* True (“red wine” is a type of “wine”)

**Example 2.**

- *Input pair:* (“wine”, “red wine”)
- *Output:* False

Table 1: Complexity Score Analysis.

(a) Robustness analysis summary by test type (Spearman  $\rho$  vs original ranking)

Test type	N variants	Min $\rho$	Mean $\rho$	Interpretation
Weight perturbation	20	0.9990	0.9996	Extremely stable to moderate reweighting
Sigmoid parameters	24	1.0000	1.0000	Ranking invariant to tested $(a, b)$ choices
Normalization	2	0.8557	0.8564	Largest sensitivity source
Equal weighting	1	0.9968	0.9968	Family-level priors matter little for ranking

(b) Spearman  $\rho$  under family-weight perturbations

Family	-20%	-10%	+10%	+20%
Graph	0.9992	0.9997	0.9998	0.9995
Coverage	0.9996	0.9998	0.9998	0.9994
Hierarchy	0.9999	0.9999	1.0000	0.9999
Breadth	0.9990	0.9996	0.9997	0.9991
Dataset	0.9997	0.9998	0.9998	0.9996

51 This task is typically formulated as a binary classification or ranking problem and is central to  
 52 constructing class hierarchies. It requires understanding semantic generalization and specialization  
 53 relationships.

### 54 1.3 Non-Taxonomic Relationship Extraction

55 **Definition.** This task focuses on identifying semantic relationships between entities that are not  
 56 hierarchical. Given two entities  $(e_i, e_j)$  and a set of possible relations  $\mathcal{R}$ , the objective is to predict  
 57 the most appropriate relation  $r \in \mathcal{R}$  or determine whether a relation exists.

#### 58 Example 1.

- 59 • *Input pair:* (“aspirin”, “headache”)
- 60 • *Candidate relations:* {treats, causes, located\_in, part\_of}
- 61 • *Output:* treats

#### 62 Example 2.

- 63 • *Input pair:* (“Paris”, “France”)
- 64 • *Output:* located\_in

65 Unlike taxonomy discovery, this task captures diverse semantic relations such as causality, function-  
 66 ality, spatial containment, and usage. It is critical for enriching ontologies with expressive relational  
 67 knowledge.

## 68 2 Complexity Score Analysis

69 The complexity score is computed by 19 metrics grouped into 5 families (graph, coverage, hierar-  
 70 chy, breadth, and dataset). We evaluate the robustness of an ontology complexity score across 180  
 71 ontologies using four stress tests, including:

- 72 • **Weight perturbation:**  $\pm 10\%$  and  $\pm 20\%$  for each family (20 variants).
- 73 • **Sigmoid sensitivity:**  $(a \in [0.3, 0.5])$ ,  $(b \in [5.5, 6.5])$  (24 variants).
- 74 • **Normalization alternatives:** min-max and z-score normalizations.
- 75 • **Equal-family-weight baseline:** all five families weights set to 0.20.

76 As a main robustness result, the Table 1a shows that the ranking is nearly unchanged under weight  
 77 and sigmoid changes, but shifts materially under dataset-dependent normalization (min-max/z-  
 78 score). This indicates robustness of the aggregation design, with sensitivity concentrated in the

79 normalization stage. Moreover, the weight Sensitivity analysis (see Table 1b), showed that even  
80  $\pm 20\%$  perturbations preserve rank order almost perfectly. Breadth is the most sensitive family (still  
81 highly stable), hierarchy the least.

82 Furthermore, the sigmoid and normalization effect study showed that for different sigmoid paramete-  
83 rs, all tested settings produce ( $\rho \approx 1.0$ ), confirming that sigmoid is primarily a monotonic rescaling  
84 in this regime. Nevertheless, the normalization method (log vs min-max:  $\rho = 0.8557$  and log vs  
85 z-score:  $\rho = 0.8570$ ), showed that normalization controls relative spacing across ontologies and can  
86 reorder mid-ranked items substantially. This is expected because min-max and z-score are dataset-  
87 dependent, whereas  $\log(1 + x)$  is pointwise and less corpus-sensitive.

88 The equal-family-weight baseline showed a Spearman correlation of  $\rho = 0.9968$  with original  
89 weighting, with top-10 overlaps with original ranking. This indicates that the complexity score’s  
90 discriminative structure is driven more by metric values than by fine-tuned family priors.

91 The rank stability analysis further supports the robustness of the proposed complexity score. Em-  
92 pirical evidence from the top-50 ontology rankings shows that under moderate weight perturbations  
93 (e.g., Graph metrics +20% or Coverage metrics +20%), the highest-ranked ontologies remain largely  
94 unchanged, with only minor local swaps (e.g., between SUMO [28] and ChEBI [27] ontologies).  
95 Similarly, under a uniform weighting configuration, the top-10 rankings are fully preserved, indicat-  
96 ing that the overall ordering is not sensitive to the specific weighting scheme. In contrast, alternative  
97 normalization strategies introduce more noticeable shifts: for instance, under min-max normaliza-  
98 tion, MatOnto [10] moves from rank 32 to 7, CSO [32] from 39 to 6, DBpedia [5] from 33 to 8,  
99 and PATO [20] drops from 7 to 18. These results suggest that while most perturbations preserve  
100 the relative ranking structure, the choice of normalization has a stronger impact on cross-ontology  
101 comparability.

102 In conclusion, the complexity score is highly robust to both family-weight perturbations (min  $\rho =$   
103  $0.9990$ ) and sigmoid parameter changes ( $\rho \approx 1.0$  across all settings). Equal-family weighting  
104 yields near-identical rankings ( $\rho = 0.9968$ ; top-10 overlap 100%), indicating limited dependence  
105 on manually chosen family priors. In contrast, the normalization strategy is the dominant sensitivity  
106 factor (log vs min-max  $\rho = 0.8557$ ; log vs z-score  $\rho = 0.8570$ ), suggesting that normalization  
107 should be fixed a priori and justified explicitly. Overall, the proposed complexity scoring framework  
108 is stable under parametric perturbations, while cross-dataset comparability depends primarily on  
109 normalization choice.

### 110 3 LLMs Contamination Analysis

111 The strong performance observed in high-resource domains such as *Geography*, *Finance*, and *Units*  
112 & *Measurements* suggests a higher likelihood of dataset contamination, as these domains are more  
113 likely to have been encountered by models during pre-training.

114 **Contamination Analysis Dataset.** We prepared the contamination analysis dataset by extracting  
115 structured knowledge triples from the ontologies across domains of interest. For each ontology,  
116 we loaded its extracted representation. We systematically parsed three types of semantic relations:  
117 taxonomic relations (*subClassOf* hierarchies), non-taxonomic relations (explicit relational links be-  
118 tween entities), and term typings (associations between entities and their types). Each extracted  
119 pair was standardized into a uniform schema containing the domain, head entity, relation type, and  
120 tail entity. After aggregating results from all ontologies, we merged them into a single dataset, re-  
121 moved duplicates to ensure uniqueness, and used the final collection for downstream contamination  
122 analysis.

123 The extracted dataset contains a total of **1,191** unique ontology-derived relation pairs distributed  
124 unevenly across four domains. The majority of the data comes from *Units & Measurements* (901  
125 pairs), followed by *Finance* (283 pairs), while *Geography* (7 pairs) is significantly underrepresented.  
126 This imbalance indicates a strong skew toward high-resource structured domains, which may influ-  
127 ence downstream analysis and bias contamination estimates toward domains that are more densely  
128 populated in existing knowledge sources.

129 We conducted the contamination analysis by evaluating how well large language models assign  
130 probability to ontology-derived relational statements constructed from structured triples. Each triple  
131  $(h, r, t)$  was converted into a natural language sentence  $x = f(h, r, t)$  (e.g., “ $h$  is a type of  $t$ ” or “ $h$

Table 2: Qwen family of LLMs perplexity statistics by ontology.

LLM	Ontology	Size	Mean	Median	STD	Min	Max	<50%	<100%	<200%	>400%
Qwen3-0.6B	OM	801	251.6	157.3	267.9	20.0	2635.4	7.9	29.5	58.6	18.1
Qwen3-0.6B	GeoNames	7	755.4	285.8	1262.5	73.8	3570.0	0.0	28.6	42.9	28.6
Qwen3-0.6B	QUDT	100	1133.6	310.5	2490.0	30.0	14749.4	3.0	16.0	31.0	39.0
Qwen3-0.6B	GoodRelations	283	14073.4	9730.6	13962.1	24.6	92253.2	0.7	2.8	6.4	93.6
Qwen3-4B-IT-2507	OM	801	133.1	81.2	163.1	9.5	1482.5	30.6	57.2	82.4	4.9
Qwen3-4B-IT-2507	QUDT	100	231.0	174.0	207.5	13.9	1060.5	15.0	29.0	58.0	15.0
Qwen3-4B-IT-2507	GeoNames	7	270.9	69.9	508.8	44.4	1419.7	28.6	71.4	85.7	14.3
Qwen3-4B-IT-2507	GoodRelations	283	7373.5	4552.5	7903.9	13.9	52838.7	2.8	5.3	6.7	93.3
Qwen3-8B	GeoNames	7	118.4	87.3	94.7	36.5	307.5	14.3	57.1	85.7	0.0
Qwen3-8B	OM	801	173.5	99.8	282.3	7.8	4247.4	23.0	50.2	77.8	6.7
Qwen3-8B	QUDT	100	345.7	183.5	766.3	19.8	5546.0	12.0	28.0	55.0	19.0
Qwen3-8B	GoodRelations	283	11974.3	7582.5	13205.6	15.9	94490.1	3.2	4.9	6.7	93.3
Qwen3-14B	GeoNames	7	127.6	72.2	95.0	49.8	277.0	14.3	57.1	71.4	0.0
Qwen3-14B	OM	801	166.2	113.4	192.2	14.5	1823.5	16.5	42.3	75.0	6.1
Qwen3-14B	QUDT	100	299.5	227.2	231.2	29.5	988.5	7.0	21.0	47.0	28.0
Qwen3-14B	GoodRelations	283	11244.6	7377.9	12379.9	27.0	88950.3	1.8	4.6	6.4	93.3
Qwen3-Next-80B-A3B-IT	GeoNames	7	124.0	74.8	141.7	33.9	431.4	28.6	71.4	85.7	14.3
Qwen3-Next-80B-A3B-IT	OM	801	128.7	77.8	150.2	10.4	1175.8	32.6	60.9	81.4	4.9
Qwen3-Next-80B-A3B-IT	QUDT	100	177.5	127.0	162.2	16.1	972.2	13.0	38.0	74.0	9.0
Qwen3-Next-80B-A3B-IT	GoodRelations	283	6148.1	2954.5	9586.1	16.5	64235.6	3.5	5.7	6.7	92.9

132 is an instance of  $t''$ ), depending on the relation type. These sentences were then used to compute  
 133 token-level likelihoods under each model using the causal language modeling objective.

134 **Perplexity-Based Contamination Analysis.** For analysis, we used a perplexity-based ap-  
 135 proach [14], where for a given statement  $x = (x_1, x_2, \dots, x_T)$ , the model assigns probability:  
 136  $P(x) = \prod_{t=1}^T P(x_t | x_{<t})$ , where we compute the average negative log-likelihood per token as  
 137 a  $\mathcal{L}(x) = -\frac{1}{T} \sum_{t=1}^T \log P(x_t | x_{<t})$ , and the perplexity is defined as a  $\text{PPL}(x) = \exp(\mathcal{L}(x))$ .  
 138 Lower perplexity indicates that the model assigns a higher likelihood to the statement, which we  
 139 interpret as a potential signal of memorization or dataset contamination from pretraining data. This  
 140 procedure was applied in batches across all generated statements for multiple models, and results  
 141 were aggregated across domains and ontologies to compare likelihood behavior across structured  
 142 knowledge sources.

143 **Results.** We conduct the contamination analysis using the Qwen family of models due to their strong  
 144 and consistent performance across all evaluated OL tasks, as demonstrated in our benchmark re-  
 145 sults. The contamination analysis, as shown in Table 2, reveals clear domain-dependent differences  
 146 in model likelihood behavior. Ontologies from *Units & Measurements* (OM [26], QUDT [16]) consis-  
 147 tently exhibit substantially lower perplexity across all models, with mean values decreasing from  
 148 251.6 to 128.7 (OM) and from 1133.6 to 177.5 (QUDT) as model size increases. These domains  
 149 also show a higher proportion of low-perplexity samples (e.g., up to 60.9% below 100 for OM),  
 150 suggesting that their underlying relations are more likely to be memorized or frequently encoun-  
 151 tered during pretraining. In contrast, the *Finance* domain ontology (GoodRelations [23]) exhibits  
 152 extremely high perplexity across all models (mean  $> 6000$  even for the largest model), with over  
 153 90% of samples exceeding 400 perplexity, indicating that these structured relations are not well cap-  
 154 tured by pretrained distributions and are unlikely to reflect memorized knowledge. This highlights  
 155 that the domain “resource richness” alone does not determine contamination; rather, the extent to  
 156 which knowledge is expressed in natural language corpora plays a critical role. Additionally, model  
 157 scaling consistently reduces perplexity across all domains, but preserves the relative ordering, sug-  
 158 gesting that contamination signals are stable and not solely driven by model capacity. Finally, results  
 159 for *Geography* domain remain inconclusive due to the very small sample size ( $n=7$ ), which leads to  
 160 high variance and unstable estimates.

## 161 4 Experimental Resources and Setup

### 162 4.1 Selected Datasets Statistics

163 Table 3 summarizes the statistics of the selected ontologies used for experimentation with On-  
 164 toLearner. The selected ontologies span diverse domains and vary significantly in size and structural  
 165 complexity, enabling evaluation across heterogeneous OL scenarios. We used the full ontologies for  
 166 evaluation, except for ChEBI [27], evaluation was conducted on a 10% subset of the data due to the  
 167 significantly larger combinatorial search space, while maintaining the same evaluation protocol.

Table 3: Selected ontological datasets ground truth statistics. The *Terms No.* is the number of terms extracted from the given ontology, whereas the *Types No.* is the number of classes, and *Relations No.* is the number of non-is-a relations. For ontologies that do not support the LLMs4OL tasks or do not contain terms or non-taxonomic relations, frequency values are empty.

	Domain	Ontology	Terms No.	Types No.	Relations No.	Term Typing	Taxonomy Discovery	Non-Taxonomic RE
1	Education	DoCO [12]	-	30	2	-	59	7
2	Events	Conference [29]	32	33	1	32	49	3
3	Finance	GoodRelations [23]	46	25	2	46	25	264
4	Food and Beverage	Wine [21]	161	20	-	161	47	-
5	Agriculture	AgrO [4]	71	4,109	67	71	10,931	1,699
6		FoodOn [15]	16	36,799	6	16	76,228	2,072
7	Geography	GeoNames [1]	699	11	-	699	18	-
8	Materials Science and Engineering	MaterialInformation [3]	404	537	2	404	605	30
9		MatOnto [10]	121	841	2	122	1,215	167
10		PeriodicTable [13]	150	-	-	150	-	-
11		MDSOnto [30]	-	1,332	30	-	1,648	259
12	Biology and Life Sciences	GO [2, 11]	-	40,635	1	-	156,430	30
13	Ecology and Environment	ENVO [8, 9]	44	6,531	5	46	16,175	147
14		SWEET [31]	2,095	10,073	3	2,219	16,111	515
15	General Knowledge	CCO [25]	350	1,422	1	362	1,532	21
16		DBpedia [5]	-	767	4	-	799	1,665
17		SchemaOrg [34]	-	910	1	-	1,058	635
18	Industry	AUTO [17]	57	1,105	1	58	2,731	42
19		PTO [24]	3,000	1,001	-	3,000	3,996	-
20	Medicine	DOID [35]	-	11,731	2	-	41,569	25
21		OBI [7]	285	4,934	3	286	11,843	38
22	Units and Measurements	OM [26]	1,686	798	-	1,953	1,124	-
23		QUDT [16]	23	84	1	27	400	12
24	Chemistry	PROCO [33, 22]	14	831	-	14	1,757	-
25		VIBSO [36]	40	470	3	40	599	23
26		ChEBI [27]	-	201,959	-	-	73,996	-

168 In the architectural design of the Retriever and Retrieval-Augmented Generation (RAG) modules  
169 within OntoLearner, constructing an ontology strictly from scratch requires clearly defined domain  
170 boundaries. This is essential to prevent the generation of hallucinated classes or properties that do  
171 not belong to the target domain. Standard autoregressive large language models (LLMs) inherently  
172 reflect biases from their pre-training distributions, which can introduce irrelevant or incorrect con-  
173 cepts. To address this limitation, the system constrains the model’s operational space to grounded  
174 domain-specific elements—namely Terms, Types, and Relations—extracted and indexed for each  
175 ontology. By restricting the generation process to these validated components, OntoLearner func-  
176 tions as a controlled and reliable assistive tool for ontology engineers. This design principle also  
177 motivates the modular structure of the system, where each component is responsible for a single,  
178 well-defined task.

179 In practical settings, ontologies often contain a large number of terms and types, along with a  
180 smaller but significant set of non-taxonomic relations. The combinatorial space of possible associa-  
181 tions among these elements is substantial and can require extensive manual effort over long periods.  
182 Specifically, the term typing task ( $\text{Term} \rightarrow \text{Type}$ ), for a dataset of size  $|\text{Terms}|$ , the total complexity  
183 is  $O(|\text{Terms}| \times |\text{Types}|)$ , while per-instance inference is  $O(|\text{Types}|)$ . The taxonomy discovery task  
184 ( $\text{Child Type} \rightarrow \text{Parent Type}$ ) involves identifying hierarchical relationships among types, resulting  
185 in a search space of  $O(|\text{Types}|^2)$ . The non-taxonomic relation extraction task ( $\text{Head Type} \xrightarrow{\text{Relation}} \text{Tail Type}$ )  
186 is even more complex, as it requires identifying meaningful relational triples across types.  
187 This leads to a search space of  $O(|\text{Types}|^2 \times |\text{Relations}|)$ , representing a significantly larger com-  
188 binatorial space focused on semantic interactions rather than hierarchical structure. In light of these  
189 computational challenges, selecting 26 ontologies provides a balanced and practical benchmark for  
190 evaluating OntoLearner while keeping the problem tractable.

## 191 4.2 How to Run OntoLearner for LLMs4OL Tasks

192 The Retrieval-Augmented Generation (RAG) module consists of two main components: a re-  
193 triever and an LLM-based re-ranker. For using such a technique for *Term Typing* task, consistent  
194 with the modular design of OntoLearner, Figure 2 illustrates the retriever implementation using  
195 `AutoRetrieverLearner`. Moreover, the *Type Taxonomy Discovery* tasks reuse the term typing  
196 formulation, with  $t_q \in \mathcal{T}$  as the query instead of  $\mathcal{L}$ . Thus, in both the retriever (see Figure 2 and  
197 re-ranker pipeline (see Figure 3), the only change needed is `task = 'taxonomy-discovery'`.  
198 However, for the LLM-Augmented retriever, Figure 4 illustrates the workflow, and Figure 5 shows  
199 the end-to-end pipeline usage. The OpenAI model is used as an LLM augments, but open-source  
200 models can also be applied. LLM-augmented learners are encapsulated and can be adapted for other  
201 approaches using the same code. Finally, the *Non-Taxonomic Relation Extraction* task follows a  
202 similar pattern, where it reuses the taxonomy discovery and term typing formulations with minimal

```

1 from ontolearner import AutoRetrieverLearner, MatOnto, evaluation_report
2 retriever = AutoRetrieverLearner(top_k=15, batch_size=10240)
3 retriever.load(model_id='Qwen/Qwen3-Embedding-8B')
4 ontology = MatOnto()
5 ontology.load()
6 data = ontology.extract()
7 task = 'term-typing'
8 retriever.fit(data, task=task)
9 predicts = retriever.predict(data, task=task)
10 truth = retriever.tasks_ground_truth_former(data=data, task=task)
11 metrics = evaluation_report(y_true=truth, y_pred=predicts, task=task)
12 print(metrics)

```

Figure 2: Example code for running the term-typing task using a retriever-based learner. To apply the retriever to the two other OL tasks, only the task specification in line 7 needs to be changed (taxonomy-discovery for taxonomy discovery and non-taxonomic-re for non-taxonomic relationship extraction). The script loads a pretrained embedding model, extracts data from the materials ontology, performs retrieval-based inference, and reports evaluation metrics.

```

1 from ontolearner import AutoLLMLearner, AutoRetrieverLearner,
2 StandardizedPrompting, MatOnto, LabelMapper,
3 AutoRAGLearner, evaluation_report
4 llm = AutoLLMLearner(prompting=StandardizedPrompting,
5 label_mapper=LabelMapper(),
6 batch_size=64, device='cuda')
7 retriever = AutoRetrieverLearner(top_k=15, batch_size=10240)
8 rag = AutoRAGLearner(llm=llm, retriever=retriever)
9 rag.load(llm_id='Qwen/Qwen3-14B', retriever_id='Qwen/Qwen3-Embedding-8B')
10 ontology = MatOnto()
11 ontology.load()
12 data = ontology.extract()
13 task = 'term-typing'
14 rag.fit(data, task=task)
15 predicts = rag.predict(data, task=task)
16 truth = rag.tasks_ground_truth_former(data=data, task=task)
17 metrics = evaluation_report(y_true=truth, y_pred=predicts, task=task)
18 print(metrics)

```

Figure 3: Example code for running the term-typing task using a RAG learner. To apply the retriever to other tasks, only the task specification in line 13 needs to be changed (taxonomy-discovery for taxonomy discovery and non-taxonomic-re for non-taxonomic relationship extraction). The script initializes the LLM and retriever components, loads pretrained models, extracts data from the materials ontology, performs RAG-based training and inference, and reports evaluation metrics.

203 modifications. In OntoLearner, either in the retriever (Figure 2) or reranker (Figure 3), it is enabled  
204 via `task = 'non-taxonomic-re'`, integrating pair discovery, natural language query generation,  
205 relation retrieval, and LLM verification in a unified RAG pipeline.

### 206 4.3 Standardized Prompts for LLMs

207 We employ standardized zero-shot prompts with strict output constraints (yes/no) to ensure compa-  
208 rability across tasks. Prompts were designed to minimize ambiguity and avoid task leakage. The  
209 prompts are supported via OntoLearner using AutoPrompt module and presented in the Figure 6.  
210 We enforce binary outputs to ensure consistency across models and enable reliable and comparable  
211 evaluation across tasks.

### 212 4.4 Computational Details

213 All experiments were conducted in a zero-shot setting using pretrained models without task-specific  
214 fine-tuning. Unless otherwise specified, default inference parameters were used for all large lan-  
215 guage models (LLMs).

216 **Hardware and Execution.** Experiments were performed on a multi-GPU system equipped with  
217 NVIDIA H100 GPUs. Larger LLMs (e.g., 27B–80B parameters) were executed on configurations  
218 utilizing two H100 GPUs (95 GB memory each), while smaller models were experimented on a

```

1 from ontolearner import LLMAugmenterGenerator, LLMAugmenter,
2                       LLMAugmentedRetriever, LLMAugmentedRetrieverLearner,
3                       MatOnto, evaluation_report, utils
4 augmenter = LLMAugmenterGenerator(token='openai-api-key')
5 ontology = MatOnto()
6 ontology.load()
7 data = ontology.extract()
8 task = 'taxonomy-discovery'
9 path = f"{task}-{ontology.ontology_id}.json"
10 augments = {"config": augmenter.get_config(),
11            task: augmenter.augment(data, task=task)}
12 utils.save_json(path, augments)
13 augmenter = LLMAugmenter(path=path)
14 base_retriever = LLMAugmentedRetriever()
15 retriever = LLMAugmentedRetrieverLearner(base_retriever=base_retriever,
16                                         top_k=15, batch_size=10240)
17 retriever.load(model_id='Qwen/Qwen3-Embedding-8B')
18 retriever.set_augmenter(augmenter=augmenter)
19 retriever.fit(data, task=task)
20 predicts = retriever.predict(data, task=task)
21 truth = retriever.tasks_ground_truth_former(data=data, task=task)
22 metrics = evaluation_report(y_true=truth, y_pred=predicts, task=task)
23 print(metrics)

```

Figure 4: Example code illustrating the end-to-end workflow for taxonomy discovery task using an LLM-augmented retriever, including ontology loading, data augmentation, retrieval-based learning, and evaluation. When a task other than `taxonomy-discovery` is specified in line 8 (i.e., `term-typing` or `non-taxonomic-re`), the augmenter defaults to a retriever-based model, as illustrated in Figure 2.

```

1 from ontolearner import LLMAugmentedRetriever, LLMAugmentedRetrieverLearner,
2                       LLMAugmentedRAGLearner, MatOnto, LabelMapper,
3                       StandardizedPrompting, evaluation_report
4 llm = AutoLLMLearner(prompting=StandardizedPrompting,
5                      label_mapper=LabelMapper(),
6                      batch_size=64,
7                      device='cuda')
8 base_retriever = LLMAugmentedRetriever()
9 retriever = LLMAugmentedRetrieverLearner(base_retriever=base_retriever,
10                                         top_k=15, batch_size=10240)
11 rag = LLMAugmentedRAGLearner(llm=llm, retriever=retriever)
12 rag.load(llm_id='Qwen/Qwen3-14B', retriever_id='Qwen/Qwen3-Embedding-8B')
13 ontology = MatOnto()
14 ontology.load()
15 data = ontology.extract()
16 task = 'taxonomy-discovery'
17 augmenter = LLMAugmenter(path=f"{task}/{ontology.ontology_id}.json")
18 rag.fit(data, task=task)
19 predicts = rag.predict(data, task=task)
20 truth = rag.tasks_ground_truth_former(data=data, task=task)
21 metrics = evaluation_report(y_true=truth, y_pred=predicts, task=task)
22 print(metrics)

```

Figure 5: Example code for running the RAG learner for taxonomy-discovery using an LLM-augmented retriever. The augmented retriever is activated only when `taxonomy-discovery` is specified; for other tasks (i.e., `term-typing` or `non-taxonomic-re`), the retriever falls back to standard RAG behavior without augmentation, as illustrated in Figure 3.

219 single H100 GPU. Multiple runs were executed concurrently across available GPUs. This parallel  
220 runs enabled large-scale evaluation across all ontologies and tasks within a reasonable time frame,  
221 while maintaining consistent experimental conditions.

222 **Batching and Inference.** Batch sizes were adjusted depending on the model size and task require-  
223 ments, as reflected in the code examples provided in Figure 5. Retrieval and reranking components  
224 were executed using pretrained embedding and LLM models without additional optimization or  
225 pruning. All tasks—term typing, taxonomy discovery, and non-taxonomic relation extraction—were  
226 evaluated using grouped metrics (precision, recall, and F1-score) computed over the full set of on-  
227 tology datasets.

228 **Run Time.** We evaluated the robustness and computational scalability of OntoLearner across all  
229 datasets and 12 LLMs. All models achieved a 100% completion rate (**756 runs**), demonstrating

<b>Term Typing Prompt</b>
<p>You are performing term typing.</p> <p>Determine whether the given term is a clear and unambiguous instance of the specified high-level type.</p> <p>Rules:</p> <ul style="list-style-type: none"> <li>- Answer "yes" only if the term commonly and directly belongs to the type.</li> <li>- Answer "no" if the term does not belong to the type, is ambiguous, or only weakly related.</li> <li>- Use the most common meaning of the term.</li> <li>- Do not explain your answer.</li> </ul> <p>Term: {term}  Type: {type}  Answer (yes or no):</p>
<b>Taxonomy Discovery Prompt</b>
<p>You are identifying taxonomic (is-a) relationships.</p> <p>Determine whether the first concept is a superclass (direct or indirect) of the second concept in a standard conceptual hierarchy.</p> <p>Rules:</p> <ul style="list-style-type: none"> <li>- A superclass means the second concept is a type or instance of the first.</li> <li>- Answer "yes" only if the relationship is a true is-a relationship.</li> <li>- Answer "no" for part-of, related-to, or associative relationships.</li> <li>- Use widely accepted general knowledge.</li> <li>- Do not explain your answer.</li> </ul> <p>Parent: {parent}  Child: {child}  Answer (yes or no):</p>
<b>Non-Taxonomic Relation Extraction Prompt</b>
<p>You are identifying non-taxonomic conceptual relationships.</p> <p>Determine whether the specified relation typically holds between the given conceptual types.</p> <p>Rules:</p> <ul style="list-style-type: none"> <li>- Answer "yes" only if the relation commonly and meaningfully applies.</li> <li>- Answer "no" if the relation is rare, indirect, or context-dependent.</li> <li>- Do not infer relations that require specific situations.</li> <li>- Use widely accepted general knowledge.</li> <li>- Do not explain your answer.</li> </ul> <p>Head type: {head}  Tail type: {tail}  Relation: {relation}  Answer (yes or no):</p>

Figure 6: Standardized prompt templates for LLM-based reranker of term typing, taxonomy discovery, and non-taxonomic relation extraction tasks.

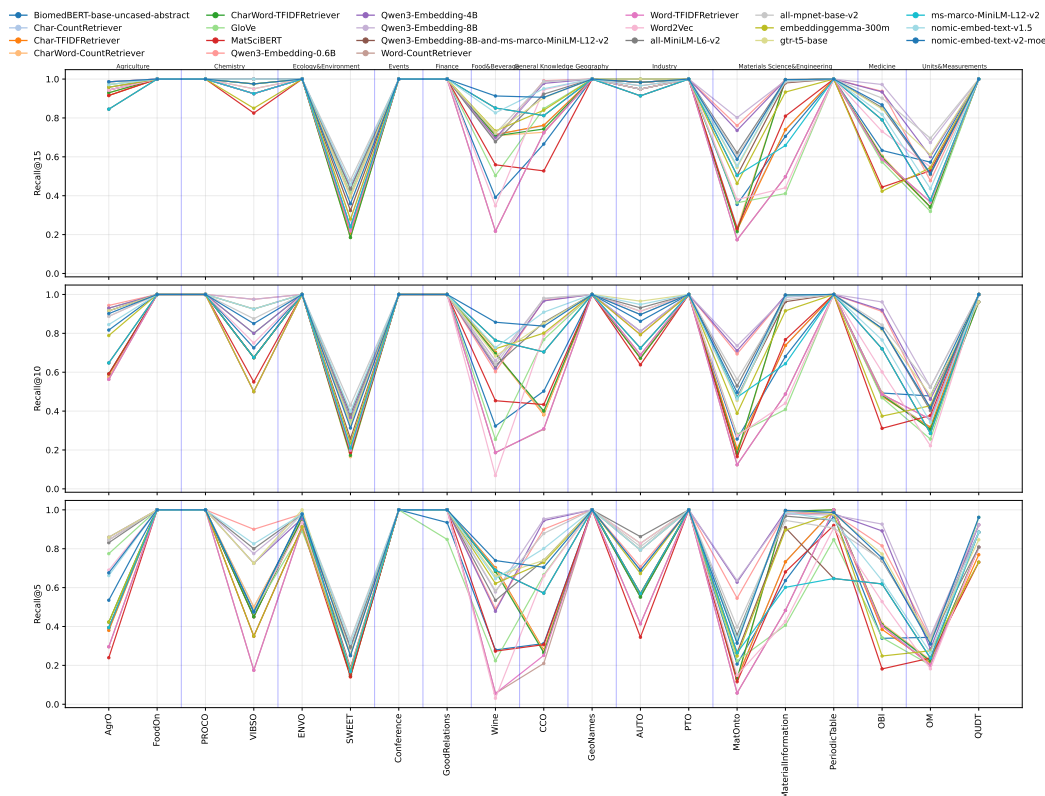


Figure 7: Performance comparison of retrieval-based methods on the **term typing** task.

230 stable and failure-free execution across tasks and domains. Each LLM was evaluated on 63 tasks,  
 231 corresponding to combinations of OL tasks and datasets. The average response time per run was ap-  
 232 proximately 21,566 seconds (~6 hours), resulting in a total runtime of approximately **378** hours per  
 233 LLM and an aggregate runtime of **4,536** hours across all models. These experiments were executed  
 234 in parallel across multiple GPUs, enabling large-scale evaluation within practical time constraints  
 235 while maintaining consistent experimental conditions. Substantially higher runtimes were observed  
 236 for large-scale ontologies such as ChEBI, where the average runtime reached approximately 387,398  
 237 seconds per run. This highlights the dominant role of ontology size, particularly the number of types,  
 238 in determining computational cost. In particular, tasks such as taxonomy discovery exhibit quadratic  
 239 complexity with respect to the number of types, leading to significantly larger search spaces in high-  
 240 dimensional ontologies. Overall, these results demonstrate that OntoLearner scales reliably across  
 241 diverse models and domains, while also revealing that computational efficiency is primarily con-  
 242 strained by the combinatorial structure of OL tasks. This suggests that future optimizations should  
 243 focus on reducing the effective search space, for example, through candidate pruning or adaptive  
 244 retrieval strategies.

245 All experiments can be reproduced on smaller subsets using the provided modular pipeline.

## 246 5 Detailed Results

### 247 5.1 Retrievers

248 The Figures 7–9 present the performance of retrieval-based approaches across the three OL tasks:  
 249 term typing, taxonomy discovery, and non-taxonomic relationship extraction. We report grouped  
 250 recall to assess how effectively each retriever supports downstream OL objectives. Overall, the  
 251 results reveal task-dependent performance variations, highlighting the importance of retrieval quality  
 252 in enabling accurate concept classification, hierarchical structuring, and relational inference.

253 Figure 7 shows that dense embedding-based retrievers consistently outperform lexical baselines  
254 across most domains, often by a substantial margin in grouped recall. This indicates that seman-  
255 tic similarity captured by dense representations is critical for accurate term typing, particularly in  
256 domains with high lexical variability or synonymy. In contrast, lexical methods struggle when sur-  
257 face forms diverge from canonical ontology labels.

258 Moreover, Figure 8 further demonstrates that LLM-Augmented retrieval yields consistent improve-  
259 ments over standard retrieval approaches in taxonomy discovery. The gains are especially pro-  
260 nounced in larger and more complex ontologies, suggesting that augmentation helps bridge the gap  
261 between surface-level similarity and deeper hierarchical reasoning. This highlights the importance  
262 of incorporating contextual or generative signals when modeling subsumption relationships.

263 Finally, Figure 9 reveals substantially higher variability in performance across domains for non-  
264 taxonomic relation extraction. This reflects the inherently diverse and context-dependent nature of  
265 such relations, which often require richer semantic understanding beyond similarity or hierarchy. As  
266 a result, retrieval performance in this setting is more sensitive to domain characteristics and relation  
267 distributions.

## 268 5.2 Reranking for the RAG Pipeline

269 The results in Figures 10–12 highlight clear trade-offs between precision and recall in different  
270 LLMs as a reranker in the RAG pipeline across OL tasks. While some LLMs achieve high precision  
271 by aggressively filtering candidates, this often comes at the cost of reduced recall, particularly in  
272 tasks with large candidate spaces such as taxonomy discovery. In contrast, recall shows that LLMs  
273 provide broader coverage but may introduce more false positives, impacting the overall F1-score.  
274 The balance between these behaviors varies across tasks: term typing benefits from higher preci-  
275 sion due to well-defined class boundaries, whereas taxonomy discovery and non-taxonomic relation  
276 extraction benefit from higher recall to capture diverse valid relationships. Overall, these results  
277 emphasize the importance of selecting LLMs that align with task-specific requirements, as no single  
278 approach uniformly performs well across all OL tasks.

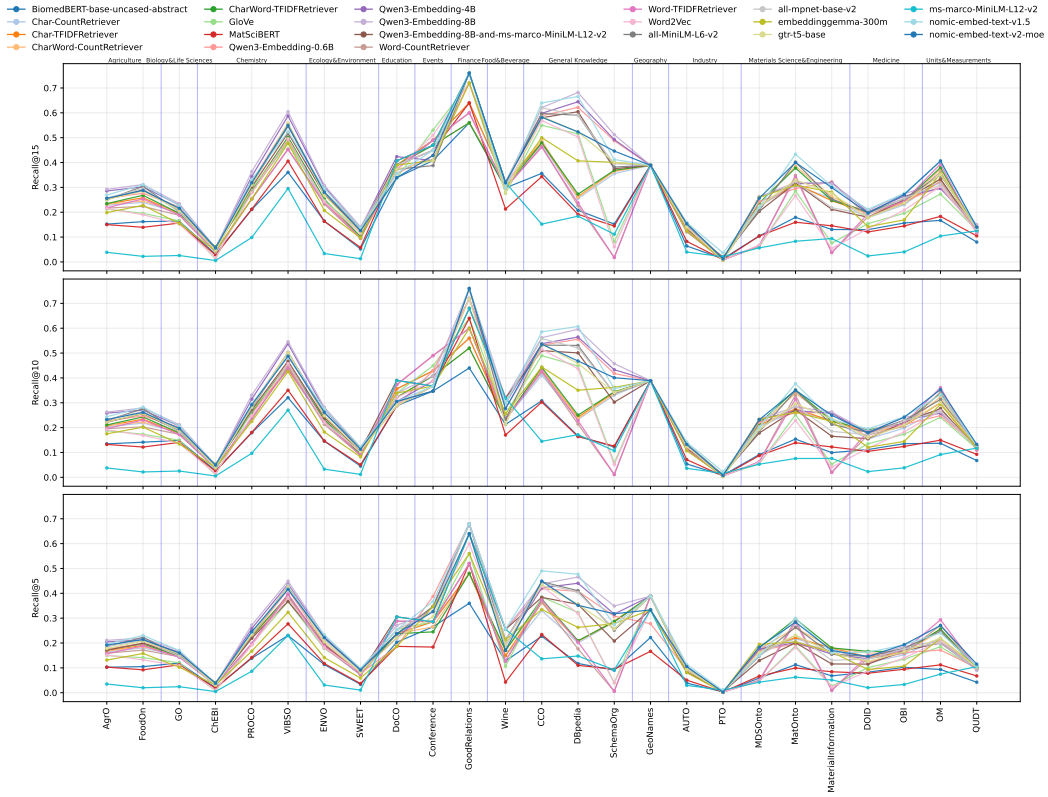
## 279 References

- 280 [1] Geonames ontology (geonames). <https://www.geonames.org/ontology>, 2022. RDF on-  
281 tology. Licensed under Creative Commons 3.0. Last updated: 2022-01-30.
- 282 [2] Michael Ashburner, Catherine A Ball, Judith A Blake, David Botstein, Heather Butler,  
283 J Michael Cherry, Allan P Davis, Kara Dolinski, Selina S Dwight, Janan T Eppig, et al. Gene  
284 ontology: tool for the unification of biology. *Nature genetics*, 25(1):25–29, 2000.
- 285 [3] Toshihiro Ashino. Materials ontology: An infrastructure for exchanging materials information  
286 and knowledge. *Data Science Journal*, 9:54–61, 2010.
- 287 [4] Céline Aubert, Marie-Angélique Laporte, Krishna Udaiwal, Pier Luigi Buttigieg, Chris  
288 Mungall, Elizabeth Arnaud, and Charles Tapley Hoyt. Agriculturalsemantics/agro: Novem-  
289 ber 2022 release, November 2022.
- 290 [5] Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary  
291 Ives. Dbpedia: A nucleus for a web of open data. In *international semantic web conference*,  
292 pages 722–735. Springer, 2007.
- 293 [6] Hamed Babaei Giglou, Jennifer D’Souza, and Sören Auer. Llms4ol: Large language models  
294 for ontology learning. In *International Semantic Web Conference*, pages 408–427. Springer,  
295 2023.
- 296 [7] Anita Bandrowski, Ryan Brinkman, Mathias Brochhausen, Matthew H Brush, Bill Bug, Mar-  
297 cus C Chibucos, Kevin Clancy, Mélanie Courtot, Dirk Derom, Michel Dumontier, et al. The  
298 ontology for biomedical investigations. *PloS one*, 11(4):e0154556, 2016.
- 299 [8] Pier Luigi Buttigieg, Norman Morrison, Barry Smith, Christopher J Mungall, Suzanna E  
300 Lewis, and Envö Consortium. The environment ontology: contextualising biological and  
301 biomedical entities. *Journal of biomedical semantics*, 4(1):43, 2013.

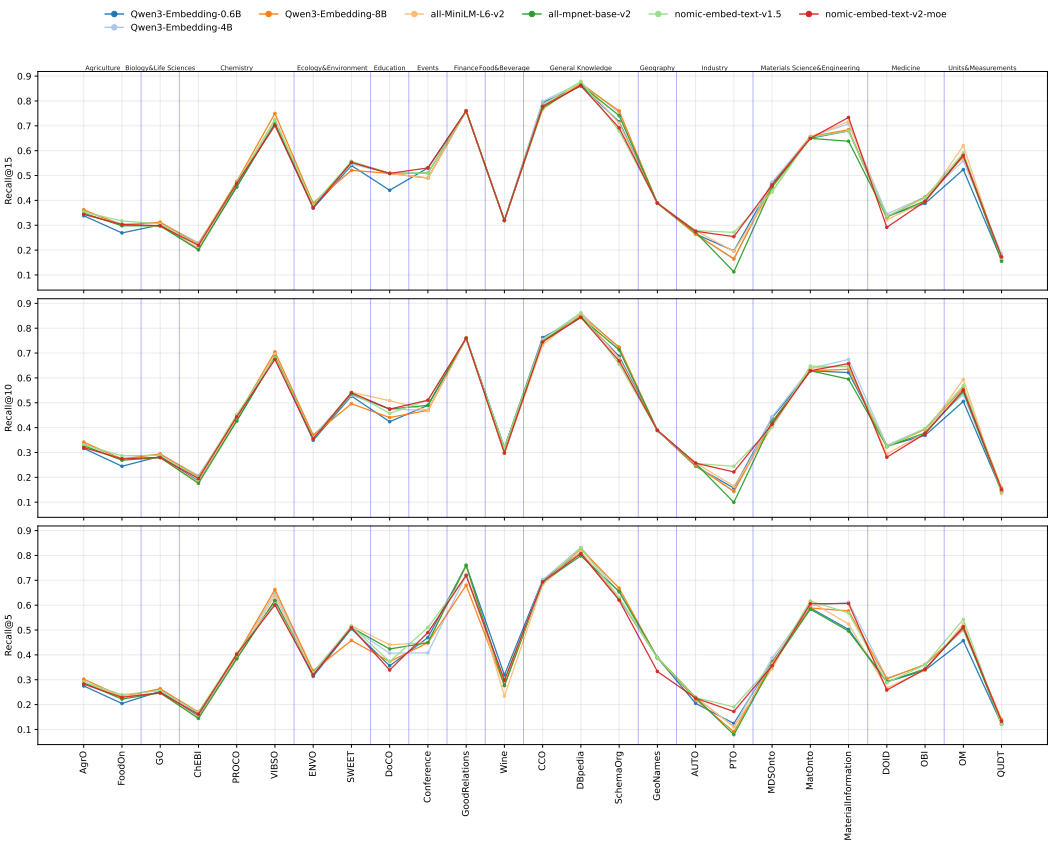
- 302 [9] Pier Luigi Buttigieg, Evangelos Pafilis, Suzanna E Lewis, Mark P Schildhauer, Ramona L  
303 Walls, and Christopher J Mungall. The environment ontology in 2016: bridging domains  
304 with increased scope, semantic density, and interoperation. *Journal of biomedical semantics*,  
305 7(1):57, 2016.
- 306 [10] Kwok Cheung, John Drennan, and Jane Hunter. Towards an ontology for data-driven discovery  
307 of new materials. In *AAAI Spring Symposium: Semantic Scientific Knowledge Integration*,  
308 pages 9–14, 2008.
- 309 [11] The Gene Ontology Consortium, Suzi A Aleksander, James Balhoff, Seth Carbon, J Michael  
310 Cherry, Harold J Drabkin, Dustin Ebert, Marc Feuermann, Pascale Gaudet, Nomi L Harris,  
311 David P Hill, Raymond Lee, Huaiyu Mi, Sierra Moxon, Christopher J Mungall, Anushya  
312 Muruganugan, Tremayne Mushayahama, Paul W Sternberg, Paul D Thomas, Kimberly  
313 Van Auken, Jolene Ramsey, Deborah A Siegele, Rex L Chisholm, Petra Fey, Maria Cristina  
314 Aspromonte, Maria Victoria Nugnes, Federica Quaglia, Silvio Tosatto, Michelle Giglio, Su-  
315 varna Nadendla, Giulia Antonazzo, Helen Attrill, Gil dos Santos, Steven Marygold, Victor  
316 Strelets, Christopher J Tabone, Jim Thurmond, Pinglei Zhou, Saadullah H Ahmed, Praoparn  
317 Asanithong, Diana Luna Buitrago, Meltem N Erdol, Matthew C Gage, Mohamed Ali Kad-  
318 hum, Kan Yan Chloe Li, Miao Long, Aleksandra Michalak, Angeline Pesala, Armalya Pri-  
319 tazahra, Shirin C C Saverimuttu, Renzhi Su, Kate E Thurlow, Ruth C Lovering, Colin Logie,  
320 Snezhana Oliferenko, Judith Blake, Karen Christie, Lori Corbani, Mary E Dolan, Harold J  
321 Drabkin, David P Hill, Li Ni, Dmitry Sitnikov, Cynthia Smith, Alayne Cuzick, James Seager,  
322 Laurel Cooper, Justin Elser, Pankaj Jaiswal, Parul Gupta, Pankaj Jaiswal, Sushma Naithani,  
323 Manuel Lera-Ramirez, Kim Rutherford, Valerie Wood, Jeffrey L De Pons, Melinda R Dwinell,  
324 G Thomas Hayman, Mary L Kaldunski, Anne E Kwitek, Stanley J F Laulerderkind, Marek A  
325 Tutaj, Mahima VEDI, Shur-Jen Wang, Peter D’Eustachio, Lucila Aimo, Kristian Axelsen, Alan  
326 Bridge, Nevila Hyka-Nouspikel, Anne Morgat, Suzi A Aleksander, J Michael Cherry, Stacia R  
327 Engel, Kalpana Karra, Stuart R Miyasato, Robert S Nash, Marek S Skrzypek, Shuai Weng,  
328 Edith D Wong, Erika Bakker, Tanya Z Berardini, Leonore Reiser, Andrea Auchincloss, Kristi-  
329 an Axelsen, Ghislaine Argoud-Puy, Marie-Claude Blatter, Emmanuel Boutet, Lionel Breuza,  
330 Alan Bridge, Cristina Casals-Casas, Elisabeth Coudert, Anne Estreicher, Maria Livia Famigli-  
331 etti, Marc Feuermann, Arnaud Gos, Nadine Gruaz-Gumowski, Chantal Hulo, Nevila Hyka-  
332 Nouspikel, Florence Jungo, Philippe Le Mercier, Damien Lieberherr, Patrick Masson, Anne  
333 Morgat, Ivo Pedruzzi, Lucille Pourcel, Sylvain Poux, Catherine Rivoire, Shyamala Sundaram,  
334 Alex Bateman, Emily Bowler-Barnett, Hema Bye-A-Jee, Paul Denny, Alexandr Ignatchenko,  
335 Rizwan Ishtiaq, Antonia Lock, Yvonne Lussi, Michele Magrane, Maria J Martin, Sandra  
336 Orchard, Pedro Raposo, Elena Speretta, Nidhi Tyagi, Kate Warner, Rossana Zaru, Alexan-  
337 der D Diehl, Raymond Lee, Juancarlos Chan, Stavros Diamantakis, Daniela Raciti, Magdalena  
338 Zarowiecki, Malcolm Fisher, Christina James-Zorn, Virgilio Ponferrada, Aaron Zorn, Sridhar  
339 Ramachandran, Leyla Ruzicka, and Monte Westerfield. The gene ontology knowledgebase in  
340 2023. *Genetics*, 224(1):iyad031, 03 2023.
- 341 [12] Alexandru Constantin, Silvio Peroni, Steve Pettifer, David Shotton, and Fabio Vitali. The  
342 document components ontology (doco). *Semantic web*, 7(2):167–181, 2016.
- 343 [13] Michael Cook. Periodic table of the elements ontology (periodictable). <https://www.daml.org/2003/01/periodictable/>, February 2004. OWL ontology. Last updated: 2004-02-05.
- 345 [14] Yihong Dong, Xue Jiang, Huanyu Liu, Zhi Jin, Bin Gu, Mengfei Yang, and Ge Li. General-  
346 ization or memorization: Data contamination and trustworthy evaluation for large language  
347 models. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages  
348 12039–12050, 2024.
- 349 [15] Damion M Dooley, Emma J Griffiths, Gurinder S Gosal, Pier L Buttigieg, Robert Hoehndorf,  
350 Matthew C Lange, Lynn M Schriml, Fiona SL Brinkman, and William WL Hsiao. Foodon:  
351 a harmonized food ontology to increase global food traceability, quality control and data inte-  
352 gration. *npj Science of Food*, 2(1):23, 2018.
- 353 [16] FAIRsharing.org. QUDT: Quantities, Units, Dimensions and Types. Online, 2025. Accessed:  
354 2026-01-06; Last edited: 2025-11-03.

- 355 [17] Michael Feld and Christian Müller. The automotive ontology: managing knowledge inside the  
356 vehicle and sharing it between cars. In *Proceedings of the 3rd International Conference on*  
357 *Automotive User Interfaces and Interactive Vehicular Applications*, pages 79–86, 2011.
- 358 [18] Hamed Babaei Giglou, Jennifer D’Souza, and Sören Auer. Llms4ol 2024 overview: The 1st  
359 large language models for ontology learning challenge. *arXiv preprint arXiv:2409.10146*,  
360 2024.
- 361 [19] Hamed Babaei Giglou, Jennifer D’Souza, Nandana Mihindukulasooriya, and Sören Auer.  
362 Llms4ol 2025 overview: The 2nd large language models for ontology learning challenge. In  
363 *Open Conference Proceedings*, volume 6, 2025.
- 364 [20] Georgios V Gkoutos, Paul N Schofield, and Robert Hoehndorf. The anatomy of phenotype  
365 ontologies: principles, properties and applications. *Briefings in bioinformatics*, 19(5):1008–  
366 1021, 2018.
- 367 [21] João Graça, Márcio Mourao, Orlando Anunciação, Pedro Monteiro, H Sofia Pinto, and Virgílio  
368 Loureiro. Ontology building process: the wine domain. In *Proc. of the 5th Conf. of EFITA*,  
369 2005.
- 370 [22] Oliver He and Wes Schafer. Extending the allotrope framework: An ontological representation  
371 and analysis of process chemistry. 2020 Fall Allotrope Connect Virtual Conference (Oral  
372 presentation), September 2020. Virtual conference, September 30, 2020.
- 373 [23] Martin Hepp. Goodrelations: An ontology for describing products and services offers on  
374 the web. In *International conference on knowledge engineering and knowledge management*,  
375 pages 329–346. Springer, 2008.
- 376 [24] Martin Hepp, Katharina Siorpaes, and Daniel Bachlechner. Harvesting wiki consensus: Us-  
377 ing wikipedia entries as vocabulary for knowledge management. *IEEE Internet Computing*,  
378 11(5):54–65, 2007.
- 379 [25] Mark Jensen, Giacomo De Colle, Sean Kindya, Cameron More, Alexander P Cox, and John  
380 Beverley. The common core ontologies. *arXiv preprint arXiv:2404.17758*, 2024.
- 381 [26] Jan Martin Keil and Sirko Schindler. Comparison and evaluation of ontologies for units of  
382 measurement. *Semantic Web*, 10(1):33–51, 2018.
- 383 [27] Adnan Malik, Muhammad Arsalan, Carlos Moreno, Juan Mosquera, Eloy Félix, Tefvik  
384 Kizilören, Venkatesh Muthukrishnan, Barbara Zdrazil, Andrew R Leach, and Noel M O’Boyle.  
385 Chebi: re-engineered for a sustainable future. *Nucleic Acids Research*, page gkaf1271, 2025.
- 386 [28] Ian Niles and Adam Pease. Towards a standard upper ontology. In *Proceedings of the in-*  
387 *ternational conference on Formal Ontology in Information Systems-Volume 2001*, pages 2–9,  
388 2001.
- 389 [29] Andrea Giovanni Nuzzolese, Anna Lisa Gentile, Valentina Presutti, and Aldo Gangemi. Se-  
390 mantic web conference ontology - a refactoring solution. In Harald Sack, Giuseppe Rizzo,  
391 Nadine Steinmetz, Dunja Mladenić, Sören Auer, and Christoph Lange, editors, *The Semantic*  
392 *Web*, pages 84–87, Cham, 2016. Springer International Publishing.
- 393 [30] Balashanmuga Priyan Rajamohan, Alexander C Harding Bradley, Van D Tran, Jonathan E Gor-  
394 don, Hayden W Caldwell, Redad Mehdi, Gabriel Ponce, Quynh D Tran, Ozan Dernek, Jarod  
395 Kaltenbaugh, et al. Materials data science ontology (mds-onto): Unifying domain knowledge  
396 in materials and applied data science. *Scientific Data*, 12(1):628, 2025.
- 397 [31] Rob Raskin and Michael Pan. Semantic web for earth and environmental terminology (sweet).  
398 In *Proc. of the Workshop on Semantic Web Technologies for Searching and Retrieving Scientific*  
399 *Data*, volume 25, 2003.
- 400 [32] Angelo A Salatino, Thiviyan Thanapalasingam, Andrea Mannocci, Francesco Osborne, and  
401 Enrico Motta. The computer science ontology: a large-scale taxonomy of research areas. In  
402 *International Semantic Web Conference*, pages 187–205. Springer, 2018.

- 403 [33] Wes Schafer, Oliver He, Anna L. Dunn, and Zachary E. X. Dance. Ontology for process  
404 chemistry – giving context to instrument data structured by the allotrope data model. Allotrope  
405 Connect Virtual Conference, April 2021. Virtual conference, April 19–26, 2021.
- 406 [34] Schema.org Community Group. Schema.org. <https://schema.org/>, 2011.
- 407 [35] Lynn M Schriml, Elvira Mitra, James Munro, Becky Tauber, Mike Schor, Lance Nickle, Vic-  
408 tor Felix, Linda Jeng, Cynthia Bearer, Richard Lichenstein, Katharine Bisordi, Nicole Cam-  
409 pion, Brooke Hyman, David Kurland, Connor Patrick Oates, Siobhan Kibbey, Poorna Sreeku-  
410 mar, Chris Le, Michelle Giglio, and Carol Greene. Human disease ontology 2018 update:  
411 classification, content and workflow expansion. *Nucleic Acids Research*, 47(D1):D955–D962,  
412 11 2018.
- 413 [36] Philip Strömert, Giacomo Lanza, Johannes Hunold, Steffen Neumann, et al. The Vibrational  
414 Spectroscopy Ontology (VIBSO), 2023.



(a) Traditional Retriever Approach



(b) LLM-Augmented Retriever Approach

Figure 8: Performance comparison of retrieval-based methods on the **taxonomy discovery** task.

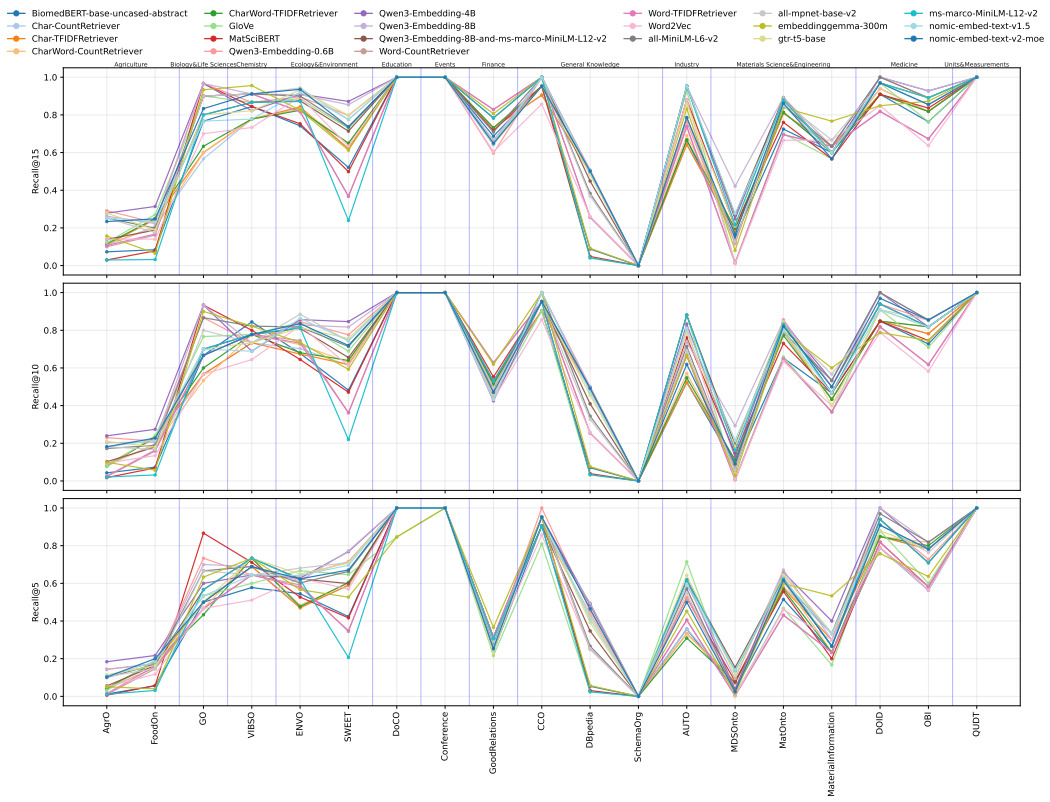


Figure 9: Performance comparison of retrieval-based methods on the **non-taxonomic relationship extraction** task.



Figure 10: Comparison of reranking methods across the three OL tasks—term typing, taxonomy discovery, and non-taxonomic relationship extraction—measured using grouped **F1-score**. The figure summarizes the overall balance between precision and recall achieved by each reranker across tasks and domains.



Figure 11: Comparison of reranking methods across the three OL tasks using grouped **precision**. The results highlight the ability of each reranker to return highly relevant candidates while minimizing false positives in different ontology learning settings.

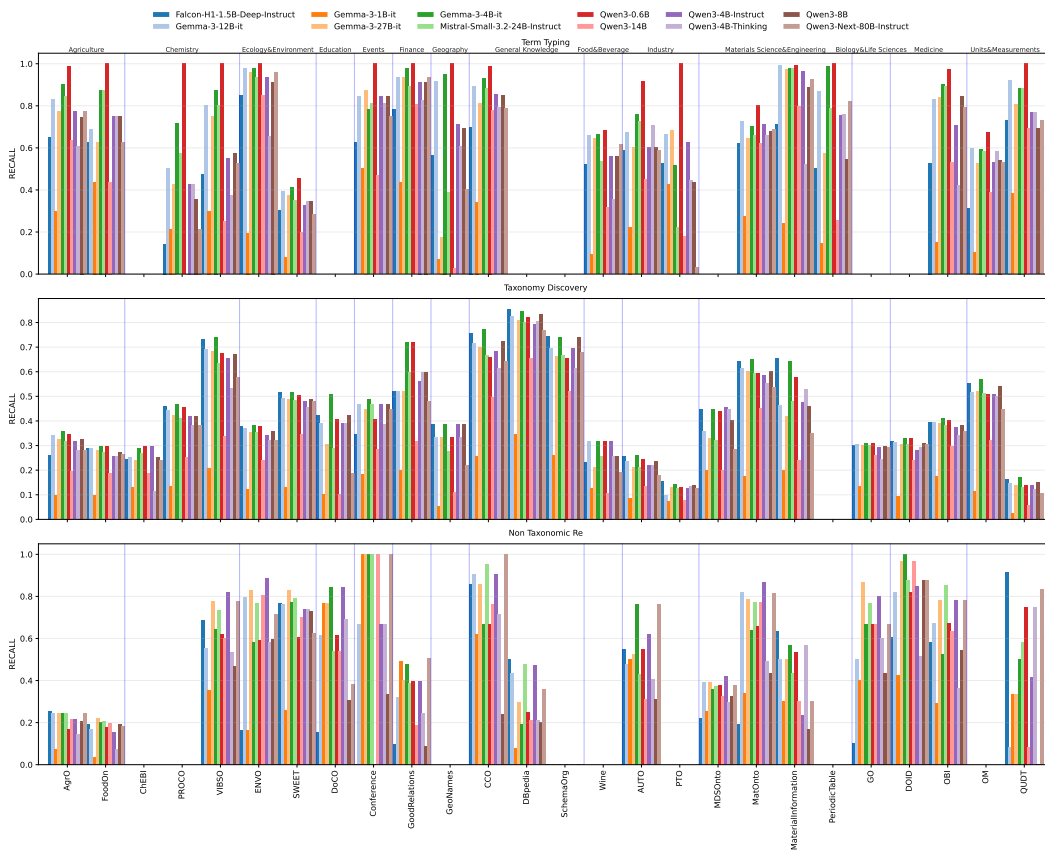


Figure 12: Comparison of reranking methods across the three OL tasks using grouped **recall**. The figure reflects how effectively each approach retrieves relevant candidates, emphasizing coverage across term typing, taxonomy discovery, and relation extraction tasks.