

# KROMA: Ontology Matching with Knowledge Retrieval and Large Language Models

Lam Nguyen, Erika Barcelos, Roger French, and Yinghui Wu

Case Western Reserve University, Cleveland OH 44106, USA  
ltn18@case.edu

**Abstract.** Ontology Matching (OM) is a cornerstone task of semantic interoperability, yet existing systems often rely on handcrafted rules or specialized models with limited adaptability. We present KROMA, a novel OM framework that harnesses Large Language Models (LLMs) within a Retrieval-Augmented Generation (RAG) pipeline, to dynamically enrich the semantic context of OM tasks with structural, lexical, and definitional knowledge. To optimize both accuracy and efficiency, KROMA integrates a bisimilarity-based concept matching and a lightweight ontology refinement step, which prune candidate concepts and substantially reduce the communication overhead from invoking LLMs. Using a variety of benchmark datasets, we experimentally verify that knowledge retrieval and context-augmented LLMs effectively improve ontology matching outperforms both traditional OM systems and state-of-the-art LLM-based OM methods, with comparable communication overhead. Our study highlights the feasibility and benefit of the proposed optimization techniques (targeted knowledge retrieval, prompt enrichment, and ontology refinement) for ontology matching at scale. Code is available at: <https://anonymous.4open.science/r/kroma/>

**Keywords:** Ontology Matching · Large Language Models · Retrieval Augmented Generation

## 1 Introduction

Ontologies have been routinely developed to unify and standardize knowledge representation to support data-driven applications. They allow researchers to harmonize terminologies and enhance knowledge and sharing in their fields. Ontologies can be classified according to their level of specificity, ranging from more abstract, general-purpose ontologies such as Basic Formal Ontology (BFO)[37] or DOLCE[5] down to more domain-oriented ‘mid-level’ ones, such as CheBi [13] in chemistry, Industrial Ontology Foundry (IOF) [16] or Common Core Ontology (CCO) [26]. Domain ontologies are data-driven, task-specific ‘low-level’ ontologies, containing concepts from domain-specific data, such as Materials Data Science Ontology (MDS-Onto) [45]. To achieve broad usability, ontologies need to be aligned for better interoperability via ontology matching.

Ontology matching has been studied to find correspondence between terms that are semantically equivalent [51]. It is a cornerstone task to ensure semantic interoperability among terms originated from different sources. Ontology

matching methods can be categorized to rule-based or structural-based (graph pattern or path-based) matching [9], matching with semantic similarity, machine learning-based approaches and hybrid methods. Linguistic (terminological) methods are often used for ontology matching tasks, ranging from simple string matching or embedding learning to advance counterparts based on natural language processing (NLP). Nevertheless, conventional rule- or structural-based methods are often restricted to certain use cases and hard to be generalized for new or unseen concepts. Learning-based approaches may on the other hand require high re-training process, for which abundant annotated or training data remains a luxury especially for *e.g.*, data-driven scientific research.

Meanwhile, the emerging Large Language Models (LLMs) have demonstrated remarkable versatility for NL understanding. LLMs are trained on vast and diverse corpora, endowing them with an understanding of language nuances and contextual subtleties. Extensive training enables them to capture semantic relationships and abstract patterns that are critical for aligning concepts across different data sources. A missing yet desirable opportunity is to investigate whether and how LLMs can be engaged to automate and improve ontology matching.

This paper introduces KROMA, a novel framework that exploits LLMs to enhance ontology matching. Unlike existing LLM-based methods that typically “outsource” ontology matching to LLMs with direct prompting (which may have low confidence and risk of hallucination), KROMA maintain a set of conceptually similar groups that are co-determined by concept similarity and LLMs, both guided by their “context” obtained via a runtime knowledge retrieval process. Moreover, the groups are further refined by a global ontological equivalence relation that incorporate structural equivalence.

**Contributions.** Our main contributions are summarized below.

- (1) We propose a formulation of semantic equivalence relation in terms of a class of bisimilar equivalence relation, and formally define the ontology structure, called concept graph, to be maintained (**Sections 2 and 3**). We justify our formulation by showing the existence of an “optimal” concept graph with minimality and uniqueness guarantee, subject to the bisimilar equivalence.
- (2) We introduce KROMA, an LLM-enhanced ontology matching framework (**Section 4**). KROMA fine-tunes LLMs with prompts over enriched semantic contexts. Such contexts are obtained from knowledge retrieval, referencing high-quality, external knowledge resources.
- (3) KROMA supports both offline and online matching, to “cold-start” from scratch, and to digest terms arriving from a stream of data, respectively. We introduce efficient algorithms to correctly construct and maintain the optimal concept graphs (**Section 5**). (a) The offline refinement performs a fast grouping process guided by the bisimilar equivalence, blending concept equivalence tests co-determined by semantic closeness and LLMs. (b) The online algorithm effectively incrementalizes its offline counterpart with fast delay time for continuous concept streams. Both algorithms are in low polynomial time, with optimality guarantee on the computed concept graphs.

(4) Using benchmarking ontologies and knowledge graphs, we experimentally verify the effectiveness and efficiency of KROMA (**Section 6**). We found that KROMA outperforms existing LLM-based methods by 10.95% on average, and the optimization of knowledge retrieval and refinement improves its accuracy by 6.65% and 2.68%, respectively.

**Related Work.** We summarize related work below.

*Large Language Models.* Large language models (LLMs) have advanced NLP by enabling parallel processing and capturing complex dependencies [52,47], which have scaled from GPT-1’s 117M [43], GPT-2’s 1.5B [44] to GPT-3’s 175B [34] and GPT-4’s 1.8T parameters [35]. Open-source models like Llama have grown to 405B [31], with Mistral Large (123B) [32] and DeepSeek V3 (671B) [11] also emerging. Recent advances in specialized reasoning LLMs (RLLMs) such as OpenAI’s O1 and DeepSeek R1 have further propelled Long Chain-of-Thought reasoning—shifting from brief, linear Short CoT to deeper, iterative exploration, and yielded substantial gains in multidisciplinary tasks [36,12,46,49,7,41,55,29].

*Ontology Matching with LLMs.* Several methods have been developed to exploit LLMs for ontology matching. Early work focused on direct prompting LLMs for ontology matching. For example, [38] frame product matching as a yes/no query, and [33] feed entire source and target ontologies into ChatGPT—both achieving high recall on small OAEI conference-track tasks but suffering from low precision. Beyond these “direct-prompt” approaches, state-of-the-art LLM-OM systems fall into two main categories: (1) retrieval-augmented pipelines, which first retrieve top- $k$  candidates via embedding-based methods and then refine them with LLM prompts (e.g. LLM4OM leverages TF-IDF and SBERT retrievers across concept, concept-parent, and concept-children representations [19], while MILA adds a prioritized depth-first search step to confirm high-confidence matches before any LLM invocation [48]); and (2) prompt-engineering systems, which generate candidates via a high-precision matcher or inverted index and then apply targeted prompt templates to LLMs in a single step (e.g. OLaLa embeds SBERT candidates into MELT’s prompting framework with both independent and multi-choice formulations [24], and LLMMap uses binary yes/no prompts over concept labels plus structural context with Flan-T5-XXL or GPT-3.5 [22]).

## 2 Ontologies and Ontology Matching

**Ontologies.** An ontology  $O$  is a pair  $(C, E)$ , where  $C$  is a finite set of concept names, and  $E \subset C \times C$  is a set of relations. An ontology has a graph representation with a set of concept nodes  $C$ , and a set of edges  $E$ . We consider ontologies as directed acyclic graphs (DAGs).

In addition, each concept node (or simply “node”)  $c$  in  $O$  carries the following auxiliary structure. (1) The *rank* of a node  $c \in C$  is defined as: (a)  $r(c) = 0$  if  $c$  has no child, otherwise, (b)  $r(c) = \max(r(c') + 1)$  for any child  $c'$  of  $c$  in  $O$ . (2) A *ground set*  $c.\mathcal{I}$ , refers to a set of entities from *e.g.*, external ontologies or knowledge bases that can be validated to belong to the concept  $c$ .

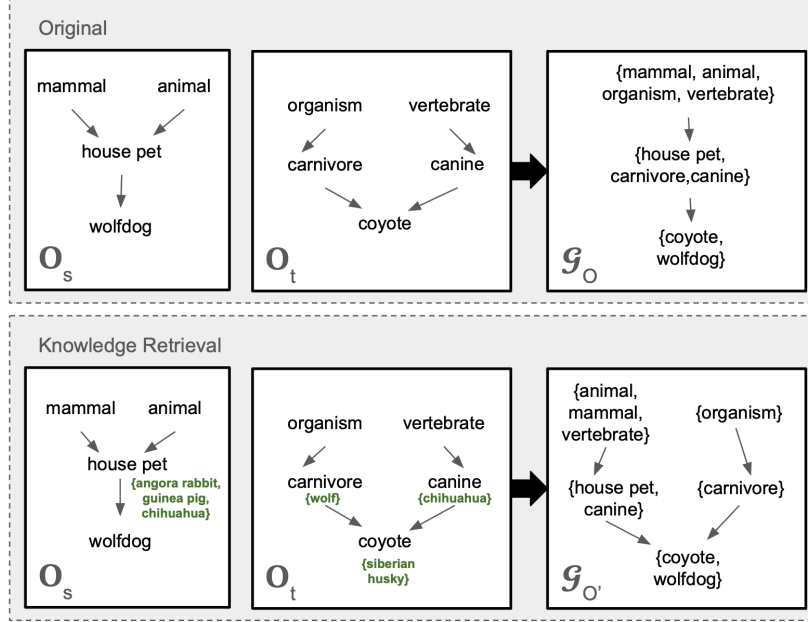


Fig. 1: Ontologies with ground sets, Ontology Matching and Concept Graphs

**Ontology Matching.** Given a source ontology  $O_s = (C_s, E_s)$  and a target ontology  $O_t = (C_t, E_t)$ , an *ontological equivalence relation*  $R_{\sim} \subseteq C_s \times C_t$  is an equivalence relation that contains pairs of nodes  $(c, c')$  that are considered to be “semantically equivalent”. The relation  $R_{\sim}$  induces an equivalence partition  $\mathcal{C}$  of the concept set  $C_s \cup C_t$ , such that each partition is an equivalence class that contains a set of pairwise equivalent concepts in  $C_s \cup C_t$ .

Consistently, we define a *concept graph*  $\mathcal{G}_O = (V_O, E_O)$  as a DAG with a set of nodes  $V_O$ , where each node  $[c] \in V_O$  is an equivalence class in  $\mathcal{C}$ , and there exists an edge  $([c], [c']) \in E_O$  if and only if there exists an edge  $(c, c')$  in  $E_s$  or in  $E_t$ . A concept graph  $\mathcal{G}_O$  can be a multigraph: there may exist multiple edges of different relation names between two equivalent classes.

Given  $O_s$  and  $O_t$ , the task of ontology matching is to properly formulate  $R_{\sim}$  and compute  $\mathcal{C}$  induced by  $R_{\sim}$  over  $O_s$  and  $O_t$ ; or equivalently, compute the concept graph  $\mathcal{G}_O$  induced by  $R_{\sim}$ .

*Example 1.* Figure 1 depicts two ontologies  $O_s$  and  $O_t$  as DAGs, involving in total 9 concept nodes. An equivalence relation may suggest that “mammal”, “animal”, “organism” and “vertebrate” are pairwise similar; and similarly for the sets {“house pet”, “carnivora” and “canine”}, and {“wolfdog”, “coyote”}. This induces a concept graph  $\mathcal{G}_o$  as shown on the top right as a result of ontology matching, with three equivalence classes.

**Language Models for Ontology Matching.** A language model  $M$  takes as input a prompt query  $q$ , and generate an answer  $q(M)$ , typically an NL statement, for downstream processing. Large language models (LLMs) are foundation

models that can effectively learn from a handful of in-context examples, included in a prompt query  $q$ , that demonstrate input–output distribution [53].

*Prompt query.* A prompt query  $q$  is in a form of NL statements that specifies (1) a task definition  $\mathcal{T}$  with input and output statement; (2) a set of in-context examples  $\mathcal{D}$  with annotated data; (3) a statement of query context  $Q$ , which describe auxiliary query semantics; and optionally (4) specification on output format, and (5) a self-evaluation of answer quality such as confidence. An evaluation of a prompt query  $q$  invokes an LLM  $M$  to infer a query result  $q(M)$ .

We specify a prompt query  $q$  and LLMs for ontology matching. A prompt query  $q$  is in the form of  $q(c, c')$ , which asks “are  $c$  and  $c'$  semantically equivalent?” An LLM  $M$  can be queried by  $q(c, c')$  and acts as a Boolean “oracle” with “yes/no” answer. An LLM is *deterministic*, if it always generate a same answer for the same prompt query. We consider deterministic LLMs, as in practice, such LLMs are desired for consistent and robust performance.

### 3 Ontology Matching with LLMs

In this section, we provide a pragmatic characterization for the ontological equivalence relation  $R_{\simeq}$ . We then formulate the ontology matching problem.

#### 3.1 Semantic Equivalence: A Characterization

**Concept similarity.** A variety of methods have been proposed to decide if two *concepts* are equivalent [6]. KROMA by default uses a Boolean function  $F$  defined by a weighted combination of a semantic closeness metric  $\text{sim}^1$  and the result from a set of LLMs  $\mathcal{M}$  (see Section 4).

$$F(c, c') = \gamma \text{sim}(c, c') + (1 - \gamma)q(c, c', \mathcal{M})$$

where  $q$  is a prompt query that specifies the context of concept equivalence for LLMs. KROMA supports a built-in library of semantic similarity functions  $\text{sim}$ , including (a) string similarity, feature and information measure [39], or normalized distances (NGDs) [27]; and (b) a variety of LLMs such as GPT-4o Mini (OpenAI) [35], LLaMA-3.3 (Meta AI) [30], and Qwen-2.5 (Qwen Team) [42].

**Ontological Bisimilarity.** We next provide a specification of the ontological equivalence relation, notably, *ontological bisimilarity*, denoted by the same symbol  $R_{\simeq}$  for simplicity. Given a source ontology  $O_s = (C_s, E_s)$ , and a target ontology  $O_t = (C_t, E_t)$ , we say a pair of nodes  $c_s \in C_s$  and  $c_t \in C_t$  are *ontologically bisimilar*, denoted as  $(c_s, c_t) \in R_{\simeq}$ , if and only if there exists a non-empty binary relation  $R_{\simeq}$ , such that: (1)  $c_s$  and  $c_t$  are concept similar, *i.e.*,  $F(c_s, c_t) \geq \alpha$ , for a threshold  $\alpha$ ; (2) for every edge  $(c'_s, c_s) \in E_s$ , there exists an edge  $(c'_t, c_t) \in E_t$ , such that  $(c'_s, c'_t) \in R_{\simeq}$ , and vice versa; and (3) for every edge  $(c_s, c''_s) \in E_s$ , there exists an edge  $(c_t, c''_t) \in E_t$ , such that  $(c''_s, c''_t) \in R_{\simeq}$ .

One can verify the following result.

<sup>1</sup> We adopt similarity metrics that satisfy transitivity, *i.e.*, if concept  $c$  is similar to  $c'$ , and  $c'$  is similar to  $c''$ , then  $c$  is similar to  $c''$ . This is to ensure transitivity of ontology equivalence, and is practical for representative concept similarity measures.

**Lemma 1.** *The ontological bisimilar relation  $R_{\sim}$  is an equivalence relation.*

We can prove the above results by verifying that  $R_{\sim}$  is reflexive, symmetric and transitive over the concept set  $C_s \cup C_t$ , ensured by the transitivity of concept similarity and by definition. Observe that two nodes that are concept similar may *not* be ontologically bisimilar. On the other hand, two ontologically bisimilar entities must be concept similar, by definition.

### 3.2 Problem Statement

We now formulate our ontological matching problem. Given a *configuration*  $W = (O_s, O_t, F, \mathcal{M}, \alpha)$  that specifies as input a source ontology  $O_s$ , a target ontology  $O_t$ , a Boolean function  $F$  and threshold  $\alpha \in (0, 1]$  that asserts concept similarity, and a set of LLMs  $\mathcal{M}$ , the problem is to compute a smallest concept graph  $\mathcal{G}_O$  induced by the ontologically bisimilar equivalence  $R_{\sim}$ .

We justify the above characterization by proving that there exists an “optimal”, invariant solution encoded by a concept graph  $\mathcal{G}_O$ . To see this, we provide a *minimality* and *uniqueness* guarantee on  $\mathcal{G}_O$ .

**Lemma 2.** *Given a configuration  $W$  and semantic equivalence specified by the ontological bisimilar relation  $R_{\sim}$ , there is a unique smallest concept graph  $\mathcal{G}_O$  that capture all semantically equivalent nodes in terms of  $R_{\sim}$ .*

**Proof sketch:** We show that the above result holds by verifying the following. (1) There is a unique, maximum ontological bisimilar relation  $R_{\sim}$  for a given configuration  $W = (O_s, O_t, F, \mathcal{M})$ , where any LLM in  $\mathcal{M}$  is a deterministic model for the same prompt query  $q$  generated consistently from  $W$ . This readily follows from Lemma 1, which verifies that  $R_{\sim}$  is an equivalence relation. (2) Let the union of  $O_s$  and  $O_t$  be a graph  $O_{st} = \{C_s \cup C_t, E_s \cup E_t\}$ . By setting  $\mathcal{G}_O$  as the quotient graph induced by the largest ontological bisimilar relation  $R_{\sim}$  over  $O_{st}$ ,  $\mathcal{G}_O$  contains the smallest number of equivalent classes (nodes) and edges. This can be verified by proof by contradiction. (3) The uniqueness of the solution can be verified by showing that  $R_{\sim}$  induces only one unique partition  $\mathcal{C}$  and results in a concept graph  $\mathcal{G}_O$  up to graph isomorphism. In other words, for any two possible concept graphs induced by  $R_{\sim}$ , they are isomorphic to each other.  $\square$

The above analysis suggests that for a configuration  $W$ , it is desirable to compute such an optimal concept graph as an invariant, stable result with guarantees on sizes and uniqueness on topological structures. We next introduce KROMA and efficient algorithms to compute such optimal concept graphs.

## 4 KROMA Framework

### 4.1 Framework Overview

Given a *configuration*  $W = (O_s, O_t, F, \mathcal{M})$  where  $\mathcal{M}$  is a set of pre-trained LLMs  $\mathcal{M}$ , KROMA has the following major functional modules that enables a multi-session ontology matching process.

*Concept Graph Initialization* Upon receiving two ontologies  $O_s = (C_s, E_s)$  and  $O_t = (C_t, E_t)$ , KROMA initializes the concept graph  $\mathcal{G}_O = (V, E)$  with  $V =$

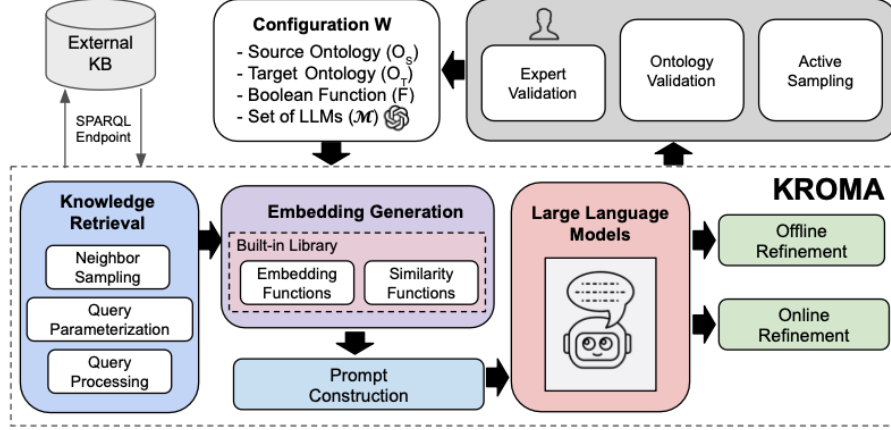


Fig. 2: KROMA Framework Overview: Major Components and Dataflow

$C_S \cup C_T$  and  $E = E_S \cup E_T$ . For each concept  $c \in V$ , KROMA computes the *rank*  $r(c)$  as described in Section 2.

*Knowledge Retrieval.* To assemble a rich, yet compact, context for each concept  $c \in V$ , we perform: (1) *Neighborhood Sampling*: traverse up to its two hops in  $\mathcal{G}_O$  to collect parents, children, and “sibling” concepts of  $c$ , creating a node induced subgraph centered at  $c$ ; (2) *subgraph parameterization*: Sample and substitute constant values from the subgraph with variables to create SPARQL queries  $S_q$ ; and (3) *ground set curation*: process  $S_q$  onto external knowledge bases to augment its ground set  $c.\mathcal{I}$  with auxiliary information (e.g., relevant entities, definition, labels, etc.).

*Embedding Generation.* In this phase, KROMA consults a built-in library of semantic similarity functions  $\text{sim}(\cdot, \cdot)$  and embedding functions  $\text{embd}(\cdot)$ . For each concept  $c \in V$ , KROMA computes:

$$z_c = \alpha \text{embd}_{\text{graph}}(c) + (1 - \alpha) \text{embd}_{\text{text}},$$

where  $\text{embd}_{\text{graph}}$  (e.g. node2vec [21]) captures the  $c$ ’s topology information and  $\text{embd}_{\text{text}}$  (e.g. SciBERT [4]) encodes  $c$ ’s textual context. After obtaining the necessary embeddings ( $z_c \forall c \in V$ ), KROMA computes pairwise concept similarity between source and target ontologies using  $\text{sim}$  functions:

$$\mathbb{S} = \{(c_s, c_t, \text{score}_{s,t}) \mid c_s \in C_S, c_t \in C_T, \text{score}_{s,t} = \text{sim}(z_{c_s}, z_{c_t})\}$$

From  $\mathbb{S}$ , we select the top- $k$  pairs with the highest similarity scores (i.e. in descending order of  $\text{score}_{s,t}$ ), yielding the candidate list  $\mathbb{C}$ .

*Example 2.* We revisit Example 1. (1) A knowledge retrieval for node “house pet” samples its neighbors and issues a set of “star-shaped” SPARQL queries centered at “house pet” to query an underlying knowledge graph  $\mathcal{KG}$ . This enriches its ground set with a majority of herbivore or omnivorous pets that are not “carnivore”. Similarly, the ground set of “carnivore” is enriched by “wolf”, unlikely a

<b>Task Description</b> Given two concepts: <b>source</b> and <b>target</b> . Determine if they are semantically equivalent or not.	<b>Task Description</b> Given two ontology concepts: <b>heatwave</b> and <b>drought</b> . Determine if they are semantically equivalent or not.
<b>Examples</b> <ul style="list-style-type: none"> <li>• <b>Example 1:</b> <b>source1</b> (+ metadata) is <b>semantically equivalent</b> to <b>target1</b> (+ metadata)</li> <li>• <b>Example 2:</b> <b>source2</b> (+ metadata) is <b>not semantically equivalent</b> to <b>target2</b> (+ metadata)</li> <li>• 3 more examples ...</li> </ul>	<b>Examples</b> <ul style="list-style-type: none"> <li>• <b>Example 1:</b> <b>urban flooding</b> (parents = [hydrological hazard], ...) is <b>semantically equivalent</b> to <b>city flood</b> (parents = [flood event], ...)</li> <li>• <b>Example 2:</b> <b>soil erosion</b> (parents = [land degradation], ...) is <b>not semantically equivalent</b> to <b>landslide</b> (parents = [mass wasting], ...)</li> <li>• 3 more examples ...</li> </ul>
<b>Query Context</b> <ul style="list-style-type: none"> <li>• Source context: <b>source</b> (parents = [...], children = [...], synonyms=[...], definition = "...")</li> <li>• Target context: <b>target</b> (parent=[...], children = [...], synonyms=[...], definition = "...")</li> </ul>	<b>Query Context</b> <ul style="list-style-type: none"> <li>• Source context: <b>heatwave</b> (parents = [extreme weather], children = [public health emergency], synonyms=[extreme heat], definition = "...")</li> <li>• Target context: <b>drought</b> (parent=[climatic anomaly], children = [crop failure], synonyms=[prolonged dryness], definition = "...")</li> </ul>
<b>Output Format</b> <ans> Yes/No </ans>	<b>Output Format</b> <ans> Yes/No </ans>
<b>Confidence</b> <conf> 0-10 </conf>	<b>Confidence</b> <conf> 0-10 </conf>

Fig. 3: KROMA Prompt Query Template and Query Example

house pet. Despite “coyote” and “wolfdog” (house pet) alone are less similar, the ground set of “coyote” turns out to be a set of canine pets *e.g.*, “husky” that are “coyote-like”, hence similar with “wolfdog”. (2) The embedding generation phase incorporates enriched ground sets and generate embeddings accordingly, which scores that distinguishes “house pet” from “carnivore” due to embedding difference, and assert “coyote” and “wolfdog” to be concept similar, hence a candidate pair to be “double checked” by LLMs in the next phase.

*Prompt Querying LLMs.* For each candidate pair  $(c_s, c_t) \in \mathbb{C}$ , KROMA automatically generates an NL prompt that includes: (1) Task description  $\mathcal{T}$  (e.g. “*Given two ontology concepts and their contextual metadata, decide if they are related or not.*”), (2) In-context examples  $\mathcal{D}$  containing both positive and negative matches, (3) Query context for  $c_s$  and  $c_t$  including their ground sets, (4) Output format and confidence (e.g., “Answer Yes or No, and provide a confidence score between 0 and 10.”). It then calls (a set of) LLMs  $\mathcal{M}$  to obtain a match decision with confidence. Low-confidence or conflicting outputs are routed to validation module. A query template and a generated example is illustrated in Figure 3.

*Ontology Refinement & Expert Validation.* KROMA next invokes a refinement process, *OfflineRefine* (Algorithm 1), to “cold-start” the construction of  $\mathcal{G}_O$ , or *OnlineRefine* (Algorithm 2), to incrementally refine  $\mathcal{G}_O$  for any unseen, newly arrived concept nodes or edges. Any pair of nodes whose structural ranks remain in “conflict” is routed into a set of queries for expert validation (see Algorithm 2, Section 4); once approved, are integrated back into  $\mathcal{G}_O$ . An active sampling strategy is applied, to select pairs of nodes that have low confidence from LLMs for expert validation, to reduce the manual effort.



**Algorithm 1:** Offline Refinement

---

**Input:** Source ontology  $O_S = (C_S, E_S)$ , target ontology  $O_T = (C_T, E_T)$   
**Output:** Concept graph  $\mathcal{G}_O$ .

```

1 set  $V \leftarrow C_S \cup C_T$ ; set  $E \leftarrow E_S \cup E_T$ ;
2 Initialize concept graph  $\mathcal{G}_O = (V, E)$ ;
3 foreach  $c \in V$  do
4    $\lfloor$  compute rank  $r(c)$  as in Section 2;
5  $\rho \leftarrow \max_{c \in V} r(c)$ ;
6 for  $i \leftarrow 0$  to  $\rho$  do
7    $B_i \leftarrow \{c : r(c) = i\}$ ;
8  $P \leftarrow \{B_0, \dots, B_\rho\}$ ;
9 for  $i \leftarrow 0$  to  $\rho$  do
10   $D_i \leftarrow \{X \in P : X \subseteq B_i\}$ ;
11  foreach  $X \in D_i$  do
12     $G \leftarrow \text{collapse}(G, X)$ ;
13  foreach  $c \in B_i$  do
14    foreach  $C \in P$  with  $C \subseteq \bigcup_{j>i} B_j$  do
15      Split  $C$  into  $C_1, C_2$  by adjacency to  $c$ ;
16       $P \leftarrow (P \setminus \{C\}) \cup \{C_1, C_2\}$ ;
17 return  $\mathcal{G}_O$ ;
```

---

*Example 3.* Continuing with Example 2, as “golden retriever” and “coyote” are asserted by the function  $F$  that combines the descision of semantic similarity function  $\text{sim}$  and LLMs, an equivalence class is created in  $\mathcal{G}'_O$ . As “house pet” and “carnivore” has quite different embedding considering the features from themselves and their ground sets, “carnivore” is separated from the group of “house pet”. This suggests further that “organism” now has a different context that distinguish it from the group  $\{\text{mammal}, \text{animal}, \text{vertebrate}\}$ , by the definition of bisimilarity equivalence. This leads to a finer-grained concept graph  $\mathcal{G}'_O$ .

## 5 Ontology Refinement

We next describe our ontology refinement algorithms. KROMA supports ontology refinement in two modes. The offline mode assumes that the source ontology  $O_s$  and the target ontology  $O_t$  are known, and performs a batch processing to compute the concept graph  $\mathcal{G}_O$  from scratch. The online refinement incrementally maintains  $\mathcal{G}_O$  upon a sequence of (unseen) triples (edges) from external resources.

**Offline Refinement.** The offline refinement algorithm is outlined as Algorithm 1. (1) It starts by initializing  $\mathcal{G}_O$  (lines 1-6) as the union of  $O_s$  and  $O_t$ , followed by computing the node ranks. At each rank, it initializes a “bucket”  $B_i$  (as a single node set) that simply include all the concept nodes at rank  $i$  (lines 7-8), and initialize a partition  $\mathcal{P}$  with all the buckets (line 9). It then follows a “bottom-up” process to refine the buckets iteratively, by checking if two concepts  $c$  and  $c'$  in a same bucket  $B_i$  are concept similar (as asserted by LLM and embedding similarity), and have all the neighbors that satisfy the requirement

of ontological bisimilar relation by definition (lines 10-13). If not, a procedure `collapse` is invoked, to (1) split the bucket  $B_i$  into three fragments:  $B_i^1 = B_i \setminus \{c, c'\}$ ,  $B_i^2 = \{c\}$ , and  $B_i^3 = \{c'\}$ , followed by a “merge” check to test if  $c$  and  $c'$  can be merged to  $B_i^1$ ; and (2) propagate this change to further “split-merge” operators to affected buckets at higher ranks (lines 14-17). This process continues until no change can be made.

*Correctness.* Algorithm 1 correctly terminates at obtaining a maximum bisimilar ontological equivalence relation  $R_{\simeq}$ , with two variants. (1) As ontologies are DAGs, it suffices to perform a one-pass, bottom-up splitting of equivalence classes following the topological ranks; (2) the `collapse` operator ensures the “mergable” cases to reduce unnecessary buckets whenever a new bucket is separated. This process simulates the correct computation of maximum bisimulation relation in Kripke structures (a DAG) [14], optimized for deriving ontology matching determined by LLM-based concept similarity and bisimilar equivalence.

*Time Cost.* The initialization of concept graph  $\mathcal{G}_O$  is in  $O(|O_s| + |O_t|)$ . Here  $|O_s| = |V_s| + |E_s|$ ; and  $|O_t|$  is defined similarly. The iterative collapse (lines 10-17) takes in  $O(|O_s| + |O_t|)$  time as the number of buckets (resp. edges) is at most  $|C_s| + |C_T|$  (resp.  $|E_s| + |E_t|$ ). The overall cost is in  $O(|O_s| + |O_t|)$ .

*Overall Cost.* We consider the cost of the entire workflow of KROMA. (1) The knowledge checking takes  $O((|C_s| + |C_T|)|KG|)$  time, where  $|KG|$  refers to the size of the external ontology or knowledge graph that is referred to by the knowledge retrieval via *e.g.*, SPARQL access. Note here we consider SPARQL queries with “star” patterns, hence the cost of query processing (to curate ground sets) is in quadratic time. (2) The total cost of LLM inference is in  $O(|C_s||C_T|T_I)$ , for a worst case that any pair of nodes in  $C_s \times C_t$  are concept similar in terms of  $\sim$ . Here  $T_I$  is the unit cost of processing a prompt query. Putting these together, the total cost is in  $O(|C_s||C_T|T_I + (|C_s| + |C_T|)|KG| + (|O_s| + |O_t|))$  time.

**Online Refinement.** We next outline the online matching process. In this setting, KROMA receives new ontology components as an (infinite) sequence of triples (edges), and incrementally maintain a concept graph  $\mathcal{G}_O$  by processing the sequence input in small batched updates  $\Delta\mathcal{G}$ . For each newly arrived concept (node)  $c$  in  $\Delta\mathcal{G}$ , KROMA conducts knowledge retrieval to curate  $c.\mathcal{I}$ ; and consult LLMs to decide if  $c$  is concept similar to any node in  $\mathcal{G}_O$ . It then invokes online refinement algorithm to enforce the ontological bisimilar equivalence.

The Algorithm 2 first updates the buckets in  $\mathcal{G}_O$  by incorporating the nodes from  $\Delta\mathcal{G}$  that are verified to be concept similar, as well as their ranks. It then incrementally update the buckets to maintain the bisimilar equivalence consistently via a “bottom-up” split-merge process as in Algorithm 1 (lines 6-15). Due to the unpredictability of the “unseen” concepts, the online refinement defers the processing of two “inconsistent” cases for experts’ validation: (1) When a concept  $c$  is determined to be concept similar by function  $\sim$  but not LLMs with high confidence; or (2) whenever for a new edge  $(c, c') \in \Delta\mathcal{G}$ ,  $(c, [c_1]) \in R_{\simeq}$ ,  $(c', [c_2]) \in R_{\simeq}$ , yet  $r(c_1) < r(c_2)$  in  $\mathcal{G}_O$ . Both require domain experts’ feedback

**Algorithm 2:** Online Ontology Refinement

---

**Input:** a concept graph  $\mathcal{G}_O=(V_O, E_O)$ ; a sequence of triples  $\Delta\mathcal{G}$ ;  
**Output:** Updated concept graph  $\mathcal{G}_O \oplus \Delta\mathcal{G}$

---

```

1 set  $\mathcal{Q} \leftarrow \emptyset$ ;
2 update  $r(c)$  for all  $c \in V_O$  w.r.t. new edges in  $\Delta\mathcal{G}$ ;
  // Process new insertion in ascending order of rank
3 Sort  $\Delta\mathcal{G}$  with ascending ranks ;
4 foreach  $e = (u \rightarrow u') \in \Delta\mathcal{G}$  do
  // Maintain bisimilar equivalence classes  $[\cdot]_{R_e}$ 
5   Split( $u, u', [u]_{R_e}, [u']_{R_e}$ );
6   if  $r([u]_{R_e}) > r([u']_{R_e})$  then
7     foreach  $v \in \mathcal{B}([u]_{R_e})$  do
8       Merge( $\{u\}, v$ );
9     foreach  $v' \in \mathcal{B}([u']_{R_e})$  do
10      Merge( $\{u'\}, v'$ );
11   else if  $r([u]_{R_e}) = r([u']_{R_e})$  then
12     foreach  $v \in \mathcal{P}([u']_{R_e})$  do
13       Merge( $\{u\}, v$ );
14     foreach  $v' \in \mathcal{C}([u]_{R_e})$  do
15       Merge( $\{u'\}, v'$ );
16   else
17     // for inconsistent ranks, defer to further
      validation
       $\mathcal{Q} \leftarrow \mathcal{Q} \cup \{(u, u')\}$ ;
18 return  $\mathcal{G}_O$  with  $\mathcal{Q}$  to be resolved;
```

---

to resolve. These cases are cached into a query set  $\mathcal{Q}$  to be further resolved in the validation phase (see Section 4).

*Analysis.* The correctness of online refinement carries over from its offline counterpart, and that it correctly incrementalize the split-merge operator for each newly arrived edges. For time cost, for each batch, it takes a delay time to update  $\mathcal{G}_O$  in  $O(|\mathcal{G}_O| + |\Delta\mathcal{G}| \log |\Delta\mathcal{G}| + |\Delta\mathcal{G}| \log |V_O|)$  time. This result verifies that online refinement is able to response faster than offline maintenance that recomputes the concept graph from scratch.

The Online Ontology Refinement procedure runs in near-linear time in the size of the compressed graph and the update batch. Recomputing all node ranks on  $\mathcal{G}_r$  costs  $O(N + M)$  (with  $N = |V_r|$ ,  $M = |E_r|$ ), and sorting the  $\Delta$  new edges by rank takes  $O(\Delta \log \Delta)$ . Each insertion is then processed via union-find splits and merges in  $O((M + \Delta) \alpha(N))$  overall—since every edge participates in at most one split or merge and each such operation is  $O(\alpha(N))$ , where  $\alpha$  is the inverse Ackermann function. Finally, dispatching  $\mathcal{Q}$  deferred pairs to the human verifier incurs an extra  $O(Q)$ . Altogether, the time complexity is

$$O(N + M + \Delta \log \Delta + (M + \Delta) \alpha(N) + Q),$$

with space overhead  $O(N + M + \Delta + Q)$ .

## 6 Experimental Study

We investigated the following research questions. **[RQ1]**: *How well KROMA improves state-of-the-art baselines with different LLMs?* **[RQ2]**: *How can knowledge retrieval and ontology refinement enhance matching accuracy?* and **[RQ3]**: *What are the impact of alternative LLM reasoning strategies?*

### 6.1 Experimental Setting

**Benchmark Datasets.** We selected five tracks from the OAEI campaign [40], covering various domain tracks. For each track, we selected two representative ontologies as a source ontology  $O_s$  and a target ontology  $O_t$ . The selected tracks include Anatomy [15] (Mouse-Human), Bio-LLM [23] (NCIT-DOID), CommonKG (CKG) [18] (Nell-DBpedia, YAGO-Wikidata), BioDiv [28] (ENVO-SWEET), and MSE [25] (MI-MatOnto).

**LLMs selection.** To underscore KROMA’s ability to achieve strong matching performance even with smaller or lower-accuracy LLMs, we selected models with relatively modest Chatbot Arena MMLU scores [8]: Gemma-2B (51.3%) and Llama-3.2-3B (63.4%), compared to the baseline systems Flan-T5-XXL (55.1%) and MPT-7B (60.1%). We have selected a diverse set of LLMs, ranging from ultra-lightweight to large-scale—to demonstrate KROMA’s compatibility with models that can be deployed on modest hardware without sacrificing matching quality. Our core evaluations use DeepSeek-R1-Distill-Qwen-1.5B [50] (1.5B), GPT-4o-mini [35] (8B), and Llama-3.3 [30] (70B), each chosen in a variant smaller than those employed by prior OM-LLM baselines. To further benchmark our framework, we include Gemma-2B [10] (2B), Llama-3.2-3B [2] (3B), Mistral-7B [3] (7B), and Llama-2-7B [1] (7B) in our ablation studies.

**Confidence Calibration.** Our selection of LLMs is justified by a calibration test with their confidence over ground truth answers. The self-evaluated confidence by LLMs align well with accuracy: over 80% of correct matches fall in the top confidence bins (9–10) (Figure 4), while fewer than 5% of errors are reported. Gemma-2B shows almost no errors above confidence 8, and both Llama-3.2-3B and Llama-3.3-70B maintain  $\geq 95\%$  precision at confidence thresholds of 9 or greater. We thus choose a confidence threshold to be 8.5 for all LLMs to accept their output.

**Test sets generation.** Following [22,20], for each dataset, we randomly sample 20 matched concept pairs from the ground truth mappings. For each source ontology concept, we select an additional 24 unmatched target ontology concepts based on their cosine similarity scores, thereby creating a total of 25 candidate mappings (including the ground truth mapping). Finally, we randomly choose 20 source concepts that lack a corresponding target concept in the ground truth and generate 25 candidate mappings for each. Each subset consists of 20 source ontology concepts with a match and 20 without matches, with each concept paired with 25 candidate mappings, totaling 1000 concept pairs per configuration.

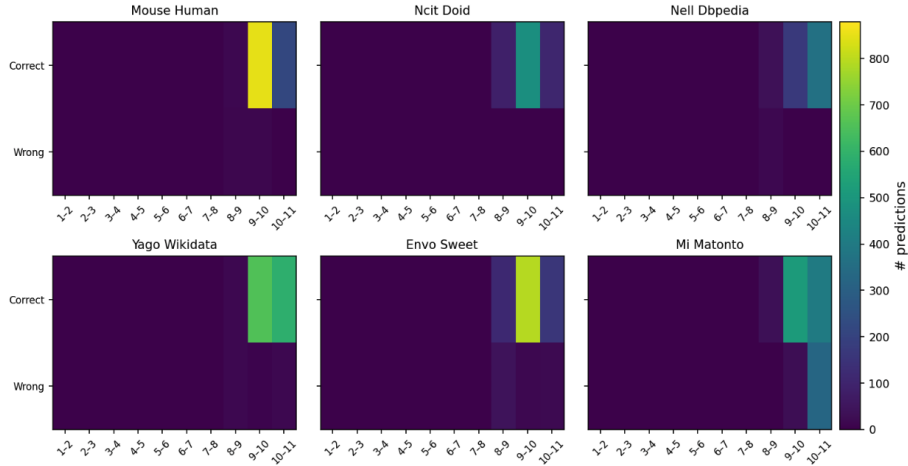


Fig. 4: LLMs Confidence Calibration

Concepts and entities not selected as test sets are treated as external knowledge base for **knowledge retrieval**. All models operate with SciBERT [4] for **Embedding Generation**, leveraging its strength in scientific data embedding.

**Baselines.** Our evaluation considers **four** state-of-the-art LLM-based ontology matching methods: LLM4OM [19], MILA [48], OLaLa [24], and LLMMap [22]. For RQ2, we also developed KROMA-NR, which skip knowledge retrieval, and KROMA-NB, which disables the bisimilarity-based clustering (hence clusters are determined by concept similarity alone).

*Naming Convention.* Each configuration uses a pattern [Method][Optional Suffix][LLM Version], where the initials (K, M, O, L) specify the method (KROMA, MILA, OLaLa, LLM4OM). Suffixes “NKR” and “NR” denote “no knowledge retrieval” and “no ontology refinement”, respectively, and the trailing version number (e.g., 3.3, 2.0, 4mini) specifies the underlying LLM release.

## 6.2 Experimental Results

**Exp-1: Effectiveness (RQ1).** In this set of tests, we evaluate the accuracy of KROMA compared with baselines, and the impact of factors such as test sizes.

*KROMA vs. Baselines: Overall Accuracy.* Across all six datasets in Table 1 (abbreviated by their first capitalized letters), the full KROMA configuration (KL3.3) achieves the highest  $F_1$  on every task, substantially outperforming competing baselines. For Mouse–Human, KL3.3 delivers 94.94  $F_1$ , eclipsing MILA’s best (ML3.1) at 92.20, OLaLa (OL2.0) at 90.20, and LLM4OM (L4G3.5) at 89.11. On NCIT–DOID, KL3.3 reaches 98.63 versus 94.80 for MILA, 83.01 for LLM4OM, and 72.10 for Flan-T5-XXL. Similar gaps appear on Nell–DBpedia (97.08 vs. 96.00/94.26), YAGO–Wikidata (95.54 vs. 92.19/93.33), ENVO–SWEET (93.18 vs. 92.98/85.06), and MI–MatOnto (61.25 vs. 59.05/32.97). These verify the effectiveness of KROMA over representative benchmark datasets. Table 9 shows details, including Precision and Recall, of the comparison.

Table 1: KROMA Performance vs. Baselines.

Model	MH	ND	NDB	YW	ES	MM
KL3.3	94.94	98.63	97.08	95.54	–	61.25
KNR3.3	91.25	95.01	94.26	91.98	–	59.95
KNB3.3	86.59	90.91	93.58	89.74	–	55.00
KL3.1	94.50	98.24	–	–	91.43	–
KL2.0	93.24	–	96.02	–	85.06	–
KG2	–	85.53	–	–	–	–
KM7	–	–	–	–	92.98	–
ML3.1	92.20	94.80	–	–	83.70	32.97
OL2.0	90.20	–	96.00	–	51.10	–
L4G3.5	89.11	83.01	94.26	–	–	–
K4mini	–	–	–	–	93.18	–
KNR4mini	–	–	–	–	92.12	–
KNB4mini	–	–	–	–	84.89	–
LLFT	–	72.10	–	–	–	–
L4L2	–	–	–	92.19	–	–
L4M7	–	–	–	–	55.09	–
L4MPT	–	–	–	–	–	32.97

*Impact of different LLMs.* Across six ontology-matching tracks, KROMA equipped with Qwen2.5-1.5B outperforms the best existing baseline on five out of the six datasets (see Table 2). In the Anatomy’s Mouse–Human track, Qwen2.5-1.5B achieves  $F_1 = 83.58$  (versus 92.20 for the OM-LLM baseline), while GPT-4o-mini reaches  $F_1 = 91.96$ . On NCIT–DOID both models surpass the baseline with  $F_1$  of 97.44 and 97.56 (baseline: 94.80), and similar gains appear on Nell–DBpedia (95.02, 95.67 vs. 96.00), YAGO–Wikidata (95.45, 95.44 vs. 92.19), and MI–MatOnto (61.25, 60.88 vs. 32.97). Only on ENVO–SWEET does the smallest model dip below the baseline (79.80 vs. 83.70), whereas GPT-4o-mini (93.18) and Llama-3.3-70B (91.95) still lead. These results confirm that KROMA with knowledge retrieval and ontology refinement can well exploit even relatively “smaller” LLMs with desirable accuracy, and further benefit from larger LLMs.

Table 2: KROMA Performance with different LLMs.

Dataset	Qwen2.5	GPT-4o-mini	Llama-3.3
Mouse–Human	83.6	92.0	<b>94.9</b>
NCIT–DOID	97.4	97.6	<b>98.6</b>
Nell–DBpedia	95.0	95.7	<b>97.1</b>
YAGO–Wikidata	95.5	95.4	<b>95.5</b>
ENVO–SWEET	79.8	<b>93.2</b>	92.0
MI–MatOnto	61.3	60.9	<b>62.2</b>

*Impact of Test sizes.* We next report the impact of ontology sizes to the performance of KROMA in accuracy (detailed results reported in Table 3). For each dataset, we varied test sizes from 200 (xsmall) to 1,000 (full) pairs. KROMA’s accuracy is in general insensitive to the change of test sizes. For example, its  $F_1$  stays within a 1.90 variance on NCIT–DOID and a 1.10-point range on YAGO–Wikidata. This verifies the robustness of KROMA in maintaining desirable accuracy for large-scale ontology matching tasks.

Table 3: KROMA Performance with different ontology sizes.

Dataset	Metric	xsmall	small	medium	large	full
NCIT-DOID	$P$	97.50	97.88	97.92	98.00	97.59
	$R$	97.50	98.12	97.50	98.13	99.69
	$F_1$	97.50	98.00	97.71	98.06	98.63
Nell-DBpedia	$P$	95.00	95.19	95.24	95.63	94.32
	$R$	96.67	98.33	99.17	99.58	100.00
	$F_1$	95.83	96.74	97.18	97.58	97.08
YAGO-Wikidata	$P$	92.00	92.50	92.67	92.75	91.80
	$R$	98.00	98.75	99.17	99.38	99.60
	$F_1$	94.90	95.58	95.89	96.04	95.54

**Exp-2: Ablation Analysis (RQ2).** In this test, we perform ablation analysis, comparing KROMA with its two variants, KROMA-NKR and KROMA-NR, removing knowledge retrieval and ontology refinement, respectively.

Table 4: KROMA Performance with/without Ontology Refinement (Left); and with/without Knowledge Retrieval (Right).

Dataset	Model	P	R	F <sub>1</sub>	Dataset	Model	P	R	F <sub>1</sub>
Mouse-Human	KL3.3	90.78	99.50	<b>94.94</b>	Mouse-Human	KL3.3	90.78	99.50	<b>94.94</b>
	KNR3.3	92.34	90.16	91.25		KNKR3.3	100.00	76.36	86.59
NCIT-DOID	KL3.3	97.59	99.69	<b>98.63</b>	NCIT-DOID	KL3.3	97.59	99.69	<b>98.63</b>
	KNR3.3	93.20	96.90	95.01		KNKR3.3	83.33	100.00	90.91
Nell-DBpedia	KL3.3	94.32	100.00	<b>97.08</b>	Nell-DBpedia	KL3.3	94.32	100.00	<b>97.08</b>
	KNR3.3	97.46	91.27	94.26		KNKR3.3	91.17	96.12	93.58
YAGO-Wikidata	KL3.3	91.80	99.60	95.54	YAGO-Wikidata	KL3.3	91.80	99.60	95.54
	KNR3.3	93.50	90.50	91.98		KNKR3.3	90.50	89.00	89.74
ENVO-SWEET	K4mini	87.23	100.00	<b>93.18</b>	ENVO-SWEET	K4mini	87.23	100.00	<b>93.18</b>
	KNR4mini	86.90	98.00	92.12		KNKR4mini	82.00	88.00	84.89
MI-MatOnto	KL3.3	45.74	92.69	<b>61.25</b>	MI-MatOnto	KL3.3	45.74	92.69	<b>61.25</b>
	KNR3.3	44.80	91.40	59.95		KNKR3.3	42.00	85.00	55.00

*Impact of Ontology Refinement.* Table 4 (Left) shows that by incorporating Ontology Refinement (KROMA vs. KROMA-NR), KROMA improves  $F_1$  score across 6 datasets: Mouse-Human (+3.69), NCIT-DOID (+3.62), Nell-DBpedia (+2.82), YAGO-Wikidata (+3.56), ENVO-SWEET (+1.06), MI-MatOnto (+1.30). By pruning false candidate pairs via offline and online refinement strategies on the concept graph, ontology refinement retains true semantic connections, by filtering out noise while preserving true alignments.

*Impact of Knowledge Retrieval.* Table 4 (Right) shows that ablating Knowledge Retrieval (KROMA-KNR vs. KROMA) results in significant  $F_1$  drop across all six benchmark datasets: Mouse-Human (-8.35), NCIT-DOID (-7.72), Nell-DBpedia (-3.50), YAGO-Wikidata (-5.80), ENVO-SWEET (-8.29), MI-MatOnto (-6.25). By enriching concepts with external, useful semantic context that are overlooked by other baselines, knowledge retrieval helps improving the accuracy.

We have also evaluated the efficiency of KROMA in terms of communication cost and performance gain from online processing. We found that KROMA is

feasible for large-scale ontology matching, and the online refinement response fast to continuous matching tasks.

**Exp-3: Efficiency (RQ3).** This tests evaluate the efficiency of KROMA. (1) Communication Cost. KROMA reduces API calls by about 26% on average: 25.9% with Qwen2.5-1.5B, 25.6% with GPT-4o-mini, and 26.4% with Llama-3.3-70B, while maintaining  $F_1$  (Table 5 (Left)). Ontologies with deeper hierarchies *e.g.*, NCIT-DOID has the largest cut (54.5%), MI-MatOnto drops  $\approx 42\%$ , and Nell-DBpedia  $\approx 26\%$ ; whereas “flatter” graphs such as Mouse-Human only shed  $\approx 3\%$ . This reflects KROMA’s ability to prune implausible pairs in rich taxonomies. (2) Offline vs. Online Refinement. We evaluate the delay time of online refinement with its offline counterpart. We found that compared with “computing from scratch” offline processing, the online refinement takes on average 9.97% of the time, with average speedup rate being  $9.03\times$  (Table 5 (Right)).

Table 5: KROMA Efficiency on Communication Cost (Left); and Offline vs. Online Refinement Delay speedup (Right).

Dataset	Qwen2.5	GPT-4o-mini	Llama-3.3	Dataset	Offline $\div$ Online
Mouse-Human	3.06%	2.88%	6.76%	Mouse-Human	$7.17\times$
NCIT-DOID	54.48%	54.63%	54.48%	NCIT-DOID	$2.33\times$
Nell-DBpedia	26.09%	25.76%	26.09%	Nell-DBpedia	$5.63\times$
YAGO-Wikidata	21.86%	21.93%	22.01%	YAGO-Wikidata	$18.76\times$
ENVO-SWEET	6.99%	6.82%	6.99%	ENVO-SWEET	$7.00\times$
MI-MatOnto	43.04%	41.66%	41.81%	MI-MatOnto	$13.31\times$
<b>Overall</b>	<b>25.92%</b>	<b>25.61%</b>	<b>26.36%</b>	<b>Overall</b>	<b><math>9.03\times</math></b>

Token size. Token footprint of KROMA remains modest even after refinement: Deep reasoning DR1D-Llama-3 (70B) issues roughly 620K total tokens (130K outputs) across all calls, versus 480K (15K outputs) for vanilla Llama-3.3 (70B), translating to an average of 605 tokens per API call compared to 470 tokens (Table 6). This  $\approx 29\%$  per-call increase incurs only a mild cost for the substantial gains in alignment quality.

Table 6: KROMA Token Usage.

Model	Total tokens	Average tokens per call
Llama-3.3-70B	480	470
DR1D-Llama-3-70B	620	605

**Exp-4: Alternative LLM Reasoning Strategies.** We evaluate KROMA’s performance with two alternative LLM reasoning strategies: “Debating” and “Deep reasoning” to understand their impact to ontology matching.

Can “Debating” help? We implemented an LLM Debate ensemble [17], where multiple agents propose alignments with their chain-of-thought and then iteratively critique one another. To bound context size, we drop 50% of historic turns and keep only each agent’s latest reply. Due to its expense, we tested this on



Table 7: KROMA Performance with/without “Debating” (Left); and with/without “Deep reasoning” (Right).

Dataset	Model	P	R	F <sub>1</sub>
Mouse–Human	D2A3R	100	70	82.35
	D2A5R	100	74	<b>85.06</b>
	D4A3R	100	66	79.52
	D4A5R	100	68	80.95
Nell–DBpedia	D2A3R	97.83	90	93.75
	D2A5R	95.91	94	<b>94.95</b>
	D4A3R	95.55	86	90.53
	D4A5R	93.88	92	92.93
YAGO–Wikidata	D2A3R	100	92	95.83
	D2A5R	100	96.9	<b>98.41</b>
	D4A3R	100	86	92.47
	D4A5R	100	90	94.73

Dataset	Model	P	R	F <sub>1</sub>
NCIT–DOID	L3.3Short	97.59	99.69	<b>98.63</b>
	L3.3Long	97.66	96.31	96.98
Nell–DBpedia	L3.3Short	94.32	100.00	<b>97.08</b>
	L3.3Long	94.67	95.48	95.07
YAGO–Wikidata	L3.3Short	91.80	99.60	<b>95.54</b>
	L3.3Long	92.31	94.80	93.53

Mouse–Human, Nell–DBpedia, and YAGO–Wikidata using four permutation of configurations: 2 or 4 agents over 3 or 5 debate rounds, denoted as  $D[\text{Number Of Agent}][\text{Number Of Round}]\text{R}$ . From Table 7 (Left), on Mouse–Human dataset, extending rounds in the 2-agent setup raised F1 from 82.35 to 85.06, whereas adding agents degraded performance (4 agents: 79.52 at 3 rounds, 80.95 at 5). Similar trends follow for Nell–DBpedia and YAGO–Wikidata. This interestingly indicates that “longer debates” help small ensembles converge to accurate matches, but larger groups introduce too much conflicting reasoning. In contrast, a single-agent, single-round KROMA pass attains  $F1 = 90.78$ , underscoring that a well-tuned solo model remains the most efficient and reliable choice.

*Can “Deep reasoning” help?* Building on DeepSeek-R1’s “Aha Moment” [12], we extended KROMA to include a forced long chain-of-thought for self-revision. Noting that vanilla DeepSeek-R1 often emits empty “<think>\n” tags (i.e. “<think>\n\n</think>”), we altered our prompt so every response must start with “<think>\n”, ensuring the model spells out its reasoning. We then compared standard Llama-3.3 (70B) (“short” reasoning, noted as L3.3Short) to DeepSeek-R1-Distill-Llama-3.3 (70B) (“long” reasoning, noted as L3.3Long), finding that the added “<think>\n” tag inflated inputs by 20% and outputs by 600% but delivered only tiny precision gains (e.g., NCIT–DOID: 97.59%  $\rightarrow$  97.66%) while slashing recall (99.69%  $\rightarrow$  96.31%), dropping F1 by 1.65–2.01 points, based on Table 7 (Right). This suggests that verbose self-reflection may improve transparency but consumes crucial context and can overly filter valid matches—especially in binary tasks, where lengthy chains of thought have been shown to harm recall, which has been consistently observed in [54].

*Can “Active Sampling” help?* In semi-supervised experiments on the Bio-LLM NCIT-DOID benchmark, KROMA paired with Llama-3.3-70B matches or exceeds the performance of MILA + Llama-3.1-70B. Under an active-learning regime with few-shot demonstrations, our method scores 1.90 F1 points higher than MILA and surpasses the previous state of the art by over 6 points. As Table 8 shows, combining bisimilarity-guided refinement with selective querying delivers a substantial boost in alignment accuracy.

Table 8: Active Learning on NCIT-DOID.

Dataset	Model	$P$	$R$	$F_1$
NCIT-DOID	KROMA + Llama-3.1 (70B)	97.85	98.15	<b>98.00</b>
	MILA + Llama-3.1 (70B)	96.70	92.80	94.67

## 7 Conclusion

We have presented KROMA, an ontology matching framework that exploits the semantic capabilities of LLMs within a bisimilar-based ontology refinement process. We show that KROMA computes a provably unique minimized structure that captures semantic equivalence relations, with efficient algorithms that can significantly reduces LLMs communication overhead, while achieving state-of-the-art accuracy across multiple OAEI benchmarks. Our evaluation has verified that LLMs, when guided by context and optimized prompting, can rival or surpass much larger models in accuracy. A future topic is to extend KROMA with more LLM reasoning strategies and ontology engineering tasks.

**Supplemental Material Statement.** Our code and experimental data has been made available at <https://anonymous.4open.science/r/kroma/>.

## Appendix

Table 9: KROMA vs Baselines.

Dataset	Model	$P$	$R$	$F_1$
Mouse-Human	KL3.3	90.78	99.50	<b>94.94</b>
	KNR3.3	92.34	90.16	91.25
	KNB3.3	100.00	76.36	86.59
	KL3.1	90.42	99.00	94.50
	KL2.0	100.00	87.34	93.24
	ML3.1	100.00	87.50	92.20
	OL2.0	91.40	89.10	90.20
	L4G3.5	90.82	87.46	89.11
NCIT-DOID	KL3.3	97.59	99.69	<b>98.63</b>
	KNR3.3	93.20	96.90	95.01
	KNB3.3	83.33	100.00	90.91
	KL3.1	97.43	99.08	98.24
	KG2	84.76	86.32	85.53
	ML3.1	96.40	93.20	94.80
	LLFT	86.10	62.00	72.10
	L4G3.5	86.19	80.06	83.01
Nell-DBpedia	KL3.3	94.32	100.00	<b>97.08</b>
	KNR3.3	97.46	91.27	94.26
	KNB3.3	91.17	96.12	93.58
	KL2.0	95.50	96.55	96.02
	OL2.0	100.00	92.20	96.00
	L4G3.5	100.00	89.14	94.26
YAGO-Wikidata	KL3.3	91.80	99.60	95.54
	KNR3.3	93.50	90.50	91.98
	KNB3.3	90.50	89.00	89.74
	KL3.2	89.09	98.00	93.33
	L4L2	100.00	85.52	92.19
ENVO-SWEET	K4mini	87.23	100.00	<b>93.18</b>
	KNR4mini	86.90	98.00	92.12
	KNB4mini	82.00	88.00	84.89
	KL3.1	87.11	96.19	91.43
	KL2.0	92.50	78.72	85.06
	KM7	87.19	99.59	92.98
	ML3.1	80.30	59.50	83.70
	OL2.0	43.10	61.30	51.10
	L4M7	59.00	51.67	55.09
MI-MatOnto	KL3.3	45.74	92.69	<b>61.25</b>
	KNR3.3	44.80	91.40	59.95
	KNB3.3	42.00	85.00	55.00
	KL3.2	56.36	62.00	59.05
	L4MPT	89.70	20.19	32.97

## References

1. AI, M.: Llama 2 7b chat model card. <https://huggingface.co/meta-llama/Llama-2-7b-chat> (2023)
2. AI, M.: Llama 3.2 3b instruct model card. <https://huggingface.co/meta-llama/Llama-3.2-3B-Instruct> (2024)
3. AI, M.: Announcing mistral 7b. <https://mistral.ai/news/announcing-mistral-7b> (2023)
4. Beltagy, I., Lo, K., Cohan, A.: Scibert: A pretrained language model for scientific text. In: EMNLP-IJCNLP (2019)
5. Borgo, S., Ferrario, R., Gangemi, A., Guarino, N., Masolo, C., Porello, D., Sanfilippo, E.M., Vieu, L.: Dolce: A descriptive ontology for linguistic and cognitive engineering. *Applied Ontology* **17**(1), 45–69 (Mar 2022)
6. Chandrasekaran, D., Mago, V.: Evolution of semantic similarity—a survey. *Acm Computing Surveys (Csur)* **54**(2), 1–37 (2021)
7. Chen, A., Song, Y., Zhu, W., Chen, K., Yang, M., Zhao, T., zhang, M.: Evaluating o1-like llms: Unlocking reasoning for translation through comprehensive analysis (2025), <https://arxiv.org/abs/2502.11544>
8. Chiang, W.L., Zheng, L., Sheng, Y., Angelopoulos, A.N., Li, T., Li, D., Zhang, H., Zhu, B., Jordan, M., Gonzalez, J.E., Stoica, I.: Chatbot arena: An open platform for evaluating llms by human preference (2024), <https://arxiv.org/abs/2403.04132>
9. Cruz, I.F., Sunna, W.: Structural Alignment Methods with Applications to Geospatial Ontologies. *Transactions in GIS* **12**(6), 683–711 (2008)
10. DeepMind, G.: Gemma 2b model card. <https://huggingface.co/google/gemma-2-2b> (2024)
11. DeepSeek-AI: Deepseek-v3 technical report (2023)
12. DeepSeek-AI: Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning (2025)
13. Degtyarenko, K., Hastings, J., Matos, P., Ennis, M.: ChEBI: An open bioinformatics and cheminformatics resource. *Current Protocols in Bioinformatics* **26**(1) (Jun 2009)
14. Dovier, A., Piazza, C., Policriti, A.: A fast bisimulation algorithm. In: Berry, G., Comon, H., Finkel, A. (eds.) *Computer Aided Verification (CAV 2001)*. Lecture Notes in Computer Science, vol. 2102 (2001)
15. Dragisic, Z., Ivanova, V., Li, H., Lambrix, P.: Experiences from the anatomy track in the ontology alignment evaluation initiative. In: *Journal of Biomedical Semantics* (2017)
16. Drobnjakovic, M., Ameri, F., Will, C., Smith, B., Jones, A.: The Industrial Ontologies Foundry (IOF) Core Ontology
17. Du, Y., Li, S., Torralba, A., Tenenbaum, J.B., Mordatch, I.: Improving factuality and reasoning in language models through multiagent debate (2023), <https://arxiv.org/abs/2305.14325>
18. Fallatah, O., Zhang, Z., Hopfgartner, F.: A gold standard dataset for large knowledge graphs matching. In: *ISWC 2020* (2020)
19. Giglou, H.B., D’Souza, J., Engel, F., Auer, S.: Llms4om: Matching ontologies with large language models (2024), <https://arxiv.org/abs/2404.10317>
20. Giglou, H.B., D’Souza, J., Engel, F., Auer, S.: Llms4om: Matching ontologies with large language models. In: *ESWC* (2024)
21. Grover, A., Leskovec, J.: node2vec: Scalable feature learning for networks (2016), <https://arxiv.org/abs/1607.00653>

22. He, Y., Chen, J., Dong, H., Horrocks, I.: Exploring large language models for ontology alignment (2023), <https://arxiv.org/abs/2309.07172>
23. He, Y., Chen, J., Dong, H., Horrocks, I.: Exploring large language models for ontology alignment. In: ISWC 2023 Posters and Demos: 22nd International Semantic Web Conference (2023)
24. Hertling, S., Paulheim, H.: Olala: Ontology matching with large language models. In: Proceedings of the 12th Knowledge Capture Conference 2023. p. 131–139 (Dec 2023). <https://doi.org/10.1145/3587259.3627571>, <http://dx.doi.org/10.1145/3587259.3627571>
25. Huschka, M., Nasr, E.: Mse benchmark (2023), <https://tinyurl.com/MSE-Benchmark>
26. Jensen, M., Colle, G.D., Kindya, S., More, C., Cox, A.P., Beverley, J.: The common core ontologies (arXiv:2404.17758) (Aug 2024)
27. Jiang, Y., Wang, X., Zheng, H.T.: A semantic similarity measure based on information distance for ontology alignment. *Information Sciences* **278**, 76–87 (2014)
28. Karam, N., Khiat, A., Algergawy, A., Sattler, M., Weiland, C., Schmidt, M.: Matching biodiversity and ecology ontologies: challenges and evaluation results. *The Knowledge Engineering Review* **35**, e9 (2020)
29. Li, X.: A survey on llm test-time compute via search: Tasks, llm profiling, search algorithms, and relevant frameworks (2025), <https://arxiv.org/abs/2501.10069>
30. Llama Team, A..M.: The llama 3 herd of models (2024)
31. Llama Team, A.: Llama 2: Open foundation and fine-tuned chat models (2023)
32. Mistral-AI: Mistral large 2 (2023)
33. Norouzi, S.S., Mahdaviinejad, M.S., Hitzler, P.: Conversational ontology alignment with chatgpt (2023), <https://arxiv.org/abs/2308.09217>
34. OpenAI: Language models are few-shot learners. In: NeurIPS 2020 (2020)
35. OpenAI: Gpt-4 technical report (2024)
36. OpenAI: Openai o1 system card (2024), <https://arxiv.org/abs/2412.16720>
37. Otte, J.N., Beverley, J., Ruttenberg, A.: Bfo: Basic formal ontology. *Applied Ontology* **17**(1), 17–43 (Jan 2022)
38. Peeters, R., Bizer, C.: Using chatgpt for entity matching (2023), <https://arxiv.org/abs/2305.03423>
39. Pirró, G., Euzenat, J.: A feature and information theoretic framework for semantic similarity and relatedness. In: ISWC (2010)
40. Qiang, Z., Taylor, K., Wang, W., Jiang, J.: Oaei-llm: A benchmark dataset for understanding large language model hallucinations in ontology matching. arXiv preprint arXiv:2409.14038 (2024)
41. Qin, L., Chen, Q., Zhou, Y., Chen, Z., Li, Y., Liao, L., Li, M., Che, W., Yu, P.S.: Multilingual large language model: A survey of resources, taxonomy and frontiers (2024), <https://arxiv.org/abs/2404.04925>
42. Qwen, Yang, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Li, C., Liu, D., Huang, F., Wei, H., Lin, H., Yang, J., Tu, J., Zhang, J., Yang, J., Yang, J., Zhou, J., Lin, J., Dang, K., Lu, K., Bao, K., Yang, K., Yu, L., Li, M., Xue, M., Zhang, P., Zhu, Q., Men, R., Lin, R., Li, T., Tang, T., Xia, T., Ren, X., Ren, X., Fan, Y., Su, Y., Zhang, Y., Wan, Y., Liu, Y., Cui, Z., Zhang, Z., Qiu, Z.: Qwen2.5 technical report (2025), <https://arxiv.org/abs/2412.15115>
43. Radford, A., Narasimhan, K., Salimans, T., Sutskever, I.: Improving language understanding by generative pre-training (2018)

44. Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I.: Language models are unsupervised multitask learners (2019)
45. Rajamohan, B.P., Bradley, A.C.H., Tran, V.D., Gordon, J.E., Caldwell, H.W., Mehdi, R., Ponon, G., Tran, Q.D., Dernek, O., Kaltenbaugh, J., Pierce, B.G., Wieser, R., Yue, W., Lin, K., Kambo, J., Lopez, C., Nihar, A., Savage, D.J., Brown, D.W., Sharma, H., Giera, B., Tripathi, P.K., Wu, Y., Li, M., Davis, K.O., Bruckman, L.S., Barcelos, E.I., French, R.H.: Materials data science ontology(mds-onto): Unifying domain knowledge in materials and applied data science. *Scientific Data* **12**(1), 628 (Apr 2025)
46. Sun, J., Zheng, C., Xie, E., Liu, Z., Chu, R., Qiu, J., Xu, J., Ding, M., Li, H., Geng, M., Wu, Y., Wang, W., Chen, J., Yin, Z., Ren, X., Fu, J., He, J., Yuan, W., Liu, Q., Liu, X., Li, Y., Dong, H., Cheng, Y., Zhang, M., Heng, P.A., Dai, J., Luo, P., Wang, J., Wen, J.R., Qiu, X., Guo, Y., Xiong, H., Liu, Q., Li, Z.: A survey of reasoning with foundation models (2024), <https://arxiv.org/abs/2312.11562>
47. Sutskever, I., Vinyals, O., Le, Q.V.: Sequence to sequence learning with neural networks. In: *NeurIPS 2014* (2014)
48. Taboada, M., Martinez, D., Arideh, M., Mosquera, R.: Ontology matching with large language models and prioritized depth-first search (2025), <https://arxiv.org/abs/2501.11441>
49. Team, K.: Kimi k1.5: Scaling reinforcement learning with llms (2025), <https://arxiv.org/abs/2501.12599>
50. Team, Q.: Introducing qwen1.5 (2024), <https://qwenlm.github.io/blog/qwen1.5/>
51. Trojahn, C., Vieira, R., Schmidt, D., Pease, A., Guizzardi, G.: Foundational ontologies meet ontology matching: A survey. *Semantic Web* **13**(4), 685–704 (2022)
52. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: *(NeurIPS 2017)*
53. Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., Yogatama, D., Bosma, M., Zhou, D., Metzler, D., Chi, E.H., Hashimoto, T., Vinyals, O., Liang, P., Dean, J., Fedus, W.: Emergent abilities of large language models (2022), <https://arxiv.org/abs/2206.07682>
54. Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q., Zhou, D.: Chain-of-thought prompting elicits reasoning in large language models (2023), <https://arxiv.org/abs/2201.11903>
55. Xu, Y., Hu, L., Zhao, J., Qiu, Z., Xu, K., Ye, Y., Gu, H.: A survey on multilingual large language models: corpora, alignment, and bias. *Frontiers of Computer Science* **19**(11) (2025)