# 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,eib14,rxf131,yxw1650}@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 performance 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. Through experiments on multiple benchmark datasets, we show that integrating knowledge retrieval with context-augmented LLMs significantly enhances ontology matching—outperforming both classic OM systems and cutting-edge LLM-based approaches—while keeping communication overhead comparable. 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. Our code and experimental dataset has been made available at: https://github.com/lamng3/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) [35] or DOLCE [5] down to more domain-oriented 'mid-level' ones, such as CheBi [14] in chemistry, Industrial Ontology Foundy (IOF) [17] 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) [42]. To achieve broad usability, ontologies need to be effectively aligned for better interoperability using ontology matching.

Ontology matching has been studied to find correspondence between terms that are semantically equivalent [50]. 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 [10], 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 expensive 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 using direct prompting (which may result in low confidence and risk of hallucination), KROMA maintains a set of conceptually similar groups that are co-determined by concept similarity and language models, both guided by their "augmented 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 technical 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 ontology matching framework leveraging the power of LLMs (**Section 4**). KROMA fine-tunes LLMs with prompts over enriched semantic contexts. Such contextual information are obtained from the knowledge retrieval process, 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) Offline refinement algorithm performs a fast grouping process guided by the bisimilar equivalence property, blending the concept equivalence tests co-determined by semantic closeness and language models.

(b) Online refinement 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 methods using large languge models by 10.95% on average, and the optimization of utilizing knowledge retrieval and refinement processes improves KROMA's accuracy by 6.65% and 2.68%, respectively.

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

*Large Language Models*. Large language models (LLMs) have advanced NLP by enabling parallel processing and capturing complex dependencies [44,51], which have scaled from GPT-1's 117M [40], GPT-2's 1.5B [41] to GPT-3's 175B [6] and GPT-4's 1.8T parameters [34]. Open-source models like Llama have grown to 405B [49], with Mistral Large (123B) [30] and DeepSeek V3 (671B) [13] 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 [8,12,29,33,38,43,46,54].

*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, [36] frame product matching as a yes/no query, and [31] 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 [20], while MILA adds a prioritized depth-first search step to confirm high-confidence matches before any LLM invocation [45]); 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 LLMap uses binary yes/no prompts over concept labels plus structural context with Flan-T5-XXL or GPT-3.5 [23]).

## 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 between the concepts. An ontology has a graph representation with a set of concept nodes $C$, and a set of edges $E$. In this paper, 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)

| Notation | Description |
|---|---|
| $\mathcal{O} = (C, E)$ | ontology $\mathcal{O}$, $C$: set of concepts, $E$: set of relations |
| $|\mathcal{O}|$ | size of ontology $\mathcal{O}$; $|\mathcal{O}| = |C| + |E|$ |
| $\mathrm{r}(c)$ | rank of concept node $c$ |
| $c.I$ | ground set of concept node $c$ |
| $\mathcal{O}_s = (C_s, E_s), \mathcal{O}_t = (C_t, E_t)$ | source and target ontology, respectively |
| $R_\simeq$ | ontological equivalence relation |
| $\mathcal{C}$ | equivalence partition of concept set $C_s \cup C_t$ |
| $\mathcal{G}_O = (V_\mathcal{O}, E_\mathcal{O})$ | concept graph $\mathcal{G}_O$, $V_O$: set of nodes, $E_O$: set of edges |
| $\Delta\mathcal{G}$ | newly arrive edges; edge updates in $\mathcal{G}$ |
| $[c] \in V_O$ | an equivalence class in $\mathcal{C}$ |
| $M$ | a language model |
| $\mathcal{M}$ | set of Large Language Model(s) |
| $q$ | a prompt query |
| $q(M)$ | a natural language answer from language model $M$ |
| $F(c, c')$ | concept similarity between two concepts $c$ and $c'$ |
| $q(c, c', \mathcal{M})$ | a natural language answer to $\mathcal{M}$ asking if $c \simeq c'$ |
| $\alpha \in (0, 1]$ | threshold for asserting concept similarity |
| $W = (O_s, O_t, F, \mathcal{M}, \alpha)$ | configuration input for ontology matching |
| $z_c$ | embedding of concept $c$ |
| $\mathbb{S}$ | set of candidate concept pairs with similarity scores |
| $\mathbb{C}$ | top-k pairs with highest similarity scores |

Table 1: Summary of Notations.

A *ground set* $c.\mathcal{I}$, refers to a set of auxilary entities from *e.g.*, external ontologies or knowledge bases, that can be validated to belong to the concept $c$.

**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* $R_\simeq \subseteq C_s \times C_t$ is defined as an equivalence relation that contains pairs of nodes $(c, c')$ that are considered to be "semantically equivalent". The relation $R_\simeq$ 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 finite 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 formulate $R_\simeq$ and $\mathcal{C}$ induced by $R_\simeq$ and compute the concept graph $\mathcal{G}_O$ accordingly.

*Example 1.* Figure 1 depicts two ontologies $O_s$ and $O_t$ as DAGs with 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 a result of ontology matching, with three equivalence classes.
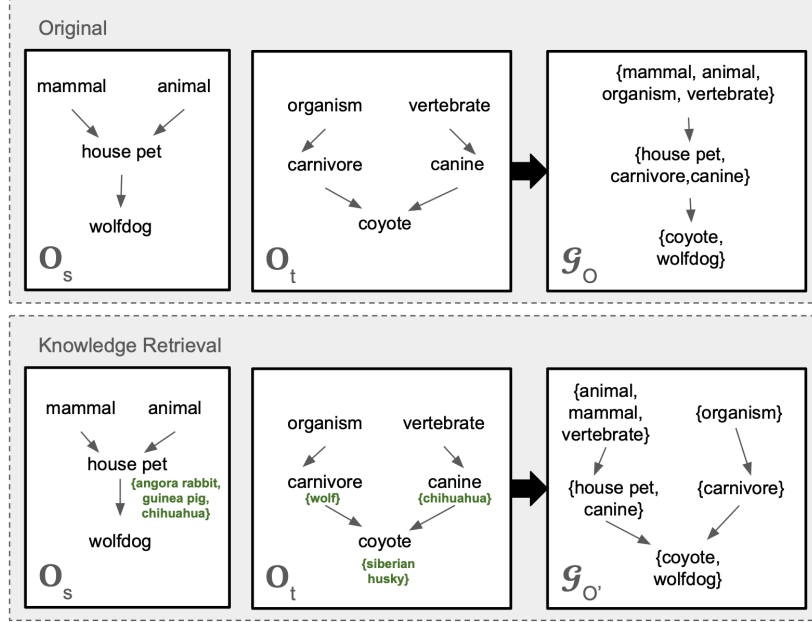
Fig. 1: Ontologies with ground sets, Ontology Matching and Concept Graphs

**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 [52].

*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 a large language model $M$ 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 determine whether two *concepts* are semantically equivalent [7]. KROMA by default uses a Boolean function $F$ defined by a weighted combination of a semantic closeness metric $\mathsf{sim}$[1] and the result from a set of LLMs $\mathcal{M}$ (see Section 4).

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

where $q$ is a prompt query that specifies the context of concept equivalence for LLMs, and $\gamma$ be a configurable parameter. KROMA supports a built-in library of semantic similarity functions $\mathsf{sim}$, including (a) string similarity, feature and information measure [37], or normalized distances (NGDs) [27]; and (b) a variety of LLMs such as GPT-4o Mini [34], LLaMA-3.3 [21], and Qwen-2.5 [39].

**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 conceptually 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 $(c_s'', c_t'') \in R_{\simeq}$.

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

We can prove the above results by verifying that $R_{\simeq}$ 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 entities that are conceptually similar may *not* be ontologically bisimilar. On the other hand, two ontologically bisimilar entities must be conceptually similar, following the 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 determines concept similarity, and a set of LLMs $\mathcal{M}$, the problem becomes computing a smallest concept graph $\mathcal{G}_O$ induced by the ontologically bisimilar equivalence $R_{\simeq}$.

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

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

---

[1] 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.
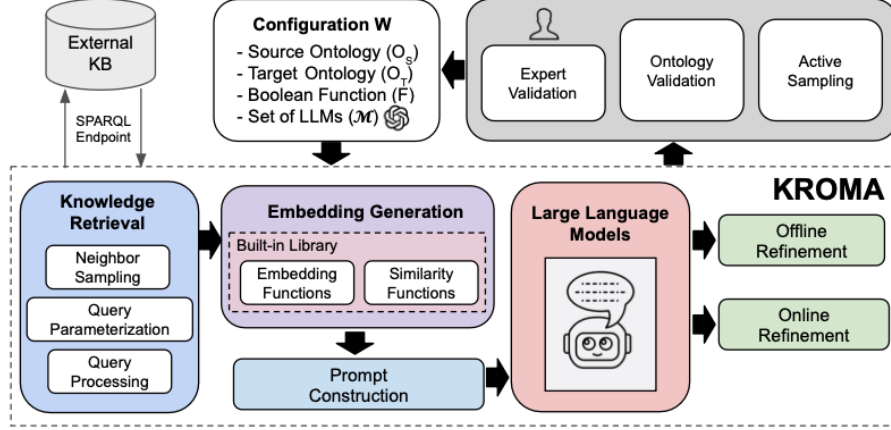
Fig. 2: KROMA Framework Overview: Major Components and Dataflow

**Proof sketch:** We show that the above result holds by verifying the following. (1) There is a unique, maximum ontological bisimilar relation $R_\simeq$ for a given configuration $W = (O_s, O_t, F, \mathcal{M}, \alpha)$, 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 $R_\simeq$ 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 $R_\simeq$ 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 $R_\simeq$ 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 $R_\simeq$, they are isomorphic to each other. □

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 the aforementioned 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 large language models $\mathcal{M}$, KROMA has the following major functional modules that enables a robust and effective multi-session ontology matching process.

*Concept Graph Initialization* Upon receiving two ontologies $O_s = (C_s, E_s)$ and $O_t = (C_t, E_t)$ as the inputs, KROMA initializes the concept graph $\mathcal{G}_O = (V, E)$ with $V = C_S \cup C_T$ and $E = E_S \cup E_T$. For each concept $c \in V$, KROMA then globally computes the structural *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*: Apply the generated SPARQL queries $S_q$ onto external knowledge bases to augment the concept's ground set $c.\mathcal{I}$ with auxiliary information (*e.g.,* relevant entities, definition, labels, synonyms, etc.).

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

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

where $\mathsf{embd}_{graph}$ (e.g. node2vec [22]) captures the $c$'s topology information and $\mathsf{embd}_{text}$ (e.g. SciBERT [4]) encodes $c$'s textual context. After obtaining the necessary embeddings ($z_c \; \forall c \in V$), KROMA then computes pairwise concept similarity between source and target ontologies using the $\mathsf{sim}$ functions:

$$\mathbb{S} = \{(c_s, c_t, score_{s,t}) \mid c_s \in C_S, c_t \in C_T, score_{s,t} = \mathsf{sim}(\mathsf{z}_{\mathsf{c}_\mathsf{s}}, \mathsf{z}_{\mathsf{c}_\mathsf{t}})\}$$

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

*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 $KG$. 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 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 conceptually similar, allowing the candidate pair to be "double checked" by LLMs in the following phase.

*Prompt Querying LLMs.* For each candidate pair $(c_s, c_t) \in \mathbb{C}$, KROMA generates an NL prompt that includes: (1) Task description $\mathcal{T}$ (e.g. "*Given two ontology concepts and their 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."). KROMA then calls (a set of) LLMs $\mathcal{M}$ to obtain a matching 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

Fig. 3: KROMA Prompt Query Template and Query Example

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, with an aim to minimalize the manual effort needed.

*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 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 {mammal, animal, 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.

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

---

**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** $\quad$ 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** $\quad$ $B_i \leftarrow \{c : r(c) = i\}$;

**8** $P \leftarrow \{B_0, \ldots, B_\rho\}$;
**9 for** $i \leftarrow 0$ **to** $\rho$ **do**
**10** $\quad$ $D_i \leftarrow \{X \in P : X \subseteq B_i\}$;
**11** $\quad$ **foreach** $X \in D_i$ **do**
**12** $\quad\quad$ $G \leftarrow \mathsf{collapse}(G, X)$;

**13** $\quad$ **foreach** $c \in B_i$ **do**
**14** $\quad\quad$ **foreach** $C \in P$ *with* $C \subseteq \bigcup_{j>i} B_j$ **do**
**15** $\quad\quad\quad$ Split $C$ into $C_1, C_2$ by adjacency to $c$;
**16** $\quad\quad\quad$ $P \leftarrow (P \setminus \{C\}) \cup \{C_1, C_2\}$;

**17 return** $\mathcal{G}_O$;

---

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) [15], 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 knowl-

edge 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 (with pseudoscope reported in [1]) 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 to resolve. These cases are cached into a query set $\mathcal{Q}$ to be further resolved in the validation phase (see Section 4). We cache these cases into a query set by using an auxiliary data structure (*e.g.,* using a priority queue ranked by LLM confidence scores) to effectively manage their processing.

*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 $O(|\mathcal{G}_O| + |\Delta G| \log |\Delta G| + |\Delta G| \log |V_O|)$ time. This result verifies that online refinement is able to response faster than offline maintenance that recomputes the concept graph from scratch. We present detailed analysis in [1].

## 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 performance?* and **[RQ3]**: *What are the impact of using alternative* LLM *reasoning strategies?*

### 6.1    Experimental Setting

**Benchmark Datasets**. We selected five tracks from the OAEI campaign [32], 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 [16] (Mouse-Human), Bio-LLM [23] (NCIT-DOID), CommonKG

(CKG) [19] (Nell-DBpedia, YAGO-Wikidata), BioDiv [28] (ENVO-SWEET), and MSE [25] (MI-MatOnto). To ensure a fair and comprehensive evaluation of KROMA, we adopted the standard benchmarks from the Ontology Alignment Evaluation Initiative (OAEI), enabling direct comparison with prior LLM-based methods. Despite the high cost of LLM inference, we tested KROMA across five diverse tracks to demonstrate its robustness and effectiveness across domains.

**LLMs selection.** To underscore KROMA 's ability to achieve strong matching performance even with smaller or lower-performance LLMs, we selected models with relatively modest Chatbot Arena MMLU scores [9]: 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 [48] (1.5B), GPT-4o-mini [34] (8B), and Llama-3.3 [21] (70B), each chosen in a variant smaller than those employed by prior OM-LLM baselines. To further benchmark our framework, we include Gemma-2B [11] (2B), Llama-3.2-3B [47] (3B), Mistral-7B [3] (7B), and Llama-2-7B [2] (7B) in our ablation studies. To run inference on the aforementioned models, we used TogetherAI and OpenAI APIs, treating them as off-the-shelf inference services on hosted, pretrained models. We are not performing any fine-tuning or weight updates to the models.

*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 performance: over 80% of correct matches fall in the top confidence bins (9–10), 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 [20,23], for each dataset, we arbitrarily designate one concept as the "source" for sampling and its target counterpart. We remark that the source and target roles are interchangeable w.l.o.g given our theoretical analysis, algorithms and test results. 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.

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 [20], MILA [45], OLaLa [24], and LLMap [23]. For RQ2, we also developed KROMA-NR, which skips the knowledge retrieval process, and KROMA-NB, which disables the use of bisimilarity-based clustering, allowing only the use of concept similarity to determine concept clusters.

*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 performance of KROMA compared with baselines, and the impact of factors such as test sizes.

Table 2: 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 | – |
| LLFT | – | 72.10 | – | – | – | – |
| L4L2 | – | – | – | 92.19 | – | – |
| L4M7 | – | – | – | – | 55.09 | – |
| L4MPT | – | – | – | – | – | 32.97 |

KROMA *vs. Baselines: Overall Performance.* Across all six datasets in Table 2 (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 cumulative results verify the effectiveness of KROMA over representative benchmark datasets. More details (*e.g.,* Precision and Recall) are reported in [1].

*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 3). 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 show that KROMA (with both knowledge retrieval and ontology refinement) delivers strong performance even on relatively "small" LLMs, and scales further with larger LLMs.

Table 3: KROMA Performance with different LLMs.

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

*Impact of Test sizes.* We next report the impact of ontology sizes to the performance of KROMA in performance (detailed results reported in [1]). For each dataset, we varied test sizes from 200 (xsmall) to 1,000 (full) pairs. KROMA's performance 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 and effectiveness of KROMA in maintaining desirable performance for large-scale ontology matching tasks.

**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 under different scenarios: with/without Ontology Refinement (Left); and with/without Knowledge Retrieval (Right).

| Dataset | Model | **P** | **R** | **F$_1$** | Dataset | Model | **P** | **R** | **F$_1$** |
|---|---|---|---|---|---|---|---|---|---|
| Mouse–Human | KL3.3 | 90.78 | 99.50 | **94.94** | Mouse–Human | KL3.3 | 90.78 | 99.50 | **94.94** |
| | KNR3.3 | 92.34 | 90.16 | 91.25 | | KNKR3.3 | 100.00 | 76.36 | 86.59 |
| NCIT-DOID | KL3.3 | 97.59 | 99.69 | **98.63** | NCIT-DOID | KL3.3 | 97.59 | 99.69 | **98.63** |
| | KNR3.3 | 93.20 | 96.90 | 95.01 | | KNKR3.3 | 83.33 | 100.00 | 90.91 |
| Nell-DBpedia | KL3.3 | 94.32 | 100.00 | **97.08** | Nell-DBpedia | KL3.3 | 94.32 | 100.00 | **97.08** |
| | 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 | **93.18** | ENVO-SWEET | K4mini | 87.23 | 100.00 | **93.18** |
| | KNR4mini | 86.90 | 98.00 | 92.12 | | KNKR4mini | 82.00 | 88.00 | 84.89 |
| MI-MatOnto | KL3.3 | 45.74 | 92.69 | **61.25** | MI-MatOnto | KL3.3 | 45.74 | 92.69 | **61.25** |
| | 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 using offline and online strategies on the concept graph, KROMA is able to retain true semantical connections between concepts.

*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 performance.

**Exp-3: Alternative LLM Reasoning Strategies**. We evaluated KROMA's performance with two alternative LLM reasoning strategies: "Debating" and "Deep reasoning" to further understand their impact to ontology matching.

Table 5: KROMA Performance with/without "Debating" (Left); and with/without "Deep reasoning" (Right).

| Dataset | Model | P | R | $F_1$ |
|---|---|---|---|---|
| Mouse–Human | D2A3R | 100 | 70 | 82.35 |
| | D2A5R | 100 | 74 | **85.06** |
| | D4A3R | 100 | 66 | 79.52 |
| | D4A5R | 100 | 68 | 80.95 |
| Nell–DBpedia | D2A3R | 97.83 | 90 | 93.75 |
| | D2A5R | 95.91 | 94 | **94.95** |
| | D4A3R | 95.55 | 86 | 90.53 |
| | D4A5R | 93.88 | 92 | 92.93 |
| YAGO–Wikidata | D2A3R | 100 | 92 | 95.83 |
| | D2A5R | 100 | 96.9 | **98.41** |
| | D4A3R | 100 | 86 | 92.47 |
| | D4A5R | 100 | 90 | 94.73 |

| Dataset | Model | P | R | $F_1$ |
|---|---|---|---|---|
| NCIT–DOID | L3.3Short | 97.59 | 99.69 | **98.63** |
| | L3.3Long | 97.66 | 96.31 | 96.98 |
| Nell–DBpedia | L3.3Short | 94.32 | 100.00 | **97.08** |
| | L3.3Long | 94.67 | 95.48 | 95.07 |
| YAGO–Wikidata | L3.3Short | 91.80 | 99.60 | **95.54** |
| | L3.3Long | 92.31 | 94.80 | 93.53 |

*Can "Debating" help?* We implemented an LLM Debate ensemble [18], 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 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[Number Of Agent]A[Number Of Round]R. From Table 5 (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% → 97.66%) while slashing recall (99.69% → 96.31%), dropping F1 by 1.65–2.01 points, based on Table 5 (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 [53].

*Can "Active Sampling" help?* In the 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 6 shows, combining bisimilarity-guided refinement with selective querying delivers a substantial boost in KROMA's alignment performance.

Table 6: Active Learning on NCIT-DOID.

| Dataset | Model | $P$ | $R$ | $F_1$ |
|---|---|---|---|---|
| NCIT-DOID | KROMA + Llama-3.1-70B | 97.85 | 98.15 | **98.00** |
| | 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 performance 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 performance. A future topic is to extend KROMA with more LLM reasoning strategies and ontology engineering tasks.

**Supplemental Material Statement**. Our code and experimental dataset has been made available at https://github.com/lamng3/kroma, including an extended version of the paper [1], providing further details for KROMA.

# References

1. Full version (2025), https://github.com/lamng3/kroma/full.pdf
2. AI, M.: Llama 2 7B Chat Model Card. Technical Blog; https://www.llama.com/llama2/ (2023)
3. AI, M.: Mistral 7b. Technical Blog; https://mistral.ai/news/announcing-mistral-7b (2023)
4. Beltagy, I., Lo, K. and Cohan, A.: SciBERT: A Pretrained Language Model for Scientific Text. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing. pp. 3615–3620 (2019)
5. Borgo, S., Ferrario, R., Gangemi, A., Guarino, N., Masolo, C. et al.: DOLCE: A Descriptive Ontology for Linguistic and Cognitive Engineering. Applied Ontology **17**, 45–69 (2022)
6. Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D. et al.: Language Models are Few-Shot Learners. In: Proceedings of the 33rd Neural Information Processing Systems. pp. 1877–1901 (2020)
7. Chandrasekaran, D. and Mago, V.: Evolution of Semantic Similarity—A Survey. ACM Computing Surveys **54**, 1–37 (2021)
8. Chen, A., Song, Y., Zhu, W., Chen, K., Yang, M., Zhao, T. and zhang, M.: Evaluating o1-like llms: Unlocking reasoning for translation through comprehensive analysis (2025)
9. Chiang, W.L., Zheng, L., Sheng, Y., Angelopoulos, A.N., Li, T. et al.: Chatbot Arena: An Open Platform for Evaluating LLMs by Human Preference. In: Proceedings of the 41st International Conference on Machine Learning. pp. 1–13 (2024)
10. Cruz, I.F. and Sunna, W.: Structural Alignment Methods with Applications to Geospatial Ontologies. Transactions in GIS **12**, 683–711 (2008)
11. DeepMind, G.: Gemma 2b model card. Technical Blog; https://blog.google/technology/developers/gemma-open-models/ (2024)
12. DeepSeek-AI, Guo, D., Yang, D., Zhang, H., Song, J. et al.: Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. Online; https://arxiv.org/abs/2501.12948 (2025)
13. DeepSeek-AI, Liu, A., Feng, B., Xue, B., Wang, B. et al.: DeepSeek-V3 Technical Report. Technical Report; https://arxiv.org/abs/2412.19437 (2023)
14. Degtyarenko, K., Hastings, J., Matos, P. and Ennis, M.: ChEBI: An Open Bioinformatics and Cheminformatics Resource. Current Protocols in Bioinformatics **14**, 14.9.1–14.9.20 (2009)
15. Dovier, A., Piazza, C. and Policriti, A.: A Fast Bisimulation Algorithm. In: Berry, G., Comon, H. and Finkel, A. (eds.) Proceedings of the 13th International Conference on Computer Aided Verification. vol. 2102, pp. 423–437 (2001)
16. Dragisic, Z., Ivanova, V., Li, H. and Lambrix, P.: Experiences from the Anatomy Track in the Ontology Alignment Evaluation Initiative. Journal of Biomedical Semantics **8**, 56 (2017)
17. Drobnjakovic, M., Ameri, F., Will, C., Smith, B. and Jones, A.: The Industrial Ontologies Foundry (IOF) Core Ontology. In: Proceedings of the 12th International Workshop on Formal Ontologies Meet Industry (2022)
18. Du, Y., Li, S., Torralba, A., Tenenbaum, J.B. and Mordatch, I.: Improving factuality and reasoning in language models through multiagent debate (2023)
19. Fallatah, O., Zhang, Z. and Hopfgartner, F.: A gold standard dataset for large knowledge graphs matching. In: Proceedings of the 19th International Semantic Web Conference (2020)

20. Giglou, H.B., D'Souza, J., Engel, F. and Auer, S.: LLMs4OM: Matching Ontologies with Large Language Models. In: Proceedings of the 21st European Semantic Web Conference (2024)
21. Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A. et al.: The Llama 3 Herd of Models. Technical Report; https://arxiv.org/abs/2407.21783 (2024)
22. Grover, A. and Leskovec, J.: node2vec: Scalable Feature Learning for Networks. In: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. pp. 855–864 (2016)
23. He, Y., Chen, J., Dong, H. and Horrocks, I.: Exploring large language models for ontology alignment. In: Proceedings of the Posters and Demos Track of the 22nd International Semantic Web Conference (2023)
24. Hertling, S. and Paulheim, H.: Olala: Ontology matching with large language models. In: Proceedings of the 12th Knowledge Capture Conference. p. 131–139 (2023)
25. Huschka, M. and Nasr, E.: Evaluation of Automatic Ontology Matching for Materials Sciences and Engineering. Master's Thesis; https://ad-publications.cs.uni-freiburg.de/theses/Master_Engy_Nasr_2020.pdf (2020)
26. Jensen, M., Colle, G.D., Kindya, S., More, C., Cox, A.P. and Beverley, J.: The Common Core Ontologies. In: Frontiers in Artificial Intelligence and Applications (2024)
27. Jiang, Y., Wang, X. and 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. and 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. Transactions on Machine Learning Research (2025)
30. Mistral-AI: Mistral Large 2. Technical Blog; https://mistral.ai/news/mistral-large-2407 (2023)
31. Norouzi, S.S., Mahdavinejad, M.S. and Hitzler, P.: Conversational Ontology Alignment with ChatGPT. In: Proceedings of the 18th International Workshop on Ontology Matching collocated with the 22nd International Semantic Web Conference ISWC (2023)
32. OAEI: Ontology Alignment Evaluation Initiative. https://oaei.ontologymatching.org (2024)
33. OpenAI, :, Jaech, A., Kalai, A., Lerer, A. et al.: Openai o1 system card. Technical Report; https://arxiv.org/abs/2412.16720 (2024)
34. OpenAI: GPT-4 Technical Report. Technical Report; https://cdn.openai.com/papers/gpt-4.pdf (2024)
35. Otte, J.N., Beverley, J. and Ruttenberg, A.: Bfo: Basic formal ontology. Applied ontology **17**, 17–43 (2022)
36. Peeters, R. and Bizer, C.: Using ChatGPT for Entity Matching. In: Proceedings of the 27th International Conference on Advances in Databases and Information Systems (2023)
37. Pirrò, G. and Euzenat, J.: A Feature and Information Theoretic Framework for Semantic Similarity and Relatedness. In: Proceedings of the 9th International Semantic Web Conference. Lecture Notes in Computer Science, vol. 6496, pp. 615–630 (2010)

38. Qin, L., Chen, Q., Zhou, Y., Chen, Z., Li, Y. et al.: Multilingual large language model: A survey of resources, taxonomy and frontiers (2024)
39. Qwen, Yang, A., Yang, B., Zhang, B., Hui, B. et al.: Qwen2.5 technical report. Technical Report; https://arxiv.org/abs/2412.15115 (2025)
40. Radford, A., Narasimhan, K., Salimans, T. and Sutskever, I.: Improving Language Understanding by Generative Pre-Training. Technical Report; https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf (2018)
41. Radford, A., Wu, J., Child, R., Luan, D., Amodei, D. and Sutskever, I.: Language Models are Unsupervised Multitask Learners. Technical Report; https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf (2019)
42. Rajamohan, B.P., Bradley, A.C.H., Tran, V.D., Gordon, J.E., Caldwell, H.W. et al.: Materials Data Science Ontology (MDS-Onto): Unifying Domain Knowledge in Materials and Applied Data Science. Scientific Data **12**,  628 (2025)
43. Sun, J., Zheng, C., Xie, E., Liu, Z., Chu, R. et al.: A Survey of Reasoning with Foundation Models (2024)
44. Sutskever, I., Vinyals, O. and Le, Q.V.: Sequence to Sequence Learning with Neural Networks. In: Proceedings of the 28th Conference on Neural Information Processing Systems. pp. 3104–3112 (2014)
45. Taboada, M., Martinez, D., Arideh, M. and Mosquera, R.: Ontology matching with large language models and prioritized depth-first search. Information Fusion **123**, 103254 (2025)
46. Team, K., Du, A., Gao, B., Xing, B., Jiang, C. et al.: Kimi k1.5: Scaling reinforcement learning with llms. Technical Report; https://arxiv.org/abs/2501.12599 (2025)
47. Team, M.A.: Llama 3.2 3b instruct model card. https://huggingface.co/meta-llama/Llama-3.2-3B-Instruct (2024)
48. Team, Q.: Introducing qwen1.5. Technical Report; https://qwenlm.github.io/blog/qwen1.5 (2024)
49. Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A. et al.: Llama 2: Open Foundation and Fine-Tuned Chat Models. Technical Report; https://arxiv.org/abs/2307.09288 (2023)
50. Trojahn, C., Vieira, R., Schmidt, D., Pease, A. and Guizzardi, G.: Foundational Ontologies Meet Ontology Matching: A Survey. Semantic Web **13**, 685–704 (2022)
51. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L. et al.: Attention Is All You Need. In: Proceedings of the 31st Conference on Neural Information Processing Systems. pp. 6000–6010 (2017)
52. Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B. et al.: Emergent Abilities of Large Language Models. Transactions on Machine Learning Research **1**, 1–21 (2022)
53. Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B. et al.: Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In: Proceedings of the 36th Conference on Neural Information Processing Systems (2022)
54. Xu, Y., Hu, L., Zhao, J., Qiu, Z., Xu, K., Ye, Y. and Gu, H.: A Survey on Multilingual Large Language Models: Corpora, Alignment, and Bias. Frontiers of Computer Science **19**, 1–23 (2025)