

## Highlights

### **Ontology Matching with Large Language Models and Prioritized Depth-First Search**

Maria Taboada, Diego Martinez, Mohammed Arideh, Rosa Mosquera

- Propose a novel retrieve-identify-prompt pipeline for ontology matching.
- Achieve the highest F-Measure score in five of the seven tasks in the unsupervised setting.
- Exhibit task-agnostic high performance, remaining stable across all tasks and settings.
- Significantly reduce the number of requests to the LLM.

# Ontology Matching with Large Language Models and Prioritized Depth-First Search

Maria Taboada<sup>a,\*</sup>, Diego Martinez<sup>b</sup>, Mohammed Arideh<sup>a</sup>, Rosa Mosquera<sup>c</sup>

<sup>a</sup>Department of Electronics and Computer Science, University of Santiago de Compostela, Santiago de Compostela, 15701, Spain

<sup>b</sup>Department of Applied Physics, University of Santiago de Compostela, , Santiago de Compostela, 15701, Spain

<sup>c</sup>Department of Functional Biology, University of Santiago de Compostela, , Santiago de Compostela, 15701, Spain

---

## Abstract

Ontology matching (OM) plays a key role in enabling data interoperability and knowledge sharing. Recently, methods based on Large Language Model (LLMs) have shown great promise in OM, particularly through the use of a *retrieve-then-prompt* pipeline. In this approach, relevant target entities are first retrieved and then used to prompt the LLM to predict the final matches. Despite their potential, these systems still present limited performance and high computational overhead. To address these issues, we introduce MILA, a novel approach that embeds a *retrieve-identify-prompt* pipeline within a prioritized depth-first search (PDFS) strategy. This approach efficiently identifies a large number of semantic correspondences with high accuracy, limiting LLM requests to only the most borderline cases. We evaluated MILA using three challenges from the 2024 edition of the Ontology Alignment Evaluation Initiative. Our method achieved the highest F-Measure in five of seven unsupervised tasks, outperforming state-of-the-art OM systems by up to 17%. It also performed better than or comparable to the leading supervised OM systems. MILA further exhibited task-agnostic performance, remaining stable across all tasks and settings, while significantly reducing runtime. These findings highlight that high-performance LLM-based OM can be achieved through a combination of programmed (PDFS), learned (embedding vectors), and prompting-based heuristics, without the need of domain-specific heuristics or fine-tuning.

**Keywords:** Ontology Matching, Retrieval Augmented Generation, Greedy Search, Large Language Models, Zero-Shot Setting

---

## 1. Introduction

In the field of information management, ontologies play a key role in semantic interoperability and knowledge exchange by providing a shared vocabulary that promotes common understanding within a domain. They act as valuable resources in various AI applications, such as autonomous communication in manufacturing environments [1], automation of information flows [2], smart contract creation [3], or 3D scene graph generation [4]. However, the proliferation of incomplete and overlapping ontologies within the same domain is an obstacle to smooth communication between applications. In this context, Ontology Matching (OM) becomes essential for integrating distributed knowledge [5]. OM has broad applications, from linking entities across public ontologies to facilitating knowledge integration in collaborative business environments [6], and merging disparate data warehouses for transactional or analytical purposes in private corporations [7]. Furthermore, OM shares many similarities with knowledge graph alignment [8], and database schema matching [9]. Therefore, advances in OM can be leveraged in a wide range of fields beyond ontologies.

OM refers to the process of identifying semantic correspondences between entities across multiple ontologies [5]. Since 2004, the Ontology Alignment Evaluation Initiative (OAEI) has been organizing annual evaluation campaigns to evaluate and benchmark OM technologies, significantly advancing the field [10]. These evaluation campaigns have

---

\*Corresponding author

Email address: maria.taboada@usc.es (Maria Taboada)

been instrumental in increasing the performance of existing approaches, particularly those focused on identifying equivalence correspondences between entities of a pair of ontologies - referred to as simple pairwise equivalence OM [6]. Consequently, current OM systems have achieved performance improvements, with techniques based on lexical, structural and semantic matching [11], as well as mapping repair techniques [12, 13]. Despite these advances, OM systems continue to have difficulty distinguishing between entities that are semantically similar and those that are merely frequently co-occurring [14, 15]. Moreover, the scalability problem prevents widespread implementation [11].

### 1.1. Motivation and main contributions

Recent OM systems often require fine-tuning with large domain-specific training datasets to achieve optimal performance [16, 17, 18, 19]. In contrast, Large Language Model (LLM)-based methods have emerged as a promising alternative for OM. They leverage pre-trained knowledge for finding correspondences across ontologies and do not require fine-tuning. However, querying LLMs for all entity pairs results in a quadratic time complexity of  $O(n^2)$ , where  $n$  is the number of entities, making this approach impractical for large datasets. To mitigate this issue, state-of-the-art LLM-based systems apply a *retrieve-then-prompt* pipeline. In this method,  $k$  relevant target entities are first retrieved and then used to prompt an LLM to predict mapping correspondences [20, 21, 22, 23]. This approach reduces the time complexity to  $O(n \cdot k)$ , significantly improving scalability [23]. Despite these advancements, these systems still face challenges, particularly in complex tracks where performance in terms of F-Measure is limited. Furthermore, while LLM-based methods are technically feasible, they face scalability issues due to the long execution times required to query the LLM, particularly when applied to large-scale ontologies in real-world scenarios and when using resource-intensive LLMs [20].

To address these challenges, we propose MILA (Minimizing LLM Prompts in Ontology Mapping), a framework designed to improve F-Measure performance and eliminate unnecessary queries to LLMs. MILA introduces a novel *retrieve-identify-prompt* pipeline, which adds an intermediate step to identify high-confidence bidirectional (HCB) correspondences with high precision. This step involves simple heuristics that operate in constant time,  $O(1)$ , eliminating  $k$  LLM queries for each identified HCB correspondence, then reducing LLM interactions to only borderline cases. For these edge correspondences, MILA applies a prioritized depth-first search (PDFS) strategy, which iteratively queries the LLM until a definitive match is confirmed. This approach also minimizes the number of queries by ending early when a valid match is found. Although the overall time complexity remains comparable to existing LLM-based OM systems, MILA significantly reduces execution time, especially when the retrieval system ranks the most relevant candidates first. To achieve efficiency gains, MILA’s retrieval system leverages the SBERT embedding model [24], prioritizing correspondences between entities that have the most semantically similar names.

Our framework was evaluated using the biomedical challenge proposed in the 2024 edition of the Ontology Alignment Evaluation Initiative (OAEI) [10]. We selected the biomedical domain due to its particular nature, which makes OM a challenge [14]: the presence of a rich and constantly evolving terminology, the significant variability in language usage, and the high frequency of rare terms, which are difficult to learn. Moreover, the limited performance of LLM-based solutions in this domain highlights the need for new approaches [23]. To further validate MILA’s robustness and its applicability beyond the biomedical domain, we also assessed its performance on two tasks from the anatomy and biodiversity challenges from the 2024 edition. The results of our evaluation demonstrate the ability of MILA to outperform state-of-the-art OM systems in terms of task-agnostic and high performance in terms of F-measure. They also demonstrate reduced execution times compared to state-of-the-art LLM-based OM systems. These excellent results corroborate the strength of our proposal.

## 2. Related work

OM technology requires the use of external background knowledge to work effectively, as most ontologies are designed in specific contexts that are not explicitly modeled [7]. The most recent OM systems consume existing pre-trained neuronal models [25] as sources of external knowledge [16, 17, 18, 19]. These systems show significant performance improvements over traditional feature engineering and machine learning strategies. They benefit from the capabilities that pre-trained neural models have to automatically interpret the textual descriptions embedded in the labels, comments and definitions of ontologies. However, most models need to be fine-tuned with large training data to perform properly, and they can only process short textual descriptions [20].

To overcome the aforementioned drawbacks, several studies have explored the promising capabilities of LLMs for OM [20, 21, 22, 23, 26, 27]. All of these studies focus on comparing the effect of different prompt inputs to LLMs [28] on OM. In [26], the LLM is provided with complete ontologies in a single prompt, along with detailed instructions on the problem definition and the query goal [26]. However, the most commonly used strategy is to include only a pair of ontology entities in each individual prompt [20, 21, 22, 23, 27]. Moreover, the performance of LLMs has been studied in both zero-shot and few-shot settings [21]. In zero-shot scenarios, where LLMs are queried without providing in-context examples, the performance increases when the prompt contains a set of explicit matching rules. Surprisingly, this zero-shot scenario is almost as effective as providing examples that are textually close to the entities to be matched (few-shot setting) [22]. In addition, prompting LLMs with a brief description of the OM task followed by positive and negative examples has achieved the best results in the anatomy domain [20]. Moreover, the use of multiple choice prompts notably reduces the query execution time, but degrades results. Furthermore, the inclusion of structural context in LLM prompting does not improve OM [27].

To meet the challenge of large-scale OM [12, 29], LLM-based OM systems integrate Retrieval-Augmented Generation (RAG) [30] to effectively reduce the problem of incorrect content generation [20, 22, 23]. They follow a naive methodology based on a *retrieve-then-prompt* pipeline that includes three sequential steps [30]: indexing the target ontology in a vector database, retrieving relevant target entities, and prompting the LLM with mapping candidates. Although the results show that the integration of RAG with LLM in OM is in some cases comparable or even better than current OM systems, the F-Measure performance of LLM-based approaches still requires improvement [20, 22, 23]. Although these systems have a high candidate recovery rate, which can reach 100%, the results can decrease by up to 30% after prompting. Even this reduction can reach 50% in the biomedical domain, where LLM-based systems perform weakly [23]. Moreover, although RAG-based approaches significantly reduce time complexity from quadratic to linear [23], they still require a long runtime when applied to large ontologies [20]. This underscores the need for alternative RAG strategies that improve F-Measure performance while minimizing the number of requests to the LLM [23]. One suggested solution is to initially use a fast and highly accurate matcher to find correspondences for straightforward cases, reserving LLM prompting only for more complex or ambiguous cases [20]. This proposal would optimize efficiency by reducing the dependence on the LLM for simple matching tasks. Therefore, the challenge now is to design the matcher so that it can effectively and reliably identify correspondences with minimal computational overhead.

### 3. Preliminaries

#### 3.1. Problem formulation

Ontology matching (OM) is the process of identifying semantic correspondences among entities of overlapping ontologies [6]. A simple pairwise OM can be defined by a function that takes as input a source ontology  $O_S$ , a target ontology  $O_T$ , an input alignment  $A$ , a set of parameters  $p$  and resources  $r$  (such as external background knowledge), and return an alignment  $A'$  (i.e., a set of correspondences) between entities (or classes) of the input ontologies [5]. A semantic correspondence is a quintuple  $\langle id, e_{O_S}, r, e_{O_T}, c \rangle$ , such as:

- $id$  identifies the correspondence,
- $e_{O_S}$  and  $e_{O_T}$  are entities of  $O_S$  and  $O_T$ , respectively,
- $r$  identifies the semantic relation between  $e_{O_S}$  and  $e_{O_T}$ ,
- $c$  is a confidence value that reflects the degree of correctness of the correspondence, which is usually a real value in the interval  $[0,1]$ .

In our approach,  $r$  is an equivalence relation ( $\equiv$ ) that links two entities through a bidirectional correspondence. Therefore,  $\langle id_i, e_{O_S}, r, e_{O_T}, c \rangle$  is the inverse of  $\langle id_j, e_{O_T}, r, e_{O_S}, c \rangle$  [6]. An example of a simple and bidirectional pairwise correspondence between the ontologies NCI Thesaurus (NCIT) [31] and Disease Ontology (DOID) [32], expressed in Description Logic (DL), is the following:

$$O_{NCIT} : \text{clear cell sarcoma of soft tissue} \equiv O_{DOID} : \text{clear cell sarcoma}.$$

### 3.2. Prioritized Depth-First Search (PDFS)

State-space search algorithms aim to find solutions to problems represented by a set of states. They organize the search space into a tree and evaluate the best path based on certain criteria, typically optimizing the cost to reach a goal. Search algorithms are divided into uninformed and informed categories. Uninformed algorithms, such as depth-first search (DFS), explore the state space without knowing how promising a state is. Informed algorithms, such as the greedy best-first search (GBFS), use a heuristic function to guide the search toward the goal by prioritizing nodes that appear closest to the goal. The term *greedy* usually denotes that the decision is never revised. However, sometimes this term is also used to describe an algorithm that backtracks when a dead end is reached, combining elements of both DFS and GBFS. Sometimes, this combined strategy is called *prioritized depth first search* (PDFS) to clearly distinguish it from a pure greedy approach.

### 3.3. Main OM components in RAG-based approaches

In RAG systems, domain knowledge is stored in vector databases, which are specifically designed to store and index individual text units based on their corresponding embedding vectors. These embedding representations enable the retrieval of relevant information when a query encoded in the same latent space is processed. The retrieval process is usually supported by the semantic similarity between the query and the indexed vectors, facilitating the extraction of contextually relevant text units to query the LLM.

The workflow of most RAG-based OM systems typically involves several sequential steps [20, 22, 23]: 1) *Vector knowledge base (KB) construction*, where the target ontology is extracted, pre-processed, and indexed; 2) *Mapping prediction*, where mapping candidates are retrieved by computing the semantic similarity between the vector representation of a source entity and the indexed target entities; 3) *Mapping refinement*, where candidates are either confirmed or discarded through prompting an LLM; 4) *Candidate filtering*, where mappings are filtered out based on predefined heuristics.

#### 3.3.1. Vector KB construction

Based on the assumption that the matched entities are likely to have labels with overlapping subtokens, traditional methods implement the inverted word-level index [12, 17]. In these approaches, the initial set of correspondences for a source entity is built by selecting target labels that share at least one sub-word token with some label of the source entity. Unlike these approaches, RAG-based OM methods encode and index complete labels and definitions, rather than their sub-words [20]. Label indexing aims at an efficient approach to retrieval, whereas definition indexing is intended to retrieve entities that are not similar in appearance [22]. Hierarchical contexts (parent and child labels) can also be extracted and encoded, although they show worse performance [23]. Therefore, in line with previous works [20, 23], our approach only indexes labels (preferred terms and synonyms), with the goal of maximizing retrieval efficiency.

Moreover, some OM approaches index only the target ontology [22, 23], while others index both ontologies, aiming to increase the initial set of candidate correspondences and thus the recall of the approach [17, 20]. MILA also indexes both ontologies, but with the goal of properly handling all the search space across them.

#### 3.3.2. Mapping prediction

An embedding-based retrieval model predicts the most similar candidates for the correspondences in OM. By computing the cosine similarity between a vector representing the source entity and each vector in the target database, the most similar alignment candidates are predicted. From these, the top k candidates per entity are selected [20, 23].

#### 3.3.3. Mapping refinement

In RAG-based approaches [20, 22, 23], an LLM filters the candidates proposed by the retriever. Most approaches verbalize each alignment candidate in text, which is used to populate an LLM prompt template, following some prompting technique [28]. Depending on the type of information provided to the LLM, these techniques can be classified as zero-shot or few-shot. In zero-shot scenarios, LLMs are prompted with no contextual examples provided [20, 23, 27], while a few contextual examples are provided in few-shot settings [20, 22]. In the latter case, it may also be appropriate to provide both examples of correct correspondences and missing correspondences [20].

Based on the amount and type of information provided, prompts can be categorized into two types: simple prompts and contextual prompts. Simple prompts provide minimal ontology context, including only entity names [20]. In contrast, contextual prompts include additional and relevant information, such as definitions or hierarchical relationships [23, 27]. This contextual information can be retrieved directly from the ontologies themselves or from external resources. Alternatively, contextual data can be selectively extracted using graph search algorithms that identify and focus on the most relevant ontology information, leading to more accurate and effective LLM answers [33]. Furthermore, prompt templates can involve binary decisions, where the LLM must decide whether a candidate is correct or not [20, 23, 27], or multiple decisions, where the LLM must select among several candidate proposals [20, 22].

In approaches prior to RAG [12, 17, 34], mapping refinement aims to discover new correspondences from predicted mappings. These approaches are based on the locality principle [12], which assumes that entities semantically related to those in a predicted correspondence are likely to be matched in new mappings. LogMap [12] computes new mappings by expanding the hierarchical contexts of each entity in a mapping and finding correspondences between classes of these two hierarchical contexts. BERTMap [17] also improves the performance of a BERT classifier with a reasonable time cost by restricting this principle to correspondences that have been predicted with a high score.

#### 3.3.4. Candidate filtering

Mapping refinement is usually finished with a post-processing step mainly aimed at filtering out incorrect correspondences. OM systems often use methods based on heuristics, such as confidence score thresholds (confidence filters) or criteria to achieve unambiguous alignments (cardinality filters) [18, 20, 23]. More sophisticated methods rely on logical reasoning [12] to remove correspondences that cause logical inconsistencies after integrating ontologies [17, 35, 36], or on probabilistic reasoning to resolve conflicts [37]. However, our approach does not apply any post-processing step.

### 4. Methodological framework

Our approach MILA aims to find simple and bidirectional pairwise correspondences. To achieve this, MILA solves the OM task through a novel *retrieve-identify-prompt* pipeline, enhanced by a PDFS strategy. The following subsections outline the key components of our approach and illustrate them with examples. An overview of the main steps is presented in Fig. 1.

#### 4.1. Vector KB construction

Given an ontology  $O$ , an entity (or class)  $e \in O$  can have multiple labels (terms) defined by annotation properties, typically including preferred and alternative labels (synonyms). MILA parses the input ontologies and indexes their entities in the *Entity-Term Index* (see Fig. 1). This index comprises three tables: entity-term relations, preferred term-entity relations, and term-entity relations. This structured indexing enables fast access and efficient searching for MILA’s other components.

Let  $\Omega(e)$  be the set of labels of an entity  $e$ . Using an embedding model, such as SBERT [24], the encoder maps each label  $\omega \in \Omega(e)$  to a vector representation, denoted by  $v(\omega) \in R^d$ , where  $v(\omega)$  is the vector representation of the label  $\omega$  in a  $d$ -dimensional vector space (see Fig. 1, *Term Encoders*). Following prior work [20, 22, 23], we encode full labels to generate the vector KBs. However, we do not incorporate ontology contexts in these KBs. This decision is motivated by the fact that test datasets do not provide definitions for most ontologies [10], and that using only the ontology’s terminology is an efficient retrieval strategy in terminology-intensive domains, such as the biomedical field [22].

#### 4.2. Mapping prediction

For each source and target entity, MILA pre-computes the most promising correspondence candidates, taking into account the vector representations of their labels. Let  $O_S$  and  $O_T$  represent source and target ontologies, respectively. Let  $e_S$  and  $e_T$  denote source and target entities, respectively. For a given source label  $\omega_S \in \Omega(e_S)$  and a target label  $\omega_T \in \Omega(e_T)$ , the function  $f_w$  computes the similarity score *sim* between their respective vector representations,  $v(\omega_S)$  and  $v(\omega_T)$ , using a similarity metric. Specifically, the function  $f_w$  is defined as

$$f_w(\omega_S, \omega_T) = \text{sim}(v(\omega_S), v(\omega_T)).$$

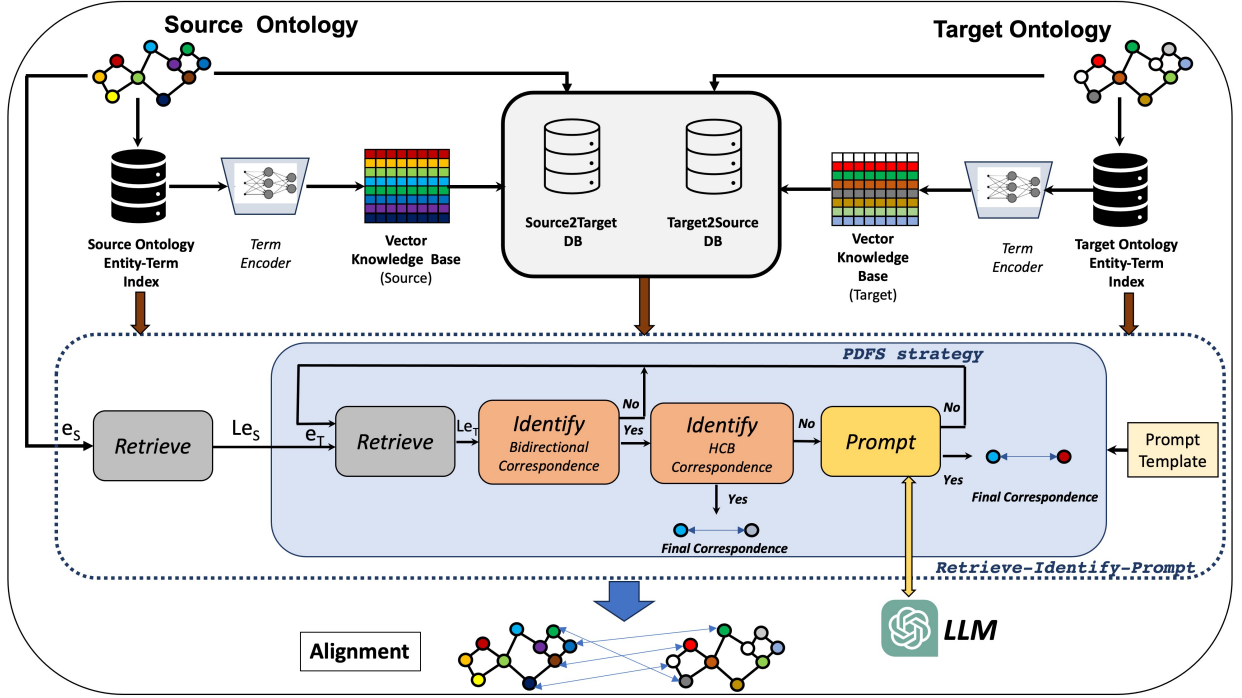


Figure 1: Overview of MILA.

We used cosine similarity in this step of the RAG model since it guarantees the retrieval of contextually relevant candidates, even when their labels differ, which is essential to achieve high-quality correspondences [38]. Fig. 2 shows the predicted correspondence candidates for the source entity  $ncit:C3745$  (Clear Cell Sarcoma of Soft Tissue) using SBERT. In this figure, the similarity score  $f_w$  is depicted in brown, representing the similarity between source labels (in green) and target labels (in gray). Specifically, for the terms  $\omega_{NCIT} = \text{clear cell sarcoma of soft tissue}$  and  $\omega_{DOID} = \text{clear cell sarcoma}$ ,  $f_w$  is

$$f_w(\omega_{NCIT}, \omega_{DOID}) = 0.80521.$$

Let  $\Omega(O_T)$  be the set of all labels in  $O_T$ . Given  $k$ , for each source label  $\omega_S \in \Omega(e_S)$ , MILA generates the  $k$  most promising correspondence candidates by selecting a subset  $C_{\omega_S} = \{\omega_{T_1}, \omega_{T_2}, \dots, \omega_{T_k}\} \subseteq \Omega(O_T)$  that maximizes the score function  $f_w$  with respect to  $\omega_S$ . To refine this selection, we introduce a threshold  $\tau$  that filters out candidates whose score function values fall below it. Therefore, the similarity scores satisfy

$$f_w(\omega_S, \omega_{T_i}) \geq \tau \geq f_w(\omega_S, \omega_{T_j}), \quad \forall (\omega_{T_i}, \omega_{T_j}) \in C_{\omega_S} \times \Omega(O_T) \setminus C_{\omega_S}.$$

For example, as illustrated in Fig. 2, when considering the term  $\omega_{NCIT} = \text{clear cell sarcoma of soft tissue}$  with  $\tau = 0.75$ , the resulting subset  $C_{\omega_{NCIT}}$  is the following

$$C_{\omega_{NCIT}} = \{\text{adult soft part clear cell sarcoma}, \text{clear cell sarcoma}, \text{clear cell chondrosarcoma}\}.$$

Given a source entity  $e_S$  described by  $n$  labels,  $n = |\Omega(e_S)|$ , MILA generates  $n$  subsets  $C_{\omega_S}^i$ , one for each label  $\omega_S^i \in \Omega(e_S)$ , with  $i = 1, 2, \dots, n$ . Let  $C_{e_S}$  be the union of all generated subsets  $C_{\omega_S}^i$ :

$$C_{e_S} = \bigcup_{i \in \{1, 2, \dots, n\}} C_{\omega_S}^i.$$

The function  $f_e$  computes the similarity score between a source entity  $e_S \in O_S$  and each target entity  $e_T \in O_T$  verifying  $C_{e_S} \cap \Omega(e_T) \neq \emptyset$  by calculating the maximum value of the similarity scores between each source label

$\omega_S \in \Omega(e_S)$  and each target label  $\omega_T \in C_{e_S} \cap \Omega(e_T)$ :

$$f_e(e_S, e_T) = \max (f_w(\omega_S, \omega_T) \mid \omega_S \in \Omega(e_S), \omega_T \in C_{e_S} \cap \Omega(e_T)).$$

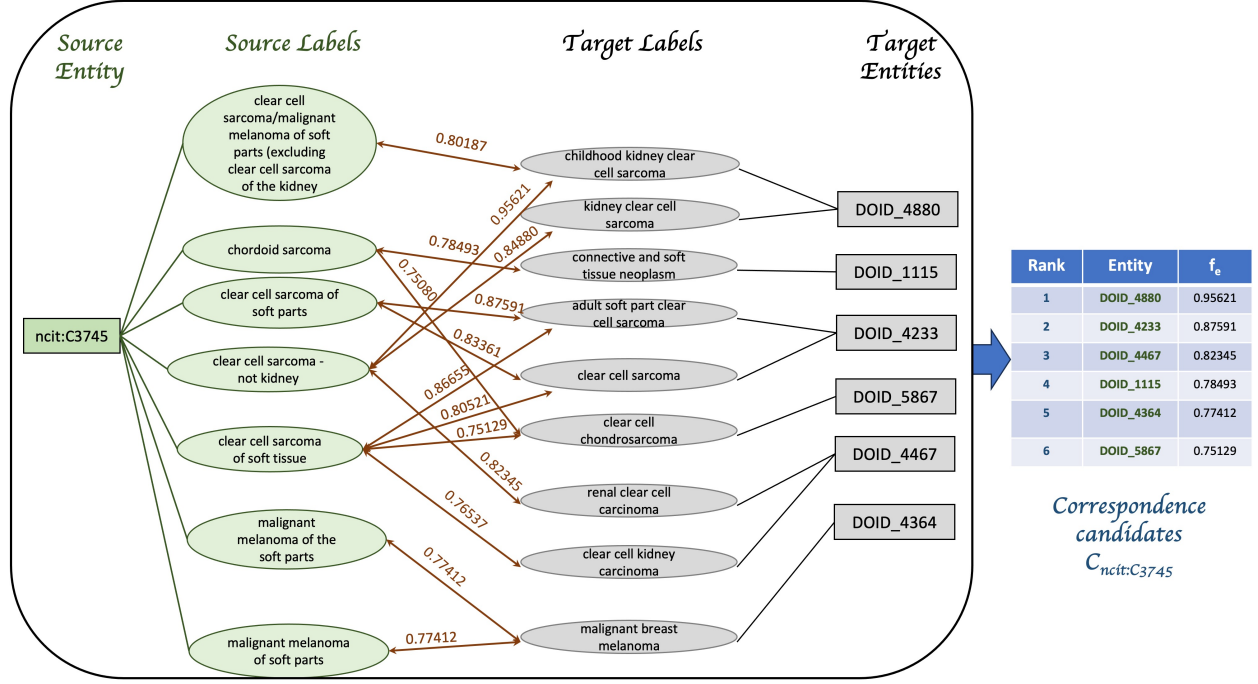


Figure 2: Prediction of correspondence candidates for the entity *ncit:C3745* (clear cell sarcoma of soft tissue) using SBERT.

Therefore, using the function  $f_e$ , MILA generates the set of the  $m := |C_{e_S}|$  most promising target correspondence entities for each source entity  $e_S$ :

$$C_{e_S} = \{e_{T_1}, e_{T_2}, \dots, e_{T_m}\} \subseteq O_T, 0 \leq m \leq n * k$$

verifying

$$f_e(e_S, e_{T_i}) \geq \tau \geq f_e(e_S, e_{T_j}), \quad \forall (e_{T_i}, e_{T_j}) \in C_{e_S} \times O_T \setminus C_{e_S}.$$

The set of source correspondence candidates  $C_{e_S}$  is transformed into an ordered sequence of elements  $L_{e_S} = [e_{T_{i_1}}, e_{T_{i_2}}, \dots, e_{T_{i_m}}] \subseteq O_T$  such that

$$f_e(e_S, e_{T_{i_1}}) \geq f_e(e_S, e_{T_{i_2}}) \geq \dots \geq f_e(e_S, e_{T_{i_m}})$$

where  $\{i_1, i_2, \dots, i_m\}$  is the permutation of indices that sorts the  $C_{e_S}$  elements by the values of the function  $f_e(e_S, e_T)$ . All sequences  $L_{e_S}$  are stored in the *Source2Target Database* in Fig. 1. Similarly, for each target entity  $e_T$ , MILA generates the set of the most promising source entities  $C_{e_T}$ , and the corresponding ordered sequence of elements,  $L_{e_T}$ .

For example, in the table on the right side of Fig. 2, the most promising correspondence candidates for *ncit:C3745* are listed along with their similarity scores,  $f_e$ . Although the entity *DOID\_4880* is described by four labels in the ontology, only two of them have similarity scores that exceed  $\tau$ , *childhood kidney clear cell sarcoma* and *kidney clear cell sarcoma*. The maximum similarity value between these labels and some label of *ncit:C3745* is achieved between the NCIT label *clear cell sarcoma - not kidney* and the DOID label *childhood kidney clear cell sarcoma*. Specifically, the value is 0.95621. Therefore, the similarity between *ncit:C3745* and *DOID\_4880* is set to 0.95621.



### 4.3. The Retrieve-Identify-Prompt pipeline

In this subsection, we first define the two types of correspondences identified by MILA. We then outline the design of the prompt template used to query the LLM. Finally, we provide a detailed description of the algorithm.

#### 4.3.1. Types of correspondences

Given a source entity  $e_S \in O_S$  and a target entity  $e_T \in O_T$ , a *bidirectional correspondence* between them exists if the following condition holds:

$$e_S \in C_{e_T} \wedge e_T \in C_{e_S}.$$

In other words, both entities must appear in each other’s candidate sets, allowing traversal in both directions. This property aligns with the fundamental principle of *symmetry* in equivalence relations. Thus, for the equivalence relation  $e_S \equiv e_T$  to hold, it is necessary that  $e_S$  and  $e_T$  are mutually accessible within their respective candidate sets.

It is important to note that the vector KBs index only the top-k most promising candidates. As a result, while a correspondence from an entity  $e_S$  to another entity  $e_T$  may be retrieved, the reverse correspondence from  $e_T$  to  $e_S$  is not always guaranteed. This limitation arises as the KB may prioritize other, more relevant candidates over  $e_T$ . Therefore, even if  $e_T \in C_{e_S}$  it does not necessarily imply that  $e_S \in C_{e_T}$  and vice versa.

Given a *bidirectional correspondence* between a source entity  $e_S \in O_S$  and a target entity  $e_T \in O_T$ , a *high confidence bidirectional (HCB) correspondence* between both entities exists if they are the most prioritized entities in each other’s correspondence candidate sets:

$$f_e(e_S, e_T) = \max(f_e(e_{S_i}, e_T) \mid e_{S_i} \in C_{e_T}) = \max(f_e(e_S, e_{T_j}) \mid e_{T_j} \in C_{e_S}).$$

#### 4.3.2. Prompt Template Design

In order to query an LLM, MILA uses a simple structured prompt with minimal ontology context to ensure clarity and focus [21]. The prompt includes only the ontology names and the preferred names of the source and target entities involved in the correspondence. This choice is based on studies reporting that a zero-shot scenario is almost as effective as a few-shot setting [20, 22]. The prompt is specifically structured to facilitate the LLM’s binary decision-making. The prompt template is as follows:

You are a helpful expert in ontology matching, which involves determining equivalence correspondences between concepts from different ontologies. The source ontology is called `[src_onto_name]` and the target ontology is called `[tgt_onto_name]`.  
 Classify whether the following concepts are equivalent:  
 Source concept: `[source_entity]`  
 Target concept: `[target_entity]`  
 If so, answer 'Yes', without adding any type of explanation. Otherwise, answer 'No'.

#### 4.3.3. The algorithm

The *retrieve-identify-prompt* pipeline, as outlined in Algorithm 1, is designed to identify valid mappings between two ontologies. For each source entity  $e_S$ , the algorithm first retrieves the ordered sequence of target candidates  $L_{e_S}$ . Then iterates over this sequence using a PDFS strategy, looking for the first target entity  $e_T \in L_{e_S}$  that is a valid match to  $e_S$ . Specifically, for each pair  $(e_S, e_T)$ , the algorithm checks whether the pair is a bidirectional correspondence. If so, it further verifies whether the pair is an HCB correspondence; if it is, the pair  $(e_S, e_T)$  is added to the final mapping. If the pair is not an HCB correspondence, the algorithm queries an LLM to confirm the potential mapping. If the LLM confirms the mapping, it is added to the final mapping, and the search stops for that source entity. If the LLM does not confirm the potential mapping, the search continues with the next candidate in  $L_{e_S}$ , iterating until a valid mapping is identified or all candidates are exhausted.

### 4.4. Examples

This subsection presents two representative examples illustrating the *retrieve-identify-prompt* pipeline. In the first example, MILA identifies an HCB correspondence in the first iteration of the pipeline, while in the second example, it iteratively applies the *retrieve-identify-prompt* pipeline until a valid correspondence is found. In both cases, the source ontology is the NCI Thesaurus (NCIT) [31] and the target ontology is the Disease Ontology (DOID) [32].

---

**Algorithm 1** Retrieve-Identify-Prompt pipeline

---

```
1: Input:  $S$  (source_entities),  $LLM$ ,  $PT$  (prompting_template)
2: Output:  $M$  (mapping)
3: Initialize  $M \leftarrow \{\}$ 
4: for all  $e_S \in S$  do
5:    $L_{e_S} \leftarrow \text{retrieve}(e_S)$ 
6:   not_found  $\leftarrow$  True
7:   while not_found and  $L_{e_S}$  is not empty do                                 $\triangleright$  PDFS strategy
8:      $e_T \leftarrow \text{pop}(L_{e_S})$ 
9:      $L_{e_T} \leftarrow \text{retrieve}(e_T)$                                  $\triangleright$  1. Retrieve
10:    if  $(e_S, e_T)$  is a bidirectional correspondence then                                 $\triangleright$  2. Identify
11:      if  $(e_S, e_T)$  is an HCB correspondence then
12:         $M \leftarrow M \cup \{(e_S, e_T)\}$ 
13:        not_found  $\leftarrow$  False
14:      else
15:         $LLM\_answer \leftarrow \text{prompt}(e_S, e_T, LLM, PT)$                                  $\triangleright$  3. Prompt
16:        if  $LLM\_answer$  is 'Yes' then
17:           $M \leftarrow M \cup \{(e_S, e_T)\}$ 
18:          not_found  $\leftarrow$  False
```

---

#### 4.4.1. Example 1

Fig. 3 illustrates the alignment of the source entity *ncit:C99383* to a corresponding entity in DOID, without requiring an LLM query, as an HCB correspondence is identified. In this case, *ncit:C99383* is uniquely labeled as *autoimmune nervous system disorder*.

- **Step 1 - Retrieve:** MILA first retrieves the ordered sequence of target candidates,  $L_{ncit:C99383}$ , for the entity *ncit:C99383*. These candidates are the most promising target candidates for the label *autoimmune nervous system disorder*:
  1. DOID:438 (autoimmune disease of the nervous system)
  2. DOID:0060004 (autoimmune disease of central nervous system)
  3. DOID:417 (autoimmune disease)
  4. DOID:11465 (autonomic nervous system disease)
- **Step 2 - Retrieve:** MILA retrieves the ordered sequence of the source candidates for the highest-ranked candidate in  $L_{ncit:C99383}$ , DOID:438.
- **Step 3 - Identify bidirectional correspondence:** Next, MILA checks whether there is a bidirectional correspondence between the highest ranked candidate, DOID:438, and the source entity.
- **Step 4 - Identify HCB correspondence:** As there is a bidirectional correspondence, MILA checks if it is an HCB correspondence. Since *ncit:C99383* and DOID:438 are each other's top-ranking candidates, MILA confirms the HCB correspondence and adds it to the set of final mappings:

$$O_{NCIT} : \text{autoimmune nervous system disorder} \equiv O_{DOID} : \text{autoimmune disease of the nervous system}$$

#### 4.4.2. Example 2

Fig. 4 illustrates how the source entity *ncit:C3745* (*clear cell sarcoma of soft tissue*) aligns with the target entity DOID:4233, requiring only two queries to the LLM - compared to the five queries needed by the baseline pipeline.

- **Step 1 - Retrieve:** MILA first retrieves the ordered sequence of target candidates  $L_{ncit:C3745}$ . The left side of Fig. 4 shows only the four most promising of these candidates:

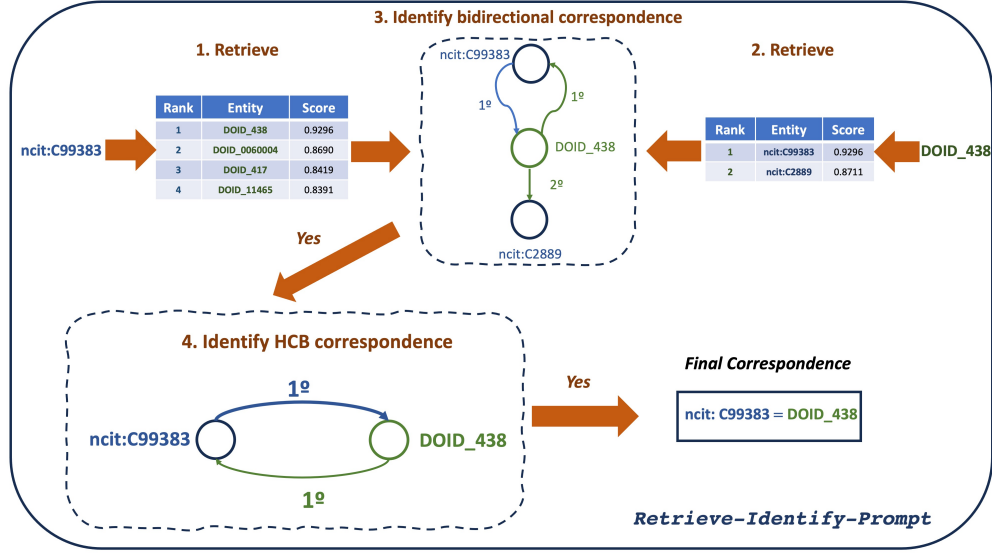


Figure 3: Example showing the identification of an HCB correspondence between the source entity ncit:C99383 and the target entity DOID\_438.

1. DOID:4880 (kidney clear cell sarcoma)
  2. DOID:4233 (clear cell sarcoma)
  3. DOID:4467 (renal clear cell carcinoma)
  4. DOID:1115 (sarcoma)
- **Step 2 - Retrieve:** MILA retrieves the ordered sequence of source candidates for the most promising target entity,  $L_{DOID:4880}$ .
  - **Step 3 - Identify bidirectional correspondence:** Next, MILA checks whether there is a bidirectional correspondence between ncit:C3745 and DOID:4880.
  - **Step 4 - Identify HCB correspondence:** As there is a bidirectional correspondence, it checks whether this is an HCB correspondence, by verifying whether ncit:C3745 is the top-ranked candidate for DOID:4880.
  - **Step 5 - Prompt:** Since there is no bidirectional correspondence between ncit:C3745 and DOID:4880, MILA prompts the LLM to confirm the match. Since the LLM does not confirm the match between ncit:C3745 and DOID:4880, MILA proceeds to the next most similar candidate, DOID:4233 (clear cell sarcoma), and repeats the process.
  - **Steps 6 and 7 - Retrieve and Identify bidirectional correspondence:** MILA retrieves the ordered sequence of source candidates for the second most promising target entity, DOID:4233, and then verifies that the pair (ncit:C3745, DOID:4233) is a bidirectional correspondence, but not an HCB correspondence.
  - **Step 8 - Prompt:** MILA prompts the LLM to confirm the pair (ncit:C3745, DOID:4233). This time, the LLM confirms the match, allowing MILA to finalize the alignment and return the valid correspondence.

## 5. OM Evaluation

In this section, we report the experimental work we have carried out with the biomedical evaluation benchmark proposed by the OEAI in the 2024 edition [10]. In addition to its relevance, the biomedical domain is characterized by its terminology richness, which can be automatically processed by pre-trained language models and LLMs. Currently, the OEAI offers large OM datasets that cover the high variability in term expression typical of this domain [14]. In

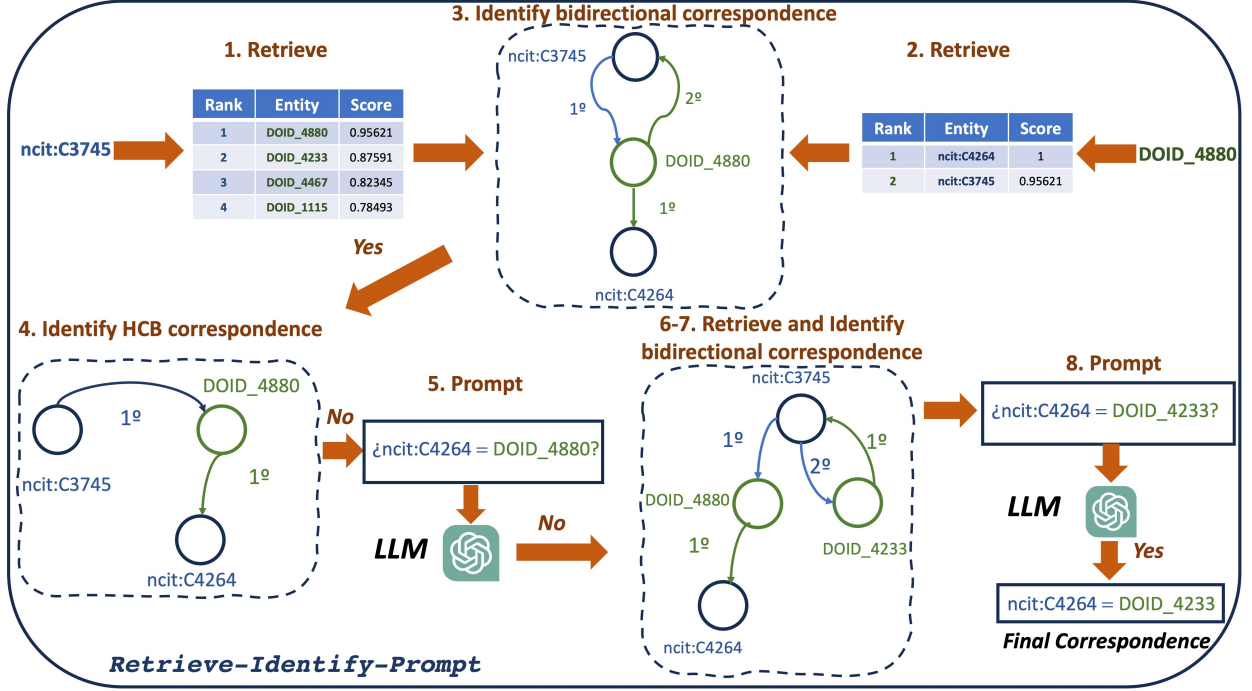


Figure 4: Example showing the *retrieve-identify-prompt* pipeline embedded into the PDFS strategy.

particular, the low performance achieved by LLM-based solutions in this domain encourages the development of new approaches [23]. To further assess the robustness and broader applicability of MILA, we also report experimental results on other well-known tracks, including the Anatomy (mouse-human task) and Biodiversity (envo-sweet task) benchmarks.

### 5.1. Evaluation Metrics

All OM systems were evaluated using the traditional information retrieval metrics of Precision (P), Recall (R) and F-Measure. Precision reflects the correctness of an alignment achieved by an OM system ( $\mathcal{A}$ ) with respect to a reference alignment ( $\mathcal{R}$ ). It is usually defined as the ratio between the number of matches correctly detected by the OM system and the total number of matches identified by the OM system. Recall reflects the completeness of an alignment achieved by an OM system ( $\mathcal{A}$ ) with respect to a reference alignment ( $\mathcal{R}$ ). It is usually defined as the ratio of the number of matches correctly detected by the OM system to the total number of matches identified in the reference alignment. F-Measure combines P and R in a unique measurement.

$$P(\mathcal{A}, \mathcal{R}) = \frac{|\mathcal{A} \cap \mathcal{R}|}{|\mathcal{A}|} \quad R(\mathcal{A}, \mathcal{R}) = \frac{|\mathcal{A} \cap \mathcal{R}|}{|\mathcal{R}|} \quad F - Measure(\mathcal{A}, \mathcal{R}) = 2 * \frac{P * R}{P + R}$$

### 5.2. Dataset and Tasks

The datasets used in our experiments cover an anatomy task, a biodiversity task, and five biomedical tasks. This diverse selection of datasets ensures a comprehensive evaluation of the proposed approach across multiple domains, highlighting its adaptability to different challenges of OM.

- **Anatomy OM Task:** This real-world scenario involves aligning the Adult Mouse Anatomy ontology with the part of the NCI Thesaurus (NCIT) that describes human anatomy.

- **Biodiversity OM Task:** This task focuses on finding correspondences between the Environment Ontology (ENVO) and the Semantic Web for Earth and Environment Technology (SWEET) Ontology, facilitating interoperability in environmental sciences.
- **Biomedical OM Tasks:** These tasks involve ontology matching across six widely used biomedical ontologies: Online Mendelian Inheritance in Man (OMIM) [39], Orphanet Rare Disease Ontology (ORDO) [40], SNOMED CT [41], Foundational Model of Anatomy (FMA) [42], Disease Ontology (DOID) [32] and NCI Thesaurus (NCIT) [31].

Each task in the biomedical benchmark includes both equivalence matching and subsumption matching, although our experiments are focused only on equivalence OM. The quality of reference mappings is ensured both by human curation and by the use of several automated techniques, such as ontology pruning or enrichment with locality modules [12, 13]. The benchmark provides two different sets of test mappings [15]. There is one set for unsupervised OM systems that do not use training mappings and another set for so-called semi-supervised systems, i.e. those that do use training mappings. Specifically, the first set contains the full set of reference mappings, while the second test set includes 70% of the set of reference mappings (excluding 30% of the training mappings). Even if an OM system is an unsupervised approach, such as MILA, it can report performance on this second set for comparison with the supervised OM systems.

### 5.3. Final Configuration

For all OM tasks, MILA used the SBERT retriever model [24], which was set to multi-qa-distilbert-cos-v1. This choice was guided by the findings of the ablation study reported in [20], where this model achieved the highest recall at  $k=1$ , highlighting its outstanding ability to identify the most relevant correspondences — a crucial factor for accurately identifying HCB correspondences. Although the model did not achieve the best recall at  $k=5$ , it still provides a solid foundation for retrieving high-quality matches, making it a reliable choice for this task. Following [23], the value  $k$  during the search for the top  $k$  neighbors ( $top_k$ ) was set to 5, as a compromise between the number of candidates to be generated and the recovery to be achieved. The candidates generated by SBERT with a confidence score lower than 0.75 were discarded and the scores were rounded to five decimal places. On the other hand, owlready2 [43] is used to extract information from ontologies, such as parent and child entities, textual labels, and annotations.

For language model-based processing, MILA used LLaMa-3.1 (70B) [44]. The experiments were conducted on a desktop computer equipped with an Intel Core i5 (4 cores, 3.50 GHz) processor, 32 GB RAM and an AMD Radeon R9 M295X (4GB) graphic card. LLaMa-3.1 (70B) was executed via an inference endpoint, leveraging 4× NVIDIA L40S GPUs. All LLM parameters were kept at their default values, with a temperature of 0.7 to maintain a balance between creativity and logical coherence. Each experiment using LLaMa-3.1 (70B) was executed in a minimum of 20 executions to obtain a reliable estimate of variance. Although a higher number of executions per experiment would provide a more stable and accurate measure of variability [23], the high associated costs prevented this approach. For runtime comparisons between the baseline and MILA pipelines, MILA used LLaMa-3.3 (8B), as it could be executed on our desktop computer. Running LLaMa-3.1 (70B) via an inference endpoint for the baseline pipeline was unaffordable due to the associated high costs.

### 5.4. Experimental Results on the biomedical evaluation benchmark

In this study, we compared the performance of MILA to fourteen OM systems on the evaluation benchmark: AMD [45], BERTMap [17], BERTMapLt [46], BioGITOM [47], BioSTransMatch [48], HybridOM [49], LLMs4OM [23], LogMap [12], LogMapBio [13], LogMapLt [13], Matcha [18], Matcha-DL [18], OLaLa [20] and SORBETMatcher [19]. Data were compiled from results published in the OAEI BIO-ML track (2023 and 2024 editions), as well as relevant literature, with no OM system intentionally excluded to the best of our knowledge. This comparison highlights the strengths and limitations of current systems, providing a benchmark to evaluate the improvements achieved by MILA in terms of performance and computational overhead.

In short, some state-of-the-art OM systems (e.g., BertMap, Matcha, Olala or MILA) use only textual knowledge to predict mapping candidates, whereas others (e.g., LogMap, SORBETMatcher or LLMs4OM) also use structure knowledge. Additionally, most OM systems use pre-trained neuronal models (such as SBERT or BERT) to encode

ontology entities. In some cases, they also include a fine-tuning stage (e.g., BertMap or Matcha-DL), or domain-specific knowledge (e.g., BioGITOM, BioSTransMatch or LogMapBio). Finally, mapping refinement is mainly based on heuristic and logical reasoning, with the exception of Olala, LLMs4OM and MILA, which leverage LLMs for mapping refinement. Compared to other current approaches, MILA’s main distinction is its use of a PDFS strategy, supported by LLMs, to solve mapping refinement. For more detailed information on the characteristics of comparative systems, the reader is referred to Appendix A.

#### 5.4.1. Results in the Unsupervised Setting

Tables 1 and 2 show the performance of all OM systems in the unsupervised setting of the evaluation dataset. MILA is the best performing algorithm in four of the five OM tasks. Even for the task in which it does not achieve the best results, its performance is comparable to that of the leading system. Below we provide a detailed interpretation of these results for each of the evaluated tasks.

- In the **OMIM-ORDO mapping task**, MILA’s F-Measure score is outstanding compared to the other approaches. It outperforms the second best OM system in this task, LogMapBio, by 17%. Specifically, MILA achieves high recall, compared to the rest of the approaches, due to the high recall of both the retrieval module and the *retrieve-identify-prompt* pipeline. Compared to the other tasks, the OMIM-ORDO task achieves the lowest F-Measure.
- In the **NCIT-DOID mapping task**, MILA outperforms the second best baseline, HybridOM, in terms of F-Measure by 3%. In addition, LogMapBio and SORBETMatcher also achieve similar results. In short, this is the task with the highest average F-Measure score of all approaches.
- In the **SNOMED-FMA mapping task**, MILA’s F-Measure score is outstanding compared to the other approaches. It outperforms the second-best OM systems in this task, BERTMap and HybridOM, by 13%. Again, MILA achieves a high recall, compared to the rest of the approaches, and a precision comparable to the best proposals (family BERTMap).
- In the **SNOMED-NCIT (Pharmacology) mapping task**, MILA is the second-best approach with 4% below the outstanding approach, HybridOM, in terms of F-Measure. Note that MILA outperforms the third baseline (i.e. AMD) in terms of F-Measure by 9%.
- In the **SNOMED-NCIT (Neoplasm) mapping task**, MILA’s F-Measure score is also outstanding compared to the other approaches. It outperforms the second best OM system in this task, LogMapBio, by 17%. As in the other tasks, MILA achieves high recall, compared to the rest of the approaches.

#### 5.4.2. Results in the Semi-Supervised Setting

Although MILA is an unsupervised system, we also show its performance against systems that use data training to improve their results. In the semi-supervised setting, MILA is the best performing algorithm in two of the five ontology mapping tasks (see Tables 3 and 4). Specifically, in the NCIT-DOID task, MILA outperforms the second-best OM system in this task, BioGITOM, by 6%, achieving an F-Measure of 0.97. In the SNOMED-NCIT task, it outperforms the second-best OM system, Matcha-DL, by 18%. In addition, for the rest of the tasks, MILA is the second or third best approach with an F-Measure between 3% and 4% below the leading approaches. Please note that no information is available for LLMs4OM, and therefore it does not appear in the tables.

#### 5.5. Experimental Results on the anatomy and biodiversity evaluation benchmarks

In this section, we report the performance of MILA on the Anatomy evaluation benchmark (MOUSE-HUMAN task) and the Biodiversity evaluation benchmark (ENVO-SWEET task). First, we compared MILA with the four top performing systems from the OAEI 2024 Anatomy track: Matcha [18], MDMapper [50], LogMap [12], and LogMapBio [13]. The evaluation results, compiled from the official benchmark dataset, are presented in Table 5. MILA achieved the second-best performance, with 2% below the outstanding approach, Matcha, in terms of the F-Measure. Next, we assessed MILA’s performance on the ENVO-SWEET task of the OAEI 2024 biodiversity track,

OM System	OMIM-ORDO			NCIT-DOID		
	P	R	F-Measure	P	R	F-Measure
AMD	0.664	0.508	0.576	0.885	0.691	0.777
BERTMap	0.734	0.576	0.646	0.888	0.878	0.883
BERTMapLt	0.834	0.497	0.623	0.919	0.772	0.839
BioSTransMatch	0.312	0.586	0.407	0.657	0.833	0.735
HybridOM	0.690	0.679	0.685	0.924	0.913	0.918
LLMs4OM	0.718	0.580	0.641	0.862	0.801	0.830
LogMap	0.876	0.448	0.593	0.934	0.668	0.779
LogMapBio	0.866	0.609	0.715	0.860	<b>0.962</b>	0.908
LogMapLt	<b>0.940</b>	0.252	0.397	<b>0.983</b>	0.575	0.725
Matcha	0.781	0.509	0.617	0.882	0.756	0.814
Matcha-DL	0.745	0.513	0.607	0.847	0.586	0.693
OLaLa	0.735	0.582	0.649	0.913	0.864	0.888
SORBETMatcher	0.693	0.635	0.663	0.920	0.907	0.913
MILA	0.879	<b>0.778</b>	<b>0.831</b>	0.964	0.932	<b>0.948</b>

Table 1: Performance comparison for the the OMIM-ORDO and NCIT-DOID tasks in the unsupervised setting.

OM System	SNOMED-FMA			SNOMED-NCIT (Pharm)			SNOMED-NCIT (Neopl)		
	P	R	F-Measure	P	R	F-Measure	P	R	F-Measure
AMD	0.890	0.633	0.740	0.962	0.670	0.790	0.836	0.481	0.610
BERTMap	<b>0.979</b>	0.662	0.790	0.971	0.585	0.730	0.557	0.762	0.643
BERTMapLt	<b>0.979</b>	0.655	0.785	0.981	0.574	0.724	0.831	0.687	0.752
BioSTransMatch	0.128	0.384	0.192	0.584	0.844	0.690	0.289	0.663	0.402
HybridOM	0.870	0.722	0.790	0.916	<b>0.889</b>	<b>0.902</b>	0.807	0.710	0.755
LLMs4OM	0.211	0.326	0.256	0.818	0.582	0.680	0.470	0.530	0.495
LogMap	0.744	0.407	0.526	0.966	0.607	0.746	0.870	0.586	0.701
LogMapBio	0.827	0.577	0.680	0.928	0.611	0.737	0.748	0.795	0.771
LogMapLt	0.970	0.542	0.696	<b>0.996</b>	0.599	0.748	<b>0.951</b>	0.517	0.670
Matcha	0.887	0.502	0.641	0.987	0.607	0.752	0.838	0.551	0.665
Matcha-DL	0.960	0.602	0.740	0.904	0.616	0.733	0.811	0.514	0.629
OLaLa	0.270	0.348	0.304	—	—	—	0.540	0.546	0.543
SORBETMatcher	0.618	0.749	0.677	0.973	0.607	0.748	0.626	0.642	0.634
MILA	0.964	<b>0.834</b>	<b>0.894</b>	0.981	0.772	0.864	0.928	<b>0.880</b>	<b>0.903</b>

Table 2: Performance comparison for the SNOMED-FMA (Body), SNOMED-NCIT (Pharmacology) and SNOMED-NCIT (Neoplasm) tasks in the unsupervised setting.

OM System	OMIM-ORDO			NCIT-DOID		
	P	R	F-Measure	P	R	F-Measure
AMD	0.601	0.567	0.583	0.858	0.770	0.811
BERTMap	0.645	0.592	0.617	0.831	0.883	0.856
BERTMapLt	0.782	0.507	0.615	0.888	0.770	0.825
BioGITOM	0.951	0.773	<b>0.853</b>	0.944	0.884	0.913
BioSTransMatch	<b>0.973</b>	0.278	0.432	0.698	0.741	0.719
HybridOM	0.611	0.683	0.645	0.895	0.913	0.904
LogMap	0.834	0.456	0.589	0.908	0.664	0.767
LogMapBio	0.821	0.614	0.703	0.811	<b>0.959</b>	0.879
LogMapLt	0.919	0.261	0.407	<b>0.976</b>	0.575	0.723
Matcha	0.718	0.519	0.602	0.839	0.750	0.792
Matcha-DL	0.745	0.732	0.738	0.847	0.834	0.841
OLaLa	0.655	0.570	0.610	0.880	0.861	0.870
SORBETMatcher	0.568	0.652	0.607	0.925	0.882	0.903
MILA	0.874	<b>0.784</b>	0.827	0.967	0.928	<b>0.970</b>

Table 3: Performance comparison for the OMIM-ORDO and NCIT-DOID tasks in the semi-supervised setting.

where all competing systems belonged to the LogMap family [13]. The results, summarized in Table 5, indicate that MILA outperforms all other approaches in this benchmark. In particular, it outperforms LogMap, the second best system, by 17%.

### 5.6. Runtime analysis

This section provides the runtimes of MILA for the datasets used in our experiments. Since MILA consists of three major steps, we report the runtimes for each of them: KB building, mapping prediction, and refinement (the *retrieve-identify-prompt* pipeline). Table 6 shows the runtimes associated with the KB building process, while 7 presents the runtimes for mapping prediction and refinement.

#### 5.6.1. KB building

Table 6 shows the processing time required to build the source and target KBs with SBERT. This runtime increases with the number of labels in the ontology. This observation is consistent with theoretical expectations, since the time complexity of the KB construction,  $O(n \cdot L^2)$ , depends on both the total number  $n$  of labels and the quadratic complexity of processing each label in transformer-based models - where  $L$  denotes the average token length per label. Scalability challenges arise primarily when  $L$  is large due to the quadratic cost of self-attention in transformers [25]. Therefore, ontologies with longer entity names, such as the pair OMIM-ORDO, show longer times compared to other pairs with higher number of labels, such as NCIT-DOID. However, since MILA only indexes entity names, which are usually short, this complexity remains tractable. As a result, KB construction in MILA scales efficiently, even for large ontologies, maintaining tractability across a variety of use cases.

#### 5.6.2. Mapping prediction

Table 7 shows that, in some cases, the processing time spent to predict correspondences can be significantly higher than the time required to build the KBs (see Table 6). For example, SNOMED-FMA, the largest dataset with 7,256 source labels in the reference dataset and 142,984 target labels, required more than 2 hours, which is more than six times the runtime of the first step (KB building). In contrast, for other tasks, such as NCIT-DOID, both stages (KB building KB and mapping prediction) were completed in approximately the same amount of time. This result is consistent with theoretical expectations, since the computation of similarities using cosine similarity between a query vector and all vectors stored in the target KB follows a complexity of  $O(n_T \cdot d)$ , where  $n_T$  is the number of target labels and  $d$  is the embedding dimension. Additionally, retrieving the top- $k$  most similar vectors adds a logarithmic term,  $O(n_T \cdot d + n_T \cdot \log k)$ . However, since  $k$  is typically much smaller than  $n_T$ , the complexity simplifies to  $O(n_T \cdot d)$ .



OM System	SNOMED-FMA			SNOMED-NCIT (Pharm)			SNOMED-NCIT (Neopl)		
	P	R	F-Measure	P	R	F-Measure	P	R	F-Measure
AMD	0.861	0.709	0.778	0.952	0.746	0.836	0.792	0.528	0.633
BERTMap	0.970	0.669	0.792	0.898	0.715	0.796	0.562	0.771	0.650
BERTMapLt	0.970	0.662	0.787	0.973	0.569	0.718	0.775	0.688	0.729
BioGITOM	0.962	<b>0.886</b>	<b>0.923</b>	0.983	0.713	0.827	—	—	—
BioSTransMatch	0.357	0.661	0.464	0.845	0.860	0.852	0.700	0.607	0.650
HybridOM	0.825	0.725	0.772	0.884	<b>0.886</b>	0.885	0.747	0.718	0.732
LogMap	0.673	0.411	0.511	0.952	0.603	0.738	0.823	0.583	0.683
LogMapBio	0.770	0.577	0.660	0.899	0.606	0.724	0.675	0.793	0.729
LogMapLt	0.958	0.542	0.693	<b>0.994</b>	0.594	0.743	<b>0.931</b>	0.514	0.662
Matcha	0.846	0.502	0.630	0.982	0.601	0.746	0.782	0.545	0.642
Matcha-DL	0.959	0.825	0.887	0.903	0.872	<b>0.888</b>	0.806	0.714	0.757
OLaLa	0.202	0.339	0.253	—	—	—	0.451	0.545	0.493
SORBETMatcher	0.794	0.704	0.746	0.876	0.604	0.715	0.731	0.605	0.662
MILA	<b>0.967</b>	0.815	0.884	0.979	0.764	0.858	0.926	<b>0.863</b>	<b>0.893</b>

Table 4: Performance comparison for the SNOMED-FMA (Body), SNOMED-NCIT (Pharmacology) and SNOMED-NCIT (Neoplasm) tasks in the semi-supervised setting.

OM System	MOUSE-HUMAN			ENVO-SWEET		
	P	R	F-Measure	P	R	F-Measure
Matcha	0.951	<b>0.931</b>	<b>0.941</b>	—	—	—
MDMapper	0.926	0.881	0.903	—	—	—
LogMap	0.917	0.848	0.881	0.776	0.659	0.713
LogMapBio	0.888	0.898	0.908	—	—	—
LogMapKG	—	—	—	0.775	0.658	0.711
LogMapLt	0.962	0.728	0.828	0.803	0.595	0.683
MILA	<b>0.974</b>	0.875	0.922	<b>0.951</b>	<b>0.748</b>	<b>0.837</b>

Table 5: Performance comparison for the the MOUSE-HUMAN and ENVO-SWEET tasks.

per query, leading to a total query complexity of  $O(n_S \cdot n_T \cdot d)$ , where  $n_S$  is the number of source labels to match. Consequently, as the size of both ontologies grows, the time required increases at a rate proportional to their product, resulting in a quadratic time complexity of the correspondence prediction step.

### 5.6.3. The retrieve-identify-prompt pipeline

The runtime analysis in Table 7 shows that the *retrieve-identify-prompt* pipeline remains efficient, with execution times consistently below 13 minutes in most tasks, except for SNOMED-NCIT (Pharm). Although the reference dataset for SNOMED-NCIT (Pharm) contains fewer entities (5,803) compared to the largest task, SNOMED-FMA (7,256), the pipeline for SNOMED-NCIT (Pharm) takes more than twice as long as for SNOMED-FMA. This discrepancy is mainly attributed to MILA identifying a lower proportion of high-confidence bidirectional (HCB) matches in SNOMED-NCIT (Pharm), resulting in a higher number of LLM queries, which in turn increases the execution time. Furthermore, the low average variability in LLM responses across all tasks suggests stable performance, regardless of the size and complexity of the ontology.

Table 8 presents a comparison of the runtimes between the *retrieve-identify-prompt* and *retrieve-then-prompt* pipelines across various OM tasks, using LLaMa-3.3-8B-Instruct. The results show that the *retrieve-identify-prompt* pipeline (column 3) requires more than three times the execution time when querying LLaMa-3.3-8B-Instruct compared to querying LLaMa-3.1-70B (column 5, Table 7). This discrepancy is mainly due to the fact that the LLaMa-3.3-8B-Instruct version operates on a machine with fewer computational resources. Furthermore, the *retrieve-identify-*

OM Task	Number of source labels	Number of target labels	Total number of labels	KB building time
OMIM-ORDO	25,890	21,556	47,446	00:04:57
NCIT-DOID	48,619	14,936	63,555	00:04:56
SNOMED-FMA	44,927	142,984	187,911	00:23:23
SNOMED-NCIT (Pharm)	33,456	64,980	98,436	00:11:56
SNOMED-NCIT (Neoplasm)	38,946	60,597	99,543	00:17:06
MOUSE-HUMAN	2,737	5,096	7,833	00:00:35
ENVO-SWEET	9,430	4,365	13,795	00:01:02

Table 6: Runtime KB construction in MILA.

OM Task	Number of target labels	Number of entities in reference	Mapping prediction time	Retrieve-identify-prompt pipeline time	Average Variance in LLM
OMIM-ORDO	21,556	3,721	00:21:32	00:12:41	0.0138
NCIT-DOID	14,936	4,686	00:26:39	00:04:45	0.0086
SNOMED-FMA	142,984	7,256	02:08:34	00:12:35	0.0085
SNOMED-NCIT (Pharm)	64,980	5,803	00:36:51	00:26:28	0.0817
SNOMED-NCIT (Neopl)	60,597	3,804	00:40:59	00:11:32	0.0106
MOUSE-HUMAN	5,096	1,516	00:01:47	00:01:38	0.0274
ENVO-SWEET	4,365	805	00:02:02	00:00:49	0.0055

Table 7: Runtime of Mapping Prediction and Refinement Using MILA with LLaMa-3.1-70B.

*prompt* pipeline significantly reduces execution time compared to the baseline pipeline, demonstrating how MILA minimizes unnecessary LLM queries. Even for larger ontologies such as SNOMED-FMA, MILA completes in approximately 45 minutes, while the baseline takes almost 47 hours. In the case of SNOMED-NCIT (Pharm), the differences in execution time are smaller since, as we have already mentioned, MILA identifies a lower proportion of HCB matches. These results highlight MILA’s efficiency in reducing LLM queries and accelerating OM, particularly for large-scale datasets.

OM Task	Number of entities in reference dataset	Retrieve-identify-prompt Time	Retrieve-then-prompt Time
OMIM-ORDO	3,721	00:43:51	44:53:05
NCIT-DOID	4,686	00:20:07	32:49:56
SNOMED-FMA	7,256	00:44:59	46:46:53
SNOMED-NCIT (Pharm)	5,803	02:03:31	35:32:36
SNOMED-NCIT (Neoplasm)	3,804	00:34:05	26:11:24
MOUSE-HUMAN	1,516	00:11:49	09:27:16
ENVO-SWEET	805	00:04:10	04:13:19

Table 8: Comparison of runtimes between the *retrieve-identify-prompt* and *retrieve-then-prompt* pipelines, using LLaMa-3.3-8B-Instruct.

These experimental results align with theoretical expectations. The time complexity in both baseline and MILA pipelines is  $O(n \cdot k)$ , where  $n$  represents the number of source labels and  $k$  the number of target candidates retrieved per source label. Therefore, the time complexity scales linearly with the number of source labels [23]. However, MILA adds a layer of decision-making that saves time by reducing unnecessary LLM queries. First, MILA incorporates a step where the first candidate can be quickly identified as an HCB correspondence. This step involves simple heuristics that can be done in constant time  $O(1)$ , eliminating  $k$  LLM queries for each entity for which an HCB correspondence

is identified. Second, if the first candidate is not identified as an HCB correspondence, then MILA proceeds with the PDFS strategy, which can also minimize the number of queries by ending early when a valid match is found. In Example 2, MILA completes the alignment in just two queries to the LLM, whereas the basic pipeline requires three additional queries.

#### 5.6.4. Performance on using different configurations

Finally, Table 9 presents the performance scores for all tasks across different MILA configurations: MILA-HCB, which works without prompting any LLM and is based only on the HCB correspondences; MILA-LLaMa-3 (8B), which queries LLaMa-3.3 (8B-Instruct); and MILA-LLaMa-3 (70B), which leverages LLaMa-3.1 (70B). The results indicate that MILA-HCB consistently achieves high precision across all datasets but tends to have a lower recall. MILA-LLaMa-3 (70B) tends to provide superior F-Measure results. In some cases, MILA-LLaMa-3 (8B) equals or improves the F-Measure, although its impact varies depending on the dataset. Although the gain achieved by MILA-LLaMa3-70B is not substantially higher compared to MILA-HCB (except for SNOMED-NCIT-Pharm), it is important to note that MILA-LLaMa3-70B is built on top of MILA-HCB to address the borderline cases where MILA-HCB fails to find a HCB correspondence. Therefore, the benefit of MILA-LLaMa3-70B is in enhancing the overall system’s performance by solving the challenging mappings that MILA-HCB cannot resolve.

OM Test	MILA-HCB			MILA-LLaMa3-8B			MILA-LLaMa3-70B		
	P	R	F-Measure	P	R	F-Measure	P	R	F-Measure
OMIM-ORDO	<b>0.911</b>	0.738	0.816	0.827	<b>0.782</b>	0.804	0.879	0.778	<b>0.831</b>
NCIT-DOID	<b>0.968</b>	0.907	0.936	0.943	0.929	0.936	0.964	<b>0.932</b>	<b>0.948</b>
SNOMED-FMA	<b>0.975</b>	0.799	0.878	0.951	<b>0.844</b>	0.894	0.964	0.834	<b>0.894</b>
SNOMED-NCIT(P)	<b>0.988</b>	0.625	0.766	0.958	<b>0.821</b>	<b>0.884</b>	0.891	0.772	0.864
SNOMED-NCIT(N)	<b>0.932</b>	0.802	0.862	0.907	0.873	0.890	0.928	<b>0.880</b>	<b>0.903</b>
MOUSE-HUMAN	<b>0.985</b>	0.835	0.904	0.920	0.859	0.889	0.974	<b>0.875</b>	<b>0.922</b>
ENVO-SWEET	<b>0.957</b>	0.718	0.820	0.935	<b>0.749</b>	0.832	0.951	0.748	<b>0.837</b>

Table 9: Performance using different MILA configurations.

## 6. Discussion and Future Work

OM plays a key role in enabling smooth communication between heterogeneous applications and ensuring data interoperability and integration across diverse domains. Despite its importance, OM remains an evolving field that continues to benefit from advances in machine learning and language modeling to improve matching performance and scalability. Our experimental results show that MILA significantly improves current LLM-based OM systems in both F-Measure performance and execution efficiency. By introducing an intermediate step that identifies HCB correspondences, MILA enhances significantly the F-Measure with regard the current LLM-based OM systems. Additionally, identification of such matches is performed in constant time, eliminating the need to perform  $k$  LLM queries for each HCB correspondence. The PDFS strategy also minimizes the number of LLM queries by early ending when a valid match is found. Therefore, although time complexity remains consistent with current LLM-based OM systems, MILA offers substantial practical benefits by reducing the number of LLM queries, which can lead to faster matching for large-scale ontologies, as shown in the results achieved in our study. The exact efficiency gain depends on the quality of the retrieval system (how well it ranks relevant candidates), the distribution of correspondence candidates (how early in the retrieval list the correct correspondence typically appears), and the scalability of the LLM (how expensive each query is in terms of time). To achieve efficiency gains, MILA incorporates a retrieval system that prioritizes correspondences between entities with the most semantically similar synonyms. Therefore, although MILA maintains the same theoretical time complexity as current LLM-based OM systems, its ability to minimize unnecessary LLM interactions offers substantial practical benefits, as evidenced by the results in our study.

Moreover, MILA shows a substantial improvement over state-of-the-art OM systems, in terms of its F-Measure performance. First, our work reports that MILA outperforms leading OM systems in four of the five tasks in the

unsupervised setting, and in two of the five tasks in the semi-supervised setting, despite being a zero-shot approach. Second, our approach exhibits task-agnostic performance, remaining nearly stable across all tasks and settings, unlike other approaches, whose performance tends to vary much more depending on the specific task. This uniformity highlights the robustness of our method, as it is either the best-performing or very close to the best in all cases, regardless of task type or setting (unsupervised or semi-supervised). In the following, we analyze in detail the results of the different tasks of the Bio-ML track.

The OMIM-ORDO task exhibits the lowest F-Measure across all tasks and approaches, highlighting the significant challenges faced in this domain. Specifically, linking disease subtype concepts is particularly difficult in the biomedical field, especially when dealing with rare diseases [22]. A key factor contributing to this performance gap may be the limited structural information embedded in the ontologies used for the task [15]. This factor may play a role in making knowledge graph embedding techniques, such as AMD [45], less effective at predicting mappings. Another factor contributing to the difficulty of this task may be the mismatch between the entities’ names and the standard biomedical nomenclature [51]. Specifically, the entity names in these ontologies tend to be excessively long, deviate from common syntactic structures used in medical terminology, and are infrequent in the relevant literature, as they focus on rare diseases. As a result, traditional matching methods that apply efficient string matching techniques, such as LogMapLt [13], are less effective. In contrast, approaches that leverage domain-specific knowledge, such as LogMapBio [13] or BioGITOM, perform better for this task. Our findings suggest that pre-trained language models, such as SBERT [24], can efficiently interpret these unconventional terminologies without the need for additional training or domain knowledge. Despite many current OM systems using the same retriever model (SBERT) as MILA, their performance does not match the results achieved by our approach. The key challenge lies in distinguishing between SBERT candidates that are semantically similar and those that merely overlap statistically [14]. By combining RAG techniques with advanced search strategies, MILA provides a more effective solution to this problem.

Similarly, the SNOMED-NCIT (Neoplasm) task involves ontologies that, while containing limited structural information [15], are more aligned with established standards. As a result, the overall performance of the OM systems tends to improve. On the other hand, the NCIT-DOID task achieves the highest F-Measure score and is therefore the least challenging task [15]. The adoption of a standard biomedical terminology and rich structural information within the involved ontologies provides a favorable environment for all OM systems. Despite this, MILA still outperforms current methods, achieving the best results in both unsupervised and semi-supervised environments. Although LogMapBio [13] and SORBETMatcher [19] achieve similar results, MILA has several advantages. Specifically, LogMapBio makes use of knowledge specific to the biomedical domain, whereas MILA is domain-agnostic, so it could be applied to align ontologies from other domains. Moreover, although SORBETMatcher is also supported by SBERT, it requires a fine-tuning stage, whereas MILA is a zero-shot approach that does not require any training data.

The SNOMED-FMA and SNOMED-NCIT (Pharm) tasks require the recognition of lexical patterns within the target ontologies. For example, patterns such as *set of <something>* are usual in FMA or *<something>-containing product* in NCIT. As a result, these tasks are well suited to automated learning-based methods. For example, MatchaDL leverages a linear neural network to effectively classify candidate mappings [52]. However, when MILA was applied to these tasks, we observed a higher proportion of ambiguous mappings compared to those in more successful tasks. This increased ambiguity led to a higher number of edge cases and then a greater dependence on LLM responses, which were less effective. Despite this, MILA’s performance is in line with that of supervised approaches, highlighting its potential even in scenarios more suited to supervised methods.

### 6.1. Limitations and Future Work

Although MILA present improvements in performance, there are still areas for future work. Currently, MILA involves the use of a greedy search strategy based exclusively on prediction values, which may not achieve the best solution. Future work will explore other informed search strategies that combine these values with structural information from ontologies [35]. Moreover, the use of simple prompts limits the full potential of LLMs. We plan to conduct further experiments that test innovative techniques for generating prompts that incorporate contextual ontology information [33] or relevant examples [20]. Furthermore, although MILA significantly reduces the number of LLM queries, compared to existing LLM-based OM systems, it could still benefit from further optimization in terms of parallelization by integrating the framework proposed by [34] into our pipeline, especially for very large ontologies. In addition, our current experiments focus on simple pairwise OM [5]. However, more sophisticated tasks, such as

predicting subsumption relations [46] or complex mappings [53], present additional challenges. In the near future, we intend to explore how these tasks can benefit from the integration of LLMs and state-space search algorithms.

Finally, a key challenge with embedding-based methods is their scalability. This challenge goes beyond LLM-based OM systems, as currently most OM systems are using or combining these methods [45, 52, 19]. As we have seen in our study, the process of generating embeddings and performing similarity searches can become computationally expensive, particularly for large ontologies with numerous entities. As the size of the ontologies grows, both the generation of embeddings and, more critically, the search for corresponding entities become progressively resource-intensive. This growing demand for computational resources can restrict the effectiveness and applicability of embedding-based methods in large-scale OM tasks, where high efficiency and scalability are needed. To improve scalability for large ontologies, MILA will test several optimizations. A recent solution proposes replacing the brute-force search with an approximate k-nearest neighbor search, using cosine similarity as the distance metric [49]. This option significantly reduces the retrieval time complexity, although at the cost of losing some accuracy. We will also test other optimizations that include the combination of logic-based modules to manage large ontologies more efficiently [34].

## **7. Conclusion**

In summary, MILA represents a significant step forward in the development of scalable, high-performance LLM-based OM systems. By combining RAG and search strategies, MILA offers an effective solution to the challenges of LLM-based OM systems by improving computational overhead and performance in critical domains, such as the biomedical domain. MILA also offers an effective solution to the challenges of OM in general, by exhibiting task-agnostic performance without the need for training data, making it a promising approach for LLM-based OM applications. Future research could focus on further enhancing the scalability of MILA and expanding its applicability to a wider range of domains and tasks.

## **8. CRediT authorship contribution statement**

Maria Taboada: Writing – original draft, Implementation, Validation, Methodology. Diego Martinez: Writing – review & editing, Methodology. Mohammed Arideh: Implementation, Validation. Rosa Mosquera: Writing – review & editing, Funding acquisition.

## **9. Declaration of competing interest**

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this article.

## **10. Acknowledgments**

The authors appreciate the support and training from the University of Santiago de Compostela, as well as the support from the projects AF4EU (101086563) and SUS-SOIL (101157560), funded by the European Union’s Horizon Europe program. The authors also express their gratitude to Diego Martinez-Taboada for insightful conversations.

## **11. Data availability**

<https://github.com/mariatab/MILA>

## Appendix A. State-of-the-Art OM Systems

Table A.10 provides a detailed overview of the state-of-the-art OM systems used for comparison in the MILA evaluation. The selected systems span a range of methodologies, including machine learning-based approaches, RAG frameworks, and traditional heuristic methods. Specifically, Table A.10 focuses on the following characteristics:

- The *type of ontology knowledge* used to predict the mapping candidates. Examples are textual knowledge and structure knowledge. Textual knowledge can include preferred names, synonyms, annotations, etc. Structure knowledge can include either any type of relationship defined in the ontology or only hierarchical relationships.
- The *prediction model* used to generate the mapping candidates. Most OM systems include pre-trained learning models and string-based matchers.
- The inclusion of a *fine-tuning* stage to adapt the pre-trained learning model to mapping prediction.
- The type of mapping refinement. Examples include heuristic-based refinement, logical reasoning-based refinement and LLM-based refinement.

## References

- [1] P. Sapel, L. Molinas Comet, I. Dimitriadis, C. Hopmann, S. Decker, A review and classification of manufacturing ontologies, *Journal of Intelligent Manufacturing* (2024) 1–25.
- [2] B. Esteves, V. Rodríguez-Doncel, Analysis of ontologies and policy languages to represent information flows in GDPR, *Semantic Web* 15 (3) (2024) 709–743.
- [3] J. A. Dominguez, S. Gonnet, M. Vegetti, The role of ontologies in smart contracts: A systematic literature review, *Journal of Industrial Information Integration* (2024) 100630.
- [4] J. Strader, N. Hughes, W. Chen, A. Speranzon, L. Carbone, Indoor and outdoor 3D scene graph generation via language-enabled spatial ontologies, *IEEE Robotics and Automation Letters* (2024).
- [5] J. Euzenat, P. Shvaiko, et al., *Ontology matching*, Vol. 18, Springer, 2007.
- [6] I. Osman, S. B. Yahia, G. Diallo, Ontology integration: approaches and challenging issues, *Information Fusion* 71 (2021) 38–63.
- [7] J. Portisch, M. Hladik, H. Paulheim, Background knowledge in ontology matching: A survey, *Semantic Web* (2022) 1–55.
- [8] X. Zhao, W. Zeng, J. Tang, W. Wang, F. M. Suchanek, An experimental study of state-of-the-art entity alignment approaches, *IEEE Transactions on Knowledge and Data Engineering* 34 (6) (2020) 2610–2625.
- [9] Y. Zhang, A. Floratou, J. Cahoon, S. Krishnan, A. C. Müller, D. Banda, F. Psallidas, J. M. Patel, Schema matching using pre-trained language models, in: *2023 IEEE 39th International Conference on Data Engineering (ICDE)*, IEEE, 2023, pp. 1558–1571.
- [10] E. Jiménez-Ruiz, O. Hassanzadeh, C. Trojahn, S. Hertling, H. Li, P. Shvaiko, J. Euzenat (Eds.), *Proceedings of the 19th International Workshop on Ontology Matching co-located with the 23rd International Semantic Web Conference (ISWC 2024)*, no. 3897 in *CEUR Workshop Proceedings*, CEUR-WS.org, 2025.
- [11] J. Chen, O. Mashkova, F. Zhapa-Camacho, R. Hoehndorf, Y. He, I. Horrocks, Ontology embedding: A survey of methods, applications and resources, *arXiv preprint arXiv:2406.10964* (2024).
- [12] E. Jiménez-Ruiz, B. Cuenca Grau, Logmap: Logic-based and scalable ontology matching, in: *The Semantic Web–ISWC 2011: 10th International Semantic Web Conference, Bonn, Germany, October 23–27, 2011, Proceedings, Part I 10*, Springer, 2011, pp. 273–288.
- [13] E. Jiménez-Ruiz, B. Grau, V. Cross, Logmap family participation in the OAEI 2017, in: *CEUR Workshop Proceedings*, Vol. 2032, CEUR Workshop Proceedings, 2017, pp. 1–5.
- [14] P. Kolyvakis, A. Kalousis, B. Smith, D. Kiritsis, Biomedical ontology alignment: an approach based on representation learning, *Journal of biomedical semantics* 9 (2018) 1–20.
- [15] Y. He, J. Chen, H. Dong, E. Jiménez-Ruiz, A. Hadian, I. Horrocks, Machine learning-friendly biomedical datasets for equivalence and subsumption ontology matching, in: *International Semantic Web Conference*, Springer, 2022, pp. 575–591.
- [16] J. Chen, E. Jiménez-Ruiz, I. Horrocks, D. Antonyrajah, A. Hadian, J. Lee, Augmenting ontology alignment by semantic embedding and distant supervision, in: *The Semantic Web: 18th International Conference, ESWC 2021, Virtual Event, June 6–10, 2021, Proceedings 18*, Springer, 2021, pp. 392–408.
- [17] Y. He, J. Chen, D. Antonyrajah, I. Horrocks, Bertmap: a BERT-based ontology alignment system, in: *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 36, 2022, pp. 5684–5691.
- [18] D. Faria, M. C. Silva, P. Cotovio, L. Ferraz, L. Balbi, C. Pesquita, Results for Matcha and Matcha-DL in OAEI 2023., in: *OM@ ISWC, 2023*, pp. 164–169.
- [19] F. Gosselin, A. Zouaq, Sorbet: A siamese network for ontology embeddings using a distance-based regression loss and BERT, in: *International Semantic Web Conference*, Springer, 2023, pp. 561–578.
- [20] S. Hertling, H. Paulheim, Olala: Ontology matching with large language models, in: *Proceedings of the 12th Knowledge Capture Conference 2023*, 2023, pp. 131–139.
- [21] R. Peeters, C. Bizer, Using ChatGPT for entity matching, in: *European Conference on Advances in Databases and Information Systems*, Springer, 2023, pp. 221–230.

- [22] Q. Wang, Z. Gao, R. Xu, Exploring the in-context learning ability of large language model for biomedical concept linking, arXiv preprint arXiv:2307.01137 (2023).
- [23] H. B. Giglou, J. D’Souza, S. Auer, LLMs4OM: Matching ontologies with Large Language Models, arXiv preprint arXiv:2404.10317 (2024).
- [24] N. Reimers, Sentence-bert: Sentence embeddings using siamese BERT-Networks, arXiv preprint arXiv:1908.10084 (2019).
- [25] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, I. Polosukhin, Attention is all you need, in: Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17, Curran Associates Inc., Red Hook, NY, USA, 2017, p. 6000–6010.
- [26] S. S. Norouzi, M. S. Mahdavi, P. Hitzler, Conversational ontology alignment with ChatGPT, arXiv preprint arXiv:2308.09217 (2023).
- [27] Y. He, J. Chen, H. Dong, I. Horrocks, Exploring large language models for ontology alignment, arXiv preprint arXiv:2309.07172 (2023).
- [28] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, G. Neubig, Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing, ACM Computing Surveys 55 (9) (2023) 1–35.
- [29] M. A. Khoudja, M. Fareh, H. Bouarfa, Deep embedding learning with auto-encoder for large-scale ontology matching, International Journal on Semantic Web and Information Systems (IJSWIS) 18 (1) (2022) 1–18.
- [30] Y. Gao, Y. Xiong, X. Gao, K. Jia, J. Pan, Y. Bi, Y. Dai, J. Sun, M. Wang, H. Wang, Retrieval-augmented generation for large language models: A survey, arXiv preprint arXiv:2312.10997 (2023).
- [31] S. d. Coronado, M. W. Haber, N. Sioutos, M. S. Tuttle, L. W. Wright, NCI Thesaurus: using science-based terminology to integrate cancer research results, in: MEDINFO 2004, IOS Press, 2004, pp. 33–37.
- [32] L. M. Schriml, C. Arze, S. Nadendla, Y.-W. W. Chang, M. Mazaitis, V. Felix, G. Feng, W. A. Kibbe, Disease ontology: a backbone for disease semantic integration, Nucleic acids research 40 (D1) (2012) D940–D946.
- [33] J. Sampels, S. Efeoglu, S. Schimmler, Exploring prompt generation utilizing graph search algorithms for ontology matching, in: Knowledge Graphs in the Age of Language Models and Neuro-Symbolic AI, IOS Press, 2024, pp. 2–19.
- [34] E. Jiménez-Ruiz, A. Agibetov, J. Chen, M. Samwald, V. Cross, Dividing the ontology alignment task with semantic embeddings and logic-based modules, in: ECAI 2020, IOS Press, 2020, pp. 784–791.
- [35] E. Jiménez-Ruiz, C. Meillicke, B. C. Grau, I. Horrocks, Evaluating mapping repair systems with large biomedical ontologies., Description Logics 13 (2013) 246–257.
- [36] E. Santos, D. Faria, C. Pesquita, F. M. Couto, Ontology alignment repair through modularization and confidence-based heuristics, PloS one 10 (12) (2015) e0144807.
- [37] W. Li, S. Zhang, Repairing mappings across biomedical ontologies by probabilistic reasoning and belief revision, Knowledge-Based Systems 209 (2020) 106436.
- [38] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel, et al., Retrieval-augmented generation for knowledge-intensive nlp tasks, Advances in neural information processing systems 33 (2020) 9459–9474.
- [39] A. Hamosh, A. F. Scott, J. Amberger, D. Valle, V. A. McKusick, Online mendelian inheritance in man (OMIM), Human mutation 15 (1) (2000) 57–61.
- [40] D. Vasant, L. Chanas, J. Malone, M. Hanauer, A. Olry, S. Jupp, P. N. Robinson, H. Parkinson, A. Rath, ORDO: an ontology connecting rare disease, epidemiology and genetic data, in: Proceedings of ISMB, Vol. 30, researchgate. net, 2014, pp. 1–4.
- [41] T. Benson, G. Grieve, T. Benson, G. Grieve, SNOMED CT, Principles of Health Interoperability: FHIR, HL7 and SNOMED CT (2021) 293–324.
- [42] C. Rosse, J. L. Mejino Jr, A reference ontology for biomedical informatics: the Foundational Model of Anatomy, Journal of biomedical informatics 36 (6) (2003) 478–500.
- [43] J.-B. Lamy, Owlready: Ontology-oriented programming in python with automatic classification and high level constructs for biomedical ontologies, Artificial intelligence in medicine 80 (2017) 11–28.
- [44] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan, et al., The Llama 3 herd of models, arXiv preprint arXiv:2407.21783 (2024).
- [45] Z. Wang, AMD results for OAEI 2022., in: OM@ ISWC, 2022, pp. 145–152.
- [46] Y. He, J. Chen, H. Dong, I. Horrocks, C. Allocca, T. Kim, B. Sapkota, Deeponto: A python package for ontology engineering with deep learning, arXiv preprint arXiv:2307.03067 (2023).
- [47] S. Oulefki, L. Berkani, L. Bellatreche, N. Boudjenah, A. Mokhtari, Results for BioGITOM in OAEI 2024, in: E. Jiménez-Ruiz, O. Hassanzadeh, C. Trojahn, S. Hertling, H. Li, P. Shvaiko, J. Euzénat (Eds.), Proceedings of the 19th International Workshop on Ontology Matching co-located with the 23rd International Semantic Web Conference (ISWC 2024), Baltimore, USA, November 11, 2024, Vol. 3897 of CEUR Workshop Proceedings, CEUR-WS.org, 2024, pp. 104–109.
- [48] S. Menad, S. Abdeddaïm, L. F. Soualmia, BioSTransMatch results in OAEI 2024, in: E. Jiménez-Ruiz, O. Hassanzadeh, C. Trojahn, S. Hertling, H. Li, P. Shvaiko, J. Euzénat (Eds.), Proceedings of the 19th International Workshop on Ontology Matching co-located with the 23rd International Semantic Web Conference (ISWC 2024), Baltimore, USA, November 11, 2024, Vol. 3897 of CEUR Workshop Proceedings, CEUR-WS.org, 2024, pp. 132–137.
- [49] M. Totoian, A. Marginean, P. Blohm, M. N. Hussain, HybridOM: Ontology matching using hybrid search, in: E. Jiménez-Ruiz, O. Hassanzadeh, C. Trojahn, S. Hertling, H. Li, P. Shvaiko, J. Euzénat (Eds.), Proceedings of the 19th International Workshop on Ontology Matching co-located with the 23rd International Semantic Web Conference (ISWC 2024), Baltimore, USA, November 11, 2024, Vol. 3897 of CEUR Workshop Proceedings, CEUR-WS.org, 2024, pp. 138–145.
- [50] X. Liu, J. Grode, M. R. Hansen, Mdmapper: A framework for aligning master data models using ontology matching techniques, in: E. Jiménez-Ruiz, O. Hassanzadeh, C. Trojahn, S. Hertling, H. Li, P. Shvaiko, J. Euzénat (Eds.), Proceedings of the 19th International Workshop on Ontology Matching co-located with the 23rd International Semantic Web Conference (ISWC 2024), Baltimore, USA, November 11, 2024, Vol. 3897 of CEUR Workshop Proceedings, CEUR-WS.org, 2024, pp. 30–42.
- [51] N. Khatwani, J. Geller, What makes a good concept anyway?, arXiv preprint arXiv:2409.06150 (2024).
- [52] P. G. Cotovio, L. Ferraz, D. Faria, L. Balbi, M. C. Silva, C. Pesquita, Matcha-DL a tool for supervised ontology alignment, preprint (2024).
- [53] M. Silva, D. Faria, C. Pesquita, Complex multi-ontology alignment through geometric operations on language embeddings, in: ECAI, 2024,

pp. 1333–1340.



OM System	Ontology Knowledge in Prediction	Prediction Model	Fine-Tuning	Domain-specific Knowledge	Mapping Refinement
AMD	Textual and structural knowledge	SBERT and TransR	Yes	No	Heuristic filtering
BERTMap	Textual knowledge	Lexical indexation, String-based matching and BERT	Yes	No	Logical reasoning based extension and filtering
BERTMapLt	Textual knowledge	Lexical indexation String-based matching and BERT	No	No	No
BioGITOM	Textual and structural knowledge	BioBERT and GNNs	No	Yes	No
BioSTransMatch	Textual knowledge	BioClinicalBERT	Yes	Yes	Heuristic filtering
HybridOM	Textual and hierarchical knowledge	gtr-t5-large, k-NN search, BM25	No	No	No
LLMs4OM	Textual and hierarchical knowledge	OpenAI text-embedding-ada	No	No	LLM-based and heuristic filtering
LogMap	Textual and structural knowledge	Lexical indexation, String- and Structure-based matching	No	No	Logical reasoning based extension and repair
LogMapBio	Textual and structural knowledge	Lexical indexation, String- and Structure-based matching	No	Yes	Logical reasoning based extension and repair
LogMapLt	Textual and structural knowledge	Lexical indexation, string-based matching	No	No	No
Matcha	Textual knowledge	String-based matchers and SBERT	No	No	Heuristic filtering
Matcha-DL	Textual knowledge	String-based matchers and SBERT	Yes	No	Heuristic filtering
OLaLa	Textual knowledge	String-based matcher and SBERT	No	No	LLM-based and heuristic filtering
SORBETMatcher	Textual and hierarchical knowledge	SBERT	Yes	No	Heuristic filtering
MILA	Textual knowledge	SBERT	No	No	PDFS with retrieve-identify-prompt pipeline

Table A.10: Overview of all OM systems that have evaluated the complete biomedical dataset