

KCMF: A Knowledge-compliant Framework for Schema and Entity Matching with Fine-tuning-free LLMs

Anonymous ACL submission

Abstract

Schema matching (SM) and entity matching (EM) tasks are crucial for data integration. While large language models (LLMs) have shown promising results in these tasks, they suffer from hallucinations and confusion about task instructions. This study presents the Knowledge-Compliant Matching Framework (KCMF), an LLM-based approach that addresses these issues without the need for domain-specific fine-tuning. KCMF employs a once-and-for-all pseudo-code-based task decomposition strategy to adopt natural language statements that guide LLM reasoning and reduce confusion across various task types. We also propose two mechanisms, Dataset as Knowledge (DaK) and Example as Knowledge (EaK), to build domain knowledge sets when unstructured domain knowledge is lacking. Moreover, we introduce a result-ensemble strategy to leverage multiple knowledge sources and suppress badly formatted outputs. Extensive evaluations confirm that KCMF clearly enhances five LLM backbones in both SM and EM tasks while outperforming the non-LLM competitors by an average F1-score of 17.93%.

1 Introduction

Schema matching is the task of identifying correspondences between elements of two or more database schemas. This task plays an important role in data integration efforts. Another task called entity matching, also known as entity resolution or record linkage, aims to identify schema instances that refer to the same real-world entity. While these two tasks vary in their definitions and approaches, they share the common goal of matching database elements. This paper attempts to tackle these tasks under the umbrella term *data matching*, employing a unified methodology while still retaining their original problem settings independently, i.e., entity matching is solved based on the database records,

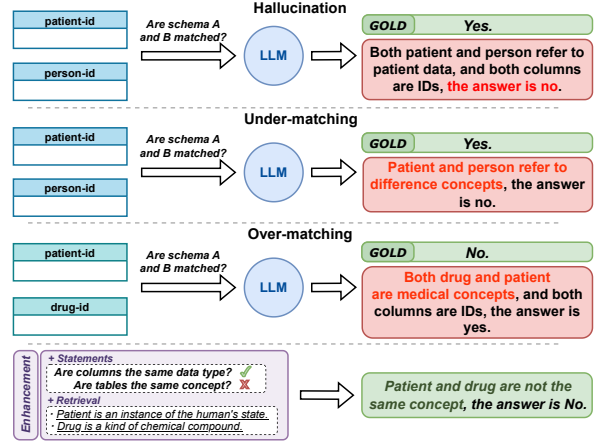


Figure 1: Three common issues in LLM-based data matching tasks and an overview of the enhancement will be discussed in this study: As demonstrated at the bottom, with pseudo-code and retrieved knowledge, by going through statements from the former, the LLM are able to reject the match between *patient-id* and *drug-id*.

whereas schema matching is solved using metadata only due to privacy considerations.

Early systems for data matching tasks rely on expert systems or traditional machine learning methods, while modern data matching approaches often employ pre-trained language models (PLMs) as the backbone (Li et al., 2020; Zhang et al., 2021; Zeakis et al., 2023), which is then fine-tuned on task-specific datasets. However, PLM-based methods usually incur efforts to collect fine-tuning data and face performance degradation when handling out-of-domain data during inference. Large Language Models (LLM) are PLMs of massive scales with billions of parameters, trained on vast and diverse datasets. With huge internal knowledge gained from pre-training and strong representation capability, LLMs show competitive performance on numerous data-wrangling benchmarks, particularly without the need for fine-tuning (Narayan et al., 2022; Peeters and Bizer, 2023; Sheerit et al., 2024; Peeters and Bizer, 2024). This remarkable ef-

ficacy liberates users from the burden of fine-tuning efforts.

Despite promising results, LLM-based methods for data matching face several performance issues. For clarity of presentation, we focus mainly on schema matching in the main text; further discussions of entity matching are provided in Appendix A. As depicted in Figure 1, these issues are categorized into three types: (1) *hallucination*, where the LLM generates incorrect deduction from correct evidence at hand; (2) *under-matching*, where the LLM rejects matching with overly strict criteria; and (3) *over-matching*, where the LLM over-generalizes the association of the input data. While hallucination is a widely recognized problem that can harm the LLM’s performance in many other tasks (Zhang et al., 2023a; Rawte et al., 2023), under-matching and over-matching, jointly termed as **confusion** in this paper, are challenges specific to the data matching tasks.

This study presents the Knowledge-compliant Matching Framework (KCMF), a fine-tuning-free and retrieval-enhanced approach to data matching tasks. To address the problem of confusion, KCMF employs a *pseudo-code-based task decomposition strategy* for LLMs. Pseudo-code here is an ordered list of task-specific conditional statements, written in natural language, that guide the matching process. The LLM can easily evaluate the validity of each condition for the data being matched, and then follow the corresponding directives. By walking through the pseudo-code, the LLM is able to reason its way to the final matching result. Unlike Chain-of-Thought (Wei et al., 2022), which relies on LLM’s internal knowledge to generate reasoning steps, KCMF uses explicit task-aware pseudo-code, eliminating the need to engineer prompts for different datasets, that is, pseudo-code is written only once for a task (such as schema matching).

To address hallucination, KCMF incorporates external knowledge in the form of natural language sentences related to the input data. To compensate for the lack of such unstructured domain knowledge, KCMF builds domain knowledge sets by leveraging readily available datasets (*Dataset as Knowledge*, DaK) and examples (*Examples as Knowledge*, EaK) from various domains. Since fine-tuning-free LLMs tend to generate improperly formatted outputs that do not match the format given in the demonstration, we employ a technique called *Inconsistency-tolerant Generation Ensemble* (IntGE) to suppress such unexpected outputs and

maintain an automated downstream workflow.

We evaluate KCMF on three SM datasets and four EM datasets. The comprehensive results show that KCMF significantly improves five LLM backbones, exceeds the current non-LLM best-performers (SMAT (Zhang et al., 2021) for SM and SUDOWOODO (Wang et al., 2023a) for EM), and achieves performance comparable to fine-tuned LLMs in certain cases, all without any fine-tuning.

1. We introduce KCMF, an LLM-based framework that incorporates a unified methodology to solve schema matching and entity matching tasks, eliminating the need for task-dependent model design and domain-specific fine-tuning.
2. We propose 1) a once-and-for-all pseudo-code-based task decomposition strategy across task types, to guide the LLM’s reasoning by obtaining task-specific conditional statements, reducing *confusion* during matching, 2) two mechanisms, namely DaK and EaK, to build knowledge sets in lack of unstructured domain knowledge, and 3) IntGE, a result-ensemble strategy designed to leverage multiple knowledge sources and suppress poorly formatted outputs.
3. We perform a comprehensive evaluation for KCMF and demonstrate its superiority against various popular SM and EM methods. Our analysis also shows that KCMF can generalize effectively across different LLMs.

2 Background and Task Definition

2.1 Related Work

Conventional *Schema Matching* (SM) approaches can generally be categorized into 1) *constraint-based methods*, which utilize attributes defined in database constraints to measure similarity among schemas (Alexe et al., 2010; Chen et al., 2018; Atzeni et al., 2019), and 2) *linguistic-based methods*, which leverage the semantic information contained in schema names or descriptions to construct mappings (Kettouch et al., 2017; Asif-Ur-Rahman et al., 2023). More recent studies have adopted deep neural networks to tackle SM. Zhang et al. (2021) use an attention-based BiLSTM with pre-trained word embeddings. Zhang et al. (2023b) leverage BERT (Kenton and Toutanova, 2019) to generate schema features, which are then used to train a linear classifier under semi-supervised learning. With the advent of LLMs, Narayan et al. (2022) were the first to apply LLMs to SM, using straightforward serialization with few-shot settings.

Recent LLM-based SM approaches, such as Zhang et al. (2024a), attempt to tackle this task using Supervised Fine-tuning (SFT). While SFT achieves strong results, it needs tremendous efforts to collect fine-tuning data and is hard to transfer to unseen domains. Another line of works, such as Sheetrit et al. (2024) and Parciak et al. (2024), utilize proprietary LLMs, which have shown promising results on several benchmarks. However, these approaches still suffer from high computational costs and produce indecisive outputs. Our approach aims to extend the scope of LLM-based SM by utilizing task-specific pseudo-code to guide LLMs’ predictions within a single-round inference. The design offers explicit criteria for the matching task, addressing the confusion problem mentioned earlier.

Entity Matching (EM) is often considered a downstream task of SM by many previous works (Nie et al., 2019; Li et al., 2020; Brunner and Stockinger, 2020; Barlaug and Gulla, 2021). Similar to SM, traditional non-LLM solutions for EM focus on computing similarity between entities (Thirumurganathan et al., 2018; Ebraheem et al., 2018; Kasai et al., 2019; Li et al., 2020). Following the initial application of LLM to EM (Narayan et al., 2022), LLM-based solutions have been actively explored (Fan et al., 2024; Peeters and Bizer, 2024; Li et al., 2024; Wang et al., 2024). From the outset, our approach is designed as a unified framework capable of addressing both SM and EM. By recognizing the deep semantic similarities and shared motivations between SM and EM — such as aligning and matching data elements — we develop a methodology that seamlessly integrates both tasks.

Retrieval Augmented Generation (RAG) incorporates retrieved documents into queries and has become a popular paradigm in mitigating hallucination in LLMs (Lewis et al., 2020; Yu et al., 2023; Shao et al., 2023; Jiang et al., 2023; Asai et al., 2023; Xu et al., 2023; Shi et al., 2024). While RAG has shown reliable results in addressing hallucination, its performance is limited by the quality of the retrieved documents and, obviously, is challenging to deploy in scenarios lacking unstructured knowledge sources. Our approach adopts the concept of introducing external knowledge to alleviate hallucination. In particular, to address the lack of unstructured domain knowledge in data matching, we present two mechanisms, DaK and EaK, for building knowledge sets by utilizing existing domain knowledge bases. Our approach showcases

the potential of external knowledge enhancement strategies on LLM-based classification tasks.

2.2 Task Definition

We denote by \mathcal{S} the *source schema* from database \mathcal{D} , and \mathcal{S}' the *target schema* from database \mathcal{D}' . The goal of **schema matching** (SM) is to identify all pairs of **attributes** $(\mathcal{A}, \mathcal{A}')$ such that $\mathcal{A} \in \mathcal{S}$ and $\mathcal{A}' \in \mathcal{S}'$, and both attributes represent the same information in their respective schemas.

In this paper, we focus on a more straightforward scenario: we enumerate all possible mappings $\mathcal{M} \subseteq \mathcal{D} \times \mathcal{D}'$ and determine whether each mapping $\{r, r'\} \in \mathcal{M}$ is correct or not. Each candidate mapping $\{r, r'\}$ consists of two items, where each item r is composed of a schema name N and a schema description C from the corresponding database \mathcal{D} .

We framework our schema matching task as an LLM generation task guided by pseudo-code instructions under a knowledge-enhanced setting. This involves 1) a list \mathcal{P} of designated task-specific pseudo-code and 2) a list K of knowledge items retrieved from all available knowledge sets \mathcal{K} .

We serialize the inputs, including the pseudo-code \mathcal{P} , candidate mapping $\{r, r'\}$, and retrieved knowledge items K , into a prompt tailored to the LLM \mathcal{L} using a function ϕ . An LLM-generated response LR is then obtained as:

$$LR \leftarrow \mathcal{L}(\phi(\mathcal{P}, \{r, r'\}, K)). \quad (1)$$

Our task objective is to classify the correctness of each candidate mapping. Specifically, we aim to obtain a binary classification result $c \in \{yes, no\}$ from the LLM’s response LR , indicating whether the mapping $\{r, r'\}$ is correct. For clarity of presentation, here we focus on defining SM, though it should be noticed that this setting can also be generalized to entity matching task (cf. Appendix A).

3 Matching Framework

3.1 Framework Overview

Figure 2 presents an overview of our proposed KCMF, which operates in four sequential stages. **S1. Pseudo-code Design:** Experts decompose the task into task-related conditional statements in natural language. A superior advantage of our approach is that this pseudo-code for a task is designed once and can be reused by anyone performing the matching task, without the need to write custom ones. **S2. Knowledge Retrieval & Construction:** KCMF constructs granular domain knowledge and database-structure knowledge via retriev-

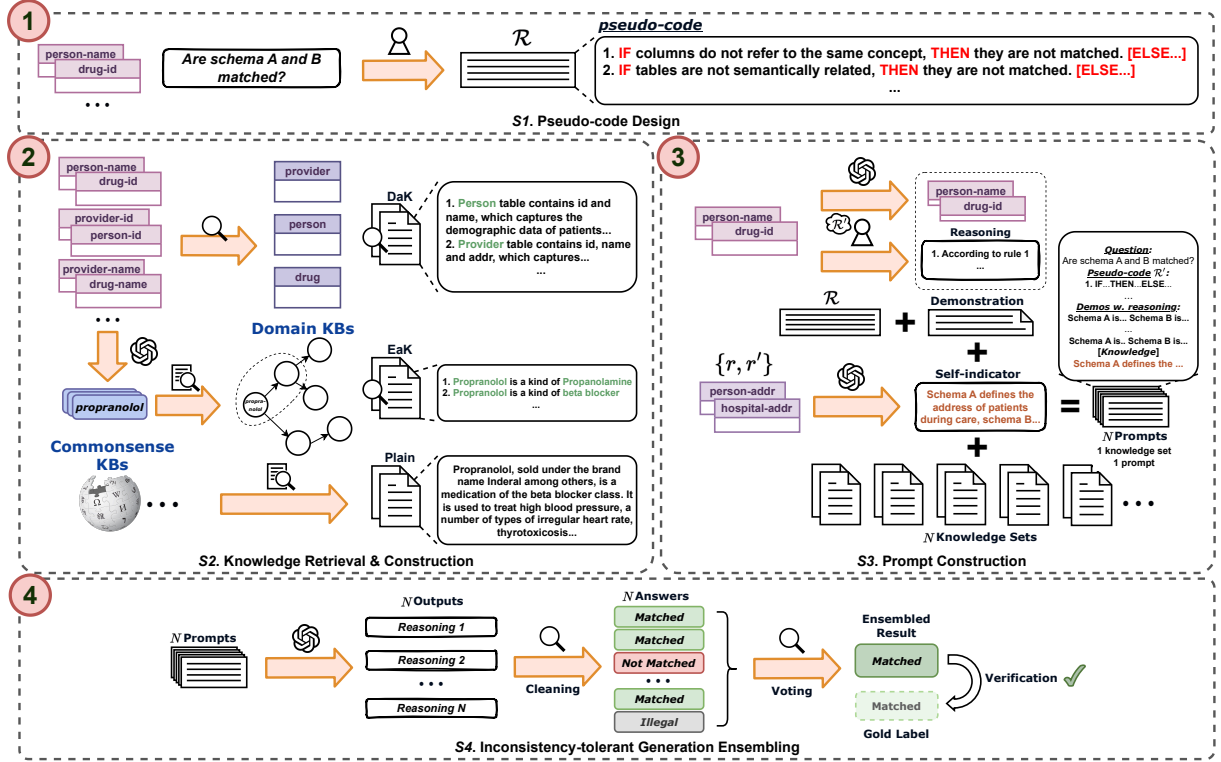


Figure 2: Overview of KCMF. Our carefully designed pseudo-code (detailed in Appendix B) offers a reusable and efficient solution for both SM and EM tasks. This eliminates the need for redesigning statements from scratch when working with new datasets, streamlining the process and enhancing adaptability.

ing information from domain knowledge bases and discovering database structures. **S3. Prompt Generation:** Utilizing the knowledge from S2 and pseudo-code from S1, KCMF generates LLM prompts that include demonstrations. **S4. Inconsistency-tolerant Generation Ensemble:** KCMF uses the constructed prompts to query the LLM. The multiple outputs generated are combined, mitigating the ill-formatting issue, to reach the final decision.

Illustrative examples of the main KCMF components in this section are provided in Appendix A.

3.2 Pseudo-code Design

As Figure 1 depicts, one of the challenges that LLMs face when performing matching is the issue of under-matching and over-matching. This issue stems from the ambiguity of the task instruction “match”. Without additional context, the word match lacks a clear definition in this setting¹. Take schema matching as an example. The term match in this context has at least three different interpretations, depending on the focus of the task:

- The *data types* of the column values are *the same*

¹As shown in Appendix F, there is a discrepancy between what GPT-3.5 understands the instruction match and the intended objective of the schema matching task.

or *convertible* to each other.

- Table definitions are *semantically related*.
- Schemas refer to *the same real-world concept*.

To address this ambiguity, we propose a task-aware strategy to decompose the task into pseudo-code composed of conditional statements directly derived from the task’s motivation. As shown in S1 of Figure 2, natural language predicates are structured into *if-then-else* constructs, each providing sufficient conditions to determine “matched” or “not matched” cases. The pseudo-code is designed to comprehensively cover all relevant task conditions. This process requires manual effort to understand the task’s motivation and design the pseudo-code \mathcal{R} (see implementation details in Appendix B). To apply this pseudo-code, we introduce a reasoning prompting strategy inspired by Chain-of-Thought (Wei et al., 2022), which will be discussed in section 3.4.

3.3 Building Domain Knowledge Set

To solve the matching task within a knowledge-enhanced setting, we first need to retrieve a domain knowledge list K . However, because answers for matching tasks cannot be explicitly derived from retrieved information and due to the scarcity

of unstructured domain knowledge, we propose two mechanisms: Dataset as Knowledge (DaK) and Examples as Knowledge (EaK). These mechanisms are designed to construct highly relevant unstructured knowledge sets tailored specifically for matching tasks. Further implementation details are provided in Appendix B.

Dataset as Knowledge (DaK) Due to privacy concerns, concrete records from the source databases \mathcal{D} and \mathcal{D}' are often inaccessible to the schema matching system (Johnson et al., 2023; Zhang et al., 2023b). This means that the matching task must be performed using only metadata. In this setting, LLMs are required to understand the structures of \mathcal{D} and \mathcal{D}' , which, however, can only be partly seen during inference through the limited metadata representations $\{r, r'\}$ (cf. Section 2.2). Hence, we propose DaK to acquire knowledge of the structure of source databases by searching metadata from the full dataset \mathcal{M} , as shown in Figure 3.

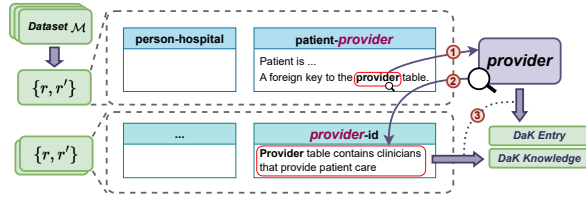


Figure 3: A toy example of DaK. Object “provider” and the description of which are identified, respectively (① & ②); then, they are integrated as an entry and a piece of DaK knowledge (③).

The procedure of DaK aims to discover database objects and their metadata from the candidate pool \mathcal{M} (see pseudo-code in Appendix A). As shown in Figure 3, DaK operates in three steps: 1) for a given candidate pair $\{r, r'\} \in \mathcal{M}$, DaK identifies an object O ; in the example, O is a table named “provider”. While this step can utilize techniques like Named Entity Recognition (NER), for rapid prototyping, we extract O by traversing all schema names N and N' associated with r, r' (see ①). 2) DaK then scans the remaining pairs $\{r, r'\}$ to find metadata K_{DaK} related to the identified object O ; in this example, K_{DaK} is the description of the table “provider”. This is accomplished by matching O with descriptions C from the other candidate pairs (see ②). 3) After obtaining an object list \mathcal{C}_{Obj} consisting of each identified object O , and a metadata list \mathcal{C}_{DaK} consisting of all matched knowledge K_{DaK} , DaK proceeds to form the name of O as an entry and K_{DaK} of O as the corresponding DaK

knowledge. As shown in Figure 3, the name of the identified object “provider” and the description of “provider” are constructed into a piece of DaK knowledge: $\{\text{provider: provider table contains clinicians that provide patient care}\}$.

Examples as Knowledge (EaK) Using text chunks from a self-built corpus has become a *de facto* approach for retrieval enhancement to reduce LLM hallucinations (Gao et al., 2024). For common-sense QA tasks, such text chunks can be easily retrieved from existing commonsense corpora. However, the knowledge required for matching tasks is domain-specific, and existing KBs are mainly structured as entity databases or thesaurus rather than natural language text. Thus, we propose EaK, as depicted in Figure 4.

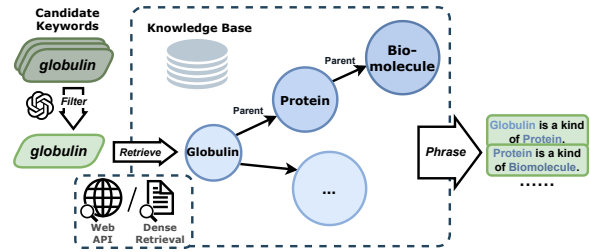


Figure 4: An example of EaK.

As its name suggests, EaK aims to explain complex concepts using examples. Given a pair $\{r, r'\}$, EaK first extracts keywords from it by querying an LLM, and then queries an LLM again to filter the domain-irrelevant ones (see pseudo-code in Appendix G). Obtaining all candidate keywords, EaK uses these candidates to search domain KBs for the top- k related records. For remote KBs, retrieval can be done through a search over the provided Web API; while for local KBs, this procedure can be implemented using stronger dense retrieval. For each retrieved record, EaK leverages its relationships and properties to form explanatory knowledge. For example, in Figure 4, EaK generates the knowledge by phrasing the parent relation.

3.4 Prompt Construction

Pseudo-code-based Reasoning Prompting In a paradigm of k -shot in-context learning, we sample k valid match pairs $\{d, d'\}$ as demonstrations together with the pseudo-code \mathcal{R} , the queried pair $\{r, r'\}$, and knowledge list K . The fundamental pattern of the prompt is illustrated in Figure 5.

There are two placeholders in the prompt, $\{\mathcal{C}_{\text{RSNG}}\}$ for the reasoning steps and $\{ans\}$ for the

Prompt
Task Instruction: Are A and B matched?
Pseudo-code \mathcal{R}: 1) If columns are not the same type of data, then the answer is no, else check the next statement...
Demonstrations $\{\{d, d'\}\}$: Schema A is drug-code, B is drug-id $\{C_{RSNG}\}\{ans\}$. Schema A is provider-id, B is drug-id $\{C_{RSNG}\}\{ans\}$.
Target Data $\{r, r'\}$: Schema A is patients-birthdate, B is person-month_of_birth.
Knowledge: Date type is a basic SQL data type.

Figure 5: A toy prompt combining all outcomes from previous sections and placeholders $\{C_{RSNG}\}$ and $\{ans\}$.

answer. To generate C_{RSNG} and ans , we first define each single statement in pseudo-code \mathcal{R} (cf. Section 3.2) as $p \rightarrow q$. Then, we check each statement sequentially: for the current statement $p \rightarrow q$, the condition p is checked if it is satisfied; if p is fulfilled, terminate the process and set the answer $ans = q$; if not, proceed to the next statement; continue this process until an answer is drawn (see the reasoning construction algorithm in Appendix A). Next, the obtained C_{RSNG} and ans will be inserted into prompt. This process is repeated for all demonstrations to obtain the full set of reasoning steps. Notably, the pseudo-code can be extended for complex cases with multiple reasoning paths by adding an "ELSE" clause for alternative statements.

Self-Indicator Extraction Inspired by the performance gain observed when generating a high-quality summary at the beginning of reasoning in our experiments, we further add a Self-Indicator Extraction pre-task. This plugin module generates a filtered text segment, called self-indicator (K_{SI}), describing key information from $\{r, r'\}$. Specifically, taking $\{r, r'\}$ and K as input, an LLM is used to generate K_{SI} (cf. Appendix B), which is then appended to the prompts created above.

Summarized Demonstrations After decomposing task instructions into pseudo-code, we leverage in-context learning to enable the LLM to learn reasoning behaviors beyond the concept of "match". To manage the prompt length, we use an LLM to summarize all demonstration pairs $\{d, d'\}$ (cf. Appendix B), and use *summarized* pairs to improve efficiency without losing essential information.

3.5 Inconsistency-tolerant Generation Ensemble (IntGE)

With multiple knowledge sets available, a straightforward utilization is to combine all retrieved knowledge within one prompt, but this practice leads to information within a prompt flooded thus

bringing unexpected outputs (Parciak et al., 2024). Therefore, we instead propose IntGE to integrate diverse information sources and improve output stability in a *fine-tuning-free* manner. We define the available knowledge sets from unique sources as $\mathcal{K} = [K_1, K_2, \dots, K_n]$. For a given pair $\{r, r'\}$, we retrieve information from each source in \mathcal{K} , resulting in n knowledge lists. Each list is used to construct a prompt following the method in Section 3.4, creating n prompts. We ensemble the n binary classification results using majority voting to determine the prediction for $\{r, r'\}$.

As noted in previous studies (Cuconasu et al., 2024; Zhang et al., 2024b), LLMs tend to prefer familiar input from their pre-training phase, and distracting context can lead to poor outputs. IntGE aims to prevent interference between different knowledge sources by separating them into distinct prompts. Inspired by Wang et al. (2023b); Shi et al. (2024), we adopt a straightforward voting strategy. This approach leverages multiple knowledge sources while enhancing the model’s robustness against poorly formatted outputs by keeping prompt lengths manageable.

4 Experiments

4.1 Experimental Setup

Dataset and Metrics Following Narayan et al. (2022), we choose the challenging Synthea, MIMIC, and CMS from the OMAP benchmark (Zhang et al., 2021) for SM. For EM, we adapt the MedMentions benchmark (Mohan and Li, 2019) to construct a new biomedical dataset MedMentions Matching (MMM), with details in Appendix B. To assess KCMF’s generalization beyond the biomedical domain, we include widely used out-of-domain EM datasets: Amazon-Google (AG), Walmart-Amazon (WA), DBLP-Scholar (DS) from DeepM (Mudgal et al., 2018).

Following Zhang et al. (2021); Narayan et al. (2022), we report F1-score across all pairs for its suitability in the *inherently imbalanced* data matching tasks and include accuracy to show precision and recall trends after applying KCMF. Detailed precision and recall results are in Appendix C.

Settings We mainly follow Section 3 to conduct experiments. However, for AG, WA and DS in EM, those derived from e-commerce and publication-profile data lacking readily available knowledge bases, we test KCMF *without* knowledge construction & retrieval. Un-

Methods	Acc	MIMIC F1	Acc	Synthea F1	Acc	CMS F1
SMAT	0.9865	0.2020	0.9902	0.3850	0.9852	0.5000
GLM-4-Flash +KCMF	0.9512 0.9816 _{+3.04%}	0.0369 0.1061 _{+6.92%}	0.9895 0.9852 _{+0.43%}	0.1143 0.1538 _{+3.95%}	0.9922 0.9836 _{-0.86%}	0.1667 0.2500 _{+8.33%}
Llama-3-8B +KCMF	0.3272 0.9774 _{+65.02%}	0.0060 0.0765 _{+7.05%}	0.8870 0.9572 _{+7.02%}	0.0233 0.0863 _{+6.30%}	0.9926 0.9481 _{-4.45%}	0.0952 0.1192 _{+2.40%}
Mistral-7B +KCMF	0.9931 0.9947 _{+0.16%}	0.1538 0.3200 _{+16.62%}	0.9686 0.9895 _{+2.09%}	0.0792 0.1143 _{+3.51%}	0.9879 0.9867 _{-0.12%}	0.2051 0.2273 _{+2.22%}
GPT-3.5 +KCMF	0.9767 0.9944 _{+1.77%}	0.0745 0.6207 _{+54.62%}	0.9865 0.9936 _{+0.71%}	0.0909 0.4242 _{+33.33%}	0.7656 0.9867 _{+22.11%}	0.0196 0.2917 _{+27.21%}
GPT-4o +KCMF	0.9964 0.9959 _{-0.05%}	0.3030 0.3500 _{+4.70%}	0.9919 0.9946 _{+0.27%}	0.0769 0.4667 _{+38.98%}	0.9922 0.9887 _{+0.20%}	0.1667 0.4082 _{+24.15%}

Table 1: The results for the schema matching task.

less otherwise noted, LLMs utilized in KCMF are GPT-3.5-turbo-1106. Detailed settings and all prompts, including pseudo-code, are in Appendix A and G. Given KCMF’s fine-tuning-free nature, we mainly compare it with (1) five *few-shot LLM baselines* including *API-styled* GPT-3.5-turbo-1106, GPT-4o, and GLM-4-Flash and *locally-deployed* Llama-3-8B-Instruct and Mistral-7B-Instruct; and (2) *non-LLM competitors*, namely SMAT (Zhang et al., 2021) achieving the best SM performance on OMAP (by classifying schema pairs just like KCMF) and SUDOWODO (Wang et al., 2023a), a self-supervised SoTA methods for EM tasks.

4.2 Main Results for SM and EM

KCMF achieves superior performance in SM As shown in Table 1, generally KCMF is significantly better than all LLM baselines and SMAT for SM. Though this superiority is less pronounced on CMS, a notable performance improvement over LLM baselines is still evident. Specifically, on Synthea, GPT-3.5-version KCMF outperforms the previous SMAT by 3.92% in F1-score, and on MIMIC, the gain increases to 41.87%. In addition, we observe that the F1-score of GPT-3.5-version KCMF outperforms all baselines by an average of 33.73% and 50.59% on Synthea and MIMIC, respectively. On MMM, the F1-score of KCMF surpasses all LLM baselines by an average of 10.62%, and on CMS this gain reaches 12.86%.

KCMF ensures positive sample recall in SM Although we see a slight drop in accuracy for GPT-4o-version KCMF, this 0.05% drop is not caused by a failure in recall, as detailed precision and recall in Appendix C show. Instead, it results from an increase in both true-positive and false-positive recalls, specifically, the drop is due to the correct

recalling of 2 extra positive samples and incorrect recalling of 5 negative samples, which looks reasonable for the 2 more true-positive raise the recall from 38.5% to 53.8% while keeping the precision around 25.0%, so we still consider this a strong result for an inherently imbalanced matching task.

KCMF generalizes well on out-of-domain EM

As shown in Table 2, KCMF consistently enhances the performance of all LLMs for EM, aligning with the SM results in Table 1. Even without knowledge retrieval, KCMF remains competitive with the semi-supervised method SUDOWODO, highlighting its strong generalizabilities across commonly used EM datasets.

KCMF generalizes across backbones for both SM and EM

We observe that after implementing KCMF, the F1-score of each backbone is improved by an average of 17.98%, 17.21%, 8.91%, 6.44%, 11.69%, 24.71%, 8.54% on MIMIC, Synthea, CMS, MMM, AG, WA, and DS, respectively. These results demonstrate that our proposed KCMF can generalize across different backbones and consistently improve the performance of LLMs on both SM and EM tasks. Notably, on MIMIC, the strongest version of KCMF achieves an F1-score of 0.6207, clearly outperforming a fine-tuned SoTA JELLYFISH (Zhang et al., 2024a) for SM, which achieves 0.4314. This result indicates the great potential of fine-tuning-free LLMs in data matching tasks on domain-specific data.

4.3 Component Analysis

Task-specific pseudo-code disambiguates task instructions

The pseudo-code is designed to prevent LLM confusion about task objectives by providing explicit conditional statements. To explore its effectiveness, we examine KCMF on Synthea

Methods	MMM		AG		WA		DS	
	Acc	F1	Acc	F1	Acc	F1	Acc	F1
SUDOWODO	-	-	0.9492	0.5045	0.9532	0.5236	0.9089	0.7706
GLM-4-Flash	0.9663	0.6841	0.6140	0.3410	0.4256	0.2332	0.9195	0.8009
+KCMF	0.9690	0.7334	0.8203	0.5164	0.8404	0.5038	0.9646	0.9035
Δ	+0.27%	+4.92%	+20.63%	+17.54%	+41.48%	+27.06%	+4.51%	+10.26%
Llama-3-8B	0.9574	0.6738	0.7553	0.2761	0.2660	0.1861	0.8377	0.6799
+KCMF	0.9730	0.7783	0.6594	0.3561	0.8023	0.4692	0.9093	0.7637
Δ	+1.56%	+10.45%	-9.59%	+8.00%	+53.63%	+28.31%	+7.16%	+8.38%
Mistral-7B	0.9671	0.7046	0.7479	0.4289	0.4300	0.2275	0.8555	0.6318
+KCMF	0.9689	0.7206	0.7771	0.4511	0.8824	0.5594	0.8948	0.6595
Δ	+0.18%	+1.60%	+2.92%	+2.22%	+45.24%	+33.19%	+3.93%	+2.77%
GPT-3.5	0.9297	0.6242	0.8731	0.2988	0.8380	0.4071	0.8915	0.6245
+KCMF	0.9727	0.7758	0.8809	0.5991	0.9517	0.7288	0.9300	0.8260
Δ	+4.30%	+15.16%	+0.78%	+30.03%	+11.37%	+32.17%	+3.85%	+20.15%
GPT-4o	0.9658	0.6738	0.9097	0.6810	0.9688	0.8469	0.9613	0.8920
+KCMF	0.9661	0.6744	0.9128	0.6875	0.9746	0.8750	0.9653	0.9035
Δ	+0.03%	+0.06%	+0.31%	+0.65%	+0.58%	+2.81%	+0.40%	+1.15%

Table 2: The results for the entity matching task.

without our proposal: all pseudo-code and corresponding reasoning steps are removed from the demonstrations, and the results are compared with the complete KCMF. We classify false-positive errors into three types: *Over-matching* (OM), as described in Figure 1; *Position Mismatching* (PM), where parts of r and r' (e.g., a schema’s table and another’s column) are incorrectly matched; and *Incorrect Reasoning* (IR), covering errors not fitting the first two categories. Table 3 shows that the pseudo-code reduces *all* three error types.

Settings	IR	# (Case) OM	PM
w/o pseudo-code	8	60	11
w. pseudo-code	5	13	1
Δ	-37.5%	-78.3%	-90.9%

Table 3: Number of erroneous cases with and without the proposed pseudo-code on Synthea.

IntGE enhances robustness against badly-formatted outputs To assess the impact of IntGE, we test a *knowledge-all-in-one* (AIO) version of KCMF, which removes IntGE and combines knowledge from all sources into a single prompt. As shown in Figure 6, results are categorized into three types: *Badly-formatted*, where the output does not follow the demonstrated format; *Well-formatted*, where it does; and *Eliminated*, where IntGE’s voting mechanism removes poorly formatted outputs, ensuring the final result is well-formatted. The results clearly indicate that the AIO setting performed poorly. In contrast, IntGE effectively reduces the number of badly-formatted output via its multi-prompt voting mechanism. The results imply a paradigm for retrieving knowledge

from different sources for inference-only LLMs, i.e., splitting retrieved knowledge into separate prompts based on the source to improve output quality while managing prompt length.

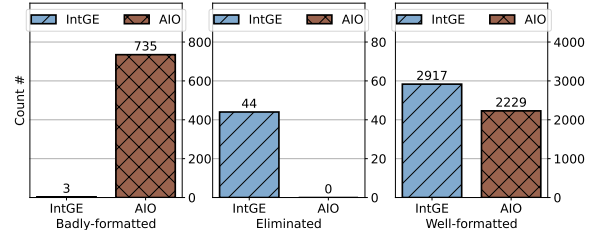


Figure 6: Ablations of IntGE on Synthea.

In addition to the above analysis, detailed ablation results in Appendix D and E further confirm the necessity and effectiveness of each component of KCMF.

5 Conclusion

We presented KCMF, a fine-tuning-free framework for solving data matching tasks under a knowledge-enhanced setting, with a pseudo-code based task instruction strategy for addressing the problem of confusion, and mechanisms DaK, EaK, and IntGE for building and utilizing knowledge from various sources in scarcity of unstructured domain knowledge. Our evaluations show that KCMF enhances five LLM backbones and outperforms the current non-LLM solutions in most cases in a completely fine-tuning-free paradigm, and our pseudo-code effectively mitigated LLMs’ confusion towards the matching tasks. Our future work involves extending KCMF to make it suitable for both metadata-based and instance-based scenarios, and further improving the efficiency of the matching workflow.

6 Limitations

One limitation of our approach is the reliance on human effort to select appropriate domain knowledge bases (KBs) as sources of information, which may pose challenges for practical implementation. Also, this study primarily focuses on matching tasks framed as predicting enumerated data pairs rather than directly identifying potential mappings within source databases. While this approach simplifies the problem formulation, it may lead to increased computational demands and thus have an environmental impact. Although cost is a concern, our cose analysis in Appendix E shows the cost-in-total remains acceptable. Nonetheless, we recognize the unexplored opportunity for efficiency optimization in future work.

References

- Bogdan Alexe, Mauricio Hernández, Lucian Popa, and Wang-Chiew Tan. 2010. Mapmerge: Correlating independent schema mappings. *Proceedings of the VLDB Endowment*, 3(1-2):81–92.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. [Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection](#). *arXiv preprint*. ArXiv:2310.11511 [cs].
- Md Asif-Ur-Rahman, Bayzid Ashik Hossain, Michael Bewong, Md Zahidul Islam, Yanchang Zhao, Jeremy Groves, and Rory Judith. 2023. [A semi-automated hybrid schema matching framework for vegetation data integration](#). *Expert Systems with Applications*, 229:120405.
- Paolo Atzeni, Luigi Bellomarini, Paolo Papotti, and Riccardo Torlone. 2019. [Meta-mappings for schema mapping reuse](#). *Proceedings of the VLDB Endowment*, 12(5):557–569.
- Nils Barlaug and Jon Atle Gulla. 2021. Neural networks for entity matching: A survey. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 15(3):1–37.
- Ursin Brunner and Kurt Stockinger. 2020. Entity matching with transformer architectures—a step forward in data integration. In *23rd International Conference on Extending Database Technology, Copenhagen, 30 March-2 April 2020*, pages 463–473. OpenProceedings.
- Chen Chen, Behzad Golshan, Alon Halevy, Wang-Chiew Tan, and AnHai Doan. 2018. BigGorilla: An Open-Source Ecosystem for Data Preparation and Integration.
- Florin Cuconasu, Giovanni Trappolini, Federico Siciliano, Simone Filice, Cesare Campagnano, Yoelle

- Maarek, Nicola Tonellotto, and Fabrizio Silvestri. 2024. The power of noise: Redefining retrieval for rag systems. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 719–729.
- Muhammad Ebraheem, Saravanan Thirumuruganathan, Shafiq Joty, Mourad Ouzzani, and Nan Tang. 2018. [DeepER – Deep Entity Resolution](#). *Proceedings of the VLDB Endowment*, 11(11):1454–1467. ArXiv:1710.00597 [cs].
- Meihao Fan, Xiaoyue Han, Ju Fan, Chengliang Chai, Nan Tang, Guoliang Li, and Xiaoyong Du. 2024. Cost-effective in-context learning for entity resolution: A design space exploration. In *2024 IEEE 40th International Conference on Data Engineering (ICDE)*, pages 3696–3709. IEEE.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. 2024. [Retrieval-Augmented Generation for Large Language Models: A Survey](#). *arXiv preprint*. ArXiv:2312.10997 [cs].
- Zhengbao Jiang, Frank Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. [Active Retrieval Augmented Generation](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7969–7992, Singapore. Association for Computational Linguistics.
- Alistair E. W. Johnson, Lucas Bulgarelli, Lu Shen, Alvin Gayles, Ayad Shammout, Steven Horng, Tom J. Pollard, Sicheng Hao, Benjamin Moody, Brian Gow, Liwei H. Lehman, Leo A. Celi, and Roger G. Mark. 2023. [MIMIC-IV, a freely accessible electronic health record dataset](#). *Scientific Data*, 10(1):1.
- Jungo Kasai, Kun Qian, Sairam Gurajada, Yunyao Li, and Lucian Popa. 2019. Low-resource deep entity resolution with transfer and active learning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5851–5861.
- Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of naacL-HLT*, volume 1, page 2. Minneapolis, Minnesota.
- Mohamed Kettouch, Cristina Luca, Mike Hobbs, and Sergiu Dascalu. 2017. [Using semantic similarity for schema matching of semi-structured and linked data](#). Pages: 133.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.

712	Huahang Li, Shuangyin Li, Fei Hao, Chen Jason Zhang,	Eitam Sheerit, Menachem Brief, Moshik Mishaeli, and	765
713	Yuanfeng Song, and Lei Chen. 2024. Booster: Lever-	Oren Elisha. 2024. ReMatch: Retrieval Enhanced	766
714	aging large language models for enhancing entity	Schema Matching with LLMs . <i>arXiv preprint</i> .	767
715	resolution. In <i>Companion Proceedings of the ACM</i>	ArXiv:2403.01567 [cs].	768
716	on <i>Web Conference 2024</i> , pages 1043–1046.		
717	Yuliang Li, Jinfeng Li, Yoshihiko Suhara, AnHai	Weijia Shi, Sewon Min, Michihiro Yasunaga, Min-	769
718	Doan, and Wang-Chiew Tan. 2020. Deep Entity	joon Seo, Richard James, Mike Lewis, Luke Zettle-	770
719	Matching with Pre-Trained Language Models . <i>Pro-</i>	moyer, and Wen-tau Yih. 2024. Replug: Retrieval-	771
720	<i>ceedings of the VLDB Endowment</i> , 14(1):50–60.	augmented black-box language models. In <i>Proceed-</i>	772
721	ArXiv:2004.00584 [cs].	<i>ings of the 2024 Conference of the North American</i>	773
722	Sunil Mohan and Donghui Li. 2019. MedMentions:	<i>Chapter of the Association for Computational Lin-</i>	774
723	A Large Biomedical Corpus Annotated with UMLS	<i>guistics: Human Language Technologies (Volume 1:</i>	775
724	Concepts . <i>arXiv preprint</i> . ArXiv:1902.09476 [cs].	<i>Long Papers)</i> , pages 8364–8377.	776
725	Sidharth Mudgal, Han Li, Theodoros Rekatsinas, An-	Saravanan Thirumuruganathan, Shameem A. Puthiya	777
726	Hai Doan, Youngchoon Park, Ganesh Krishnan, Ro-	Parambath, Mourad Ouzzani, Nan Tang, and Shafiq	778
727	hit Deep, Esteban Arcaute, and Vijay Raghavendra.	Joty. 2018. Reuse and Adaptation for Entity Res-	779
728	2018. Deep Learning for Entity Matching: A Design	olution through Transfer Learning . <i>arXiv preprint</i> .	780
729	Space Exploration . In <i>Proceedings of the 2018 Inter-</i>	ArXiv:1809.11084 [cs, stat].	781
730	<i>national Conference on Management of Data</i> , pages	Runhui Wang, Yuliang Li, and Jin Wang. 2023a. Su-	782
731	19–34, Houston TX USA. ACM.	dowoodo: Contrastive Self-supervised Learning for	783
732	Avanika Narayan, Ines Chami, Laurel Orr, and Christo-	Multi-purpose Data Integration and Preparation . In	784
733	pher Ré. 2022. Can foundation models wrangle	<i>2023 IEEE 39th International Conference on Data</i>	785
734	your data? <i>Proceedings of the VLDB Endowment</i> ,	<i>Engineering (ICDE)</i> , pages 1502–1515, Anaheim,	786
735	16(4):738–746.	CA, USA. IEEE.	787
736	Hao Nie, Xianpei Han, Ben He, Le Sun, Bo Chen,	Tianshu Wang, Xiaoyang Chen, Hongyu Lin, Xuanang	788
737	Wei Zhang, Suhui Wu, and Hao Kong. 2019. Deep	Chen, Xianpei Han, Hao Wang, Zhenyu Zeng, and	789
738	Sequence-to-Sequence Entity Matching for Heteroge-	Le Sun. 2024. Match, Compare, or Select? An Inves-	790
739	neous Entity Resolution . In <i>Proceedings of the 28th</i>	tigation of Large Language Models for Entity Match-	791
740	<i>ACM International Conference on Information and</i>	ing . <i>arXiv preprint</i> . ArXiv:2405.16884 [cs].	792
741	<i>Knowledge Management</i> , pages 629–638, Beijing	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le,	793
742	China. ACM.	Ed Chi, Sharan Narang, Aakanksha Chowdhery, and	794
743	Marcel Parciak, Brecht Vandevoort, Frank Neven,	Denny Zhou. 2023b. Self-Consistency Improves	795
744	Liesbet M. Peeters, and Stijn Vansummeren. 2024.	Chain of Thought Reasoning in Language Models .	796
745	Schema Matching with Large Language Mod-	<i>arXiv preprint</i> . ArXiv:2203.11171 [cs].	797
746	els: an Experimental Study . <i>arXiv preprint</i> .	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	798
747	ArXiv:2407.11852 [cs].	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	799
748	Ralph Peeters and Christian Bizer. 2023. Using chatgpt	et al. 2022. Chain-of-thought prompting elicits rea-	800
749	for entity matching. In <i>European Conference on Ad-</i>	soning in large language models. <i>Advances in neural</i>	801
750	<i>vances in Databases and Information Systems</i> , pages	<i>information processing systems</i> , 35:24824–24837.	802
751	221–230. Springer.	Fangyuan Xu, Weijia Shi, and Eunsol Choi. 2023. RE-	803
752	Ralph Peeters and Christian Bizer. 2024. Entity Match-	COMP: Improving Retrieval-Augmented LMs with	804
753	ing using Large Language Models . <i>arXiv preprint</i> .	Compression and Selective Augmentation . <i>arXiv</i>	805
754	ArXiv:2310.11244 [cs].	<i>preprint</i> . ArXiv:2310.04408 [cs].	806
755	Vipula Rawte, Amit Sheth, and Amitava Das. 2023. A	Wenhao Yu, Hongming Zhang, Xiaoman Pan, Kaixin	807
756	Survey of Hallucination in Large Foundation Models .	Ma, Hongwei Wang, and Dong Yu. 2023. Chain-	808
757	<i>arXiv preprint</i> . ArXiv:2309.05922 [cs].	of-Note: Enhancing Robustness in Retrieval-	809
758	Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie	Augmented Language Models . <i>arXiv preprint</i> .	810
759	Huang, Nan Duan, and Weizhu Chen. 2023. En-	ArXiv:2311.09210 [cs].	811
760	hancing Retrieval-Augmented Large Language Mod-	Alexandros Zeakis, George Papadakis, Dimitrios Sk-	812
761	els with Iterative Retrieval-Generation Synergy . In	outas, and Manolis Kourarakis. 2023. Pre-Trained	813
762	<i>Findings of the Association for Computational Lin-</i>	Embeddings for Entity Resolution: An Experimen-	814
763	<i>guistics: EMNLP 2023</i> , pages 9248–9274, Singapore.	tal Analysis . <i>Proceedings of the VLDB Endowment</i> ,	815
764	Association for Computational Linguistics.	16(9):2225–2238.	816
		Haochen Zhang, Yuyang Dong, Chuan Xiao, and Masa-	817
		fumi Oyamada. 2024a. Jellyfish: A Large Lan-	818
		guage Model for Data Preprocessing . <i>arXiv preprint</i> .	819
		ArXiv:2312.01678 [cs].	820

- Jing Zhang, Bonggun Shin, Jinho D. Choi, and Joyce C. Ho. 2021. [SMAT: An Attention-Based Deep Learning Solution to the Automation of Schema Matching](#). In Ladjel Bellatreche, Marlon Dumas, Panagiotis Karras, and Raimundas Matulevičius, editors, *Advances in Databases and Information Systems*, volume 12843, pages 260–274. Springer International Publishing, Cham. Series Title: Lecture Notes in Computer Science.
- Wei Zhang, Hao Wei, Bunyamin Sisman, Xin Luna Dong, Christos Faloutsos, and David Page. 2020. Autoblock: A hands-off blocking framework for entity matching. In *Proceedings of the 13th International Conference on Web Search and Data Mining*, pages 744–752.
- Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. 2023a. [Siren’s Song in the AI Ocean: A Survey on Hallucination in Large Language Models](#). *arXiv preprint*. ArXiv:2309.01219 [cs].
- Yuji Zhang, Sha Li, Jiateng Liu, Pengfei Yu, Yi R. Fung, Jing Li, Manling Li, and Heng Ji. 2024b. [Knowledge Overshadowing Causes Amalgamated Hallucination in Large Language Models](#). *arXiv preprint*. ArXiv:2407.08039 [cs].
- Yunjia Zhang, Avriella Floratou, Joyce Cahoon, Subru Krishnan, Andreas C. Müller, Dalitso Banda, Fotis Psallidas, and Jignesh M. Patel. 2023b. [Schema Matching using Pre-Trained Language Models](#). In *2023 IEEE 39th International Conference on Data Engineering (ICDE)*, pages 1558–1571. ISSN: 2375-026X.

A Methodology

Appendix A.1 provides examples illustrating the key concepts from the main text used to construct the final prompt. Appendix A.2 discusses the adaptations required for applying KCMF to entity matching tasks. Appendix A.3 presents algorithms formalized based on the main procedure described in Section 3.

A.1 Examples

The following examples correspond to the KCMF techniques described in Section 3 and illustrated in Figure 2.

Schema Description C
For schema A <u>patients-birthdate</u>: patient demographic data; the date the patient was born.
Schema Description C
For schema B <u>person-month_of_birth</u>: {NL description of table <u>person</u> }; {NL description of attribute <u>month_of_birth</u> }
Self-indicator K_{SI}
From the descriptions and knowledge, I know that schema A is the date of birth for the patient in the patient table, and schema B is the month of birth for the patient in the person table.
Candidate Pair $\{r, r'\}$
Schema A: patient-birthdate Description of schema A: $\{C \text{ of patient-birthdate}\}$ Schema B: person-month_of_birth Description of schema B: $\{C \text{ of person-month_of_birth}\}$

Figure 7: Examples of schema description C , self-indicator K_{SI} , and candidate pair $\{r, r'\}$ based on the schema *patients-birthdate* and *person-month_of_birth*.

In-context Demonstration
$\{r, r'\}$ K_{SI} 2. Based on rule I, I must check each rule sequentially. 3. The DATE of BIRTH and the MONTH of BIRTH are both date-related data types, I need to check rule III. 4. The PERSON table and the PATIENT table both refer to individuals in the healthcare system, I need to check rule IV. 5. The MONTH of BIRTH and the DATE of BIRTH both refer to the birth information of a patient, based on rule IV, the answer is yes.

Figure 8: An example of the in-context learning demonstration, i.e., an example of the [Examples] in Table 26, directly derived from the components from Figure 7.

A.2 Discussions on Entity Matching (EM)

Since entity matching operates at the level of records, the task of entity matching can be defined

by slightly modifying the definition of r described in Section 2.2 for schema matching:

$$r = \{N, attr_1, attr_2, \dots, attr_m\}, \quad (2)$$

where N represents the name of the entity r , but the schema description C is replaced with a sequence of attributes $\{attr_1, attr_2, \dots, attr_m\}$, which corresponds to all the other column values of the entity r .

For KCMF, the only difference between solving schema matching (SM) and entity matching (EM) tasks is the number of attributes requiring serialization. In SM, KCMF only needs to serialize schema name and description, whereas in EM, KCMF needs to serialize the entity name and potentially a variable number of attributes. In addition, since SM operates at the schema level while EM at the record level, EM datasets are often much larger in scale in practice. As a result, introducing *blocking* as a preprocessing phase, as commonly performed in conventional EM methods (Ebraheem et al., 2018; Zhang et al., 2020), could achieve higher efficiency in handling those larger EM datasets.

A.3 Algorithms

Algorithm 1 formalizes from the metadata discovering procedure introduced in Section 3.3. As discussed in Section 3.3, the functions ObjectDiscovery() and FindMetadata() can be implemented using readily available Named Entity Recognition (NER) models or LLMs.

Algorithm 1 Metadata Discovery

Input: A matching task dataset D
Output: A metadata list C_{DaK} and an object list C_{Obj}

```

1:  $C_{DaK} \leftarrow \emptyset, C_{Obj} \leftarrow \emptyset$ 
2: for  $\{r_i, r'_i\} \in D$  do
3:    $O \leftarrow \text{ObjectDiscovery}(r_i, r'_i)$ 
4:   for  $\{r_j, r'_j\} \in D, j \neq i$  do
5:      $K_{DaK} \leftarrow \text{FindMetadata}(O, r_j, r'_j)$ 
6:     if  $K_{DaK} \neq \text{NULL}$  then
7:        $C_{DaK}.\text{insert}(K_{DaK})$ 
8:        $C_{Obj}.\text{insert}(O)$ 
9:     end if
10:  end for
11: end for
12: return  $C_{DaK}, C_{Obj}$ 

```

Algorithm 2 formalizes the scan-then-check procedure over statements while constructing prompts, which has been discussed in Section 3.4.

Here, in Algorithm 2, getQ() takes a statement R as input and returns its conclusion (the q from $p \rightarrow q$).

Algorithm 2 Constructing Reasoning Steps

Input: Current demonstration $\{d, d'\}$ and the pseudo-code \mathcal{R}
Output: A reasoning steps C_{RSNG} and a corresponding final answer ans

```
1:  $C_{\text{RSNG}} \leftarrow \emptyset$ 
2: for  $R \in \mathcal{R}'$  do
3:    $C_{\text{RSNG}}.append(\text{genReasoningStep}(R, d, d'))$ 
4:   if  $\text{isConditionMet}(R, d, d')$  then or
      $\text{isLastStatement}(R, \mathcal{R})$  then
5:      $ans \leftarrow \text{getQ}(R)$ 
6:     return  $C_{\text{RSNG}}, ans$ 
7:   end if
8: end for
```

B Implementation

In this section, we present detailed information on implementing KCMF, specifically, B1. guidelines for designing the pseudo-code for a specific task (i.e., SM or EM), B2. the knowledge sources and the models used for knowledge retrieval and construction, B3. the models and parameters used for pre-tasks (e.g., self-indicator extraction) during prompt construction, B4. the models and parameters used in inference, B5. the modifications made to the MedMentions dataset to construct MMM, and B6. the accessibility datasets and code used in the experiments.

B1. Pseudo-code Design The pseudo-code used in our experiments for SM and EM is shown in Table 24 and Table 25, respectively. It was directly authored based on the objectives of the matching tasks.

As emphasized in the main text, the pseudo-code statements are defined as “IF-THEN-ELSE”. To better utilize the LLM’s natural language comprehension capability, this structure is adapted to “IF, THEN, otherwise” in the actual design of the pseudo-code. To ensure that the given pseudo-code is checked in the correct order, a statement is added at the beginning of the designed pseudo-code to explicitly specify that the subsequent statements should be checked in order. This precautionary step can be omitted when using more advanced LLMs (e.g., GPT-4o).

Additionally, to cover all the task motivations of the matching task, we include overlapped conditions in the pseudo-code statement design. For the last statement of all the pseudo-codes, we set a general condition to guarantee comprehensive coverage of all cases.

B2. Knowledge Retrieval & Construction The retrieval pipeline in this paper is based on key-

words extracted from $\{r, r'\}$, then the quality-managed keywords are used to retrieve and construct the knowledge sets following the techniques demonstrated in Section 3. Specifically, we adopt GPT-3.5 to extract domain-specific, difficult-to-understand keywords from $\{r, r'\}$ (see Table 16 for details of the prompt). After that, we use GPT-3.5 to filter the extracted keywords (see Table 17 and Table 18 for details of the prompt), and we empirically design a blacklist for further rejecting low-quality keywords.

To construct the EaK knowledge set, we select SNOMED-CT as the knowledge base (KB) and use its API² based on the above keywords to search for the associated entities in the KB, specifically here we keep the top-1 search results, and for the searched entities, we query their children entities again using the API and randomly sample up to 3 children entities. The parent and children are then serialized into EaK knowledge in the form of “One of parent is children”.

For out-of-domain knowledge sources, we based on the English WIKIDATA API³ to search for entity codes using extracted keywords. We use the searched top-1 entity code to construct SPARQL statements for knowledge query from WIKIDATA⁴ while using this entity code to query the page extraction from English WIKIPEDIA⁵. Then, for the retrieved results, we use GPT-3.5 to summarize the retrieved extractions to limit their length to no more than 1000 words under a zero-shot setting (see Table 21).

B3. Prompt Construction GPT-3.5 is used to extract the self-indicator of schema pairs (or entity pairs in MMM). The used prompts can be found in 19. To implement Demonstration-summarization-only, we also use GPT-3.5 to compress the input $\{r, r'\}$. Specifically, to control the length of the prompt, we split $\{r, r'\}$ into two queries (see Table 20 for the detailed prompt).

B4. Inference Settings For the MIMIC, Synthea and CMS datasets, we take WIKIDATA with DaK, WIKIPEDIA with EaK, and EaK as the knowledge sources to evaluate KCMF under the 4-shot setting. For MMM, we take WIKIDATA, WIKIPEDIA, and EaK as the knowledge sources to evaluate the

²<https://browser.ihtsdotools.org/snowstorm/snomed-ct>

³<https://www.wikidata.org/w/api.php>

⁴<https://query.wikidata.org/>

⁵<https://en.wikipedia.org/w/api.php>

framework under the 2-shot setting. We empirically set temperature to 0 and top_p to 0.1 for stable outputs if these two parameters are available. Also, we directly took poorly formatted outputs as negative predictions in our experiments for metrics calculation and workflow automation. The pseudo-code we designed for these two tasks can be found in Appendix G. We derived the main results from the best result of running five times.

B5. Construction of the MMM Dataset We present MedMentions Matching, or MMM, a clinical entity matching dataset modified from MedMentions (Mohan and Li, 2019). MedMentions is an entity linking dataset whose data are recognized entity mentions from PubMed documents labeled with Unified Medical Language System (UMLS) code. Inspired by the fact that entities referred to the same real-world concept are matched from the view of EM, we re-organized the MedMentions as per Figure 9: For entities belonging to the same UMLS code, we keep the pairs with lowest similarity as positive. For those belonging to different, we sample the pairs with the highest similarity as negative. Then we use sentences where sampled entities are located as context. In total, MMM contains 7,359 positive and 100,000 negative examples and an example can be found in Figure 10.

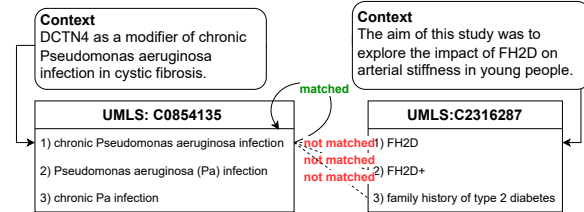


Figure 9: The process of constructing MedMentions Matching (MMM) dataset.

Entity A	Context A	Entity B	Context B	Label
chronic Pseudomonas aeruginosa	DCTN4 as a modifier of chronic Pseudomonas aeruginosa..	FH2D	The aim of this study was to explore the impact of FH2D...	0 (not matched)
...				

Figure 10: An example from MMM.

B6. Dataset and Code All datasets can be accessed through original papers cited in this study, specifically for modifying MMM, the script is given in our supplementary codes. Table 4 presents the statistics of the test sets used in the main experiment, sampled as 10% stratified subsets.

Dataset	Task	# (Instance)	# (Positive)	IR
Synthea	SM	2964	8	370.5
MIMIC	SM	6408	13	492.9
CMS	SM	2564	20	128.2
MMM	EM	10736	736	14.6
AG	EM	2293	234	9.8
WA	EM	2049	193	10.6
DS	EM	5842	1070	5.5

Table 4: Statistics of evaluation datasets. IR = # (Instance) / # (Positive) is the Imbalance Ratio. The larger the IR, the harder the matching task.

C Performance Gain Analysis

Methods	MIMIC		Synthea		CMS	
	P	R	P	R	P	R
SMAT	11.5	84.6	24.4	90.9	33.9	95.0
GLM-4	1.9	46.2	7.4	25.0	50.0	10.0
+KCMF	5.9	53.8	9.7	37.5	19.4	35.0
Δ	+4.0	+7.6	+2.3	+12.5	-30.6	+25.0
Llama-3	0.3	100.0	1.2	50.0	100.0	5.0
+KCMF	4.2	46.2	4.6	75.0	6.9	45.0
Δ	+3.9	-53.8	+3.4	+25.0	-93.1	+40.0
Mistral	10.3	30.8	4.3	50.0	21.1	20.0
+KCMF	21.6	61.5	9.1	37.5	20.8	25.0
Δ	+11.3	+30.7	+4.8	-12.5	-0.3	+5.0
GPT-3.5	4.1	46.2	5.6	25.0	1.0	30.0
+KCMF	56.3	69.2	28.0	87.5	25.0	35.0
Δ	+52.2	+23.0	+22.4	+62.5	+24.0	+5.0
GPT-4o	25.0	38.5	5.6	12.5	50.0	10.0
+KCMF	25.9	53.8	31.8	87.5	34.5	40.8
Δ	+0.9	+15.3	+26.2	+75.0	-15.5	+29.2

Table 5: Precision (P) and Recall (R) measures in SM. KCMF excels in imbalanced tasks like SM by effectively balancing precision and recall, leveraging pseudo-code to enhance recall without compromising precision.

To further explore how KCMF achieves a higher F1-score compared to both LLM baselines and the conventional method, we provide a detailed analysis of Precision (P) and Recall (R) on the three SM datasets with the highest imbalance ratios (see Table 4). The results are provided in Table 5.

In *imbalanced* tasks like SM, often struggle to balance precision and recall, a challenge also evident in LLM baselines. As shown in Table 5, different LLM baselines exhibit varying tendencies toward recall. However, KCMF demonstrates performance improvements by addressing these imbalances: for baselines that prioritize precision, KCMF enhances recall by identifying more positive samples while effectively managing false positives, ensuring precision remains stable. This aligns with the practical preference for higher recall in matching tasks, reducing the need for additional human verification. For baselines that emphasize

recall, KCMF improves precision by correctly rejecting negative samples, with less emphasis on further recall gains when unnecessary.

The performance gain of KCMF owes to its ability to enhance recall through pseudo-code, which operates independently of the model’s reliance on the exposed data distribution. This enables KCMF to accurately reject negative instances while improving the recall of positive samples.

D Ablation Study

In this section, we incrementally stack techniques onto a bare LLM, demonstrating how techniques proposed in the main text transform a backbone model into a comprehensive matching framework. The evolution process is visually summarized in Figure 11 taking the SM task on Synthea as an example, which provides an overview of how each technique incrementally enhances the F1-score. The process begins with a plain zero-shot LLM query (GPT-3.5 baseline), which achieves an initial F1-score of 9.1%. Through the addition of various techniques, the final framework, KCMF, attains an F1-score of 42.4%, significantly outperforming both the GPT-3.5 baseline (9.1%) and the SMAT method (38.5%). The following subsections detail the techniques that contribute to this substantial performance improvement.

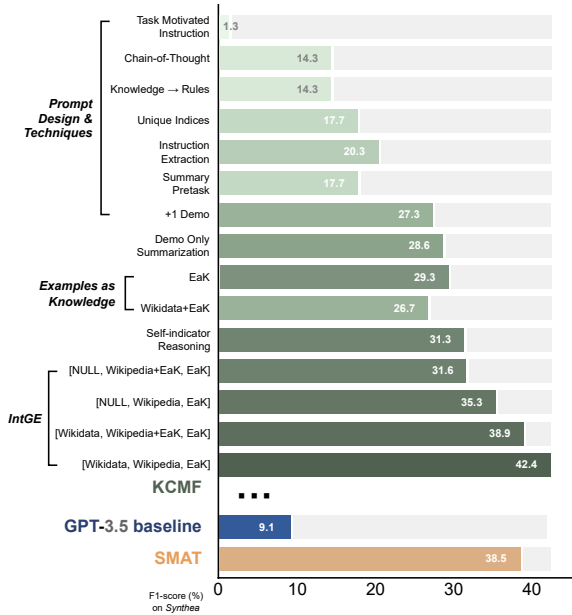


Figure 11: KCMF evolved from a plain zero-shot LLM query method. The horizontal axis depicts the F1-score for each version of KCMF on Synthea.

D.1 Prompt Design and Construction

We conducted incremental experiments using WIKIDATA+DaK with a 3-shot setting on Synthea as the baseline. The results, summarized in Table 6, illustrate the impact of each optimization.

Increment	Shot	Accuracy	F1-score
TO Inst.	3	0.5735	0.0125
CoT	3	0.9798	0.1429
Know2Rule	3	0.9798	0.1429
U-indices	3	0.9811	0.1765
Inst. Extraction	3	0.9841	0.2034
Summary Pretask	3	0.9781	0.1772
+1 demo	4	0.9892	0.2727

Table 6: Ablation studies for prompt design and construction.

- **TO Inst.** Task-Oriented Instruction: We revised the instruction from Narayan et al. (2022) to “Can records in schema B be transformed and stored into schema A?” to better align with the schema matching task’s objectives.
- **Know2Rule** (Knowledge to Rule): The “knowledge for the task” prompt element was renamed to “rules for the task” for improved clarity and alignment with task requirements.
- **U-indices** (Unique Indices): We applied unique indices for the different elements of the prompt that need to be numbered. We used Roman numerals, lowercase letter numbers, and numeric numbers separately for the row number of pseudo-code, the sequence number of retrieved knowledge, and sequence number of reasoning steps.
- **Inst. Extraction** (Instruction Extraction): We extracted the instructions and rules for each demonstration and placed them at the beginning of the prompt.
- **Summary Pretask and +1 Demo**: The Summary Pretask was designed to shorten prompt length, and +1 Demo added an extra demonstration.

All the optimizations we added are motivated by achieving higher precision while ensuring promising recalls.

As shown in Table 6, CoT, U-indices and Inst. Extraction brought 13.04%, 3.36%, 2.69% improvement in F1-score on Synthea. Although Summary Pretask caused a temporary drop in performance, the combination of Summary Pretask

and +1 Demo still shows a strong result; the F1-score is improved by 6.93% compared to the Inst. Extraction step. We examine the results of the Summary Pretask and the +1 Demo as a whole because the motivation for the Summary Pretask was to shorten the length of the prompt.

D.2 Summary Pretask

For the summary pretask, we apply two strategies: **All** and **Demo Only**. The Demo Only strategy involves summarizing only the demonstrations used in the prompt, as described in Section 3.4. The All strategy builds on Demo Only by additionally summarizing descriptions of the target schema pairs.

As shown in Table 7, the compression strategy of Demo Only outperforms the All strategy by 1.3% on F1-score under the 4-shot setting. This indicates that focusing compression on demonstration examples is more effective than extending it to target schema pairs.

Strategy	Accuracy	F1-score
All	0.9892	0.2727
Demo only	0.9899	0.2857

Table 7: Ablation study of the summary pretask’s strategies on Synthea.

D.3 Examples as Knowledge (EaK)

Following the **Demo Only** setting in Appendix D.2, we further verify the validity of the knowledge constructed through EaK. We introduce EaK as a knowledge source and form three distinct knowledge sets: WIKIDATA, EaK, and WIKIDATA+EaK, which are used for constructing the prompt.

As shown in Table 8, EaK, characterized by its concise and high-quality knowledge, achieves the best performance among the three configurations. Using EaK alone as the knowledge source outperforms WIKIDATA and WIKIDATA+EaK by 0.7% and 2.6%, respectively, in terms of F1-score. This highlights the advantage of leveraging high-quality, compact examples as a knowledge source.

Knowledge Source	Accuracy	F1-score
WIKIDATA	0.9899	0.2857
EaK	0.9902	0.2927
WIKIDATA+EaK	0.9889	0.2667

Table 8: Ablation study of EaK on Synthea.

D.4 Self-indicator Reasoning (SIR)

To evaluate the effectiveness and generalizability of Self-Indicator Reasoning (SIR), we introduce it into demonstrations across different knowledge sources for both the MMM and Synthea datasets.

Results on MMM Table 9 presents the results with and without SIR on the MMM dataset with WIKIDATA, WIKIPEDIA, and EaK as the knowledge sources, respectively. We observe that all cases with SIR outperform cases without SIR, with their F1-scores improving by an average of 7.73%. More specifically, F1-scores of SIR with WIKIDATA, WIKIPEDIA, and EaK as knowledge sources improved by 8.43%, 7.05%, and 7.70%, respectively. These results demonstrate the robustness and effectiveness of SIR in enhancing the reasoning capabilities of the model.

Knowledge Source	w. SIR	Accuracy	F1-score
WIKIDATA	No	0.9498	0.6534
WIKIDATA	Yes	0.9667	0.7387
WIKIPEDIA	No	0.9603	0.6979
WIKIPEDIA	Yes	0.9719	0.7684
EaK	No	0.9593	0.6997
EaK	Yes	0.9729	0.7767

Table 9: The ablation study of SIR on MMM (2-shot setting).

Results on Synthea Meanwhile, for Synthea, also based on the *Demo Only* setting, we introduce SIR with WIKIDATA as the knowledge source. As shown in Table 10, SIR improved the F1-score on Synthea by 2.68%, further confirming its effectiveness.

Knowledge Source	w. SI	Accuracy	F1-score
WIKIDATA	No	0.9899	0.2857
WIKIDATA	Yes	0.9926	0.3125

Table 10: The ablation study of SIR on Synthea (4-shot setting).

D.5 Inconsistency-tolerant Generation Ensemble (IntGE)

The ablation experiments for IntGE build upon the findings from Appendix D.4. Here, we evaluate the performance of IntGE by combining different knowledge sources and testing various configurations. Given that the impact of SIR is less pronounced on Synthea compared to MMM, we also include SIR as an experimental variable.

Results on Synthea Table 11 showcases the results of IntGE applied to Synthea using different combinations of knowledge sources. All ensembled configurations outperform the single-knowledge-source baseline (WIKIDATA). The F1-scores of IntGE under combinations [NULL, WIKIPEDIA+EaK, EaK], [NULL, WIKIPEDIA, EaK], [WIKIDATA, WIKIPEDIA+EaK, EaK], and [WIKIDATA, WIKIPEDIA, EaK] are improved by 7.74%, 4.04%, 11.04%, 0.33%, respectively. These results highlight the effectiveness of combining multiple knowledge sources in improving model performance.

Knowledge Source	Voting	Accuracy	F1-score
WIKIDATA	No	0.9926	0.3125
[NULL, WIKIPEDIA+EaK, EaK]	Yes	0.9926	0.3889
[NULL, WIKIPEDIA, EaK]	Yes	0.9926	0.3529
[WIKIDATA, WIKIPEDIA+EaK, EaK]	Yes	0.9936	0.4242
[WIKIDATA, WIKIPEDIA, EaK]	Yes	0.9912	0.3158

Table 11: Ablation study of IntGE on Synthea (4-shot setting). NULL denotes no knowledge source.

Results on MMM We conducted similar experiments on the MMM dataset, as shown in Table 12. The F1-scores of IntGEs based on the combinations [WIKIDATA, WIKIPEDIA, EaK], [WIKIDATA, WIKIPEDIA, EaK*], [WIKIDATA*, WIKIPEDIA, EaK], [WIKIPEDIA, EaK, NULL] increased by 3.61%, 2.88%, 3.71%, and 3.71%, respectively, as compared to the results based on the single WIKIDATA data source. Here, the asterisk (*) indicates that the prompt for the corresponding knowledge source excludes self-indicators in reasoning steps. These results confirm the robustness and adaptability of IntGE across varying configurations and datasets.

Knowledge Source	Voting	Accuracy	F1-score
WIKIDATA	No	0.9667	0.7387
[WIKIDATA, WIKIPEDIA, EaK]	Yes	0.9727	0.7748
[WIKIDATA, WIKIPEDIA, EaK*]	Yes	0.9715	0.7675
[WIKIDATA*, WIKIPEDIA, EaK]	Yes	0.9727	0.7758
[WIKIPEDIA, EaK, NULL]	Yes	0.9727	0.7758

Table 12: Ablation study of IntGE on MMM (2-shot setting).

Summary The experiments here demonstrate that IntGE consistently improves performance by leveraging diverse knowledge sources, with significant F1-score gains across both MMM and Synthea. The results emphasize the value of combining high-quality knowledge sources and adapting reasoning strategies for optimal results.

E Cost Analysis

This section presents an ablation study of the main components of KCMF, analyzing their token costs and contributions to performance improvement. The results are grouped by techniques to evaluate their trade-offs in terms of cost and gain, providing deeper insights into the effectiveness of each component.

In Figure 12, token costs and performance gain of GPT-4o-mini versioned KCMF are visualized as multiple bars, those unhatched represent the percentage of tokens the corresponding technique takes, while hatched bars represent the gain in F1-score of each technique to the baseline. It should be noted that the performance gain here and in the following text are obtained by ablations, e.g., the gain of the self-indicator is the difference between versions of KCMF with and without self-indicators. According to Figure 12, both self-indicator and knowledge show consistent and independent performance gain, among which, we observe that self-indicator is significantly cost-effective, especially for schema matching tasks, which may be due to its effect in clarifying the complexity and heterogeneity of schemata. Also, knowledge is necessary for KCMF to achieve its optimal performance.

The token cost and gain of the F1-score of each technique are listed in Table 14, and the API invoking costs are reported in Table 13. Generally, the average cost of KCMF is times the baseline, due to the pseudo-code, reasoning steps, and the retrieved knowledge included in its query, and we consider this overhead worthwhile because 1) KCMF’s total cost is still satisfactory, where around 60 cents for each dataset, and 2) it shows consistent and competitive performance gain towards different LLMs.

Further, we introduce the cost performance (CP), which is the ratio of the gain of F1-score to the token percentage (of total) for the main components of the KCMF, to measure the trade-off between cost and gain for techniques in KCMF. CP is formally defined in Equation 3, where $|\mathcal{T}'|$ and $|\mathcal{T}|$ denote numbers of tokens corresponding to the single technique and the whole prompt respectively.

$$CP = \frac{\Delta F1}{|\mathcal{T}'|/|\mathcal{T}|} \quad (3)$$

To ensure fairness, CPs are averages calculated on Synthea, MIMIC, CMS, and MMM due to the absence of knowledge in the other three datasets. According to Table 14, we find that self-indicator is more cost-effective than knowledge, but knowledge

	MIMIC	Synthea	CMS	MMM	AG	WA	DS
Baseline	0.0077¢	0.0068¢	0.0074¢	0.0048¢	0.0034¢	0.0054¢	0.0056¢
KCMF	0.0260¢	0.0278¢	0.0253¢	0.0207¢	0.0207¢	0.0271¢	0.0220¢

Table 13: API costs using GPT-4o-mini endpoint.

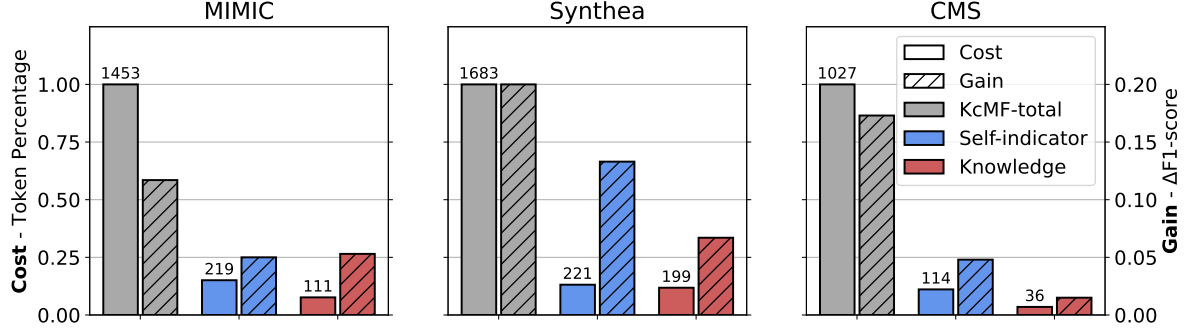


Figure 12: Cost and gain of the KCMF’s main components, using GPT-4o-mini as the backbone model.

is also necessary for the whole KCMF to achieve its best performance. Also, we observe that the self-indicator generalizes well on AG, WA, and DS, which indicates that self-indicator is a versatile method for improving LLM’s performance on matching tasks.

Dataset		Token Cost/Query			
		K	SI	Infer	Total
MIMIC	Cost	111	219	1123	1453
	Gain	5.3%	5.0%	-	11.7%
Synthea	Cost	199	221	1263	1683
	Gain	6.7%	13.3%	-	20.0%
CMS	Cost	36	114	877	1027
	Gain	1.5%	4.8%	-	17.3%
MMM	Cost	23	61	1051	1135
	Gain	0.2%	0.3%	-	8.2%
AG	Cost	-	90	795	885
	Gain	-	1.2%	-	15.6%
WA	Cost	-	109	1044	1153
	Gain	-	2.2%	-	41.0%
DS	Cost	-	94	912	1006
	Gain	-	0.2%	-	19.9%
CP		0.441	0.474	-	-

Table 14: Cost and gain ablation results on all datasets in the main experiment, with averaged cost performance (CP) for knowledge (K) and self-indicator (SI).

F Discrepancies in LLM’s Understanding of Schema Matching Instructions

Table 15 presents the prompt using vanilla zero-shot Chain-of-Thought on Synthea and the response from GPT-3.5. It can be observed that there is a discrepancy between what GPT-3.5 understands the instruction “match” and the actual objective of the schema matching task.

<p>Are schema A and B matched? Let’s think step by step.</p> <p>schema A: person-month_of_birth</p> <p>description of schema A: the person domain contains records that uniquely identify each patient in the source data who is time at-risk to have clinical observations recorded within the source systems.;the month of birth of the person. for data sources that provide the precise date of birth, the month is extracted and stored in this field.</p> <p>schema B: patients-birthdate</p> <p>description of schema B: patient demographic data.;the date the patient was born.</p>
<p>Based on the descriptions provided:</p> <p>**Schema A:**</p> <ul style="list-style-type: none"> - **Name:** person-month_of_birth - **Description:** This schema appears to be related to a person domain where records uniquely identify each patient. It specifically captures the month of birth of the person. If the source data provides the precise date of birth, only the month is extracted and stored in this field. <p>**Schema B:**</p> <ul style="list-style-type: none"> - **Name:** patients-birthdate - **Description:** This schema relates to patient demographic data and specifically captures the date when the patient was born. <p>Now, to determine if schema A and schema B are a match, we need to consider the following points:</p> <ol style="list-style-type: none"> 1. **Data Captured:** <ul style="list-style-type: none"> - Schema A captures only the month of birth if the precise date is available. - Schema B captures the full birthdate. 2. **Granularity:** <ul style="list-style-type: none"> - Schema A has a lower granularity by capturing only the month, whereas Schema B captures the full date. 3. **Context and Usage:** <ul style="list-style-type: none"> - Schema A seems more focused on situations where the precise date is not necessary, and only the month suffices (possibly for anonymization or simplification purposes). - Schema B is more comprehensive as it captures the exact birthdate, which is typically more useful for demographic analysis and age calculation. <p>**Conclusion:**</p> <p>Based on the descriptions:</p> <ul style="list-style-type: none"> - Schema A (‘person-month_of_birth’) and Schema B (‘patients-birthdate’) do not appear to be a direct match because they capture different levels of detail regarding the birthdate of patients. Schema A focuses on the month of birth, while Schema B captures the full birthdate. <p>Therefore, they serve different purposes and would likely be used in different contexts within a database or system that manages patient or person-related data.</p>

Table 15: GPT-3.5’s response on Synthea.

G Detailed Prompts

In this section, we present all prompt templates used in the experiments covered in the main text, specifically:

• Tables 16 to 21

- Prompts for pre-tasks, including keywords extraction and quality management (DaK introduced in Section 3.3), self-indicator extraction, demonstration and retrieval summarization (Section 3.4).

• Tables 22 and 23

- Prompts of baselines for comparison in the main experiment, Tables 22 and 23 are for SM and EM, separately.

1278
1279
1280
1281

- **Tables 24 and 25**
 - Task-specific pseudo-codes for SM and EM used by KCMF.
- **Tables 26 and 27**
 - Prompts obtained from Section 3.4, which are used to query LLMs directly.

<p>You need to extract all the keywords in the schema that require special domain knowledge to understand, keywords should be separated by commas.</p> <p>Example 1: Schema: provider-npi Schema description: the provider table contains a list of uniquely identified healthcare providers. these are individuals providing hands-on healthcare to patients, such as physicians, nurses, midwives, physical therapists etc.;the national provider identifier (npi) of the provider. Answer: national provider identifier, npi</p> <p>Example 2: Schema: imaging_studies-sop description Schema description: patient imaging metadata.;description of the sop code. Answer: sop</p> <p>Example 3: Schema: procedure_occurrence-modifier_concept_id Schema description: the procedure_occurrence table contains records of activities or processes ordered by, or carried out by, a healthcare provider on the patient to have a diagnostic or therapeutic purpose. procedures are present in various data sources in different forms with varying levels of standardization.;a foreign key to a standard concept identifier for a modifier to the procedure (e.g. bilateral). these concepts are typically distinguished by 'modifier' concept classes (e.g., 'cpt4 modifier' as part of the 'cpt4' vocabulary). Answer: foreign key, identifier, cpt4</p> <p>Your turn:</p>

Table 16: Prompt for keyword extraction.

<p>You need to find out which of the given keywords are relevant to the database or medical field and return them, keywords should be separated by commas.</p> <p>Example 1: birthday, home, dcm, location, primary key Answer: dcm, primary key</p> <p>Example 2: endtime, id, recognition, observation, data model, algorithm, artificial Answer: id, data model</p> <p>Your turn:</p>

Table 17: Prompt for keyword quality management for the schema matching datasets (MIMIC, Synthea, and CMS).

<p>Tell if the given word is a hard-to-understand biomedical domain term, only yes or no.</p> <p>Example 1: leaf Answer: no</p> <p>Example 2: deoxynivalenol Answer: yes</p> <p>Your turn:</p>

Table 18: Prompt for keywords quality management for the entity matching datasets (MMM, AG, WA, and DS).

Instruction: Given two schemas, you need to summarize the column and table for each considering the given knowledge.
[Examples]
Your turn:
 $\{r, r'\}$
 $[K_i]$
Answer:

Table 19: Prompt for self-indicator extraction (SIR).

Instruction: You need to summarize the given schema based on its schema name and description. The summary should be focused on retaining and explaining concepts in the database domain. Schema name is the table and column names of the schema separated by a dash. Schema description is the table and column descriptions of the schema separated by a semicolon.
[Examples]
Your turn:
 $\{r, r'\}$
Answer:

Table 20: Prompt for demonstration summarization.

Instruction: You need to summarize the given text into a paragraph less than 1000 words.
{CONTENT}
Your answer:

Table 21: Prompt for retrieval summarization.

Are schema A and B the same? ONLY yes or no.
[Examples]
Your turn:
 $\{r, r'\}$
{RESPONSE}

Table 22: Prompt template of baseline for the schema matching datasets (MIMIC, Synthea, and CMS).

Question: Do entity A and entity B refer to the same real-world concept? Only yes or no.
[Examples]
Your turn:
 $\{r, r'\}$
Answer:
{RESPONSE}

Table 23: Prompt template of baseline for the entity matching datasets (MMM, AG, WA, and DS).

I: Rules II, III, and IV MUST be checked SEQUENTIALLY until you conclude an answer.
 II: If the columns of the two schemas can not be the same type of data in the database, the answer is no, otherwise, check rule III.
 III: If the tables of the two schemas are not semantically the same, the answer is no, otherwise, check rule IV.
 IV: If the columns of the two schemas do not refer to the same concept, the answer is no, otherwise, the answer is yes.

Table 24: Pseudo-code for the schema matching datasets (MIMIC, Synthea, and CMS).

I: Rules II, III, and IV MUST be checked SEQUENTIALLY until you conclude an answer.
 II: If Entity A is an abbreviation of Entity B or vice versa, the answer is yes, otherwise check rule III.
 III: If Entity A is an alias of Entity B or vice versa, the answer is yes, otherwise check rule IV.
 IV: If Entity A and Entity B refer to the same real-world concept the answer is yes, otherwise the answer is no.

Table 25: Pseudo-code for the entity matching datasets (MMM, AG, WA, and DS).

<p>Question: Can records in schema B be transformed and stored in schema A? The task should be solved by completing the reasoning steps and concluding a final answer ONLY yes or no. Do not stop until you draw a final answer. Schema name is the table and column names of the schema separated by a dash.</p> <p>Rules for the task: [\mathcal{R}] [<i>Demonstrations</i>]</p> <p>Your turn: $\{r, r'\}$</p> <p>Knowledge for the task: [K_i]</p> <p>Reasoning: 1. {<i>Self-indicator</i>}</p> <p>Please continue the reasoning until you draw a final answer ONLY yes or no: {<i>RESPONSE</i>}</p>

Table 26: Prompt template of KCMF for the schema matching datasets (MIMIC, Synthea, and CMS).

<p>Question: Do entity A and entity B refer to the same real-world concept? You must think step by step, and finally draw an answer only yes or no.</p> <p>Rules for the task: [\mathcal{R}] [<i>Demonstrations</i>]</p> <p>Your turn: $\{r, r'\}$</p> <p>Knowledge for the task: [K_i]</p> <p>Reasoning: 1. {<i>Self-indicator</i>}</p> <p>Please continue the reasoning until you draw a final answer ONLY yes or no: {<i>RESPONSE</i>}</p>

Table 27: Prompt template of KCMF for the entity matching datasets (MMM, AG, WA, and DS).