

# LLMATCH: a Unified Schema Matching Framework with Large Language Models

No Author Given

No Institute Given

**Abstract.** Schema matching is crucial for data integration and involves finding matching columns across schemas. Traditional tools and benchmarks focus on simple single table mappings while multi-table matching remains a challenging task. We present LLMATCH, a solution for complex schema matching tasks. LLMATCH decomposes the problem into three modular components: schema preparation, table selection, and column matching, each supporting multiple configurable options to allow flexible customization of various method. We also introduce a multi-level matching optimization to the framework for complex matching problems. The optimization includes a Rollup component, which aggregates semantically related columns prior to high-level matching. Following this, a Drilldown component expands the high-level matches back into their original columns for detailed alignment. Further, we created a benchmark dataset with seven multi-table schema pairs from three industries. We used this dataset to test configurations and improve state-of-the-art performance. The experiments provided insights for schema matching and LLM applications. Tests with human data engineers showed significant productivity improvements in business scenarios.

**Keywords:** Schema Matching · LLM · Data Management

## 1 Introduction

Schema matching is an important task in database systems. It supports applications like data integration, warehousing, and exchange. The process identifies semantic correspondences between columns of source and target schemas, allowing data from various systems to be transformed into a unified format. Industries such as healthcare [39], energy [35], and agriculture [6] are adopting standardized data formats to improve interoperability. In healthcare, a unified data model [47] can facilitate cross-institutional research and speed up new treatment discovery [33,29,30]. Each country performs schema matching to convert their healthcare data into this unified format [21,17,31,15,22]. In the business world, companies match and migrate data during mergers and acquisitions (M&As). Errors in this process can cause significant disruptions. For example, Deutsche Bank’s 2023 merger with Postbank resulted in account lockouts due to data migration issues, leading to regulatory scrutiny and reputational damage [14].

Classical schema matching tools and benchmarks are designed for tabular data in CSV files, but their accuracy drops sharply on complex real-world

schemas due to inconsistent naming conventions and hierarchical structures [51]. For instance, in the toy matching task shown in Fig. 1, fields like *PersonID*, *CustomerID*, and *ContactID* differ lexically but are considered matches. Consequently, complex matching scenarios often depend on domain experts to manually create and maintain mappings [9], a process that can take up to 500 hours for two experts to map a single schema [37] in some cases. This approach is not only time-intensive but also prone to errors [2, 44].

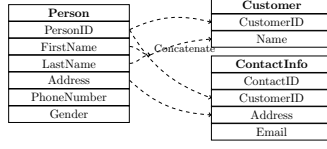


Fig. 1: Toy matching example.

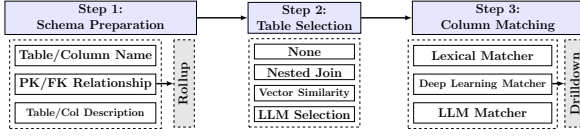


Fig. 2: Workflow design for schema matching task.

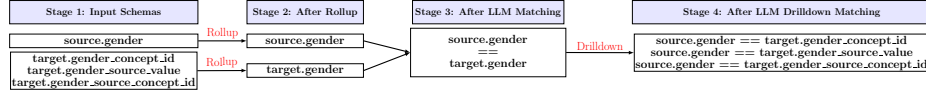


Fig. 3: Example multi-level matching with Rollup and Drilldown

Recently, Large language models (LLMs) have shown great potential in tackling complex database tasks, including entity resolution, datalake discovery, and query synthesis [49, 12]. While several LLM-based schema matching tools [43, 18] have demonstrated promising capabilities, fair comparison remains challenging due to the lack of a standardized evaluation framework and appropriate datasets. Existing methods address tasks with varying complexities: *table-to-table* matching involves aligning columns in isolated tables, while *schema-to-schema* matching requires understanding relationships across entire schemas, such as primary key-foreign key mappings. Most evaluation datasets focus on *table-to-table* matching and fail to capture the broader scope of schema-level tasks. This disparity in evaluation settings obscures the strengths and weaknesses of different approaches and emphasizes the need for a unified framework with datasets that address both simple and complex matching tasks.

In this paper, we introduce a unified framework that divides the schema matching process into three phases: *schema preparation*, *table selection*, and *column matching* (Fig. 2). Schema preparation identifies the components of the source and target data. Table selection determines the most relevant target tables for each source table, reducing the context size for the next phase. After selecting target tables, the column matching phase aligns columns between each source table and its corresponding target tables, generating the final mappings. To enhance the adaptability and efficiency of our framework, especially in complex multi-table environments, we have developed an optimization featuring the *Rollup* and *Drilldown* techniques, illustrated in Fig. 3. Rollup serves as a preprocessing step where semantically related columns from the target schema are aggregated into a single, simplified representation, such as combining `target.gender_concept_id`, `target.gender_source_value`, and `target.gender_source_concept_id` into a consolidated `target.gender` column. This

aggregation aids in abstracting the detailed data into broader categories that facilitate high-level matching, making the schema matching process more manageable and focused. Following the Rollup, the schema, the simplified schemas are fed into original framework for high-level matching. Should the initial high-level matching be successful, the Drilldown component is activated to detail the process. It reverses the Rollup by decomposing the aggregated target.gender back into its detailed components. This step ensures that each attribute within the target schema is precisely aligned with the corresponding source data, providing thorough and accurate mappings. This multi-level approach not only optimizes the matching process by reducing computational overhead but also improves precision, particularly in handling complex data schemas.

To evaluate our framework, we curated seven benchmark datasets spanning finance, healthcare, and entertainment domains. These schemas average 14 tables and 135 columns, in sharp contrast to simpler benchmarks, which typically include just one table and 15 columns (Table 1). Using the proposed framework and benchmark, we conducted extensive experiments with state-of-the-art methods and uncovered several key insights:

- **Irrelevant prompt context hurts performance:** Presenting all target tables in a single prompt with a large context size or processing them individually across multiple prompts with a small context size often leads to suboptimal performance. Instead, selecting a subset of the most relevant target tables while maintaining a moderate context size yields better results.
- **The default RAG configuration is not the optimum for this task:** Embedding vector similarity, commonly used in the RAG framework, is not ideally suited for schema matching tasks. Our experiments indicate that allowing the LLM to directly handle target selection often produces better results, highlighting the need for more adaptive and context-aware approaches.
- **Multi-level matching aids in managing complex schemas:** The Rollup and Drilldown operations are essential for managing complexity at different stages of the matching process. By first simplifying the schema through Rollup and then applying Drilldown once a general term is matched, this approach effectively improves performance on complex datasets.

These findings render the interplay between model capabilities, context quality, and schema complexity, offering insights for advancing schema matching solutions. We release our framework as well as the new benchmark datasets [5].

## 2 Problem Statement

A database schema  $S$  consists of a set of tables  $T$ , columns  $A$ , descriptions of tables and columns  $D$ , and relationships  $R$ , including primary keys (PK) and foreign keys (FK). Given two database schemas  $S_s$  (source) and  $S_t$  (target), the objective is to map tables and columns from  $S_s$  to those in  $S_t$ .

**Motivation for n:m Mapping.** In many real-world data integration scenarios, especially those involving relational databases with multi-table structures,

matching often requires mapping multiple columns from the source schema to multiple columns in the target schema. We refer to this as an n:m mapping. This pattern, while expansive, helps ensure that important relationships are captured, even if it occasionally results in higher false positives. For most industrial applications, prioritizing recall (minimizing missed matches) is more desirable than optimizing precision, as it reduces the risk of missed correspondences [36].

**Definition 1.** In an n:m match, the task is to find a mapping  $\Psi : A_s \rightarrow A_t$ , where each column  $a \in A_s$  is associated with a subset  $\Psi(a) \subseteq A_t$ .

**Metrics.** The performance of n:m matching is thereby measured by the F1 score, which is defined in terms of the precision and recall from the predicted mapping  $\Psi$  and the ground truth mapping  $\hat{\Psi}$ .

$$\text{precision}(\Psi, \hat{\Psi}) = \frac{|\{(a, b) \mid b \in \Psi(a) \cap \hat{\Psi}(a)\}|}{|\{(a, b) \mid b \in \Psi(a)\}|}, \quad \text{recall}(\Psi, \hat{\Psi}) = \frac{|\{(a, b) \mid b \in \Psi(a) \cap \hat{\Psi}(a)\}|}{|\{(a, b) \mid b \in \hat{\Psi}(a)\}|}$$

### 3 Dataset Development

Traditional schema matching benchmarks typically involve a single source and target table, as traditional NLP tools struggle to interpret PK/FK relationships in relational schemas. More recent works on complex schema matching primarily uses two healthcare datasets: *mimic-omop* and *synthea-omop* [43,18]. We expand this scope to include seven pairs of source and target schemas spanning healthcare, entertainment, and finance, creating a more diverse benchmark. For healthcare data, we leverage mappings from prior research [43,37] and official mapping [29,9]. For entertainment data, we use ground truth mappings from [50]. In the finance domain, we release two test schemas along with their ground truth mappings annotated by industry experts.

**Table 1:** Dataset statistics for traditional benchmarks (top) and the complex benchmark (bottom) we developed. Complex datasets have more mapped target tables per source table and a low ration of simple (1:1) mappings.

Mapping Task	Domain	Src # Tables/- Cols/PKs/FKs	Tgt # Tables/- Cols/PKs/FKs	# Total Map	Avg/Max # Tgt Tab. per Src Tab.	1:1 Map Ratio (%)	Src/Tgt Column Map Ratio (%)
mjs-mjt	Entertainment	1 / 13 / 0 / 0	1 / 13 / 0 / 0	6	1 / 1	100%	46.15% / 46.15%
msjs-msjt	Entertainment	1 / 14 / 0 / 0	1 / 14 / 0 / 0	8	1 / 1	100%	57.14% / 57.14%
mus-mut	Entertainment	1 / 20 / 0 / 0	1 / 20 / 0 / 0	20	1 / 1	100%	100% / 100%
mvs-mvt	Entertainment	1 / 13 / 0 / 0	1 / 13 / 0 / 0	6	1 / 1	100%	46.15% / 46.15%
imdb-sakila	Entertainment	7 / 39 / 2 / 7	16 / 90 / 12 / 22	22	1.00 / 3	26.67%	56.41% / 42.22%
bank1-bank2	Finance	9 / 27 / 5 / 11	9 / 36 / 4 / 12	11	0.44 / 1	100.00%	81.48% / 63.89%
cms-omop	Healthcare	5 / 96 / 1 / 4	39 / 432 / 12 / 58	157	7.40 / 11	3.70%	60.42% / 25.69%
synthea-omop	Healthcare	12 / 111 / 3 / 19	39 / 432 / 12 / 58	101	1.50 / 3	19.44%	49.55% / 31.02%
cprd_aurum-omop	Healthcare	8 / 76 / 5 / 21	39 / 432 / 12 / 58	42	1.75 / 3	48.00%	60.53% / 21.06%
cprd_gold-omop	Healthcare	9 / 123 / 4 / 21	39 / 432 / 12 / 58	52	2.44 / 4	25.00%	36.59% / 21.30%
mimic_iii-omop	Healthcare	26 / 324 / 6 / 55	39 / 432 / 12 / 58	189	1.54 / 5	14.73%	56.79% / 34.26%

Table 1 contrasts traditional single-table schema matching datasets from earlier works [20] with our new complex datasets that feature richer and more intricate schemas, particularly evident in the healthcare domain. Unlike traditional datasets that predominantly feature simple 1:1 mappings, our healthcare dataset showcases a significant complexity with an average number of more than one target table mapped per source table and a low percentage of 1:1 mappings. This complexity highlights the need for advanced schema matching techniques that can handle multiple and complex table relationships, reflecting more realistic scenarios found in large-scale data environments. We have released the schemas and ground truths at [5] to serve as a robust benchmark for advancing research and improving schema matching solutions.

## 4 Related Work

There are two primary types of information commonly used in the matching process: (1) schema information and (2) instance information. We first examine how each type is utilized to establish correspondences, then discuss recent approaches that leverage deep learning and LLMs.

**Schema-based Approaches.** Schema-based methods utilize three types of metadata: linguistic knowledge, constraint knowledge, and structural knowledge. Linguistic knowledge derives meaning from attribute names and descriptions to identify semantic similarities between schema elements [28,10,52,23]. Constraint knowledge uses metadata such as data types, primary keys, foreign keys, and uniqueness constraints to establish correspondences [52,7]. Structural knowledge examines table relationships, including one-to-one and one-to-many relationships, as well as their cardinalities [52,7,27]. These methods assume the availability of rich metadata, but real-world databases often include abbreviations or domain-specific acronyms, which limit their effectiveness. To address this, approaches such as those by [46,34,16] incorporate external resources like dictionaries and thesauri to enhance schema semantics.

**Instance-based Approaches.** When schema metadata is limited, instance-based methods analyze data values, column semantics, and distribution patterns to identify matches [24,28]. For example, numeric data distributions can reveal semantic relationships that schema-based methods might miss. Machine learning (ML)-based matchers, such as those in [41], leverage data patterns to predict correspondences between schemas. However, these methods require significant amounts of data and often need retraining for new domains, limiting their scalability [11]. Additionally, instance-based approaches are computationally expensive and prone to errors when applied to large datasets. Some solutions combine schema and instance information to improve correspondence accuracy [6,10,52,45].

**Deep-learning and LLM-based Approaches.** Recent methods leverage deep learning and LLMs to uncover deeper relationships between schema elements. SMAT [50] and SemProp [13] use pre-trained word embeddings to move beyond syntactic similarities [38,26], while EmbDI [8] incorporates instance-level

information, and REMA [19] applies graph embeddings to capture structural relationships. LLM-based approaches further advance schema matching. LSM [51] combines LLMs with active learning and real-time user feedback. Our approach differs by operating without human-in-the-loop feedback. Parciak et al. [36] use LLMs as standalone matching tools but limit their scope to single-source and single-target table matching, which is less applicable to multi-table scenarios. Huang et al. [18] propose a schema matching tool involving offline PK/FK identification and sample generation, followed by online SQL transformation and iterative debugging. In contrast, our work avoids query generation and does not rely on strict SQL constraints like DISTINCT for correctness evaluation. Recent multi-table methods, such as REMATCH [43], show strong results with GPT-4 [1] but struggle with models like GPT-3.5 [32], which often fail to capture deeper relational patterns.

## 5 LLMATCH Framework

**Framework Design.** Unifying different schema matching methods presents significant challenges due to the diversity in their design:

- **Task Scope Differences:** Some methods are designed for table-to-table matching, focusing on aligning columns between individual tables. Others address schema-to-schema matching, which involves capturing relationships across multiple interrelated tables. These differences in scope require a flexible framework that accommodates both granular and holistic approaches.
- **Integrating Large Language Models (LLMs):** LLMs are increasingly ubiquitous in schema matching. However, their integration into existing methods raises fairness concerns. Some methods may use LLMs extensively, while others rely minimally or not at all. A fair framework must standardize LLM usage to enable meaningful comparisons across diverse approaches.

To address these challenges, we propose a unified framework with a focus on structured evaluation and fair comparison. The processing framework includes three major steps: *Schema Preparation*, *Table Selection* and *Column Matching*.

**Schema Preparation:** In this step, we categorize schema information into table and column names (syntactic content), primary and foreign keys (PK/FK) (structural relationships), and table/column descriptions (semantic metadata), enabling our framework to assess their impact on matching quality and identify the most influential features across different contexts. Additionally, we introduce **Rollup**, a schema transformation that merges semantically related columns into a single higher-level representation to simplify matching. As shown in Fig. 3, the *Rollup phase* aggregates columns such as: `target.gender_concept_id`, `target.gender_source_value`, `target.gender_source_concept_id` into a single column, `target.gender`. This reduces schema complexity before the matching process begins.

**Table Selection:** This step narrows down the set of target tables for each source table  $T_s$ , ensuring the next stage focuses only on relevant candidates. We compare four candidate selection strategies:

- **None:** Matches all target tables at once.
- **Nested Join:** Matches each target table separately in distinct prompts.
- **Vector Similarity:** Uses embeddings to select the top- $k$  most similar tables.
- **LLM Selection:** Lets the LLM determine relevant target tables based on schema context.

**Column Matching with Drilldown:** Once relevant tables are selected, column matching is performed. To enhance alignment, we introduce **Drilldown**, which refines matches by expanding rolled-up columns back into their original components. As shown in Fig. 3, once `source.gender` is matched to `target.gender`, the *Drilldown phase* decomposes it back into: `target.gender_concept_id`, `target.gender_source_value`, `target.gender_source_concept_id`. This approach ensures a fine-grained alignment.

Our framework provides a systematic and flexible approach to evaluating schema matching methods. By categorizing schema information, standardizing LLM integration, and structuring the matching process into table candidate selection and column matching, we enable fair comparisons and deeper insights into the performance of diverse methods. This structured evaluation advances the understanding and design of LLM-enabled schema matching pipelines.

## 6 Evaluation

We evaluate the following baselines using our proposed LLMATCH framework.

- COMA [10]: A classical schema matching approach that combines multiple schema-based matchers, representing schemata as rooted graphs.
- SF [25]: Transforms schemas into directed graphs and propagates similarity scores through neighboring nodes.
- CUPID [23]: Represents schemas as hierarchical tree structures and calculates similarity as a weighted sum of linguistic (name-based) and structural (context-based) similarities.
- UNICORN [48]: A the-state-of-the-art deep learning model that learns from multiple datasets and tasks using a mixture-of-experts model.
- REMATCH [43]: An LLM-based method that serializes tables into documents, generates document embeddings, and uses vector similarity to pair the most similar target tables with the source table. The selected tables are then fed into an LLM to produce column mappings.
- LLMATCH (our approach): We leverage LLMs for both table candidate selection and column matching. Similar to REMATCH, we include schema information such as names, primary/foreign keys, and descriptions in the prompt. Details of the two prompt templates and schema serialization format are provided in the extended report [5].

**Evaluation Metric.** We use the **F1 score** as the primary metric for performance evaluation, as it is widely adopted in schema matching research [36,50,3,10]. While alternative metrics such as `recall@k` [48] and `accuracy@k` [43] are also used



Strategy	Complex Schema Matching Task							Simple Schema Matching Task			
	imdb-sakila	cms-omop	cprd_au-omop	cprd_gold-omop	mimic_iii-omop	synthea-omop	bank1-bank2	mjs-mjt	msjs-msjt	mus-mut	mvs-mvt
COMA	0.18	0.01	0.13	0.03	0.04	0.08	<u>0.76</u>	0.91	0.71	0.57	<b>0.91</b>
SF	0.21	0.04	0.08	0.03	0.09	0.03	0.50	0.71	0.70	0.89	0.59
CUPID	0.00	0.01	0.00	0.00	0.00	0.00	0.32	0.33	0.80	0.86	0.33
UNICORN	0.00	0.05	0.05	0.01	0.01	0.03	0.38	<b>0.92</b>	0.80	0.87	0.80
REMATCH GPT-3.5	0.36	0.07	0.06	0.04	0.07	0.06	0.40	0.86	0.70	0.88	0.67
REMATCH GPT-4o-mini	<u>0.64</u>	0.16	<u>0.20</u>	<u>0.19</u>	<u>0.20</u>	0.19	<u>0.74</u>	<b>0.92</b>	<u>0.84</u>	0.81	0.71
LLMATCH GPT-3.5	0.47	<u>0.20</u>	0.15	0.16	0.13	<u>0.26</u>	0.67	<b>0.92</b>	0.82	<u>0.95</u>	<u>0.83</u>
LLMATCH GPT-4o-mini	<b>0.73</b>	<b>0.37</b>	<b>0.36</b>	<b>0.30</b>	<b>0.33</b>	<b>0.37</b>	<b>0.85</b>	<b>0.92</b>	<b>0.94</b>	<b>0.97</b>	<u>0.83</u>

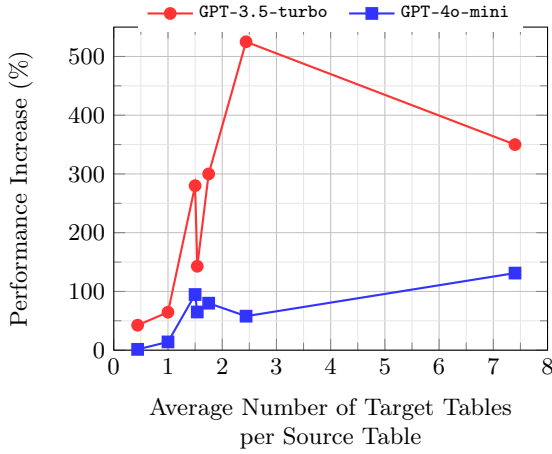
**Table 2:** F1 scores of schema matching methods on complex and simple tasks. Traditional methods (COMA, SF, CUPID, and UNICORN) perform well on simple schema matching tasks but experience a significant drop in performance on complex schema matching problems. In contrast, the two LLM-based methods, REMATCH and LLMATCH, achieve top performance across both task types, demonstrating their effectiveness in handling complex schema matching.

in related work, we prioritize the F1 score over rank-based metrics because our approach does not require the LLM to produce ranked results. Following the setup in [36], we ask the model to generate a list of mappings, treating each mapping with equal importance. In industrial applications, recall is prioritized over precision because identifying and addressing missing matches typically requires more human effort than correcting incorrect ones [36]. In this context, a high recall rate is generally more valuable than high precision, and the F1 score captures this balance effectively. For F1 score calculations, we treat mappings to foreign keys (FKs) as equivalent to their corresponding primary keys (PKs).

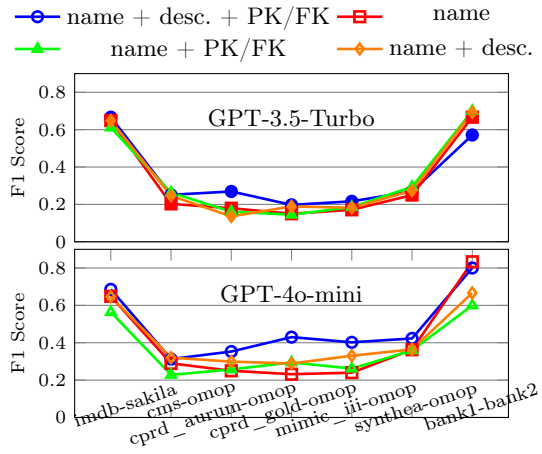
**Implementation.** We use the `gpt-3.5-turbo` and `gpt-4o-mini` models with their default settings, acknowledging the limited reproducibility of conversations due to the inherent non-determinism of these models [4]. We do not modify any model parameters or fine-tune the models. All embeddings in our experiments are generated using SBERT [40]. Since the specific wording of a prompt influences the output of an LLM [42], we include all prompt templates and schema serialization formats in the full report [5].

**Evaluation Goals.** The primary goal of our evaluation is to assess the effectiveness of LLMATCH for schema matching, particularly in complex, multi-table relational databases. To achieve this, we investigate several key questions: how different approaches perform on traditional and new benchmarks (Sec. 6.1), the impact of table selection methods on performance (Sec. 6.2), and the role of schema elements such as names, descriptions, and PK/FK relationships (Sec. 6.3). Additionally, we examine how multi-level schema matching with Rollup and Drilldown affects performance (Sec. 6.4), the implications of LLM context limits for large datasets (Sec. 6.5), and the productivity gains in schema matching with and without machine assistance (Sec. 6.6).





(a) Percentage performance increase of LLMATCH over SOTA (REMATCH) across datasets with varying complexity. Higher values indicate higher complexity.



(b) Schema element ablation study.: removing the PK/FK relationship and descriptions reduces the F1 score in the matching task, particularly for the more capable GPT-4o-mini model.

## 6.1 End-to-End Performance Evaluation

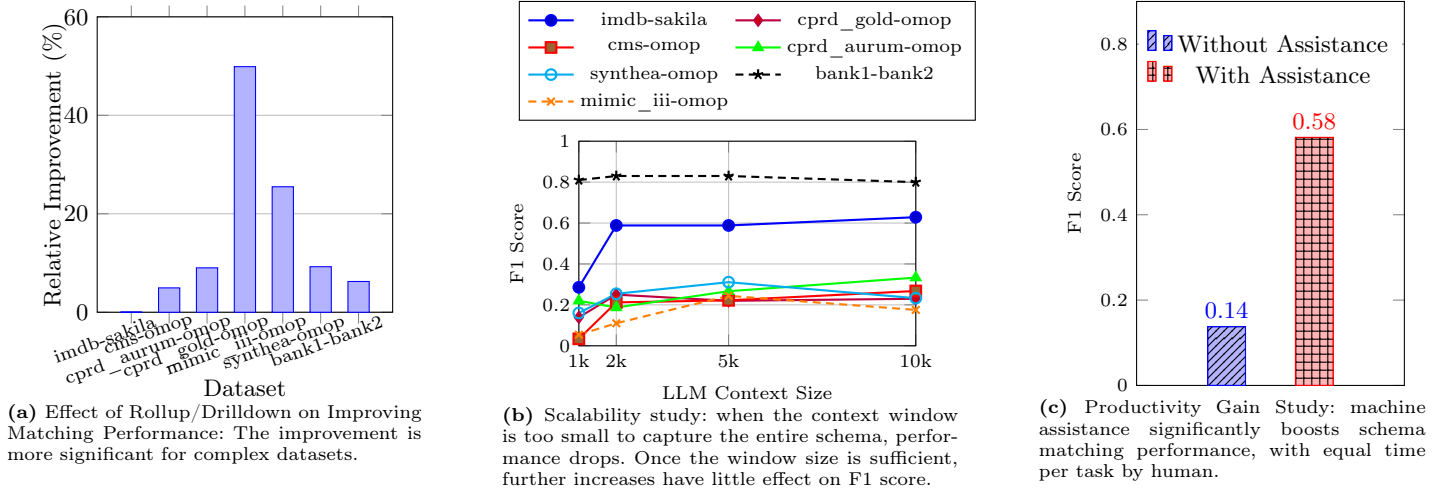
We conducted an extensive evaluation of our proposed method against several baseline strategies on both complex and simple schema matching tasks. The results are summarized in Table 2. For the simple schema matching tasks, both traditional methods and LLM-based methods performed well, achieving high scores on all datasets. However, traditional methods’ performance dropped significantly on complex tasks, a phenomenon also reported by LSM [51].

For the complex matching problems, LLMATCH significantly outperforms all baselines. For instance, on the largest dataset, `mimic_iii-omop`, LLMATCH achieved a score of 0.4, doubling the performance of REMATCH, which scored 0.2. Traditional methods like COMA and SF performed considerably worse, with scores of 0.04 and 0.09, respectively. An exception was observed with the `bank1-bank2` dataset, where traditional methods showed comparable performance. A detailed inspection revealed that the two schemas share a high degree of common vocabulary, making the task easier for lexical matching approaches.

Our method also consistently outperforms the state-of-the-art LLM matcher REMATCH across both GPT-3.5-turbo and GPT-4o-mini models. The gain becomes more substantial as the complexity of the matching task increases. Fig. 4a shows the percentage improvement relative to the average number of target tables matched with each source table. For simpler tasks, the difference in performance is minimal. However, as the number of target tables approaches 2, the improvement from our method becomes more pronounced. Notably, the performance gap is even wider for GPT-3.5-turbo, indicating that our approach is effective even with less capable models where other methods struggle.

Table Selection Strategy	imdb-sakila	cms-omop	cprd_aurum-omop	cprd_gold-omop	mimic_iii-omop	synthea-omop	bank1-bank2
	Score Count	Score Count	Score Count	Score Count	Score Count	Score Count	Score Count
None	0.39 (16.0)	0.21 (39.0)	0.28 (39.0)	0.22 (39.0)	0.13 (39.0)	0.23 (39.0)	0.74 (9.0)
nested_join	0.32 (1.0)	0.10 (1.0)	0.20 (1.0)	0.18 (1.0)	0.07 (1.0)	0.10 (1.0)	0.64 (1.0)
vector_similarity (top5)	<u>0.67</u> (5.0)	0.35 (5.0)	0.19 (5.0)	0.27 (5.0)	0.16 (5.0)	0.29 (5.0)	<u>0.85</u> (5.0)
vector_similarity (top10)	<u>0.67</u> (10.0)	<u>0.36</u> (10.0)	<u>0.33</u> (10.0)	<b>0.37</b> (10.0)	0.25 (10.0)	0.29 (10.0)	<b>0.88</b> (9.0)
vector_similarity (top15)	0.59 (15.0)	0.31 (15.0)	0.27 (15.0)	0.29 (15.0)	<u>0.26</u> (15.0)	<u>0.32</u> (15.0)	<b>0.88</b> (9.0)
llm (gpt-4o-mini)	<b>0.73</b> (4.0)	<b>0.37</b> (5.4)	<b>0.36</b> (6.3)	<u>0.30</u> (6.9)	<b>0.33</b> (6.9)	<b>0.37</b> (4.7)	<u>0.85</u> (4.6)

**Table 3:** Table Selection Strategy Study: F1 scores for various table selection methods, with the number of selected tables in parentheses. LLM-based methods select fewer tables but achieve higher quality. (Column matching uses LLM gpt-4o-mini).



**Fig. 5:** Comparison of Different Factors on Schema Matching Performance

## 6.2 Table Selection Strategy Study

In this section, we evaluate the impact of various table selection strategies (Step 2 in Figure 2) on schema matching performance. We compare four distinct table selection methods while keeping all other settings constant. For each source table, target table candidates are selected using the following strategies. Firstly, in the *None* strategy, no specific table selection is applied; all target tables are considered relevant, so the output size matches the number of tables in the target schema. Secondly, the *nested\_join* strategy matches each source table with a single target table in separate prompts during column matching. Thirdly, the *vector\_similarity* strategy, a commonly used filtering method, selects the top- $k$  target tables for each source table based on vector similarity scores, with  $k$  set to 5, 10, and 15 in our experiments. Lastly, the LLM strategy provides descriptions of the source and target tables to the LLM, which then identifies relevant tables, producing a variable number of outputs.

The experimental results are presented in Table 3. For each dataset, the first column displays the final F1 score of the column matching results, and the

second column shows the average number of target tables selected for each source table. From the results, it is evident that the quality of table selection impacts the final matching performance. The LLM-based strategies, with smaller output size, out-perform vector similarity methods.

This insight is significant because vector similarity methods have been the popular choice in retrieval-augmented generation (RAG) frameworks. Our findings suggest that LLMs possess a deeper understanding of the semantic relationships between tables, leading to better schema matching performance even with fewer tables selected.

### 6.3 Schema Elements Ablation Study

To assess the impact of schema elements—*name*, *description*, and *relationship*—on matching performance, we conducted an ablation study. Figure 4b reports F1 scores for GPT-3.5-Turbo and GPT-4o-mini across different schema element combinations. The results show that incorporating *description* and *relationship* improves performance, with the *name* + *description* + *relationship* combination yielding the highest F1 scores. The additional context helps disambiguate schema elements with similar names but different roles. GPT-4o-mini consistently outperforms GPT-3.5-Turbo, particularly when enriched with all schema elements, indicating its superior ability to capture complex semantic relationships. However, the performance gap varies across datasets, suggesting that dataset characteristics influence how effectively each model leverages schema information.

### 6.4 Effect of Rollup and Drilldown

The results, as shown in Fig. 5a, demonstrate the impact of Rollup and Drilldown on schema matching performance across multiple datasets. The relative improvement varies significantly depending on the dataset’s complexity and structure. The *cprd\_gold-omop* dataset exhibits the highest improvement, suggesting that Rollup effectively consolidates fragmented schema elements, allowing the LLM to establish stronger high-level correspondences before Drilldown refines the column-level mappings. Similarly, *mimic\_iii-omop* and *synthea-omop* show notable gains, indicating that multi-level schema matching is particularly beneficial for datasets with many-to-many schema relationships. However, datasets such as *imdb-sakila* and *bank1-bank2* show relatively minor improvements, likely due to their simpler schema structures where traditional column-level matching is already effective. This suggests that Rollup and Drilldown offer the greatest benefits when schema fragmentation is high, reinforcing the importance of hierarchical schema alignment in complex integration scenarios. These findings highlight the strengths of multi-level schema matching, demonstrating that schema abstraction (Rollup) and fine-grained alignment (Drilldown) significantly improve performance on complex datasets while providing only marginal benefits in simpler cases.

### 6.5 Scalability Study

In this experiment, we limit input context size while keeping output unrestricted due to technical constraints. To standardize evaluation across LLMs, we approximate context size using word count. Tasks exceeding a single-prompt limit are split into multiple prompts with small contexts, ensuring the largest source and target tables fit together. This constraint applies to both table selection and column matching. As shown in Fig. 5b, overly small contexts fragment the task, reducing performance. This underscores the need for LLMs with sufficient context capacity to handle large schemas effectively.

### 6.6 Productivity Gain Study

To assess the impact of LLM-generated matching results on manual schema matching tasks, we conducted an experiment with two scenarios: one without machine-generated results and one with them. Each scenario was limited to two minutes. Four data engineers participated, and the results, shown in Fig. 5c, indicate that machine assistance significantly enhances performance, even when LLM-generated matches are not entirely accurate. These findings align with prior research[51], which demonstrated that LLM assistance improves accuracy and reduces labeling costs by up to 81% compared to fully manual efforts.

## 7 Conclusion

This paper presents a framework that decomposes schema matching into three stages: schema preparation, table selection, and column matching. A key feature of our approach is the introduction of Rollup and Drilldown, a multi-level schema matching addon that simplifies complex schemas before matching and refines alignments afterward. Through this framework, we also found that vector similarity, a widely used method, is not the most effective choice for table selection, offering valuable insights for other LLM-based applications. Additionally, we developed and released a comprehensive multi-industry complex schema matching benchmark to support future research in this field.

## References

1. Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F.L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al.: Gpt-4 technical report. arXiv preprint arXiv:2303.08774 (2023)
2. Ackerman, R., Gal, A., Sagi, T., Shraga, R.: A cognitive model of human bias in matching. In: PRICAI 2019: Trends in Artificial Intelligence: 16th Pacific Rim International Conference on Artificial Intelligence, Cuvu, Yanuca Island, Fiji, August 26–30, 2019, Proceedings, Part I 16. Springer (2019)
3. Alwan, A.A., Nordin, A., Alzeber, M., Abualkashik, A.Z.: A survey of schema matching research using database schemas and instances. International Journal of Advanced Computer Science and Applications **8**(10) (2017)

4. Anadkat, S.: How to make your completions outputs consistent with the new seed parameter (2023), [https://cookbook.openai.com/examples/reproducible\\_outputs\\_with\\_the\\_seed\\_parameter](https://cookbook.openai.com/examples/reproducible_outputs_with_the_seed_parameter), openAI Cookbook
5. Anonymous Authors: Llmitch code and dataset. <https://anonymous.4open.science/r/LLMatch-ACOE/README.md> (2025)
6. Asif-Ur-Rahman, M., Hossain, B.A., Bewong, M., Islam, M.Z., Zhao, Y., Groves, J., Judith, R.: A semi-automated hybrid schema matching framework for vegetation data integration. *Expert Systems with Applications* **229** (2023)
7. Bernstein, P.A., Madhavan, J., Rahm, E.: Generic schema matching, ten years later. *Proceedings of the VLDB Endowment* **4**(11) (2011)
8. Cappuzzo, R., Papotti, P., Thirumuruganathan, S.: Embdi: generating embeddings for relational data integration. In: 29th Italian Symposium on Advanced Database Systems (SEDB), Pizzo Calabro, Italy (2021)
9. Development, J.R.: Etl lambdabuilder documentation (2024), <https://ohdsi.github.io/ETL-LambdaBuilder/>, accessed: 2025-03-01
10. Do, H.H., Rahm, E.: Coma—a system for flexible combination of schema matching approaches. In: VLDB’02: Proceedings of the 28th International Conference on Very Large Databases. Elsevier (2002)
11. Feng, J., Hong, X., Qu, Y.: An instance-based schema matching method with attributes ranking and classification. In: 2009 Sixth International Conference on Fuzzy Systems and Knowledge Discovery. vol. 5. IEEE (2009)
12. Fernandez, R.C., Elmore, A.J., Franklin, M.J., Krishnan, S., Tan, C.: How large language models will disrupt data management. *Proc. VLDB Endow.* **16**(11) (Jul 2023). <https://doi.org/10.14778/3611479.3611527>, <https://doi.org/10.14778/3611479.3611527>
13. Fernandez, R.C., Mansour, E., Qahtan, A.A., Elmagarmid, A., Ilyas, I., Madden, S., Ouzzani, M., Stonebraker, M., Tang, N.: Seeping semantics: Linking datasets using word embeddings for data discovery. In: 2018 IEEE 34th International Conference on Data Engineering (ICDE). IEEE (2018)
14. Financial Times: Deutsche Bank struggles with fallout after huge Post-bank IT migration: Regulator makes unprecedented rebuke as many clients locked out of their accounts for weeks. <https://www.ft.com/content/4138876c-5a10-46d4-b6c6-7d421cbd9df0> (2025), accessed: 2025-03-01
15. Frid, S., Pastor Duran, X., Bracons Cucó, G., Pedrera-Jiménez, M., Serrano-Balazote, P., Muñoz Carrero, A., Lozano-Rubí, R.: An ontology-based approach for consolidating patient data standardized with european norm/international organization for standardization 13606 (en/iso 13606) into joint observational medical outcomes partnership (omop) repositories: Description of a methodology. *JMIR Med Inform* **11**, e44547 (Mar 2023). <https://doi.org/10.2196/44547>, <https://medinform.jmir.org/2023/1/e44547>
16. Giunchiglia, F., Shvaiko, P., Yatskevich, M.: S-match: an algorithm and an implementation of semantic matching. In: The Semantic Web: Research and Applications: First European Semantic Web Symposium, ESWS 2004 Heraklion, Crete, Greece, May 10-12, 2004. Proceedings 1. Springer (2004)
17. Haberson, A., Rinner, C., Schoberl, A., Gall, W.: Feasibility of mapping austrian health claims data to the omop common data model. *Journal of Medical Systems* **43**(10), 314 (2019). <https://doi.org/10.1007/s10916-019-1436-9>, <https://doi.org/10.1007/s10916-019-1436-9>
18. Huang, Z., Guo, J., Wu, E.: Transform table to database using large language models. *Proceedings of the VLDB Endowment*. ISSN **2150** (2024)

19. Koutras, C., Fragkoulis, M., Katsifodimos, A., Lofi, C.: Rema: Graph embeddings-based relational schema matching. In: EDBT/ICDT Workshops (2020)
20. Koutras, C., Siachamis, G., Ionescu, A., Psarakis, K., Brons, J., Fragkoulis, M., Lofi, C., Bonifati, A., Katsifodimos, A.: Valentine: Evaluating matching techniques for dataset discovery. In: ICDE. IEEE (2021)
21. Lamer, A., Abou-Arab, O., Bourgeois, A., Parrot, A., Popoff, B., Beuscart, J.B., Tavernier, B., Moussa, M.D.: Transforming anesthesia data into the observational medical outcomes partnership common data model: Development and usability study. *J Med Internet Res* **23**(10), e29259 (Oct 2021). <https://doi.org/10.2196/29259>, <https://www.jmir.org/2021/10/e29259>
22. Lamer, A., Depas, N., Doutreligne, M., Parrot, A., Verloop, D., Defebvre, M.M., Ficheur, G., Chazard, E., Beuscart, J.B.: Transforming french electronic health records into the observational medical outcome partnership’s common data model: a feasibility study. *Applied clinical informatics* **11**(01), 013–022 (2020)
23. Madhavan, J., Bernstein, P.A., Rahm, E.: Generic schema matching with cupid. In: *vldb*. vol. 1 (2001)
24. Mehdi, O., Ibrahim, H., Affendey, L.: An approach for instance based schema matching with google similarity and regular expression. *International Arab Journal of Information Technology (IAJIT)* **14**(5) (2017)
25. Melnik, S., Garcia-Molina, H., Rahm, E.: Similarity flooding: A versatile graph matching algorithm and its application to schema matching. In: *Proceedings 18th international conference on data engineering*. IEEE (2002)
26. Mikolov, T.: Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781* (2013)
27. Mukherjee, D., Bandyopadhyay, A., Chowdhury, R., Bhattacharya, I.: Learning knowledge graph for target-driven schema matching. In: *Proceedings of the 3rd ACM India Joint International Conference on Data Science & Management of Data (8th ACM IKDD CODS & 26th COMAD)* (2021)
28. Munir, S., Khan, F., Riaz, M.A.: An instance based schema matching between opaque database schemas. In: *2014 4th International Conference on Engineering Technology and Technopreneuship (ICE2T)*. IEEE (2014)
29. OHDSI: Observational Health Data Sciences and Informatics (2019)
30. OHDSI: Omop common data model - version 5.4. <https://ohdsi.github.io/CommonDataModel/cdm54.html> (2024), accessed: 2025-03-01
31. Oja, M., Tamm, S., Mooses, K., Pajusalu, M., Talvik, H.A., Ott, A., Laht, M., Malk, M., Lõo, M., Holm, J., Haug, M., Šuvalov, H., Särg, D., Vilo, J., Laur, S., Kolde, R., Reisberg, S.: Transforming Estonian health data to the Observational Medical Outcomes Partnership (OMOP) Common Data Model: lessons learned. *JAMIA Open* **6**(4), ooad100 (12 2023). <https://doi.org/10.1093/jamiaopen/ooad100>, <https://doi.org/10.1093/jamiaopen/ooad100>
32. OpenAI: Gpt-3.5 turbo fine-tuning and api updates. <https://openai.com/index/gpt-3-5-turbo-fine-tuning-and-api-updates/> (August 22 2023), retrieved Feb 12, 2025
33. Overhage, J.M., Ryan, P.B., Reich, C.G., Hartzema, A.G., Stang, P.E.: Validation of a common data model for active safety surveillance research. *Journal of the American Medical Informatics Association* **19**(1) (2012)
34. Palopoli, L., Terracina, G., Ursino, D., et al.: The system dike: Towards the semi-automatic synthesis of cooperative information systems and data warehouses. In: *ADBIS-DASFAA Symposium* (2000)
35. Pan, Z., Yang, M., Monti, A.: Schema matching based on energy domain pre-trained language model. *Energy Informatics* **6**(Suppl 1) (2023)

36. Parciak, M., Vandevoort, B., Neven, F., Peeters, L.M., Vansummeren, S.: Schema matching with large language models: An experimental study. Joint Workshops at the 50th International Conference on Very Large Data Bases (VLDBW'24) — TaDA'24: 2nd International Workshop on Tabular Data Analysis (2024)
37. Paris, N., Lamer, A., Parrot, A.: Transformation and evaluation of the mimic database in the omop common data model: development and usability study. *JMIR Medical Informatics* **9**(12) (2021)
38. Pennington, J., Socher, R., Manning, C.D.: Glove: Global vectors for word representation. In: *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)* (2014)
39. Reich, C., Ostropelets, A., Ryan, P., Rijnbeek, P., Schuemie, M., Davydov, A., Dymshyts, D., Hripcsak, G.: Ohdsi standardized vocabularies—a large-scale centralized reference ontology for international data harmonization. *Journal of the American Medical Informatics Association* **31**(3) (2024)
40. Reimers, N., Gurevych, I.: Sentence-bert: Sentence embeddings using siamese bert-networks. In: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics (11 2019)
41. Rong, S., Niu, X., Xiang, E.W., Wang, H., Yang, Q., Yu, Y.: A machine learning approach for instance matching based on similarity metrics. In: *The Semantic Web—ISWC 2012: 11th International Semantic Web Conference, Boston, MA, USA, November 11–15, 2012, Proceedings, Part I* 11. Springer (2012)
42. Schulhoff, S., Ilie, M., Balepur, N., Kahadze, K., Liu, A., Si, C., Li, Y., Gupta, A., Han, H., Schulhoff, S., et al.: The prompt report: A systematic survey of prompting techniques. *arXiv preprint arXiv:2406.06608* (2024)
43. Sheerit, E., Brief, M., Mishaeli, M., Elisha, O.: Rematch: Retrieval enhanced schema matching with llms. *arXiv preprint arXiv:2403.01567* (2024)
44. Shrager, R., Amir, O., Gal, A.: Learning to characterize matching experts. In: *2021 IEEE 37th International Conference on Data Engineering (ICDE)*. IEEE (2021)
45. Shrager, R., Gal, A., Roitman, H.: Adnev: Cross-domain schema matching using deep similarity matrix adjustment and evaluation. *Proceedings of the VLDB Endowment* **13**(9) (2020)
46. Sorrentino, S., Bergamaschi, S., Gawinecki, M., Po, L.: Schema label normalization for improving schema matching. *Data & Knowledge Engineering* **69**(12) (2010)
47. Stang, P.E., Ryan, P.B., Racoosin, J.A., Overhage, J.M., Hartzema, A.G., Reich, C., Welebob, E., Scarnecchia, T., Woodcock, J.: Advancing the science for active surveillance: rationale and design for the observational medical outcomes partnership. *Annals of internal medicine* **153**(9) (2010)
48. Tu, J., Fan, J., Tang, N., Wang, P., Li, G., Du, X., Jia, X., Gao, S.: Unicorn: A unified multi-tasking model for supporting matching tasks in data integration. *Proceedings of the ACM on Management of Data* **1**(1) (2023)
49. Wornow, D., Suh, G., Elmore, A.J., Krishnan, S., Parameswaran, A.: Automating the enterprise with foundation models. *Proceedings of the VLDB Endowment* **17**(12) (2024)
50. Zhang, J., Shin, B., Choi, J.D., Ho, J.C.: Smat: An attention-based deep learning solution to the automation of schema matching. In: *Advances in Databases and Information Systems: 25th European Conference, ADBIS*. Springer (2021)
51. Zhang, Y., Floratou, A., Cahoon, J., Krishnan, S., Müller, A.C., Banda, D., Psalidas, F., Patel, J.M.: Schema matching using pre-trained language models. In: *International Conference on Data Engineering (ICDE)*. IEEE (2023)
52. Zhao, H., Ram, S.: Combining schema and instance information for integrating heterogeneous data sources. *Data & Knowledge Engineering* **61**(2) (2007)



16      No Author Given

## **A    Appendix**

Listing 1.1: Table Selection Prompt

```

1 You are an expert in matching database schemas. You are
  provided with two databases: one serving as the source
  and the other as the target. Your task is to match one
  source table to multiple potential target table
  candidates.
2
3 **Objective**: Identify and list all potential target tables
  that can map to columns in the given source table.
4
5 **Source Table Details**:
6 {{source_table}}
7
8 **Target Tables Details**:
9 {{target_tables}}
10
11 **Matching Criteria**:
12 - Identify target tables may potentially have columns that
  can be matched to the source table.
13 - One source column might be matched to multiple target
  columns. Redundant matches are allowed.
14
15 **Expected Output**:
16 Provide the matches in the following JSON format:
17 ```json
18 {
19   "source_table1": [
20     {
21       "target_table": "target_table1",
22       "reasoning": "..."
23     },
24     {
25       "target_table": "target_table2",
26       "reasoning": "..."
27     }
28   ]
29 }
30 ```
31 Return only the JSON object and no other text.

```

Listing 1.2: Column Matching Prompt

```

1 You are an expert in databases. Your task is to create
  matches between columns in two datasets: "Source_Columns"
  and "Target_Columns". One source column can be matched
  to multiple target columns. The matches should be based
  on the semantic similarity of the entities described by
  the columns, considering the context provided in their
  descriptions.

```

```

2
3 **Matching Criteria:**
4
5     - Entity Similarity: The matched entries should
      describe the same or very similar entities. The
      source entry can be part of the target entry and
      vice versa. e.g. full_name => first_name,
      last_name. registration_date => registration_date
      , registration_time
6     - Contextual Alignment: each column represent
      different types of entities. make sure that the
      matched columns are of the same type. Negative
      examples: bank.name != bank_branch.name, customer
      .name != staff.name
7     - Data Type Compatibility: Ensure that the data
      types of the matched columns are compatible. A
      single element can be matched with multiple
      elements and vice versa. e.g. source_table.
      language => target_table.languages.
8
9 **Instructions:**
10
11     1. Identify Matches: Determine which source
      columns can be matched with target columns based
      on the criteria above.
12     2. Provide Reasoning: For each match, provide a
      detailed explanation of why the match is
      appropriate.
13     3. Review Matches: Ensure data belongs to the same
      domain and is semantically similar.
14
15 **Output Format:**
16
17     - Provide the matches in the following JSON
      format:
18     ``json
19     {
20         "source_table1.source_column1": [
21             {
22                 "mapping": "target_table1.target_column1",
23                 "reasoning": "explanation_of_the_match"
24             },
25             {
26                 "mapping": "target_table2.target_column2",
27                 "reasoning": "explanation_of_the_match"
28             }
29         ],
30         "source_table2.source_column2": [
31             {
32                 "mapping": "None",

```

```

33         "reasoning": "...
34     }
35 ]
36 }
37 '''
38 **Source Tables:**
39 {{source_columns}}
40
41 **Target Tables:**
42 {{target_columns}}
43
44 Return only the JSON object and no other text.

```

**Listing 1.3:** Example of target table serialization during table selection. It demonstrates a concise schema representation optimized to fit within a single prompt.

```

1 {
2   "account": {
3     "description": "Represents bank accounts, specifying type
4     , balance, and associated branch.",
5     "foreign_keys": "branch_id=>branch.branch_id",
6     "non_foreign_key_columns": "account_balance, account_id,
7     account_type"
8   },
9   "branch": {
10    "description": "Represents a branch of the bank,
11    containing basic information about each branch.",
12    "foreign_keys": "",
13    "non_foreign_key_columns": "assets, branch_address,
14    branch_id, branch_name"
15  },
16  ...
17 }

```

**Listing 1.4:** Table Serialization Format. This is the full schema representation. It is used to serialize the source table during table selection and column matching. The same format is applied to the target table during column matching.

```

1 {
2   "account": {
3     "table": "account",
4     "columns": {
5       "acc_type": {
6         "name": "acc_type",
7         "description": "Type of account (e.g., savings,
8         checking)"
9       },
10      "account_no": {
11        "name": "account_no",
12        "description": "Unique identifier for each account",

```

```

12         "is_primary_key": true,
13         "foreign_keys": [
14             "hold_by.account_no",
15             "maintain.account_no"
16         ]
17     },
18     "balance": {
19         "name": "balance",
20         "description": "Current balance in the account"
21     },
22     "branch_id": {
23         "name": "branch_id",
24         "description": "Foreign key linking the account to the branch where it is maintained",
25         "is_foreign_key": true,
26         "linked_entry": "branch.branch_id"
27     }
28 },
29 "table_description": "Account stores information about accounts held by customers at the bank.";
30 },
31 "branch": {
32     "table": "branch",
33     "columns": {
34         "address": {
35             "name": "address",
36             "description": "Address of the branch"
37         },
38         "bank_code": {
39             "name": "bank_code",
40             "description": "Foreign key linking the branch to its bank",
41             "is_foreign_key": true,
42             "linked_entry": "bank.code"
43         },
44         "branch_id": {
45             "name": "branch_id",
46             "description": "Unique identifier for each branch",
47             "is_primary_key": true,
48             "foreign_keys": [
49                 "account.branch_id",
50                 "loan.branch_id",
51                 "maintain.branch_id",
52                 "offer.branch_id"
53             ]
54         },
55         "name": {
56             "name": "name",
57             "description": "Name of the branch"
58         }
59     }
60 }

```

```

59     },
60     "table_description": "Branch stores details about each
        branch of the bank.";
61   },
62   ...
63 }

```

Listing 1.5: Rollup (column merge) Prompt

```

1 Task: Database Schema Pre-Processing for Data Migration
2
3 Objective:
4 You will be provided with a database schema in JSON format.
  Your goal is to identify columns within the same table
  that hold related or similar data (e.g., first_name and
  last_name, start_date and start_datetime, end_date and
  end_datetime). Then, merge these related columns into a
  single column in your output.
5
6 Steps:
7   1. Read the Input JSON Schema
8     - Understand the structure of each table.
9     - Look for columns in the same table that are similar
      or can be combined logically.
10  2. Identify Mergable Columns
11     - If two or more columns represent "parts of a whole"
      (like first and last name, or address1 and
      address2), these can be merged.
12     - If the columns are not logically related, do not
      merge them.
13  - If two columns represent the same data but are named
      differently, they can be merged. e.g. table.address,
      table.address_id
14  3. Create a Merged Output
15     - For each table with identified merges, create an
      entry in a JSON array named "merged_columns".
16     - Inside each entry, show:
17     - table_name: The name of the table.
18     - merged_columns: An array of merged column
      definitions, each of which must include:
19       1. column_name: The new merged name.
20       2. column_description: A short explanation of
      why or how these columns are merged.
21       3. original_columns: A list of the fully
      qualified column references (e.g., "customer.
      first_name", "customer.last_name").
22  4. Output Format
23     - Your final output must be valid JSON.
24     - Use the following structure (with any necessary
      merges included):

```

```

25 {
26   {
27     "merged_columns": [
28       {
29         "table_name": "<TABLE_NAME>",
30         "merged_columns": [
31           {
32             "column_name": "<MERGED_COLUMN_NAME>",
33             "column_description": "<WHY_MERGED>",
34             "original_columns": [
35               "<table.columnA>",
36               "<table.columnB>"
37             ]
38           }
39         ]
40       }
41     ]
42   }
43
44
45
46 Important:
47   - Only list tables that actually have merged columns.
48   - If a table has no columns that need merging, skip it.
49   - Do not delete or rename any columns unless they are
50     part of a merged set.
51
52 Sample Input
53
54 (The JSON schema from your database might be large, so here
55  is a short version showing two tables as an example.)
56
57 {
58   "actor": {
59     "table": "actor",
60     "columns": {
61       "actor_id": {
62         "name": "actor_id",
63         "description": "Primary_key"
64       },
65       "first_name": {
66         "name": "first_name",
67         "description": "Actor_first_name"
68       },
69       "last_name": {
70         "name": "last_name",
71         "description": "Actor_last_name"
72       }
73     }
74   },

```



```

73  "customer": {
74    "table": "customer",
75    "columns": {
76      "customer_id": {
77        "name": "customer_id",
78        "description": "Primary_key"
79      },
80      "first_name": {
81        "name": "first_name",
82        "description": "Customer_first_name"
83      },
84      "last_name": {
85        "name": "last_name",
86        "description": "Customer_last_name"
87      }
88    }
89  }
90 }

```

#### Sample Output

Here is the kind of JSON output we expect after identifying mergable columns. Notice how first\_name and last\_name have been combined into one column called name. Adjust the language/description to fit your needs:

```

95 {
96   "tables": [
97     {
98       "table_name": "actor",
99       "merged_columns": [
100         {
101           "column_name": "name",
102           "column_description": "Actor's_full_name_(merged_
103             from_first_name_and_last_name)",
104           "original_columns": [
105             "actor.first_name",
106             "actor.last_name"
107           ]
108         }
109       ]
110     },
111     {
112       "table_name": "customer",
113       "merged_columns": [
114         {
115           "column_name": "name",
116           "column_description": "Customer's_full_name_(merged_
117             from_first_name_and_last_name)",

```

```
118         "customer.first_name",
119         "customer.last_name"
120     ]
121 }
122 ]
123 }
124 ]
125 }
126
127 Input :
```