

KNOWTRANS: Boosting Transferability of Data Preparation LLMs via Knowledge Augmentation

Yuhang Ge[†], Fengyu Li[†], Yuren Mao[†], Yanbo Yang[†], Congcong Ge[†], Zhaoqiang Chen[‡], Jiang Long[†], Yunjun Gao[†]

[†]Zhejiang University, Hangzhou, China

[‡]Huawei Cloud Computing Technologies Co., Ltd, China

[†]{yuhangge,fengyuli,yuren.mao,yanboyang,gcc,longjiangl,gaoyj}@zju.edu.cn

[‡]chenzhaoqiang1@huawei.com

Abstract—Data Preparation (DP), which involves tasks such as data cleaning, imputation and integration, is a fundamental process in data-driven applications. Recently, Large Language Models (LLMs) fine-tuned for DP tasks, i.e., DP-LLMs, have achieved state-of-the-art performance. However, transferring DP-LLMs to novel datasets and tasks typically requires a substantial amount of labeled data, which is impractical in many real-world scenarios. To address this issue, we propose a knowledge augmentation framework for data preparation, dubbed KNOWTRANS. This framework allows DP-LLMs to be transferred to novel datasets and tasks with a few data, significantly decreasing the dependence on extensive labeled data. KNOWTRANS comprises two components: *Selective Knowledge Concentration* and *Automatic Knowledge Bridging*. The first component re-uses knowledge from previously learned tasks, while the second automatically integrates additional knowledge from external sources. Extensive experiments on 13 datasets demonstrate the effectiveness of KNOWTRANS. KNOWTRANS boosts the performance of the state-of-the-art DP-LLM, Jellyfish-7B, by an average of 4.93%, enabling it to outperform both GPT-4 and GPT-4o.

Index Terms—data preparation, large language model

I. INTRODUCTION

Data Preparation (DP) aims to transform raw data into a form that is ready for analysis, which is one of the most important problems in data management [2]. It often requires extensive manual intervention and is time-consuming. To avoid the labor-intensive operations, deep learning models have been widely used to automatically solve the data preparation tasks, such as Ditto [3] for Entity Matching (EM), SMAT [3] for Schema Matching (SM), Doduo [4] for Column Type Annotation (CTA) and IPM [5] for Data Imputation (DI). However, these models are typically dataset-specific or task-specific, necessitating specialized model training on each dataset. Consequently, the transferability of these methods is constrained.

To improve the transferability, recent studies focus on unifying diverse input formats into a text-to-text form, and fine-tuning open-source Large Language Models (LLMs) to handle multiple DP tasks concurrently [6]. These methods typically implement multi-task supervised fine-tuning (SFT) on various data preparation datasets, achieving state-of-the-art performance [7, 8, 9, 10]. We refer to this multi-task SFT stage as *upstream learning*, and the fine-tuned LLMs as *upstream DP-LLMs*. Although upstream DP-LLMs can perform well on

trained datasets, they still require a significant amount of labeled data to adapt to new datasets and tasks. It is often expensive and impractical in many real-world situations. To relieve the dependence on labeled data, MELD [10] proposes a heuristic data augmentation technique in few-shot settings. However, it still requires a significant number of labeled examples, which is typically up to 10% of the dataset. In some cases, thousands of labeled samples are needed, such as for SM, CTA, and DI. Overall, the existing methods are data-greedy and have inferior few-shot transferability.

To boost the few-shot transferability of DP-LLMs, this paper investigates the reasons that necessitate a large number of labeled data and proposes a novel framework for adapting DP-LLMs to novel datasets and tasks with a few data (e.g., only twenty labeled data). We present that inferior few-shot transferability of DP-LLMs is primarily caused by **Cross-Datasets/Tasks Knowledge Distraction** and **Dataset-Informed Knowledge Gap**.

- **Cross-Datasets/Tasks Knowledge Distraction.** During upstream learning, DP-LLM simultaneously learns from multiple data preparation datasets or tasks, each containing unique knowledge. However, the fine-tuning process forces the model to update all datasets/tasks within a shared parameter space, leading to gradient conflicts. As shown in Fig. 1 (left), when learning two datasets/tasks, \mathcal{L}_1 and \mathcal{L}_2 , conflicting gradient update directions can create a “tug-of-war” effect [11, 12], preventing the model parameters from converging effectively to an optimal state. Ultimately, this conflict results in overlapping parameter representations for different datasets/tasks. When we perform cross-datasets/tasks transfer, DP-LLM is susceptible to interference from the knowledge of others stored in the model weights. We refer to this phenomenon as “knowledge distraction”.
- **Dataset-Informed Knowledge Gap.** Different data preparation datasets may have unique data-informed rules for processing, which is hard to capture from limited labeled samples. In Fig. 1 (right), we illustrate an example using the EM dataset Walmart-amazon, where the discrimination of whether two entities match typically involves the following knowledge: 1) primary identifiers are the product’s model numbers; 2) in case of missing or “NaN” values, focus on comparing other attributes; 3) product prices can

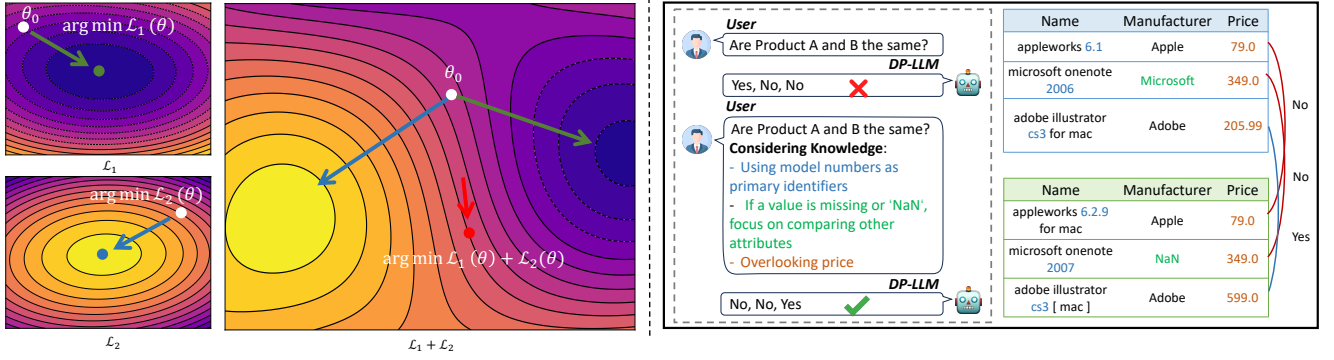


Fig. 1: (Left) Illustration of the loss landscape for \mathcal{L}_1 , \mathcal{L}_2 and $\mathcal{L}_1 + \mathcal{L}_2$. The “tug-of-war” effect arises when conducting multi-task upstream learning, such as the two tasks in $\mathcal{L}_1 + \mathcal{L}_2$ have nearly opposite gradient update directions, i.e. the angle is obtuse. (Right) Comparison without v.s. with knowledge for DP-LLM on EM task. DP-LLM cannot handle all cases without knowledge, but it succeeds when knowledge is applied. Text colors represent corresponding matched knowledge.

be disregarded. However, if the DP-LLM model has not encountered this dataset during upstream learning, it may lack this specific knowledge. Previous methods have tried to manually construct this knowledge [7], but this approach is both labor-intensive and costly. Thus, this gap hinders the DP-LLM’s performance in a few-shot setting.

To address the above issues, we propose a novel framework, **Knowledge Augmentation** for boosting the few-shot **Transferability** of DP-LLM, dubbed **KNOWTRANS**. KNOWTRANS only requires a small amount of labeled data—specifically, just twenty labeled samples. KNOWTRANS integrates two essential components: the *Selective Knowledge Concentration* (SKC) and the *Automatic Knowledge Bridging* (AKB). Specifically, SKC extracts knowledge from the upstream dataset into modular “knowledge patches” that are isolated from each other in parameter space. These knowledge patches are integrated into the upstream model and selectively re-used for few-shot fine-tuning, ensuring that the DP-LLM concentrates on valuable knowledge for new datasets and tasks. On the other hand, AKB employs a more capable closed-source LLM (e.g., GPT-4o) to incorporate knowledge derived from the dataset. Through an iterative process driven by an error feedback mechanism, it generates accurate, dataset-informed knowledge. This optimized knowledge is subsequently utilized as supplementary prompts during inference, effectively addressing the dataset-informed knowledge gap.

We conducted extensive experiments on 13 datasets, including both novel datasets and tasks in a few-shot setting (e.g., 20 samples per dataset). Our results show that the KNOWTRANS significantly improves the transferability of 4 different backbone DP-LLMs across datasets. For example, KNOWTRANS improves the state-of-the-art DP-LLM Jellyfish-7B’s performance by **4.93%** on average, surpassing other DP-LLMs of the same size and outperforming GPT-4 and GPT-4o. Furthermore, KNOWTRANS boosts Jellyfish-13B by **6.1%**, significantly outperforming GPT-4 and GPT-4o by **7.03%** and **6.07%**, respectively. Furthermore, the framework shows impressive scalability, consistently outperforming the performance of the backbone DP-LLM as the labeled sample size increases.

The remaining sections are organized as follows: Section II covers the preliminaries for this work. Section III defines the problem. Section IV provides an overview of our proposed framework KNOWTRANS. Section V details the SKC component of KNOWTRANS, while Section VI focuses on the AKB component of KNOWTRANS. Section VII presents experimental results on novel downstream datasets and tasks. Section VIII discusses related work, and Section IX concludes the paper.

II. PRELIMINARIES

A. Large Language Model

Large Language Models (LLMs), such as GPT-3 [13], LLaMA [14] and Mistral [15], are pre-trained on large-scale amounts of data, acquiring extensive world knowledge [16]. Thanks to the scale-up of parameters and data, LLMs exhibit emergent abilities [17], which refer to unexpected capabilities that arise only when the model reaches a certain scale, allowing it to perform tasks it was not explicitly trained for. As a result, these models demonstrate remarkable generalization capabilities across various natural language processing tasks [18]. Recently, some works [19, 20] begin to apply LLMs to data preparation tasks. However, LLMs are trained on general text data, which remains a gap between tabular data used in DP tasks. To bridge the gap, instruction tuning is commonly employed to adapt LLMs to DP tasks to achieve better performance.

Definition 1. (*Instruction Tuning*) A learning paradigm where an LLM is trained to follow natural language instructions for a variety of tasks by fine-tuning instruction datasets. The LLM learns to generalize across tasks by understanding and performing based on the provided instructions. For instance, the prompt below demonstrates a task where the model is asked to generate the antonym of a given word according to the instruction [21].

$$\text{LLM}(\underbrace{\text{Instruction: “Give the antonym of fast.”}}_{\text{Instruction}}) \rightarrow \underbrace{\text{slow}}_{\text{Response}}$$

B. Multi-Task Learning

Multi-task learning is a training paradigm where a model learns multiple tasks simultaneously, using different datasets

for each task. The aim is to improve the model’s performance on all tasks by sharing useful information across them. Given a set of **tasks** \mathcal{T} with associated datasets \mathcal{D} , a **multi-task model** f is trained to be capable of handling all tasks [11, 22, 23].

C. Model Fusion

Model Fusion (a.k.a. model merging) aims to integrate the parameters of multiple pre-trained models into a single model. This process combines the knowledge and capabilities of the individual models, resulting in a unified model that exhibits enhanced performance and generalization relative to the constituent models [24, 25, 26].

Vanilla Model Fusion. The vanilla model fusion method is a straightforward method where parameters from multiple models are combined by averaging. Formally, given a set of models $\{M_1, M_2, \dots, M_n\}$, where each model M_i is represented by a parameter set θ_i , the averaged merged model M^* is defined by the parameter vector θ^* , formulated as follows:

$$\theta^* = \frac{1}{n} \sum_{i=1}^n \theta_i, \quad (1)$$

where n is the number of models being merged. This equation signifies that each parameter of θ^* is the arithmetic mean of the corresponding parameters across all models.

III. PROBLEM DEFINITION

Data preparation (DP) typically involves cleaning, transforming, and organizing raw data for analysis. Formally, let T be a table (e.g., a relational or web table) with rows $\{r_1, r_2, \dots\}$ and columns $\{c_1, c_2, \dots\}$. Each cell $v_{i,j}$ represents the value in row i and column j . We define a set of data preparation tasks, each with training instances \mathcal{X} and labels \mathcal{Y} . A multi-task model f is applied to each task, taking specific inputs and producing corresponding outputs. This work focuses on the following DP tasks:

Entity Matching (EM). Binary classification task. Given two rows r_1 and r_2 from table T , determine whether they refer to the same real-world entity, i.e.,

$$f(r_1, r_2) \rightarrow \{0, 1\}, \text{ where } 1 \text{ indicates a match.}$$

Data Imputation (DI). Open-domain generation task. Given a missing cell value $v_{i,j}$ of row r , infer the missing value based on other values in the same row.

$$f(v_{i,j}, r) \rightarrow \hat{v}_{i,j}.$$

Schema Matching (SM). Binary classification task. Given a pair of column names of (c_j, c_k) with its corresponding description (d_j, d_k) , determine whether they refer to the same attribute, i.e.,

$$f((c_j, d_j), (c_k, d_k)) \rightarrow \{0, 1\}, \text{ where } 1 \text{ indicates a match.}$$

Error Detection (ED). Binary classification task. Given a row

r and a specific cell $v_{i,j}$, identify whether the value in $v_{i,j}$ is erroneous based on the other cell values in the same row, i.e.,

$$f(v_{i,j}, r) \rightarrow \{0, 1\}, \text{ where } 1 \text{ indicates an error.}$$

Data Cleaning (DC). Open-domain generation task. Given a row r and a specific erroneous cell $v_{i,j}$, correct the erroneous based on other cell values in the same row, i.e.,

$$f(v_{i,j}, r) \rightarrow \hat{v}_{i,j}.$$

Column Type Annotation (CTA). Multi-class/Multi-label classification task. Given a column c_j in table T , assign a semantic type \mathcal{C} to the entire column based on the cell values in c_j , i.e.,

$$f(c_j) \rightarrow \mathcal{C}.$$

Attribute Value Extraction (AVE). Open-domain generation task. Given a text s and a target attribute c_j , extract the corresponding value v_j from s , i.e.,

$$f(s, c_j) \rightarrow v_j.$$

To develop a versatile LLM capable of addressing diverse data preparation tasks, training a multi-task **upstream model** is essential. Let N upstream tasks $\{\mathcal{T}_1, \dots, \mathcal{T}_N\}$ be associated with a collection of labeled datasets \mathcal{D} , referred to as **upstream data**. A multi-task **upstream DP-LLM** can be trained using the upstream data. However, the upstream DP-LLM struggles to transfer to novel datasets and tasks with limited labeled data, due to cross-datasets/tasks knowledge distraction and dataset-informed knowledge gap. We aim to address these two issues and thus boost the transferability of DP-LLM with limited labeled data, which can be defined as follows.

Few-Shot DP-LLM Transfer. Building on the upstream DP-LLM, we consider a set of novel **downstream data**, comprising novel datasets $\{\mathcal{D}'_1, \mathcal{D}'_2, \dots\}$ and novel tasks $\{\mathcal{T}'_1, \mathcal{T}'_2, \dots\}$, each with few-shot labeled samples. Our goal is to boost the transferability of the upstream DP-LLM when fine-tuned on these few-shot downstream datasets and tasks.

IV. OVERVIEW

Fig. 2 illustrates the knowledge augmentation KNOWTRANS. KNOWTRANS consists of two key components: the *Selective Knowledge Concentration* (SKC) component and the *Automatic Knowledge Bridging* (AKB) component. They work at training time and inference time, respectively. In the SKC component, the upstream data are reused, enabling the upstream DP-LLM to concentrate on relevant knowledge when fine-tuning the novel downstream data and mitigating the knowledge distraction issue. The AKB component generates proper knowledge as supplementary prompts for each novel dataset during inference to alleviate the dataset-informed knowledge gap.

SKC Component. In this component, we propose a novel method to fine-tune the upstream DP-LLM to concentrate on relevant knowledge for novel downstream data, with the following three stages: ① *Upstream knowledge patches extraction*. It first dynamically extracts knowledge of the upstream datasets

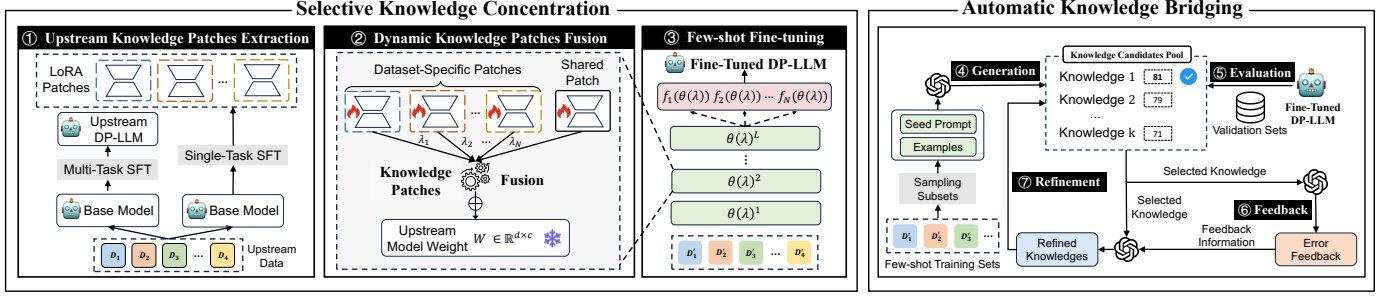


Fig. 2: The illustration of our proposed knowledge augmentation framework KNOWTRANS, which consists of two components: Selective Knowledge Concentration (SKC) component and Automatic Knowledge Bridging (AKB) component.

into independent knowledge patches, using a cross-model low-rank parameterize method. ② *Dynamic knowledge patches Fusion*. It incorporates extracted knowledge patches and an additional new knowledge patch, dynamically weighting them to integrate transferable knowledge from upstream datasets and shared knowledge for the new data. ③ *Few-shot fine-tuning*. It keeps the upstream model fixed and fine-tunes only the knowledge patches and weights on few-shot downstream data. After these stages, the upstream model effectively concentrates relevant knowledge for novel datasets and tasks.

AKB Component. This component automatically generates dataset-informed knowledge in an iterative process, which involves the following four steps: ⑤ *Generation*. A subset of the dataset is sampled and transformed into prompt-response pairs. These pairs are concatenated with a seed prompt containing initial knowledge and fed into a more capable closed-source LLM to generate the initial knowledge candidates pool. ⑥ *Evaluation*. It evaluates these knowledge candidates using the fine-tuned DP-LLM on a validation set, identifies the best knowledge, and collects error cases. ⑦ *Feedback*. It generates error feedback information based on these error cases. ⑧ *Refinement*. It refines the knowledge based on error feedback and then incorporates the refined knowledge into the pool for the next iteration. After completing the iteration, the most appropriate data-informed knowledge should be generated to closely align with the dataset characteristics.

Together, these components form an integrated framework that effectively enhances the transferability of DP-LLMs by optimizing both training and inference, resulting in better performance on novel downstream data.

V. SELECTIVE KNOWLEDGE CONCENTRATION

This section introduces the Selective Knowledge Concentration (SKC) component, which enables the DP-LLM to concentrate knowledge relevant to novel datasets and tasks. As shown in Fig. 2, it involves three key stages: 1) Upstream knowledge patches extraction (Section V-A) stage, 2) Dynamically knowledge patch fusion (Section V-B stage and 3) Few-shot fine-tuning (Section V-C) stage. Next, we provide a detailed introduction for each stage.

A. Upstream Knowledge Patches Extraction

To prevent conflicts within the parameter space, we extract knowledge from upstream datasets and incorporate it into sepa-

rate knowledge patches. These patches are isolated, inherently preventing any overlap in the parameter space. We can then reuse these isolated knowledge patches to transfer dataset-informed knowledge to downstream data. However, extracting knowledge patches from the upstream DP-LLM is non-trivial, as the model has already been fine-tuned on upstream data. Further fine-tuning cannot extract valuable information into knowledge patches.

To address this challenge, we propose reusing the upstream data and fine-tuning it on the base model of the upstream DP-LLM using Low-Rank Adaptation (LoRA). For instance, if we employ the upstream DP-LLM Jellyfish-7B [7], the corresponding base model can be Mistral-7B [15]. We refer to the LoRA module obtained after fine-tuning as a “knowledge patch”. LoRA’s plug-and-play nature makes it an ideal choice for serving as a modular knowledge patch, and its parameter efficiency allows us to maintain only the lightweight module for each dataset. This approach benefits from the shared architecture and pre-trained data between the base and upstream DP-LLMs, enabling an effective transfer of knowledge patches from the base model to the upstream DP-LLM.

Cross-Model Low-Rank Parameterize. Assume we have a base model of the upstream DP-LLM. Given a pre-trained weight matrix $\theta_0 \in \mathbb{R}^{d \times c}$ of the base model and a various of upstream tasks $\{\mathcal{T}_1, \dots, \mathcal{T}_N\}$. For each upstream task \mathcal{T}_i , we initialize a LoRA module ΔW_i on the base model and fine-tune it with task-specific data \mathcal{D}_i using a low-rank parameter update, i.e.,

$$\theta_i = \theta_0 + \alpha \Delta W_i = \theta_0 + \alpha B_i A_i. \quad (2)$$

This allows ΔW_i to extract knowledge patches from each dataset by introducing only a low-rank transformation $\Delta W = B \times A$ for effective knowledge encoding, where $B \in \mathbb{R}^{d \times r}$ and $A \in \mathbb{R}^{r \times c}$ with rank $r \ll \min(d, c)$. B and A are initialized with a random Gaussian distribution and zero respectively. The hyperparameter α represents the scaling factor that controls the strength of updates.

Knowledge Patch Training. To train knowledge patch for each upstream task \mathcal{T}_i , we keep the main transformer layers of \mathcal{M}_b fixed, while updating only the parameters in ΔW_i . The training objective can be formulated as the standard maximum-

likelihood training with cross-entropy for conditional language modeling, i.e.,

$$\min_{\Delta \mathbf{W}_i} \sum_{(x,y) \in \mathcal{D}_i} \sum_{t=1}^{|y|} -\log(P_{\theta_0 + \Delta \mathbf{W}_i}(y_t | x, y_{<t})), \quad (3)$$

where $|y|$ is the length of the output sequence y , and y_t is the target token at position t in the output sequence. Fig. 2 ① depicts the training process of our cross-model knowledge patch training. We use the dataset from the Jellyfish-Instruct [7] for training these knowledge patches, detailed in Section VII-A. Our experiments in Section VII-B show that this cross-model training method substantially enhances the performance of novel datasets and tasks.

This confirms that the analogous base model can effectively serve as a proxy for creating transferable knowledge patches. After training, these knowledge patches are then integrated into the upstream DP-LLM, providing an initialization point for subsequent dynamic fusion.

B. Dynamic Knowledge Patch Fusion

To ensure that the DP-LLM focuses on relevant knowledge for new datasets and tasks, we propose reusing the extracted knowledge patches during fine-tuning. To achieve this, we introduce a dynamic knowledge patch fusion module that seamlessly integrates upstream knowledge patches while incorporating newly added knowledge patches. By fusing these patches through a weighted ensemble approach, the module dynamically combines dataset-specific transferable knowledge with shared knowledge, enabling upstream DP-LLM effective transfer to novel datasets and tasks.

Knowledge Patches Fusion. Formally, given a pre-trained weight matrix $\hat{\theta}_0$ of upstream DP-LLM, we learn $\lambda = [\lambda_1, \dots, \lambda_N]$ that control the contribution of each upstream knowledge patch to the model, i.e.,

$$\hat{\theta}(\lambda) = \hat{\theta}_0 + \alpha \left(\sum_{i=1}^N \lambda_i \Delta \mathbf{W}_i + \Delta \mathbf{W}_{N+1} \right). \quad (4)$$

Intuitively, $\hat{\theta}_0$ is the backbone model, and each $\Delta \mathbf{W}_N$ is the upstream knowledge patch with dataset-specific knowledge while $\Delta \mathbf{W}_{N+1}$ is the newly added knowledge patch that is expected to capture the shared information across multiple downstream datasets. The knowledge patch fusion module attempts to elicit the most relevant knowledge from upstream tasks and transfer it to downstream datasets, effectively preventing the issue of knowledge distraction. The illustration of the combination process can be found in Fig. 2 ②.

C. Few-shot Fine-Tuning

After integrating the knowledge patches and fusion module into the upstream DP-LLM, we fine-tune it for new datasets and tasks. Similar to the earlier training of individual knowledge patches, we keep the parameters of the backbone model fixed and use cross-entropy to maximize the likelihood in conditional language modeling. However, in this stage, we fine-tune the

Algorithm 1 Selective Knowledge Concentration

Input: Pre-trained base model \mathcal{M}_b with parameters θ_0 , upstream tasks $\{\mathcal{T}_1, \dots, \mathcal{T}_N\}$ with respective datasets $\{\mathcal{D}_1, \dots, \mathcal{D}_N\}$, few-shot data \mathcal{D}' of novel datasets and tasks

- 1: *Stage 1: Upstream Knowledge Patches Extraction*
 - 2: **for** each upstream task \mathcal{T}_i **do**
 - 3: Initialize random $\Delta \mathbf{W}_i$ on \mathcal{M}_b
 - 4: Fine-tune $\Delta \mathbf{W}_i$ on \mathcal{D}_i using Eq. 2
 - 5: Extracting low-rank knowledge patch $\Delta \mathbf{W}_i$
 - 6: **end for**
 - 7: *Stage 2: Dynamic Knowledge Patches Fusion*
 - 8: Initialize interpolation weights $\{\lambda_1, \dots, \lambda_N\}$ on the upstream dataset-specific knowledge patches $\{\Delta \mathbf{W}_1, \dots, \Delta \mathbf{W}_N\}$
 - 9: Initialize shared knowledge patch $\Delta \mathbf{W}_{N+1}$
 - 10: Compute fused weight matrix $\hat{\theta}$ using Eq. 4
 - 11: *Stage 3: Few-shot Fine-Tuning*
 - 12: **for** each mini-batch $\xi \in \mathcal{D}'$ **do**
 - 13: Update knowledge patches and $\{\lambda_1, \dots, \lambda_N\}$ using gradient descent by optimizing Eq. 5
 - 14: **end for**
 - 15: **return** Fine-tuned model $\hat{\theta}$
-

upstream DP-LLM across all datasets. Ultimately, the objective function for multi-task learning can be described as:

$$\min_{\lambda, \Delta \mathbf{W}_1, \dots, \Delta \mathbf{W}_{N+1}} \mathbb{E}_{\xi} \left[\sum_{i=1}^N f_i(\theta(\lambda); \xi) \right], \quad (5)$$

where ξ is a mini-batch of multi-task training data $\{(x_i^1, y_i^1), \dots, (x_i^N, y_i^N)\}_{i=1}^q$ and q is the batch size. The illustration of few-shot fine-tuning can be found in Fig. 2 ③. Algorithm 1 summarized the algorithm of SKC.

VI. AUTOMATIC KNOWLEDGE BRIDGING

This section presents the Automatic Knowledge Bridging (AKB) component, which is designed to integrate dataset-informed knowledge and address the knowledge gap in DP-LLM. As shown in Fig. 2 (right), AKB formulates the search for dataset-informed knowledge as an optimization problem, implemented through a four-step workflow: *Generation*, *Evaluation*, *Feedback*, and *Refinement*. We will first give an overall introduction to the optimization problem (Section VI-A) and the workflow (Section VI-B) of AKB, followed by a detailed explanation of each step.

A. Automatic Knowledge Bridging

Manually Deriving knowledge for a target dataset is time-consuming and impractical, especially in scenarios where novel datasets evolve rapidly. Thus, we attempt to frame the knowledge search process as an optimization problem to achieve automated searching. However, due to the vast search space, finding the right knowledge can be extremely difficult, making the knowledge optimization process intractable. Recent advances in prompt engineering have shown that LLMs can write better prompts than humans [27, 28, 29]. Inspired by this, we adopt closed-source LLM (e.g., GPT-4o) with

comprehensive world knowledge to generate dataset-informed knowledge, which can be seen as a segment of the prompt.

Considering we have the fine-tuned upstream DP-LLM \mathcal{M} , a novel data preparation dataset $\mathcal{D}'_i = \{(x, y)\}$ of input-output pairs, which associated with a seed prompt P_i for the dataset. Our goal is to find the most suitable knowledge ρ from a closed-source LLM \mathcal{M}_{gpt} . When \mathcal{M} is prompted with the concatenation $[P_i(\rho)||x]$, where $P_i(\rho)$ denote knowledge insert into the prompt. Providing this as input, we expect \mathcal{M} to generate the corresponding response y . To achieve this goal, we formulate it as a prompt optimization problem. Formally, we search for the optimal knowledge that maximizes the expected per-sample score $\mathcal{S}(\rho, x, y)$ across all possible (x, y) :

$$\rho^* = \arg \max_{\rho} \mathcal{S}(\rho) = \arg \max_{\rho} \mathbb{E}_{(x, y)} [\mathcal{S}(\rho, x, y)], \quad (6)$$

where ρ^* is the founded optimal knowledge. Next, we will instantiate the optimization problem using a four-step workflow.

B. Workflow

In practice, AKB consists of four steps: *Generation*, *Evaluation*, *Feedback* and *Refinement*, as illustrated in Fig 2 (right). Specifically, given a seed prompt and a subset of examples from the few-shot training set, AKB generates an initial pool of knowledge candidates. Then, it evaluates each candidate in the pool on a validation set with the fine-tuned upstream DP-LLM and selects the top performer. Subsequently, the best knowledge is used to identify error cases in the training set. Based on various errors, the closed-source LLM is utilized to generate diverse error feedback. Finally, based on feedback information, the selected knowledge is refined several times and re-added into the knowledge pool for the next iteration. Once the convergence criteria are met, AKB will provide high-quality knowledge as a supplement to the prompt for each inference. Next, we will discuss the details of each step.

Generation. To create an initial pool of potential knowledge candidates \mathcal{K} , we use a seed prompt with a set of examples as input. The seed prompt corresponds to the task prompt of the dataset and contains initial handcrafted knowledge. The examples X_{examples} are a set of input-output pairs sampled from the dataset \mathcal{D}'_i , similar to the role of demonstrations in ICL. In this work, examples implicitly convey data-informed dataset knowledge, which is provided to the closed-source LLM to generate data-informed knowledge as text descriptions, it can be described as follows:

$$\rho = \mathcal{M}_{\text{gpt}}(P_{\text{gen}} || X_{\text{examples}}), \quad (7)$$

where P_{gen} is the knowledge generation prompt to guide the \mathcal{M}_{gpt} in generating knowledge candidates. To illustrate, we take the EM task as an example, the sketch of its generation prompt is shown in Fig. 3 (a), and the complete prompt can be found in the full version [1].

Evaluation. To obtain the rank of knowledge in the candidates, we need a score function that accurately measures the alignment between the task performance and the generated knowledge. Following previous works [27, 28], we directly choose the task

metric as the score function, since our goal is to improve the performance of the target task, making the metric a suitable measure for our purposes.

Specifically, we first integrate each piece of knowledge ρ from the set \mathcal{K} into the task prompt P_i , resulting in $P_i(\rho)$. Next, we evaluate this knowledge-enhanced prompt using a validation dataset $\mathcal{D}_{\text{valid}}$. We use the validation dataset the same as dataset \mathcal{D}'_i . For each pair (x, y) in $\mathcal{D}_{\text{valid}}$, we concatenate x with the prompt $P_i(\rho)$, i.e. $[P_i(\rho)||x]$, and feed it into the model \mathcal{M} to generate an output. Finally, we compare the output to the ground truth y using the task's specific metric. This process can be summarized as follows:

$$\mathcal{S}(\rho, x, y) = \mathcal{E}(\mathcal{M}(P_i(\rho)||x), y), \quad (8)$$

where $\mathcal{E}(\cdot)$ denotes the metric of the target task (e.g. binary-F1 for EM). Ultimately, we rank the knowledge in the pool and select the ρ with the highest score for further refinement.

Feedback. To further improve the quality of knowledge, we attempt to refine the selected knowledge by introducing an error feedback mechanism. Before that, we need to analyze and provide feedback information for knowledge refinement. We use model \mathcal{M}_{gpt} again to analyze and generate feedback information. Specifically, we first test the selected knowledge ρ_t on the training dataset $\mathcal{D}_{\text{train}}$ to obtain the error set E in the t iteration. Then, we sample a subset X_{errors} of E for producing the feedback information. Given the selected knowledge ρ_t and error subset X_{errors} , the feedback prompt P_{fb} (as shown in Fig. 3 (b)) is employed to guide \mathcal{M}_{gpt} in producing feedback. The process is formalized as follows:

$$f_{b_t} = \mathcal{M}_{\text{gpt}}(P_{\text{fb}} || P_i(\rho_t) || X_{\text{errors}}). \quad (9)$$

Here, the feedback f_{b_t} is a text response from \mathcal{M}_{gpt} that includes error feedback information summarized of error examples. Intuitively, error examples highlight deficiencies in the current knowledge ρ_t , helping the model produce useful feedback information. This feedback information will help \mathcal{M}_{gpt} generate more accurate knowledge through refinement.

Refinement. With the error feedback in hand, we refine the selected knowledge. Given the knowledge ρ_t , error set X_{errors} and error feedback information f_{b_t} , the refinement prompt P_{refine} (refer to Fig. 3 (c)) is used to guide \mathcal{M}_{gpt} in generating refined knowledge $\hat{\rho}_t$. This process is formalized as follows:

$$\hat{\rho}_t = \mathcal{M}_{\text{gpt}}(P_{\text{refine}} || P_i(\rho_t) || X_{\text{errors}} || f_{b_t}). \quad (10)$$

The knowledge refinement process alternates between feedback and refinement multiple times to cover diverse error types. Specifically, we sample different subsets from error example set E , and use them to get different knowledge refinement. Additionally, we retain the knowledge optimization trajectory $\rho_{0:t-1}$ involves all generated knowledge from the previous iterations. By considering the full knowledge optimization trajectory, AKB implicitly summarizes the common mistakes from past solutions and avoids repeating them. Thus, Eq. 10 is in fact implemented as:

$$\hat{\rho}_t = \mathcal{M}_{\text{gpt}}(P_{\text{refine}} || X_{\text{errors}} || f_{b_t} || \rho_{0:t-1}). \quad (11)$$

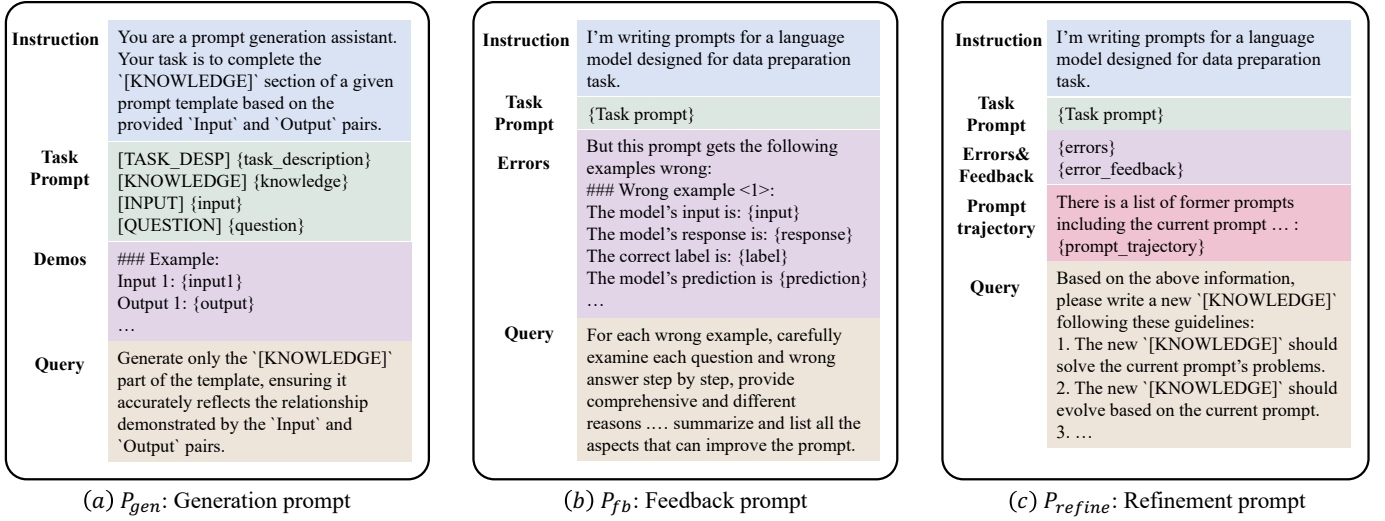


Fig. 3: The illustration of three prompts used in the AKB component.

After processing all error subsets, the refined knowledge is aggregated and incorporated into the candidate pool \mathcal{K} for re-evaluation in the next iteration. Once the specified number of iterations is completed, the highest-scoring knowledge candidate is selected as the final output. This iterative AKB process ensures progressive enhancement of ρ , optimizing task performance and effectively addressing the dataset-informed knowledge gap in the upstream DP-LLM for novel datasets.

Algorithm 2 Automatic Knowledge Bridging

Input: \mathcal{M}_{gpt} : closed-source LLM, \mathcal{M}' : knowledge fused DP-LLM, \mathcal{D}_{valid} : validation dataset

- 1: Randomly sample subsets $X_{demos} \subset \mathcal{D}_{valid}$
- 2: Generate initial candidate knowledge pool \mathcal{K} using Eq. 7
- 3: $t \leftarrow 0$
- 4: **while** not converged **do**
- 5: $\rho_t = \arg \max_{\rho \in \mathcal{K}} f(\rho, \mathcal{D}_{valid})$ // Score (Eq. 8)
- 6: $E = \{(x, y) \in \mathcal{D}_{valid} \mid \mathcal{M}'(\rho_t, x) \neq y\}$ // Error set
- 7: **for** iteration $j \in \{0, 1, \dots\}$ **do**
- 8: Choose a random errors subset $X_{errors}^j \subset E$
- 9: Generate feedback f_{b_j} using Eq. 9
- 10: Get refined knowledge ρ'_t using Eq. 10
- 11: $\mathcal{K} \leftarrow \mathcal{K} \cup \rho'_t$
- 12: **end for**
- 13: $t \leftarrow t + 1$
- 14: **end while**
- 15: **Return:** ρ_t

Algorithm 2 describes the optimization process of AKB in detail. It begins by generating an initial pool of knowledge candidates (lines 1–3). At each iteration, the algorithm selects the most promising knowledge candidate using a scoring function (line 5) and identifies the error set, which consists of instances where the model's predictions diverge from the ground truth (line 6). The knowledge pool is then updated by integrating refined knowledge derived from feedback (lines

7–11). This iterative process continues until convergence, ultimately yielding the optimized data-informed knowledge.

VII. EXPERIMENTS

In this section, we conduct extensive experiments to evaluate our proposed framework KNOWTRANS.

A. Experiments Setup

Datasets. We evaluate KNOWTRANS on publicly available datasets across four upstream data preparation tasks: **ED**: Flights, Rayyan, Beer [30]. **DI**: Flipkart [31], Phone [32]. **SM**: CMS [3]. **EM**: Abt-Buy, Walmart-Amazon [33]; and three novel tasks: **CTA**: SOTAB [34]. **AVE**: AE-110k, OAmine. [35] **DC**: Rayyan, Beer [36]. These datasets encompass a wide range of domains and present various challenges, allowing us to thoroughly assess the capabilities of KNOWTRANS.

Baselines. We compare KNOWTRANS with the following three types of methods.

(1) Open-source LLM-based Methods (DP-LLMs):

- **Mistral** [15]: A versatile model for general NLP tasks, we fine-tuned the Mistral-7B using a data preparation dataset.
- **TableLLaMA** [9]: An open-source, versatile model supports a wide range of table-based tasks.
- **MELD** [10]: A Mixture-of-Expert (MoE)-based models adapted for low-resource data preparation tasks.
- **Jellyfish** [7]: A versatile DP-LLM capable of multiple data preparation tasks. We include three different model sizes: 7B, 8B, and 13B for this model.
- **Jellyfish-ICL** [7]: An enhanced version of Jellyfish-7B with in-context learning.
- **KNOWTRANS**: Our proposed framework, KNOWTRANS, is designed to boost the transferability of DP-LLMs. In this work, we specifically focus on improving the performance of the Jellyfish-7B/8B/13B and Mistral-7B models.

(2) Closed-source LLM-based Methods:

- **GPT-3.5** [37]: Using the OpenAI API for in-context learning, with version *gpt-3.5-turbo-1106*.

TABLE I: Statistic of Datasets.

Task	Dataset	Training Set	Few-shot	Test Set
ED	Flights	12,256	20	2,000
	Rayyan	9,000	20	2,000
	Beer	10,050	20	2,000
DI	Flipkart	11,460	20	2,675
	Phone	2,547	20	1,194
SM	CMS	23,068	20	2,564
EM	Abt-buy	5,743	20	1,916
	Walmart-Amazon	6,144	20	2,049
CTA	SOTAB	356	20	250
AVE	AE-110k	4,405	20	1,495
	OA-mine	7,360	20	2,451
DC	Rayyan	9,000	20	2,000
	Beer	10,050	20	2,000

- **GPT-4** [37]: Using the OpenAI API for in-context learning, with version *gpt-4-0613*.
- **GPT-4o** [37]: Using the OpenAI API for in-context learning, with version *gpt-4o-2024-08-06*.

(3) *Non-LLM Methods*: we also include some non-LLM methods that do not use LLM: (a) ED: **Raha** [30], (b) DI: **IPM** [5], (c) SM: **SMAT** [3], (d) EM: **Ditto** [38], (e) CTA: **Doduo** [4], (f) AVE: **MAVE** [39], (g) EC: **Baran** [36].

Implementation. We implement KNOWTRANS in Pytorch. We fine-tuned the upstream model for the SKC component using the LlamaFactory framework [40]. Specifically, For the SKC component, we apply LoRA with a default rank of 32. The learning rate is set to $6e-5$, with a batch size of 4 and a gradient accumulation step of 4. The model is trained for 3 epochs, with a 9:1 split for the train/validation dataset. For the AKB component, we use the closed-source LLM GPT-4o (*gpt-4o-2024-08-06* API) for generation, feedback and refinement steps, configured with a temperature of 0.9. For evaluation, we use DP-LLM for inference with a temperature setting of 0.35, a top-k value of 10, and a top-p value of 0.9. We provide 10 examples for knowledge generation, with 4 error examples for refinement. The default setting runs 3 iterations, with 5 error random samples in each iteration. All experiments are conducted 3 times and the averaged performances are reported.

Environment. All experiments are conducted on a server equipped with 256GB RAM, Intel(R) Xeon(R) Silver 4314 CPU @ 2.40GHz and NVIDIA A40 GPU. The software environment consists of Ubuntu 20.04 LTS operating system, Python 3.9, PyTorch 2.3.1, Hugging Face Transformers Library 4.46.1 and LLM serving engine vLLM 0.5.3.

Evaluation Metrics. To evaluate these data preparation tasks, we use accuracy for DI. For EM, ED, SM, DC and AVE, we use the F1-score, while for CTA, we use the micro-F1 score. All metrics are reported on a 100-point scale.

B. Effectiveness Evaluation

1) *Performance Against Open-Source LLMs-Based and non-LLM Methods*: To evaluate the effectiveness of KNOWTRANS,

we compared it with 7B open-source DP-LLMs and non-LLM methods. As shown in Table II, KNOWTRANS achieved state-of-the-art performance on 13 novel datasets in few-shot settings, with an average performance of **79.26%**, surpassing the best DP-LLM, Jellyfish, by **4.93%**. Additionally, KNOWTRANS significantly outperforms non-LLM methods by an average of **33.7%**. Non-LLM methods often overfit in few-shot scenarios because they typically depend on feature learning or smaller pre-trained models. Although DP-LLMs use larger models to mitigate overfitting, they still face issues like cross-datasets/tasks knowledge distraction and dataset-informed knowledge gap. In contrast, our proposed KNOWTRANS concentrates on transferable knowledge selectively and can effectively bridge the knowledge gap by integrating SKC and AKB components.

For **novel datasets** in ED, DI, SM, and EM, KNOWTRANS excels by leveraging its AKB component to incorporate dataset-specific knowledge. For example, it achieves **85.37%** and **85.88%** on ED and EM tasks, improving by **9.98%** and **3.96%**, respectively. While MELD can acquire transferable knowledge by assigning weights to the top-k experts within a MoE module, it employs an instance-level expert combination approach that fails to utilize dataset-level knowledge. For **novel tasks** such as CTA, AVE, and DC, KNOWTRANS addresses knowledge distraction with its SKC component, which reuses relevant knowledge from upstream tasks. On the CTA task, KNOWTRANS achieves **83.61%**, outperforming Mistral (**80.08%**). For AVE, it achieves **63.90%**, surpassing Mistral (**62.65%**). In DC, it delivers competitive results, such as **96.27%** on Rayyan and **98.54%** on Beer datasets.

2) *Performance Against Closed-Source LLMs-Based Methods*: We also compare KNOWTRANS’s performance against open-source LLM-based methods, particularly the GPT-series models such as GPT-3.5, GPT-4, and GPT-4o. The results are summarized in Table IV. We present three model sizes—7B, 8B, and 13B—boosted with the corresponding version of Jellyfish. All these enhanced models demonstrate superior performance compared to GPT-3.5, GPT-4, and GPT-4o across multiple datasets and tasks. For example, on average, KNOWTRANS-13B outperforms GPT-4 and GPT-4o by **7.03%** and **6.07%**, respectively, and shows a substantial improvement of **13.54%** over GPT-3.5. This improvement underscores the advantage of our knowledge augmentation framework in a few-shot scenario.

GPT-series models, while powerful in general applications, face challenges when directly applied to domain-specific tasks without further adaptation. Additionally, API calls to GPT can be quite costly, and we will provide a detailed comparison in the next section. Our proposed approach mitigates these challenges by leveraging SKC and AKB, enabling the model to specialize effectively in data preparation tasks. The results confirm that the KNOWTRANS provides robust solutions for adapting data preparation LLMs to novel datasets and tasks, bridging the gap between general-purpose LLMs and specialized tasks. This makes KNOWTRANS an efficient choice for scenarios where task-specific tuning is essential but access to closed-source LLMs may be limited or costly.

TABLE II: Comparison of state-of-the-art 7B open-source DP-LLMs and non-LLM methods on thirteen datasets. Each dataset is evaluated in a few-shot setting. The best result is highlighted in **bold**, while the second-best is underlined.

Task	Dataset	Model						
		Non-LLM-Methods	Mistral	TableLLaMA	MELD	Jellyfish	Jellyfish-ICL	KNOWTRANS
ED	Flights	44.00	45.67	53.02	66.48	<u>68.65</u>	64.67	74.38
	Rayyan	62.00	45.00	36.99	<u>79.79</u>	78.89	74.17	89.40
	Beer	70.00	12.99	38.06	77.84	<u>78.62</u>	45.27	92.33
	Average	58.67	34.55	42.69	74.70	<u>75.39</u>	61.37	85.37
DI	Flipkart	2.54	81.27	42.59	79.74	78.09	<u>82.47</u>	82.88
	Phone	8.20	84.09	70.35	<u>85.09</u>	83.17	83.92	85.68
	Average	5.37	82.68	56.62	82.42	80.63	<u>83.20</u>	84.28
SM	CMS	2.10	18.75	1.86	26.67	27.59	30.30	<u>27.69</u>
EM	Abt-buy	57.14	20.09	42.58	<u>85.52</u>	77.62	74.56	87.86
	Walmart-Amazon	80.00	39.83	34.70	78.31	<u>82.74</u>	79.08	83.89
	Average	68.57	29.96	38.64	<u>81.92</u>	80.18	76.82	85.88
CTA	SOTAB	25.13	<u>80.08</u>	20.31	58.78	79.22	42.75	83.61
AVE	AE-110k	3.91	<u>65.08</u>	18.93	60.54	59.27	59.51	67.86
	OA-mine	1.63	60.22	17.01	57.16	57.57	42.76	<u>59.93</u>
	Average	2.77	<u>62.65</u>	17.97	58.85	58.42	51.14	63.90
DC	Rayyan	63.00	96.82	84.23	91.57	<u>96.37</u>	92.69	96.27
	Beer	87.00	95.83	<u>99.68</u>	99.72	98.54	95.10	98.54
	Average	75.00	96.33	91.96	91.57	97.46	93.90	97.41
Average		45.56	57.36	43.11	70.62	<u>74.33</u>	66.71	79.26

TABLE III: Comparison of input tokens, output tokens, and cost per instance for KNOWTRANS and other GPT models.

	Input Tokens	Output Tokens	Price
GPT-3.5	751.08	2.86	\$0.0004
GPT-4o	751.08	2.86	\$0.0038
GPT-4	751.08	2.86	\$0.0227
KNOWTRANS	20.41	8.21	\$0.0002

C. Efficiency and Cost Analysis of KNOWTRANS

We further explore the cost-effectiveness of KNOWTRANS-7B compared to closed-source LLMs like GPT-3.5, GPT-4o, and GPT-4, as demonstrated in Tab. III. Our analysis focuses on the number of input tokens required, the number of output tokens generated, and the cost implications for each method.

- **Input and Output Tokens:** KNOWTRANS requires significantly fewer input tokens, averaging **20.41** tokens, compared to the 751.08 tokens required by GPT-3.5, GPT-4o, and GPT-4. The reduction in token usage is primarily attributed to our approach of incorporating few-shot examples into model parameters through training, whereas GPT requires few-shot samples to be placed in the context of demonstrations. This reduction in token usage not only minimizes computational load but also enhances processing speed. Additionally, due to the additional knowledge and error feedback in the AKB component, the output token count for KNOWTRANS is **8.21**, slightly higher than the other models (**2.86** tokens), suggesting more detailed and comprehensive output generation. However, since the AKB process only requires searching once for each dataset, the increase can be essentially ignored.

- **Cost Efficiency:** The cost per instance for KNOWTRANS is \$0.0002, markedly lower than that of GPT-3.5 (**\$0.0004**), and significantly less than the costs associated with GPT-4o (**\$0.0038**) and GPT-4 (**\$0.0227**). This pricing advantage makes KNOWTRANS an economical choice for large-scale implementations where multiple queries are processed.

The substantial reduction in input token requirements coupled with cost advantages positions KNOWTRANS as a highly efficient and scalable solution for data preparation tasks. These attributes make KNOWTRANS particularly suitable for deployment in environments where computational resources are limited or where cost containment is essential. Thus, KNOWTRANS not only provides a technologically advanced solution but also offers practical benefits in real-world applications.

D. Scalability Analysis

To further assess the effectiveness of the proposed framework KNOWTRANS across varying levels of data availability, we conducted a detailed scalability analysis using four datasets: DC (Rayyan), SM (CMS), EM (Walmart-Amazon), and AVE (AE-110k). This analysis examines the performance as the number of instances increases, simulating different resource conditions ranging from low to high data availability.

The results, presented in Fig. 4, demonstrate that the KNOWTRANS consistently outperforms the baseline model, Jellyfish-7B, across all datasets. As the number of instances grows, the performance of both models converges, reducing the gap in F1 scores. However, the KNOWTRANS maintains a consistent advantage in the low-to-medium range of instances, emphasizing its robustness in scenarios where data resources

TABLE IV: Comparison of state-of-the-art closed-source LLMs on thirteen datasets. Each dataset is evaluated in a few-shot setting. The best result is highlighted in **bold**, while the second-best is underlined.

Task	Dataset	Model					
		GPT-3.5	GPT-4	GPT-4o	KNOWTRANS-7B	KNOWTRANS-8B	KNOWTRANS-13B
ED	Flights	61.39	68.65	66.22	74.38	66.91	84.53
	Rayyan	51.87	63.80	52.54	89.40	78.13	84.40
	Beer	75.91	83.89	76.90	<u>92.33</u>	98.93	83.93
	Average	63.06	72.11	65.22	85.37	81.32	84.29
DI	Flipkart	82.47	85.64	<u>85.57</u>	82.88	81.87	81.20
	Phone	85.43	85.51	86.18	<u>85.68</u>	83.33	84.34
	Average	83.95	<u>85.58</u>	85.88	84.28	82.60	82.77
SM	CMS	30.77	30.30	28.57	27.69	<u>35.9</u>	40.91
EM	Abt-buy	37.21	<u>92.49</u>	94.09	87.86	87.98	92.42
	Walmart-Amazon	80.80	84.07	<u>87.37</u>	83.89	83.55	91.15
	Average	59.01	88.28	<u>90.73</u>	85.88	85.77	91.79
CTA	SOTAB	90.76	<u>96.40</u>	98.00	83.61	81.60	81.48
AVE	AE-110k	50.99	49.72	61.10	<u>67.86</u>	66.67	72.47
	OA-mine	55.12	46.47	<u>60.47</u>	59.93	55.08	63.57
	Average	53.06	48.10	60.79	<u>63.90</u>	60.88	68.02
DC	Rayyan	91.87	90.82	86.50	<u>96.27</u>	95.65	96.82
	Beer	99.75	99.17	94.35	98.54	97.69	<u>99.52</u>
	Average	95.81	95.00	90.43	<u>97.41</u>	96.67	98.17
Average		67.85	74.76	75.32	<u>79.40</u>	77.87	81.39

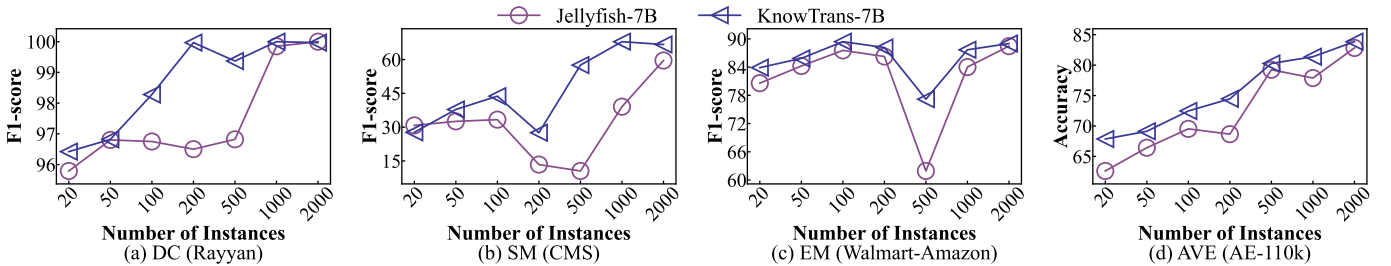


Fig. 4: Comparison of the performance of Jellyfish-7B and KNOWTRANS-7B across four tasks (datasets). The results illustrate F1-score and accuracy for varying numbers of labeled instances (20, 50, 100, 200, 1000, and 2000).

are constrained. This scalability is attributed to the framework’s SKC and AKB components, which enable efficient utilization of limited labeled data and mitigate challenges such as knowledge distraction and gap. These results underline KNOWTRANS’s scalability and effectiveness, making it a practical solution for diverse data preparation tasks, particularly in real-world applications where data availability often varies widely.

E. Ablation Study

To evaluate the effectiveness of each component in KNOWTRANS, we conducted an ablation study by removing the SKC and AKB components individually. Table V shows the performance of KNOWTRANS-7B on multiple tasks and datasets under three configurations: remove the SKC component, remove the AKB component, and remove both of them (i.e. original upstream Jellyfish-7B model).

The results demonstrate that both the SKC and AKB components contribute significantly to the overall performance of KNOWTRANS. Specifically, the absence of the SKC component leads to a drop in average performance from **83.94%** to **79.88%**,

TABLE V: Performance for Ablation Study.

Task	Dataset	Model			
		w/o SKC & AKB	w/o SKC	w/o AKB	KNOWTRANS
DI	Flipkart	78.36	81.79	81.42	82.88
	Phone	83.50	84.76	<u>85.34</u>	85.68
CTA	SOTAB	79.47	82.08	<u>82.60</u>	83.61
AVE	AE-110k	56.48	61.58	<u>65.79</u>	67.86
	OA-mine	53.20	57.89	<u>58.63</u>	59.93
DC	Rayyan	90.39	96.82	95.99	<u>96.27</u>
	Beer	93.56	97.98	<u>98.15</u>	98.54
Average		76.64	79.88	<u>80.74</u>	83.94

indicating that the selective knowledge concentration is crucial for mitigating knowledge distraction. Similarly, removing the AKB component reduces performance from **83.94%** to **80.74%**, highlighting the importance of generation dataset-informed for addressing the knowledge gap in novel tasks. These findings confirm that integrating SKC and AKB components significantly improves the upstream DP-LLM’s adaptability to novel data preparation datasets and tasks.

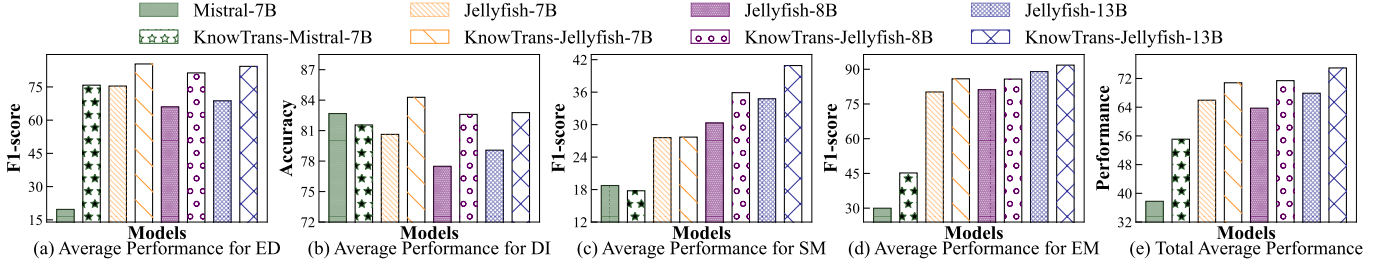


Fig. 5: Comparison of different backbone models using our proposed KNOWTRANS on novel datasets.

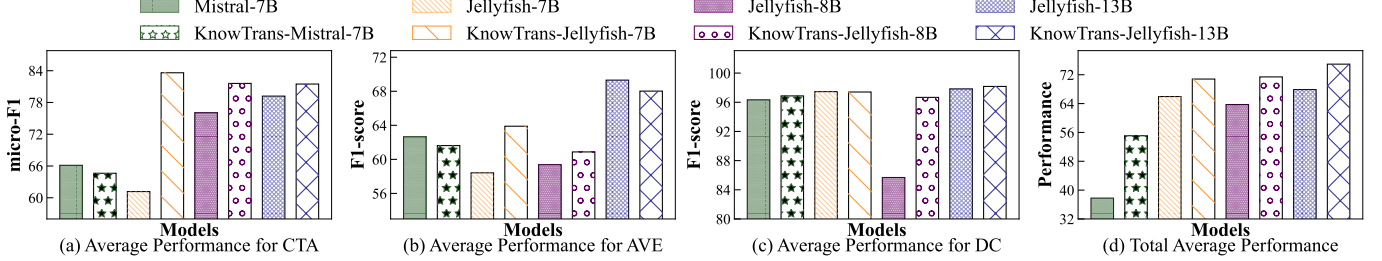


Fig. 6: Comparison of different backbone models using our proposed KNOWTRANS on novel tasks.

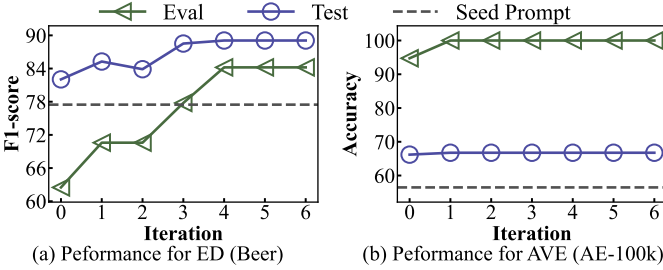


Fig. 7: Performance of different round times for refinement.

F. Hyper-parameter Analysis

Effect of Different Models. The results in Fig. 5 and Fig. 6 demonstrate the effectiveness of KNOWTRANS in boosting various backbone models, including Mistral-7B, Jellyfish-7B, Jellyfish-8B, and Jellyfish-13B, across novel datasets and tasks. KNOWTRANS consistently delivers superior performance compared to the baselines. Specifically, KNOWTRANS exhibits consistent and significant performance improvements for all model sizes, though the extent of improvement varies by model capacity. Larger models, such as Jellyfish-13B, benefit the most from the KNOWTRANS framework within the same architecture, achieving the highest scores across multiple datasets and tasks. This is attributed to their larger parameter capacity, which allows them to better utilize the augmented knowledge introduced by KNOWTRANS’s SKC and AKB components.

Smaller models, also experience substantial performance gains, indicating the versatility of KNOWTRANS. Particularly, Mistral-7B exhibits the most significant performance gain after applying the KNOWTRANS. This is particularly notable because Mistral-7B had not undergone upstream multi-task training on data preparation tasks, unlike the Jellyfish models. These findings highlight the general applicability and efficiency of KNOWTRANS in boosting the transferability of DP-

TABLE VI: Performance of using different weight strategies.

Task	Dataset	Model			
		Single	Uniform	Adaptive	KNOWTRANS
ED	Flights	65.99	69.08	<u>71.69</u>	74.38
	Rayyan	74.33	81.79	<u>87.30</u>	89.40
EM	Abt-buy	77.36	80.43	<u>81.17</u>	87.86
AVE	AE-110k	56.48	63.09	65.79	67.86
Average		69.00	73.60	<u>76.49</u>	79.90

LLMs, regardless of model size. The remarkable improvement can be attributed to KNOWTRANS’s ability to dynamically incorporate knowledge patches and refine dataset-specific knowledge through the SKC and AKB components, effectively compensating for the absence of prior upstream training.

Effect of weight strategy. To evaluate the impact of different adaptive strategies on model performance, we present the performance improvement comparison across different adaptive strategies: single (i.e. no upstream knowledge patches), uniform and adaptive (i.e. used in SKC), over Jellyfish-7B. As shown in Table VI, the overall average performance of uniform and adaptive strategies surpasses that of the single model, demonstrating that dynamically integrating knowledge from upstream tasks can effectively enhance the transferability of DP-LLMs. Furthermore, the adaptive strategy outperforms uniform, indicating that our proposed adaptive knowledge integration module, SKC, plays a crucial role in optimizing the use of dataset-specific and shared knowledge. Notably, KNOWTRANS, which incorporates SKC, achieves the best performance across all datasets, further validating the effectiveness of this approach.

Effect of round times for refinement. We then explore the effect of the iterative round times for error feedback on the performance of KNOWTRANS. We demonstrate the experiment results on two tasks: Attribute Value Extraction (AVE) and Error Detection (ED), as shown in Fig. 7.

- **Error Detection (ED):** In contrast, Fig. 7(a) shows a notable improvement in both eval and test performance for the ED task. Initially, eval performance increases sharply, indicating rapid assimilation of feedback in early iterations. After the third round, the eval performance begins to plateau, while the test performance continues to improve slightly, suggesting enhanced model transferability from iterative refinements.
- **Attribute Value Extraction (AVE):** As depicted in Fig. 7(b), the eval and test performance for AVE remain constant over 6 rounds of feedback, which suggests that the model quickly stabilizes and additional feedback does not contribute further to significant changes in model behavior for this metric. This indicates that for certain data, additional knowledge may not be helpful. Nonetheless, AKB does not impair its performance, and we can also benefit from the SKC.

This iterative feedback effectively refines the knowledge by eliciting insights from error patterns, showcasing its potential to enhance performance in diverse data preparation tasks.

VIII. RELATED WORK

In this section, we introduce related works including Non-LLM methods, LLM-based methods and model fusion.

A. Non-LLM Methods

Early works in data preparation mainly focus on rule-based method [41, 42, 43] and statistical model [44, 45, 46]. For example, employing manual rules [47], consistency checks [48], and similarity constraints [49], probabilistic models [44, 45, 46] in ED, DC and DI. Similarly, For AVE, early works also relied on rule-based methods [50, 51] and statistical model [52] to identify attribute-value pairs. For match tasks such as EM and SM, early works rely on similarity functions [53, 54] or hashing [55, 56]. For CTA, it focuses on capturing tabular semantics through feature engineering techniques, including hand-crafted [57, 58], statistical [59] and semantic features [60].

Recently, some works have shifted towards utilizing machine learning (ML) techniques. For instance, to achieve more accurate error detection, correction and imputation, few-shot learning [61], transfer learning [36], meta-learning [62], GANs [63] and pre-train language models (PLMs) [62] are adopted in ED, DC and DI tasks. Additionally, for EM and SM, some methods propose to utilize word embeddings [64], DNNs [65], gaussian mixture model [66], graph neural networks [67] and PLMs [38, 68, 69, 70, 71], aim to obtain a better tabular representation. Similarly, PLM-based methods also shows outperformed performance in CTA [4, 72, 73, 74, 75, 76, 77, 78, 79] and AVE [39, 80, 81, 82, 83] tasks. However, most of the methods mentioned above are restricted to a single data preparation task. In contrast, we propose an all-in-one multitask model that can be applied to a variety of datasets and tasks.

B. LLM-based Methods

More recent, some works employing LLMs have achieved notable progress in data preparation tasks, including closed- and open-source LLM-based methods. For closed-source LLMs-based methods, prompt engineering [33, 35, 84, 85, 86, 87, 88,

89] technique is frequently utilized on ChatGPT and GPT-4o through OpenAI API. However, due to domain gaps, relying solely on APIs makes it difficult to adapt these LLMs to tabular data, which is common for data preparation tasks. Thus, TableGPT [8] further enhances ChatGPT via instruction tuning, improving its performance on table-related tasks. Nevertheless, closed-source models still pose high costs, instability, and data privacy concerns. Thus, some works attempt to fine-tune the open-source LLMs. For instance, AnyMatch [90] fine-tunes GPT-2 for zero-shot EM, Archetype [91] improves Alpaca-7B [14] for CTA tasks. Furthermore, DP-LLMs are proposed to manage multiple data preparation tasks within a single LLM, such as Jellyfish [7], TableLLaMA [9] and MELD [10]. In contrast to these methods of training multiple tasks from scratch, we aim to boost the existing DP-LLMs' transferability to novel datasets and tasks in a few-shot scenario.

C. Model Fusion

Model fusion combines individually trained models to enhance multi-task capabilities [92, 93], often surpassing the original models in task generalization [24, 94]. Early approaches like SWA [92] and TWA [95] averaged model weights along the training trajectory, while methods such as Model Soup [93], Model-Ratatouille [25], DiWA [24], DART [96], and WEMoE [97] fused independently trained or fine-tuned models to leverage their diversity for improved generalization. However, these approaches are computationally intensive, particularly for large models. To address this, lightweight parameter-efficient fusion methods like LoRA [98] and Adapter [99] have been proposed, enabling functional transferability without merging entire models. LoRAHub [94] optimizes LoRA module fusion using black-box coefficient search, while AdapterSoup [99] averages adapters during test time for domain adaptation. LoRARetriever [100] dynamically retrieves and fuses relevant LoRAs, and MoLE [101] adapts a Mixture-of-Experts framework to learn optimal composition weights for LoRA modules. Unlike these works, our method focuses on fusing LoRAs from upstream models to boost the transferability of DP-LLM in a cross-model manner.

IX. CONCLUSION

This paper proposed KNOWTRANS, a knowledge augmentation framework that enhances the transferability of existing DP-LLMs to novel datasets and tasks in a few-shot setting. By integrating two key components—Selective Knowledge Concentration (SKC) and Automatic Knowledge Bridging (AKB)—KNOWTRANS effectively mitigates cross-datasets/tasks knowledge distraction and dataset-informed knowledge gap, enabling robust performance in few-shot settings and scalability with increasing data availability. Experimental results demonstrate the KNOWTRANS can effectively boost the DP-LLM's transferability over various DP-LLMs, highlighting its practicality for diverse and evolving data preparation tasks. Furthermore, this work provides a robust and scalable approach, setting a strong foundation for future innovations in knowledge-augmented LLMs.

REFERENCES

- [1] “Code, datasets and full version,” <https://github.com/allen0/knowtrans>, 2025.
- [2] C. Chai, J. Wang, Y. Luo, Z. Niu, and G. Li, “Data management for machine learning: A survey,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 35, no. 5, pp. 4646–4667, 2022.
- [3] J. Zhang, B. Shin, J. D. Choi, and J. C. Ho, “Smat: An attention-based deep learning solution to the automation of schema matching,” in *ADBIS*, 2021.
- [4] Y. Suhara, J. Li, Y. Li, D. Zhang, Ç. Demiralp, C. Chen, and W.-C. Tan, “Annotating columns with pre-trained language models,” in *SIGMOD*, 2022.
- [5] Y. Mei, S. Song, C. Fang, H. Yang, J. Fang, and J. Long, “Capturing semantics for imputation with pre-trained language models,” in *ICDE*, 2021.
- [6] J. Tu, J. Fan, N. Tang, P. Wang, G. Li, X. Du, X. Jia, and S. Gao, “Unicorn: A unified multi-tasking model for supporting matching tasks in data integration,” *SIGMOD*, 2023.
- [7] H. Zhang, Y. Dong, C. Xiao, and M. Oyamada, “Jellyfish: Instruction-tuning local large language models for data preprocessing,” in *EMNLP*, 2024.
- [8] P. Li, Y. He, D. Yashar, W. Cui, S. Ge, H. Zhang, D. Rifinski Fainman, D. Zhang, and S. Chaudhuri, “Table-gpt: Table fine-tuned gpt for diverse table tasks,” *Proceedings of the ACM on Management of Data*, vol. 2, no. 3, pp. 1–28, 2024.
- [9] T. Zhang, X. Yue, Y. Li, and H. Sun, “Tablellama: Towards open large generalist models for tables,” in *NAACL*, 2024.
- [10] M. Yan, Y. Wang, K. Pang, M. Xie, and J. Li, “Efficient mixture of experts based on large language models for low-resource data preprocessing,” in *KDD*, 2024.
- [11] Y. Mao, Z. Wang, W. Liu, X. Lin, and W. Hu, “Banditmtl: Bandit-based multi-task learning for text classification,” in *ACL*, 2021.
- [12] B. Liu, X. Liu, X. Jin, P. Stone, and Q. Liu, “Conflict-averse gradient descent for multi-task learning,” *NeurIPS*, 2021.
- [13] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” in *NeurIPS*, 2020.
- [14] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale *et al.*, “Llama 2: Open foundation and fine-tuned chat models,” *arXiv:2307.09288*, 2023.
- [15] A. Q. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D. S. Chaplot, D. d. l. Casas, F. Bressand, G. Lengyel, G. Lample, L. Saulnier *et al.*, “Mistral 7b,” *arXiv:2310.06825*, 2023.
- [16] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong *et al.*, “A survey of large language models,” *arXiv:2303.18223*, 2023.
- [17] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler *et al.*, “Emergent abilities of large language models,” *Transactions on Machine Learning Research*, 2022.
- [18] B. Min, H. Ross, E. Sulem, A. P. B. Veyseh, T. H. Nguyen, O. Sainz, E. Agirre, I. Heintz, and D. Roth, “Recent advances in natural language processing via large pre-trained language models: A survey,” *ACM Computing Surveys*, vol. 56, no. 2, pp. 1–40, 2023.
- [19] S. Arora, B. Yang, S. Eyuboglu, A. Narayan, A. Hojel, I. Trummer, and C. Ré, “Language models enable simple systems for generating structured views of heterogeneous data lakes,” *arXiv:2304.09433*, 2023.
- [20] R. C. Fernandez, A. J. Elmore, M. J. Franklin, S. Krishnan, and C. Tan, “How large language models will disrupt data management,” *VLDB*, 2023.
- [21] S. Zhang, L. Dong, X. Li, S. Zhang, X. Sun, S. Wang, J. Li, R. Hu, T. Zhang, F. Wu *et al.*, “Instruction tuning for large language models: A survey,” *arXiv:2308.10792*, 2023.
- [22] O. Sener and V. Koltun, “Multi-task learning as multi-objective optimization,” *NeurIPS*, 2018.
- [23] Y. Mao, W. Liu, and X. Lin, “Adaptive adversarial multi-task representation learning,” in *ICML*, 2020.
- [24] A. Rame, M. Kirchmeyer, T. Rahier, A. Rakotomamonjy, P. Gallinari, and M. Cord, “Diverse weight averaging for out-of-distribution generalization,” *NeurIPS*, 2022.
- [25] A. Ramé, K. Ahuja, J. Zhang, M. Cord, L. Bottou, and D. Lopez-Paz, “Model ratatouille: Recycling diverse models for out-of-distribution generalization,” in *ICML*, 2023.
- [26] E. Yang, L. Shen, G. Guo, X. Wang, X. Cao, J. Zhang, and D. Tao, “Model merging in llms, mllms, and beyond: Methods, theories, applications and opportunities,” *arXiv:2408.07666*, 2024.
- [27] Y. Zhou, A. I. Muresanu, Z. Han, K. Paster, S. Pitis, H. Chan, and J. Ba, “Large language models are human-level prompt engineers,” *arXiv:2211.01910*, 2022.
- [28] L. Chen, J. Chen, T. Goldstein, H. Huang, and T. Zhou, “Instructzero: Efficient instruction optimization for black-box large language models,” in *ICML*, 2024.
- [29] J. Cheng, X. Liu, K. Zheng, P. Ke, H. Wang, Y. Dong, J. Tang, and M. Huang, “Black-box prompt optimization: Aligning large language models without model training,” *arXiv:2311.04155*, 2023.
- [30] M. Mahdavi, Z. Abedjan, R. Castro Fernandez, S. Madden, M. Ouzzani, M. Stonebraker, and N. Tang, “Raha: A configuration-free error detection system,” in *SIGMOD*, 2019.
- [31] PromptCloud, “Flipkart products dataset,” <https://www.kaggle.com/datasets/PromptCloudHQ/flipkart-products>, 2017.
- [32] —, “Amazon reviews: Unlocked mobile phones,” <https://www.kaggle.com/datasets/PromptCloudHQ/amazon-reviews-unlocked-mobile-phones>, 2016.
- [33] A. Narayan, I. Chami, L. Orr, S. Arora, and

- C. Ré, “Can foundation models wrangle your data?” *arXiv:2205.09911*, 2022.
- [34] K. Korini and C. Bizer, “Column type annotation using chatgpt,” 2023.
- [35] A. Brinkmann, R. Shraga, and C. Bizer, “Product attribute value extraction using large language models,” *arXiv:2310.12537*, 2023.
- [36] M. Mahdavi and Z. Abedjan, “Baran: Effective error correction via a unified context representation and transfer learning,” *VLDB*, 2020.
- [37] H. Zhang, Y. Dong, C. Xiao, and M. Oyamada, “Large language models as data preprocessors,” in *VLDB*, 2024.
- [38] Y. Li, J. Li, Y. Suhara, A. Doan, and W.-C. Tan, “Deep entity matching with pre-trained language models,” *arXiv:2004.00584*, 2020.
- [39] L. Yang, Q. Wang, Z. Yu, A. Kulkarni, S. Sanghai, B. Shu, J. Elsas, and B. Kanagal, “Mave: A product dataset for multi-source attribute value extraction,” in *WSDM*, 2022.
- [40] Y. Zheng, R. Zhang, J. Zhang, Y. Ye, Z. Luo, Z. Feng, and Y. Ma, “Llamafactory: Unified efficient fine-tuning of 100+ language models,” 2024.
- [41] F. Geerts, G. Mecca, P. Papotti, and D. Santoro, “The llunatic data-cleaning framework,” *VLDB*, 2013.
- [42] X. Ding, H. Wang, J. Su, M. Wang, J. Li, and H. Gao, “Leveraging currency for repairing inconsistent and incomplete data,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 34, no. 3, pp. 1288–1302, 2020.
- [43] W. Fan, P. Lu, and C. Tian, “Unifying logic rules and machine learning for entity enhancing,” *Science China Information Sciences*, vol. 63, pp. 1–19, 2020.
- [44] T. Rekatsinas, X. Chu, I. F. Ilyas, and C. Ré, “Holoclean: Holistic data repairs with probabilistic inference,” *VLDB*, 2017.
- [45] S. De, Y. Hu, M. V. Vamsikrishna, Y. Chen, and S. Kambhampati, “Bayeswipe: A scalable probabilistic framework for cleaning bigdata,” *arXiv:1506.08908*, 2015.
- [46] A. Lew, M. Agrawal, D. Sontag, and V. Mansinghka, “Pclean: Bayesian data cleaning at scale with domain-specific probabilistic programming,” in *AISTATS*, 2021.
- [47] W. Fan, J. Li, S. Ma, N. Tang, and W. Yu, “Towards certain fixes with editing rules and master data,” *The VLDB journal*, vol. 21, pp. 213–238, 2012.
- [48] G. Cong, W. Fan, F. Geerts, X. Jia, and S. Ma, “Improving data quality: Consistency and accuracy,” in *VLDB*, 2007.
- [49] S. Song, Y. Sun, A. Zhang, L. Chen, and J. Wang, “Enriching data imputation under similarity rule constraints,” *IEEE transactions on knowledge and data engineering*, vol. 32, no. 2, pp. 275–287, 2018.
- [50] D. Vandic, J.-W. Van Dam, and F. Frasincar, “Faceted product search powered by the semantic web,” *Decision Support Systems*, vol. 53, no. 3, pp. 425–437, 2012.
- [51] L. Zhang, M. Zhu, and W. Huang, “A framework for an ontology-based e-commerce product information retrieval system,” *J. Comput.*, vol. 4, no. 6, pp. 436–443, 2009.
- [52] D. Putthividhya and J. Hu, “Bootstrapped named entity recognition for product attribute extraction,” in *EMNLP*, 2011.
- [53] M. Hadjieleftheriou, A. Chandel, N. Koudas, and D. Srivastava, “Fast indexes and algorithms for set similarity selection queries,” in *ICDE*, 2008.
- [54] A. Doan, P. Konda, P. Suganthan GC, Y. Govind, D. Paulsen, K. Chandrasekhar, P. Martinkus, and M. Christie, “Magellan: toward building ecosystems of entity matching solutions,” *Communications of the ACM*, vol. 63, no. 8, pp. 83–91, 2020.
- [55] G. Papadakis, G. Koutrika, T. Palpanas, and W. Nejdl, “Meta-blocking: Taking entity resolution to the next level,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 26, no. 8, pp. 1946–1960, 2013.
- [56] P. Konda, S. Das, P. S. G. C., A. Doan, A. Ardalani, J. R. Ballard, H. Li, F. Panahi, H. Zhang, J. F. Naughton, S. Prasad, G. Krishnan, R. Deep, and V. Raghavendra, “Magellan: Toward building entity matching management systems,” *VLDB*, 2016.
- [57] S. K. Ramnandan, A. Mittal, C. A. Knoblock, and P. Szekely, “Assigning semantic labels to data sources,” in *ESWC*, 2015.
- [58] M. Pham, S. Alse, C. A. Knoblock, and P. Szekely, “Semantic labeling: a domain-independent approach,” in *ISWC*, 2016.
- [59] D. Zhang, Y. Suhara, J. Li, M. Hulsebos, Ç. Demiralp, and W.-C. Tan, “Sato: Contextual semantic type detection in tables,” *arXiv:1911.06311*, 2019.
- [60] M. Hulsebos, K. Hu, M. Bakker, E. Zraggen, A. Satyanarayan, T. Kraska, Ç. Demiralp, and C. Hidalgo, “Sherlock: A deep learning approach to semantic data type detection,” in *KDD*, 2019.
- [61] A. Heidari, J. McGrath, I. F. Ilyas, and T. Rekatsinas, “Holodetect: Few-shot learning for error detection,” in *SIGMOD*, 2019.
- [62] Z. Miao, Y. Li, and X. Wang, “Rotom: A meta-learned data augmentation framework for entity matching, data cleaning, text classification, and beyond,” in *SIGMOD*, 2021.
- [63] J. Peng, D. Shen, N. Tang, T. Liu, Y. Kou, T. Nie, H. Cui, and G. Yu, “Self-supervised and interpretable data cleaning with sequence generative adversarial networks,” *VLDB*, 2022.
- [64] J. MESTS and M. Tang, “Distributed representations of tuples for entity resolution,” *VLDB*, 2018.
- [65] S. Mudgal, H. Li, T. Rekatsinas, A. Doan, Y. Park, G. Krishnan, R. Deep, E. Arcaute, and V. Raghavendra, “Deep learning for entity matching: A design space exploration,” in *SIGMOD*, 2018.
- [66] R. Wu, S. Chaba, S. Sawlani, X. Chu, and S. Thirumuranathan, “Zeroer: Entity resolution using zero labeled examples,” in *SIGMOD*, 2020.
- [67] B. Li, W. Wang, Y. Sun, L. Zhang, M. A. Ali, and

- Y. Wang, “Grapher: Token-centric entity resolution with graph convolutional neural networks,” in *AAAI*, 2020.
- [68] C. Zhao and Y. He, “Auto-em: End-to-end fuzzy entity-matching using pre-trained deep models and transfer learning,” in *WWW*, 2019.
- [69] R. Peeters and C. Bizer, “Dual-objective fine-tuning of bert for entity matching,” *VLDB*, 2021.
- [70] J. Tu, J. Fan, N. Tang, P. Wang, C. Chai, G. Li, R. Fan, and X. Du, “Domain adaptation for deep entity resolution,” in *SIGMOD*, 2022.
- [71] P. Wang, X. Zeng, L. Chen, F. Ye, Y. Mao, J. Zhu, and Y. Gao, “Promptem: prompt-tuning for low-resource generalized entity matching,” *arXiv:2207.04802*, 2022.
- [72] X. Deng, H. Sun, A. Lees, Y. Wu, and C. Yu, “Turl: Table understanding through representation learning,” *ACM SIGMOD Record*, vol. 51, no. 1, pp. 33–40, 2022.
- [73] J. Herzig, P. K. Nowak, T. Müller, F. Piccinno, and J. M. Eisenschlos, “Tapas: Weakly supervised table parsing via pre-training,” *arXiv:2004.02349*, 2020.
- [74] P. Yin, G. Neubig, W.-t. Yih, and S. Riedel, “Tabert: Pretraining for joint understanding of textual and tabular data,” *arXiv:2005.08314*, 2020.
- [75] H. Iida, D. Thai, V. Manjunatha, and M. Iyyer, “Tabbie: Pretrained representations of tabular data,” *arXiv:2105.02584*, 2021.
- [76] Y. Zhou, S. Singh, and C. Christodoulopoulos, “Tabular data concept type detection using star-transformers,” in *CIKM*, 2021.
- [77] Y. Sun, H. Xin, and L. Chen, “Reca: Related tables enhanced column semantic type annotation framework,” *VLDB*, 2023.
- [78] Z. Miao and J. Wang, “Watchog: A light-weight contrastive learning based framework for column annotation,” *SIGMOD*, 2023.
- [79] Y. Wang, H. Xin, and L. Chen, “Kglink: A column type annotation method that combines knowledge graph and pre-trained language model,” *arXiv:2406.00318*, 2024.
- [80] H. Xu, W. Wang, X. Mao, X. Jiang, and M. Lan, “Scaling up open tagging from tens to thousands: Comprehension empowered attribute value extraction from product title,” in *ACL*, 2019.
- [81] Q. Wang, L. Yang, B. Kanagal, S. Sanghai, D. Sivakumar, B. Shu, Z. Yu, and J. Elsas, “Learning to extract attribute value from product via question answering: A multi-task approach,” in *KDD*, 2020.
- [82] X. Zhang, C. Zhang, X. Li, X. L. Dong, J. Shang, C. Faloutsos, and J. Han, “Oa-mine: Open-world attribute mining for e-commerce products with weak supervision,” in *WWW*, 2022.
- [83] J. Yan, N. Zalmout, Y. Liang, C. Grant, X. Ren, and X. L. Dong, “Aadatag: Multi-attribute value extraction from product profiles with adaptive decoding,” *arXiv:2106.02318*, 2021.
- [84] R. Peeters and C. Bizer, “Using chatgpt for entity matching,” in *ADBS*, 2023.
- [85] Z. Cheng, J. Kasai, and T. Yu, “Batch prompting: Efficient inference with large language model apis,” *arXiv:2301.08721*, 2023.
- [86] N. Guan, K. Chen, and N. Koudas, “Can large language models design accurate label functions?” *arXiv:2311.00739*, 2023.
- [87] M. Fan, X. Han, J. Fan, C. Chai, N. Tang, G. Li, and X. Du, “Cost-effective in-context learning for entity resolution: A design space exploration,” in *ICDE*, 2024.
- [88] C. Fang, X. Li, Z. Fan, J. Xu, K. Nag, E. Korpeoglu, S. Kumar, and K. Achan, “Llm-ensemble: Optimal large language model ensemble method for e-commerce product attribute value extraction,” in *SIGIR*, 2024.
- [89] N. Seedat and M. van der Schaar, “Matchmaker: Self-improving large language model programs for schema matching,” *arXiv:2410.24105*, 2024.
- [90] Z. Zhang, P. Groth, I. Calixto, and S. Schelter, “Any-match - efficient zero-shot entity matching with a small language model,” *arXiv:2409.04073*, 2024.
- [91] B. Feuer, Y. Liu, C. Hegde, and J. Freire, “Archetype: A novel framework for open-source column type annotation using large language models,” *VLDB*, 2024.
- [92] P. Izmailov, D. Podoprikin, T. Garipov, D. Vetrov, and A. G. Wilson, “Averaging weights leads to wider optima and better generalization,” *arXiv:1803.05407*, 2018.
- [93] M. Wortsman, G. Ilharco, S. Y. Gadre, R. Roelofs, R. G. Lopes, A. S. Morcos, H. Namkoong, A. Farhadi, Y. Carmon, S. Kornblith, and L. Schmidt, “Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time,” in *ICML*, 2022.
- [94] C. Huang, Q. Liu, B. Y. Lin, T. Pang, C. Du, and M. Lin, “Lorahub: Efficient cross-task generalization via dynamic lora composition,” *arXiv:2307.13269*, 2023.
- [95] T. Li, Z. Huang, Q. Tao, Y. Wu, and X. Huang, “Trainable weight averaging: Efficient training by optimizing historical solutions,” in *ICLR*, 2022.
- [96] S. Jain, S. Addepalli, P. K. Sahu, P. Dey, and R. V. Babu, “Dart: Diversify-aggregate-repeat training improves generalization of neural networks,” in *CVPR*, 2023.
- [97] A. Tang, L. Shen, Y. Luo, N. Yin, L. Zhang, and D. Tao, “Merging multi-task models via weight-ensembling mixture of experts,” *arXiv:2402.00433*, 2024.
- [98] E. J. Hu, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, W. Chen *et al.*, “Lora: Low-rank adaptation of large language models,” in *ICLR*, 2022.
- [99] A. Chronopoulou, M. E. Peters, A. Fraser, and J. Dodge, “Adaptersoup: Weight averaging to improve generalization of pretrained language models,” *arXiv:2302.07027*, 2023.
- [100] Z. Zhao, L. Gan, G. Wang, W. Zhou, H. Yang, K. Kuang, and F. Wu, “Loraretriever: Input-aware lora retrieval and composition for mixed tasks in the wild,” *arXiv:2402.09997*, 2024.
- [101] X. Wu, S. Huang, and F. Wei, “Mixture of lora experts,” in *ICLR*, 2024.

X. APPENDIX

A. Upstream Datasets

Tab. VII shows the upstream datasets that are used to train upstream DP-LLMs. It includes four upstream tasks: ED, DI, SM, and EM, comprising a total of 36,304 training samples. We follow the experimental setting of Jellyfish benchmark [7].

TABLE VII: Statistic of Upstream Datasets.

Task	Dataset	#Samples	#Positives
ED	Adult	1100	70
	Hospital	3420	88
DI	Buy	586	N/A
	Restaurant	778	N/A
SM	MIMIC	7000	11
	Synthea	5000	18
EM	Amazon-Google	6874	699
	Beer	359	54
	DBLP-ACM	5000	885
	DBLP-GoogleScholar	5000	924
	Fodors-Zagats	757	88
	iTunes-Amazon	430	105

B. Example of Task Prompt Template

In this work, each sample in the datasets is converted into a prompt-response pair using predefined instruction templates. We utilize the task prompt templates provided by Jellyfish Benchmark. To illustrate, here is an example of the ED task:

Listing 1 Task Prompt for ED (Flights).

(System Message) You are an AI assistant that follows instruction extremely well. User will give you a question. Your task is to answer as faithfully as you can.

(Task Description) Your task is to determine if there is an error in the value of a specific attribute within the whole record provided. The attributes may include datasource, flight, scheduled departure time, actual departure time, scheduled arrival time, and actual arrival time.

(Seed Knowledge) Errors may include, but are not limited to, spelling errors, missing values, inconsistencies, or values that don't make sense given the context of the whole record.

(Input) Record [{attribute}: {value}...]

Attribute for Verification: [{attribute}: {value}]

(Question) Is there an error in the value of the {attribute} attribute?

(Output Format) Choose your answer from: [Yes, No]

C. Prompt Templates Used in Automatic Knowledge Bridging Component

The automatic knowledge bridging component utilizes the three prompt templates: generation prompt, feedback prompt, and refinement prompt. As shown in follows:

Listing 2 Generation Prompt.

You are a prompt generation assistant. Your task is to complete the '[KNOWLEDGE]' section of a given prompt template based on the provided 'Input' and 'Output' pairs.

Task Template:

[TASK_DESP] {task_description}

[KNOWLEDGE] {knowledge}

[INPUT] {input}

[QUESTION] {question}

Example:

Input 1: {input1}

Output 1: {output1}

...

Generate only the '[KNOWLEDGE]' part of the template, ensuring it accurately reflects the relationship demonstrated by the 'Input' and 'Output' pairs.

Listing 3 Feedback Prompt.

I'm writing prompts for a language model designed for a task. My current prompt is: {prompt} But this prompt gets the following examples wrong:

Wrong example <{index}>:

The model's input is:

input

The model's response is:

response

The correct label is: label

The model's prediction is prediction

For each wrong example, carefully examine each question and wrong answer step by step, provide comprehensive and different reasons why the prompt leads to the wrong answer. At last, based on all these reasons, summarize and list all the aspects that can improve the prompt.

Listing 4 Refinement Prompt.

I'm writing prompts for a language model designed for data preparation task. My current prompt is:

{prompt}

But this prompt gets the following examples wrong:

{error}

Based on these errors, the problems with this prompt and the reasons are:

{error_feedback}

There is a list of former prompts including the current prompt, and each prompt is modified from its former prompts:

{prompt_trajectory}

Based on the above information, please write a new [KNOWLEDGE] following these guidelines:

1. The new [KNOWLEDGE] should solve the current prompt's problems.

2. The new [KNOWLEDGE] should evolve based on the current prompt.

3. Each new [KNOWLEDGE] should be wrapped with [KNOWLEDGE] and [\KNOWLEDGE].

The new prompt is:

D. Searched Knowledge of Various Datasets

The search knowledge by our proposed automatic knowledge bridging component for each dataset is as follows:

Flights (ED)	<p>Consider the following attributes when making your decision:</p> <ol style="list-style-type: none">1. Format Consistency: Ensure that date and time attributes follow a standard format, such as "hh:mm a.m./p.m. MMM d" (e.g., "7:10 a.m. Dec 1"). Inconsistencies or deviations from this format are errors.2. Logical Consistency: Check for logical inconsistencies in time sequences. For example, if the actual arrival time is earlier than the scheduled arrival time, this might be unusual but not necessarily incorrect.3. Missing Values: Values such as 'nan' and 'N/A' are indicators of missing information and should always be considered errors.4. Contextual Errors: Assess if the value makes sense within the context of the whole record. Anomalies or unexpected patterns, like inconsistent time formats or implausible sequences, should be flagged as errors.5. Error Types: Expand your assessment to include format errors, logical inconsistencies, spelling errors, and missing values.6. Guiding Questions: Ask yourself: Does the format match the expected pattern? Does the value logically fit within the sequence of events described in the record?
Rayyan (ED)	<p>Consider the following attributes when making your decision:</p> <ol style="list-style-type: none">1. Typographical and Naming Errors: Carefully check for spelling errors and inconsistencies in attribute names and values. For example, <code>jounral_abbreviations</code> should be <code>journal_abbreviation</code>. Correct such errors to maintain consistency.2. Date Format and Consistency: Ensure date attributes like <code>article_jcreated_at</code> follow the "YYYY-MM-DD" format. Highlight and correct any dates that do not adhere to this format, such as "4/3/15", which should be "2015-04-03".3. Acceptable Numerical Values: Recognize that "0" can be a valid value for attributes such as <code>article_jissue</code> and <code>article_jvolumn</code> in specific contexts (e.g., when the article is not part of a traditional issue). Do not flag these as errors unless context suggests otherwise.4. Standard Format Verification: Validate attributes like <code>journal_issn</code> for proper format compliance (e.g., "1234-5678") and actual validity beyond format correctness.5. Contextual Coherence: Evaluate the logical consistency of attribute values in relation to the entire record. Understand that certain placeholder values like "0" or "nan" may indicate incomplete data but are not necessarily errors if contextually appropriate.6. Holistic Record Assessment: Analyze the record as a whole for logical consistency. Avoid making decisions based solely on individual attributes without considering the context provided by the full record.7. Balanced Sensitivity to Errors: Maintain balanced sensitivity to error detection. Be cautious not to overinterpret placeholder values as errors, and ensure that format-related inconsistencies are accurately identified and addressed.8. Illustrative Examples: Use examples to illustrate correct and incorrect data entries, demonstrating both acceptable deviations and errors to guide better judgment.9. Ambiguity Clarification: Address any ambiguity in attribute values. Provide necessary clarifications to eliminate misinterpretations, ensuring clarity in data representation.10. Character and Encoding Checks: Be alert to unusual characters or encoding issues, especially in text-based attributes, and correct them to prevent data misinterpretation.

Beer (ED)	<p>Consider the following attributes when making your decision: 1. Emphasize the No-Percent Rule for ABV: eliminate misinterpretations, ensuring clarity in data representation.</p> <ul style="list-style-type: none"> - The ABV attribute must be a decimal value between 0 and 1, without a "%" symbol. Highlight this rule using bold and color text to ensure it's memorable. - Examples: - Correct: ABV = "0.05" - Incorrect: ABV = "0.05%" <p>2. Visual Cues for Correct and Incorrect Formats:</p> <ul style="list-style-type: none"> - Use bold, underlining, or color coding to clearly differentiate correct formats from incorrect ones. - Examples for IBU: - Correct: IBU = "20" - Incorrect: IBU = "nan" <p>3. Clarify Handling of Non-Numeric Fields:</p> <ul style="list-style-type: none"> - Minor spelling variations in fields like city names are acceptable unless they cause confusion. Use examples to show what constitutes a significant error. - Acceptable Variations: - "NYC" vs. "New York City" - Significant Errors: - "San Fransico" for "San Francisco" <p>4. Consistency in Rule Application:</p> <ul style="list-style-type: none"> - Reinforce rules consistently across all records. Provide repeated examples to illustrate consistent application. <p>5. Strengthen Feedback Mechanisms:</p> <ul style="list-style-type: none"> - Implement feedback loops to provide corrective feedback, aiding the model in learning from errors and improving accuracy over time. <p>6. Differentiate Error Types:</p> <ul style="list-style-type: none"> - Clearly distinguish between format errors (e.g., symbols in numeric fields) and logical/contextual errors (e.g., inconsistent names). <p>7. Expand on Edge Cases:</p> <ul style="list-style-type: none"> - Offer examples covering edge cases and contextual variations, explaining the rationale for classifying certain scenarios as errors. <p>8. Enhance Contextual Understanding:</p> <ul style="list-style-type: none"> - Add explanations and examples to help the model understand when variations are errors or not, especially for non-numeric fields.
Flipkart (DI)	<p>To infer the manufacturer of the product from the given record, make your judgment based on the following information:</p> <ol style="list-style-type: none"> 1. Product Name Analysis: The brand is often mentioned at the beginning or within the product name. <ul style="list-style-type: none"> - Example: "Trinketbag Tasli Green Alloy Necklace" ->Brand: Trinketbag 2. Description Examination: The brand name can sometimes be repeated within the product description. <ul style="list-style-type: none"> - Example: "Buy Allure Auto CM 1254 Car Mat Tata Safari for Rs.1599 online. Allure Auto CM 1254 Car Mat Tata Safari at best prices..." ->Brand: Allure Auto 3. Product Naming Pattern: Look for distinct naming patterns where the brand name stands out among generic product descriptors. <ul style="list-style-type: none"> - Example: "Naisha Casuals" ->Brand: Naisha 4. Repetitive Clues: The brand might be highlighted multiple times within the record for emphasis. <ul style="list-style-type: none"> - Example: "Gift Studios Laxmi Ganesh Ji Stone Showpiece - 17.6 cm" ->Brand: Gift Studios 5. Unique Brand Names: If the product name includes a unique, non-descriptive term, it is likely to be the brand. <ul style="list-style-type: none"> - Example: "FRENEMY MPAD389 Mousepad" ->Brand: FRENEMY
Phone (DI)	<p>To determine the brand from the product name, look for the first recognizable and distinct brand name within the product name. This is typically a well-known manufacturer or company associated with cellphones or electronics. The brand name often appears at the beginning of the product name and can be recognized through familiar company names or trademarks.</p>

CMS (SM)	<p>Consider the following attributes when making your decision:</p> <ol style="list-style-type: none"> 1. Focus on the semantic meaning and context provided in the descriptions, not just the attribute names. 2. Check if the attributes refer to the same type of information or data, such as dates, codes, or locations. 3. Consider if the attributes are used in similar contexts, even if they are in different tables or have different names. 4. If attributes describe related but not identical concepts, they are not semantically equivalent. 5. Attributes describing the same information but with slight differences in context can be considered equivalent. 6. Pay attention to the specific terminology used, such as "code," "date," "location," etc., and how they relate to each other. 7. Consider whether the attributes could logically be merged without losing meaningful information. 8. Different coding systems (e.g., ICD9, ethnicity codes) usually imply non-equivalence unless explicitly stated otherwise. 9. Attributes referring to different stages or aspects of an event (e.g., start vs. end) are not semantically equivalent. 10. Ensure that the attributes serve a similar purpose within their respective datasets to determine equivalence.
Abt-Buy (EM)	<p>Consider the following attributes when making your decision:</p> <ul style="list-style-type: none"> - Prioritize the recognition of model numbers as key identifiers. If the model numbers match, this is a strong indication that the products may be the same, even if other attributes like color differ. - Handle cases with incomplete data carefully: if one product lacks a detailed description but shares key identifiers like model numbers or significant attributes with the other, consider them potentially the same. - Be aware of synonymous terms and different naming conventions that may refer to the same product. Check for alternative names or terms that are commonly used interchangeably. - Recognize that variations such as color, minor design elements, or finish (e.g., black vs. red) may not signify different products if they share the same model number or core specifications. - Focus on the overall functionality and core features of the products. If the descriptions imply they serve the same purpose, they may be the same product despite minor descriptive differences. - When dealing with different model numbers, verify if they belong to a series or product line that shares key characteristics, which could still imply the same fundamental product with variations. - Emphasize implicit matching: understand that an incomplete description does not necessarily indicate a different product if the available information aligns significantly with the other product's details.
Walmart- Amazon (EM)	<p>Consider the following attributes when making your decision:</p> <ol style="list-style-type: none"> 1. Prioritize matching key identifiers like brand names, model numbers, and capacities as these are critical for determining if two products are the same. 2. If a description is missing or labeled as "nan," focus on comparing the available product names and identifiers rather than assuming a lack of similarity. 3. Pay special attention to unique identifiers such as model numbers and capacities that can confirm product sameness, even when additional descriptive details are missing. 4. When names indicate the same core product (e.g., titles or model identifiers), treat these as strong indicators of sameness, regardless of extra descriptive differences. 5. Ensure that decisions are based on the strongest identifiers available, with model numbers and capacity being primary indicators, while taking into account the broader context provided by the available details.

SOTAB (CTA)	<p>Consider the following attributes when making your decision:</p> <ol style="list-style-type: none"> 1. Country Identification: When the column consists of repeated codes or abbreviations (e.g., "BE BE BE"), it likely represents a country. 2. Schema.org Identifiers: If the column contains repeated URLs from the Schema.org vocabulary, they typically relate to specific types, such as "event status" or "event attendance mode," depending on the context of the URLs. 3. Event Descriptions: Columns with detailed, narrative-style text usually describe events. 4. Locality or Address Components: Repeated city or locality names (e.g., "Monza and Brianza") in a column indicate a "locality of address." 5. Coordinates: Numerical values that resemble geographic coordinates are classified as "coordinate." 6. Price Ranges: Symbols like "\$\$\$" signify a "price range." <p>When classifying a new column, look for these patterns to determine the most appropriate category.</p>
AE (AVE)	<p>To accurately extract the attribute value from the product title, follow these guidelines:</p> <ol style="list-style-type: none"> 1. Case Sensitivity: Retain the original case of keywords or capitalize them appropriately when extracting attributes. Ensure that extracted attributes match the format of the correct labels. 2. Single Attribute Extraction: Focus on extracting only one attribute value, even if multiple relevant features are present. Choose the most prominent or relevant attribute based on the attribute type requested. 3. Attribute Prioritization: If multiple attributes could potentially be extracted, prioritize the one that appears first in the title or that aligns best with typical expectations for the requested attribute type. 4. No Normalization: Do not normalize or change the format (such as case) of the extracted keywords unless explicitly instructed to do so. 5. Example Reference: Use examples of correct output formatting to guide the extraction process. For instance, if extracting "Gender," ensure the output matches the expected format like "Men" or "Women." 6. Attribute Type Focus: Pay attention to the attribute type requested and extract information directly related to it, such as "Running" for "Sport Type" or "Breathable" for "Feature." 7. Fallback to 'n/a': If the specific attribute is not clearly indicated in the product title, default to 'n/a'.
OA (AVE)	<p>Consider the following attributes when making your decision:</p> <ol style="list-style-type: none"> 1. Prioritize Descriptive Terms: Focus on identifying key descriptive terms such as flavors, scents, or colors in the product title. These terms should be given precedence as they offer specific information about the product's characteristics. 2. Refined Hierarchical Attribute Evaluation: Establish a hierarchy where descriptive terms are prioritized first, followed by unique characteristics like "Decaf," and finally brand names. Ensure this hierarchy is clearly communicated and adhered to. 3. Contextual Term Evaluation: Encourage the model to assess the context within the product title, distinguishing between primary descriptors and secondary or less relevant terms. This understanding is crucial for accurate attribute extraction. 4. Guidance on Brand Names: Clarify that brand names should only be used as attributes when they indicate a unique or distinguishing feature that aligns closely with the product's specific characteristics. 5. Examples and Common Mistakes: Integrate illustrative examples within the KNOWLEDGE, showing both correct and incorrect attribute identification. This will guide understanding and help the model avoid common errors. 6. Robust 'N/A' Criteria: Define clear conditions for when to default to 'N/A', ensuring all potential attributes are thoroughly considered before making this choice. 7. Validation of Extracted Attributes: Implement guidelines for validating the output, focusing on both descriptive accuracy and proper formatting, to ensure that the extracted attribute meets the expected criteria. 8. Handling Multiple Attributes: Emphasize the importance of selecting the most relevant attribute when faced with multiple options. Use the established hierarchy and context cues to make an informed decision, resulting in concise and informative output. 9. Iterative Improvement: Consider establishing a feedback loop where the language model can learn from incorrect outputs, facilitating ongoing refinement of both the KNOWLEDGE and the model's performance.

**Rayyan
(DC)**

Consider the following attributes when making your decision:

1. Inference for Missing Numerical Data:

- When numerical fields such as "article_jissue" are missing (e.g., "nan"), consider using contextual clues from the record like "article_jvolumn," "article_pagination," or "article_jcreated_at" to infer the most likely value. If no context is available, and the inference is not feasible, default to "-1" or other specified conventions.

2. Handling Placeholder Values:

- Distinguish between legitimate "nan" values and those requiring inference. If "nan" is used as a placeholder but an inference can be made from other data points, do so. Use "-1" only when a correction can't be logically inferred.

3. Utilizing Contextual Clues:

- Leverage all available information in the record to make educated guesses. This includes cross-referencing related attributes, such as "journal_title" and "journal_abbreviation," to deduce values.

4. Examples and Patterns:

- Pay attention to similar records or known patterns to make educated inferences about missing or incorrect values. Use common journal conventions or past records to guide decisions when specific data is missing.

5. Comprehensive Record Analysis:

- Always analyze the entire record for interrelated information that can provide insight into missing or incorrect attributes. Consider how attributes may inform each other and reflect upon common data entry errors or patterns.

6. Prioritization of Correction:

- When a specific attribute is given for correction, prioritize finding a correct value for that attribute over correcting other potential errors unless those errors directly impact the interpretation of the target attribute.

**Beer
(DC)**

Consider the following attributes when making your decision:

1. Spelling Verification Emphasis: Pay close attention to spelling errors in categorical attributes like beer_name, brewery_name, style, city, and state. Always verify the correctness of spelling using a reference list or contextual clues from the entire record.

2. Numerical Formatting Clarity: Explicitly state the rules for handling numerical precision. For attributes like ounces, simplify the value by removing non-essential decimal places unless they are critical for the context. Conversely, ensure that ABV values fall within a realistic range, and if necessary, correct by adding or adjusting decimal places.

3. Contextual Consistency: Leverage the full context of the record to make informed decisions. For instance, ensure that style corresponds logically with beer_name and other associated attributes. Use this context to validate if corrections are needed.

4. Learning from Examples: Incorporate learning from past examples by recognizing patterns in corrections. Highlight typical error types encountered previously and the correct approach to resolving them, such as common misspellings or numerical anomalies.

5. Error Type Awareness: Clearly list potential error types that the model should be aware of, such as spelling mistakes, missing values, unrealistic numerical values, and inconsistencies among related attributes. Provide guidance on how to address each type.

6. Instruction Precision and Clarity: Provide precise and clear instructions for processing each attribute type. Use sub-points or examples to illustrate how to handle specific errors, ensuring there is no ambiguity in the corrections needed.

7. Typical Value Guidelines: Include guidelines for typical or expected values for specific attributes. For instance, ABV should generally be within a realistic range (e.g., 0.02 to 0.12), and any deviations should be carefully considered.

TABLE VIII: The knowledge searched by automatic knowledge bridging component for each downstream dataset.