Can an LLM find its way around a Spreadsheet?

Abstract—Spreadsheets are routinely used in business and scientific contexts, and one of the most vexing challenges is performing data cleaning prior to analysis and evaluation. The ad-hoc and arbitrary nature of data cleaning problems, such as typos, inconsistent formatting, missing values, and a lack of standardization, often creates the need for highly specialized pipelines. We ask whether an LLM can find its way around a spreadsheet and how to support end-users in taking their free-form data processing requests to fruition. Just like RAG retrieves context to answer users' queries, we demonstrate how we can retrieve elements from a code library to compose data preprocessing pipelines. Through comprehensive experiments, we demonstrate the quality of our system and how it is able to continuously augment its vocabulary by saving new codes and pipelines back to the code library for future retrieval.

Index Terms—LLMs, code generation, data cleaning, end-user programming

I. INTRODUCTION

Pre-trained large language models (LLMs) have demonstrated significant proficiency in generating code from natural language prompts, heralding a new era in software development [1, 2, 38, 39]. In addition to generating code, modern integrated development environments (IDEs) also incorporate LLMs to assist with error correction, code refactoring, and multilingual programming knowledge. By engaging in English conversations, LLMs can function as coding assistants, enabling users with limited proficiency to produce accurate and executable code to accomplish their tasks [3, 37].

Nevertheless, LLMs have not yet attained the maturity to address all or most challenges programmers face in software development and machine learning. One of the vexing aspects for programmers involves processing tabular data, e.g., in the form of spreadsheets. Such datasets (Fig. 1 for example) frequently suffer from issues like typographical errors, inconsistent formatting, missing values, and lack of standardization, necessitating the creation of highly specialized pipeline [4, 5].

Data preprocessing is not merely a matter of resolving inconsistencies; it can significantly impact downstream results. For example, correcting typographical errors in a dataset may inadvertently lead to over-clustering or under-clustering of the relevant column, while inaccurately resolving entities can distort the true distribution of data points. Therefore, data preprocessing must be approached with care and deliberation, with careful attention to sequential transformations. The need for automation support is widely acknowledged [35, 44, 45], particularly for a "human-in-the-loop" pipeline to guide the cleaning process and prevent spurious results.

LLMs have demonstrated impressive programming abilities; however, their capabilities in processing tabular data remain underexplored. Therefore, in this paper, we investigate whether



Fig. 1. Missing values, inconsistent formatting, misspellings, and other similar issues are frequently encountered in large spreadsheets. Our approach, TradeSweep, proposes the use of an LLM to systematize data cleaning and transformations.

an LLM can effectively find its way around a spreadsheet and how it can support end-users in fulfilling their free-form data preprocessing requests. Our approach, TradeSweep (named for its focus on tabular trade datasets), interprets English requests for data preprocessing and generates code proposals that can be composed and applied to targeted datasets, achieving high performance. Similar to how Retrieval-Augmented Generation (RAG) retrieves context to answer users' queries, we demonstrate how elements from a code library can be stored and retrieved to compose complex pipelines.

The contributions of this paper are as follows:

• TradeSweep utilizes English conversations to understand and respond to users' requests with Python code encompassing preprocessing functions, supporting three main capabilities:

- 1) It produces new code when requested for a completely new task.
- 2) It can precisely modify its proposed code based on users' feedback in English conversations.
- The proposed code is automatically tailored to the target data, including accurate column names and suitable algorithms.
- TradeSweep develops and continuously expands a library of fundamental data preprocessing codes by storing executable functions that have been successfully deployed previously. The augmentation of the code library supports the composition and creation of elaborate data pipelines.
- To incorporate feedback from users who lack programming expertise, TradeSweep offers both code proposals and execution results on example data. This allows users to determine whether to accept TradeSweep's proposal or to make modifications based on output data visualizations, rather than focusing on code specifics.
- We perform extensive experiments on three trade datasets. Results show that TradeSweep is capable of generating executable and efficient code for data preprocessing, significantly reducing time and effort for data analysts.

II. RELATED WORK

A. Generating Formulas and SQL queries

Formulas and SQL queries are the lingua franca of spreadsheets. "Formula Language Model for Excel" (FLAME) [6] is a T5-based model trained exclusively on Excel formulas and used for code generation tasks such as formula repair, formula completion, and similarity-based formula retrieval. (TradeSweep is designed to address various formats of tabular data, including Excel and CSV.)

A significant amount of research focuses on converting natural language questions into executable SQL queries (often referred to as Text-to-SQL) [25, 26, 27, 28]. Most existing benchmarks primarily focus on small databases, which do not accurately reflect the challenges of working with large databases in real-world situations. BIRD (Big Bench for Large-scale Databases) [7] is a Text-to-SQL benchmark designed to narrow the gap between experimental and practical scenarios.

B. LLMs with Information Retrieval

While LLMs are highly effective when fine-tuned for particular NLP tasks, their ability to access and manipulate knowledge is intrinsically constrained [29, 30, 31, 41]. RAG [8] has achieved great success by integrating retrieval and generation techniques, leading to the development of various related approaches and variations. This line of work has facilitated the convergence between information retrieval and information generation.

Language models have shown substantial improvements in code completion tasks by learning from "internal" source code contexts [32, 33, 34]. Several retrieval-augmented frameworks have been proposed, e.g., Retrieval-Augmented Code Completion framework (ReACC) [19], RedCoder [20], and RepoCoder [21], These systems however are intended to support developers rather than end-users.

C. LLMs with Data Interaction

Tian et al. presented SpreadsheetLLM [12], a framework designed to enable LLMs to comprehend and analyze spreadsheet data, aiming to create an effective encoding approach that leverages the powerful comprehension and reasoning abilities of LLMs. Although SpreadsheetLLM focuses on spreadsheet comprehension rather than preprocessing, it addresses the challenges associated with handling large and complex datasets with LLMs. The compression strategies established in this study provide valuable insights for enhancing TradeSweep's data preprocessing capabilities, particularly for handling massive spreadsheets.

D. LLMs and Code Generation

Li et al. developed SKCoder [13], a sketch-based code generation strategy designed to mimic the code reuse behaviors of software engineers. To avoid hallucination [36, 40], works such as ALGO [15] and CleanAgent [18] implement guardrails to structure code generation outputs. In general, while LLMs have strong capabilities in generating code based on user prompts, they still encounter difficulties in handling algorithmic complexities and often require human verification [42, 43] to ensure correct and/or efficient code. For example, a minor typo in column names can cause the model to generate code that applies to the wrong column, and failing to specify an algorithm in the user's request may result in functions with low efficiency. For instance, this could involve generating code that uses a brute-force approach to sorting numbers instead of a more efficient algorithm like bubble sort. Table I presents a comparison between TradeSweep and the systems surveyed above.

III. APPROACH

We present TradeSweep, an LLM-based tool that leverages the LLM's code-writing abilities to produce executable programs for data preprocessing tasks. Users are only required to provide a dataset in either CSV or Excel format and submit a preprocessing request. As illustrated in Fig. 2, TradeSweep comprises three primary components:

- **Prompt augmentation**: This component employs information retrieval techniques to select the top-k relevant codes from the code library. (Section III-B)
- Code generation: The LLM generates a code proposal and visualized samples of execution results, awaiting user feedback. (Section III-C)
- Code library: We build a reference document that includes classic Python scripts for data preprocessing tasks. (Section III-D)

A. Problem Definition

Let \mathcal{D} represent a table with n rows and m columns, with column names $c_1, c_2, ..., c_m$ denoted as \mathcal{C} . Each element x_{ij} , where $i \in \{1, ..., n\}$ and $j \in \{1, ..., m\}$, represents the

	SkCoder	Jigsaw	ALGO	CleanAgent	Devin	TradeSweep
English conversations	 ✓ 		\checkmark	√	 ✓ 	✓
Simple user request				\checkmark	\checkmark	✓
Feedback enabled		√			\checkmark	✓
Novel code generation	√	√	√	\checkmark	\checkmark	✓
Dynamic code library						✓
Visualized execution results		\checkmark				\checkmark

 TABLE I

 Feature Comparison: Existing Tools vs. TradeSweep



Fig. 2. An overview of TradeSweep with three key components: code retrieval for prompt augmentation, LLM-based code generation followed by human evaluation, and continuous updates to the code library

value of the j-th feature of the i-th data record. In practical applications, it is often necessary to perform data cleaning (e.g., handling inconsistent formats, outliers, missing values) and preprocessing (e.g., normalization, label encoding) on the dataset \mathcal{D} before training a machine learning model.

Let \mathcal{A} be a subset of $\{c_1, c_2, ..., c_m\}$, representing the collection of columns that require processing. TradeSweep (\mathcal{M}) is designed to receive a preliminary description r of the user's data preprocessing requirements in English and to automatically generate Python code f and a processed dataset $\hat{\mathcal{D}}$. To achieve this goal, TradeSweep (\mathcal{M}) begins by constructing a prompting curriculum denoted as $\mathcal{P} = [r, \mathcal{C}]$. This curriculum is then inputted into an LLM as a prompt to develop Python code that meets the specified requirements r. Using the code generated by TradeSweep (code proposal \mathcal{G}), the system executes and evaluates \mathcal{G} on a subset of \mathcal{D} to obtain an execution outcome \mathcal{O} that can be presented to the user. It is important to note that, TradeSweep operates without requiring programming expertise from the user, as the specific columns \mathcal{A} necessitating modification are automatically determined by the model, and the entire process from input description to execution outcome is managed by the system.

One of the commonly reported challenges with using LLMs is its vulnerability to hallucinations. In the context of code generation, this can result in code that includes non-existent functions or incorrect syntax, among other issues. To address hallucination challenges, we expanded the existing curriculum $\mathcal{P} = [r, \mathcal{C}]$ into $\mathcal{P} = [r, \mathcal{C}, \mathcal{F}]$. Here, \mathcal{F} denotes a finite set of closely interconnected fundamental functions that serve as a reference for the LLM. This approach offers a dual advantage:

- 1) **Decreased Occurrence of Hallucinations:** By incorporating a set of fundamental functions, the likelihood of hallucinations is reduced.
- Enhanced Code Generation: LLMs can develop new code in addition to leveraging the fundamental reference functions.

To implement the expanded curriculum $\mathcal{P} = [r, \mathcal{C}, \mathcal{F}]$, TradeSweep employs an adaptable code library \mathcal{L} . First, the request r is transformed into an embedding representation emb_r . A set of k relevant fundamental functions $\mathcal{F} = \{f_1, ..., f_k\}$ is then extracted from the code database $\langle emb_i, function_i \rangle$ based on the minimum cosine similarity $\min_{i \in L} \text{cosine}(emb_r, emb_i)$. Once a new code proposal \mathcal{G} is successfully developed to meet complex requirements r, typically through interaction with users, it is saved in the library \mathcal{L} . This allows the library \mathcal{L} to continually learn and provide more accurate and advanced reference functions in the future.

Furthermore, incorporating human feedback in the code generation process ensures that TradeSweep produces code meeting user expectations. Specifically, we present the code proposal \mathcal{G} along with an execution result \mathcal{O} . This allows users to inspect and determine if the outcome meets their

expectations. If satisfied, the code is then executed on the entire dataset, enabling users to verify if the processed data aligns with their domain knowledge rather than examining the code in detail. This makes the tool accessible for non-expert programmers. If the intended outcome is not achieved, TradeSweep will initiate a continuous conversation $\mathcal{P} = [r', \mathcal{C}, \mathcal{F}]$, where r' denotes a revision request. In the Results section (Section V), we demonstrate how TradeSweep achieves a high probability of meeting users' requests on the LLM's first attempt and adapts more rapidly than baselines when revisions are requested.

B. Prompt Augmentation

Following the submission of the user's request for data preprocessing, TradeSweep examines all functions contained within the code library to learn from and utilize them as references. The library comprises code functions commonly applied in various preprocessing tasks. Along with these functions, data information is also provided to the LLM to aid its understanding of the dataset structure.

However, inputting both code functions and data information can result in a lengthy prompt containing redundant information. This not only reduces the LLM's performance in accurately identifying the most relevant function but also significantly delays the LLM's response time. To address this issue, we aim to shorten the input context. Instead of presenting all code functions in the prompt, we utilize information retrieval (IR) techniques to perform a more precise selection on code functions, providing the LLM with the most relevant and useful codes.

We implemented this approach by representing our code library as a vector database. Each vector includes an embedded description of a code function, with its corresponding function code serving as the payload of the vector. During each round of code generation, we query the vector database to retrieve k code functions deemed most relevant to the user's request. These selected functions are then included in the LLM prompt. This process not only effectively shortens the prompt length and reduces the LLM's response time but also helps the LLM focus on functions most pertinent to fulfilling the user's request (see Fig. 3).



Fig. 3. Identifying relevant code functions for prompt augmentation

C. Code Generation

Once the top-k relevant codes have been retrieved as candidate codes, a prompt for the LLM is created by combining the candidate codes with data information and the user's request. The LLM examines the provided information and produces a code proposal designed to accomplish the requested task. If the LLM does not find any candidate code that aligns with the user's request, it generates a novel code function.

The process of code generation includes iterative enhancements under the following scenarios after the initial code proposal is generated:

- Incorrect Code Proposal Format: We provide the LLM with a response template specifying that the proposal format should include a Python code function that reads a dataframe and returns it at the end. The proposal should also include a function call for applying the code to the data. If either of these components is missing in the generated code, the LLM automatically modifies the proposal to meet the required format.
- 2) Execution Error: Once a code proposal is generated in the expected format, it is tested on sample data. If execution fails due to bugs, syntax errors, exceptions, or other issues, both the code and its corresponding error message are returned to the LLM for revision.
- 3) User Feedback: Users are provided with a visualized execution result on sample data, showing examples of input values and their corresponding output values when the code is applied. Based on the user's review of the code proposal and execution outcomes, they can request code revisions by providing feedback describing necessary fixes to the LLM.

Finally, once the user confirms that both the code proposal and execution examples are correct, the code is applied to the target dataset, as shown in Fig. 4. For novel code or when users request saving multiple functions into a single pipeline, the newly generated code or series of codes are added to our code library.



Fig. 4. Code generation and enhancement based on execution results and user feedback

D. Code Library

Due to the complexity of data preprocessing tasks - such as those requiring specific algorithms or involving multiple datasets - a code library is beneficial as a reference document for the LLM to learn from and follow when generating code proposals. In TradeSweep's code library, each function is designed to handle a particular data preprocessing task and includes comments describing its usage. To enhance the efficacy of TradeSweep, we have developed a code library capable of supporting the addition of new functions as interactions progress:

- Novel Code: When the LLM does not identify any function that corresponds to the user's request after analyzing the retrieved codes, it generates a novel code proposal. Since this newly generated code is not initially included in the library, we incorporate it back into the library to improve efficiency and accuracy for future code generations.
- 2) **Pipeline Creation Request**: After a series of code functions have been generated, the user may request that the entire procedure be stored as a pipeline. During this process, the several functions previously applied to the dataset are combined into a single code function, which is then added to our code library.

As illustrated in Fig. 5, the process of integrating the novel code function into the code library involves a sequence of steps. First, we use an LLM to generate a description of the novel code, which is then added as comments. Subsequently, both the newly generated code function and its associated description are then added into our vector database code library. A new query vector is constructed using the function description, and its corresponding payload is generated by the code.



Fig. 5. Code library dynamically updated with each novel LLM-generated code or user pipeline request.

IV. EXPERIMENTS

In this section, we present the design of our experiments, including the choice of datasets, large language model, discussion about the baselines, and evaluation methodology.

A. Experimental Setup

Data: We use shipment-level bill of lading data [22] that captures business-to-business international trade and highlights the complexity of supply chains. The effectiveness of identifying shipments potentially circumventing economic sanctions, high tariffs, or engaging in suspicious activities largely depends on the quality of data initially cleaned and preprocessed. Specifically, we use three trade datasets involving imports and exports across multiple nations for commodities that have been subject to recent import prohibitions, high tariff rates, and sanctions, and all may contain risks associated with origin fraud [23, 24]. The datasets are as follows:

1) **Teak**: This dataset includes 69,134 teakwood transactions exporting from 116 countries to the United States, spanning from July 1, 2007, to August 10, 2023 (5,885 days in total). The data source is Panjiva¹.

- Grain: This dataset comprises 145,217 grain transactions from Russia to 118 global destinations, covering the period from May 20, 2021, to November 30, 2022 (560 days in total). The data sources are ExportGenius² and ImportGenius³.
- Timber: This dataset contains 3,087,822 timber exports from Russia to 173 countries, starting from October 20, 2021, to March 31, 2023 (528 days in total). The data sources are ExportGenius and ImportGenius.

These datasets span significant periods and were manually entered, making them susceptible to errors such as typos, inconsistent formatting, and missing information. Consequently, data analysis becomes a laborious and time-consuming task for analysts, highlighting the urgency and importance of effective data preprocessing.

Vector database: We use Qdrant⁴, a widely recognized and rapidly expanding vector database, as our information retrieval system. Qdrant is employed to efficiently store and retrieve code functions based on its vector embeddings.

Large Language Model: For this study, we employed CodeLlama-13b-Instruct⁵. This model is used to learn from retrieved codes, generate executable Python functions, and modify codes in response to user feedback. Unlike API-based LLMs, CodeLlama-Instruct allows for local execution, offering greater flexibility and control.

Code Library: Our initial code library comprises 12 widely recognized and commonly used data preprocessing functions. These functions cover a range of tasks, including standardizing date formats, removing punctuation marks, correcting misspellings, and filling in missing values based on other columns. Table II presents details on some of the functions from the initial code library.

Hardware Environment: The experiments were conducted using a Tesla P40 GPU with 8 cores, 38 GB of RAM, and 500 GB of disk memory. This hardware setup ensures efficient processing and management of computational tasks associated with code generation and evaluation.

B. Baselines

To evaluate the performance of TradeSweep in generating effective and relevant code, we developed three baselines for comparison:

 Baseline 1 (B1): State-of-the-Art (SOTA) Simulation -Code Generation Using Only LLM. For B1, we abstracted the code generation components of state-of-theart tools such as FLAME and Jigsaw. These tools' codewriting capabilities focus solely on LLM and do not utilize any code libraries or additional augmentations. The LLM must independently generate code based on

⁵https://github.com/meta-llama/codellama

¹https://panjiva.com/

²https://www.exportgenius.in/

³https://www.importgenius.com/

⁴https://qdrant.tech/

 TABLE II

 PREPROCESSING FUNCTIONS FROM THE INITIAL CODE LIBRARY

Task	Explanation		
Remove columns	Delete unwanted column(s).		
Remove Delete rows that satisfy a certain condition.			
Clean num- bers	Remove non-numeric symbols and convert the value to numbers.		
Clean strings	Remove punctuation marks, quotation marks, and any extra spaces.		
Standardize dates	Standardize all date values to a YYYY-mm-dd format.		
Correct mis- spellings	Apply word-embedding and clustering to a column to cluster similar values. Then, in each cluster group, find the most frequent value and correct others to that.		
Compare columns and clean	Between a to-clean column and a reference column, group the two columns. Then, for all to-clean values that have the same reference value, find the most frequent to- clean value and update others on this.		
Lookup doc- ument	Given a to-clean column in the dataset and an external CSV/Excel document, map the to-clean values to a reference column in the document, then create a new column with the mapped values.		

the user's request without external references. This setup evaluates the LLM's capability to produce relevant code purely from the textual description provided by the user, which may involve significant effort and may result in less effective data cleaning outcomes.

- 2) Baseline 2 (B2): LLM Prompted with Candidate Codes and User's Request. For Baseline 2, we provided the LLM with a set of top-k candidate codes retrieved from the code library, along with the user's request. Unlike TradeSweep, which also includes data information, this baseline excludes specific data details from the prompt. This design assesses the impact of not providing data context on the LLM's performance. The LLM must generate code based solely on the provided candidate codes and user request, potentially leading to less accurate code generation due to the lack of data-specific guidance.
- 3) Baseline 3 (B3): Providing the Entire Code Library Without Descriptions. In this baseline, the LLM receives the entire code library, but without any accompanying descriptions or comments. The candidate codes are provided in their raw form, with no explanatory notes. This setup explores the effect of removing code descriptions on the LLM's ability to generate relevant code. Additionally, B3 does not utilize the vector database for prompt augmentation; instead, it provides the full code library as part of the prompt. This approach helps understand the influence of having complete access to code functions without context or descriptions on code generation performance.

C. Evaluation Methodology

For each dataset, we defined a set of data preprocessing tasks based on the dataset's specific content and characteristics,

 TABLE III

 Data preprocessing tasks used in our experiments

	Teak	Grain	Timber
Remove columns	✓		✓
Standardize dates	✓	√	√
Clean numbers	✓		 ✓
Clean strings		\checkmark	 ✓
Correct misspellings	✓		 ✓
Compare columns and clean	✓	\checkmark	√
Lookup documents		\checkmark	 ✓
Others (not in library)	√	\checkmark	√

The code for each task is generated using TradeSweep and the three baselines, and evaluated based on several key metrics. We record both the initial and final versions of the generated codes, noting any revisions made, and measure the time taken to generate the code. After generating the code, we apply it to the initial, unprocessed dataset and compare the execution outputs to manually preprocessed data to assess the accuracy and effectiveness of each method.

When multiple columns are assigned to perform the same task within the same dataset, we generate separate codes for each column using the same algorithm but with slight modifications to tailor the code to the specific column requirements. This approach allows us to evaluate how well each method adapts to different columns and their unique attributes, ensuring a comprehensive comparison of code generation efficiency and performance.

V. RESULTS

This section outlines the results and presents a discussion on them. Fig. 6 illustrates the user interface of TradeSweep. First, users upload a spreadsheet (e.g., Timber.csv), which is displayed on the left half of the screen. They then enter their data preprocessing request, such as *Only show records where trading country is US*, into the text box. The interface also provides access to all previously completed data preprocessing tasks. Once the LLM generates a code proposal, it appears in the box below, where users can test it on sample data and request revisions if necessary.

We applied the codes generated by TradeSweep and the three baselines to the trade datasets for evaluation. In our analysis, we refer to the initial data as *Init* and the manually cleaned data as GT (Ground Truth) for simplicity. The Ground Truth data represents the fully preprocessed and cleaned version of the dataset, achieved through meticulous manual adjustments by data analysts. Our experiments are designed to answer the following questions:

- 1) Can TradeSweep preprocess data accurately? (Section V-A and V-D)
- 2) Does TradeSweep generate high-quality data preprocessing functions? (Section V-B)
- How effective can TradeSweep pass assertion checks (i.e. generating code that functions properly)? (Section V-B)

ē		Trade									
				Q https://trade-sweep.ai							⊻ 💿 ᢓ 🛎 🚔
Tra	de	Sw	/eep				Search CSV		٩		TS
#	10	D	Date	DeclarationNumber	SenderTaxID	SenderNameEng		SenderAddressEng	Recipient	File Timber.csv	×
1	1		2020-01-01	10210200/010120/0000017	1121009024	"LLC ""SYKTYVKAR PLYWOOD PLANT		"167026, KOMI, SYKTYVKAR city, UKHTINSKOE SHOSS	E, PRICE-G	query Remove the CustomsCode Column	
2	1		2020-01-01	10210200/010120/0000019	1121009024	"LLC "SYKTYVKAR PLYWOOD PLANT		"167026, KOMI, SYKTYVKAR city, UKHTINSKOE SHOSS	E, PRICE-G	query: Convert the Dates in the data column to	use - instead of /
3	1		2020-01-01	10210200/010120/0000021	1121009024	"LLC ""SYKTYVKAR PLYWOOD PLANT	***	"167026, KOMI, SYKTYVKAR city, UKHTINSKOE SHOSS	E, MULTIPL	query: Remove AOKLM KO from SenderNameE	ng
4	1		2020-01-01	10210200/010120/0000022	1121009024	"LLC ""SYKTYVKAR PLYWOOD PLANT	***	"167026, KOMI, SYKTYVKAR city, UKHTINSKOE SHOSS	E, MULTIPL	query: Only show TradingCountryCode for US	
5	1		2020-01-01	10210200/010120/0000023	1121009024	"LLC ""SYKTYVKAR PLYWOOD PLANT	***	"167026, KOMI, SYKTYVKAR city, UKHTINSKOE SHOSS	E, MULTIPLY		
6	1		2020-01-01	10210200/010120/0000024	1121009024	"LLC ""SYKTYVKAR PLYWOOD PLANT	***	"167026, KOMI, SYKTYVKAR city, UKHTINSKOE SHOSS	E, MULTIPLY	Enter command	>
7	1		2020-01-01	10210200/010120/0000026	1121009024	"LLC ""SYKTYVKAR PLYWOOD PLANT	***	"167026, KOMI, SYKTYVKAR city, UKHTINSKOE SHOSS	E, MULTIPLY		
8	1		2020-01-01	10210200/010120/0000027	1121009024	"LLC ""SYKTYVKAR PLYWOOD PLANT	***	"167026, KOMI, SYKTYVKAR city, UKHTINSKOE SHOSS	E, MULTIPLY	det filter_us_trading_countries(d return df[df['TradingCountries(d df = filter us_trading_countries(d	<pre>ide'].str.contains('USA', case=False)] if)</pre>
9	1		2020-01-01	10317120/010120/0000011	107008490	"LLC ""TORNADO""		"385200, REPUBLIC OF ADYGEY, city of ADYGEYSK, st	tr GVOL HA		
10	1		2020-01-02	10610080/020120/0000029	2407009152			"660049, KRASNOYARSKY KRAI, KRASNOYARSK city, s	stre CHENGD	TEST	MODIFY
11	1		2020-01-02	10610080/020120/0000029	2407009152			"660049, KRASNOYARSKY KRAI, KRASNOYARSK city, s	stre CHENGD	GENERA	TE SCRIPT
12	1		2020-01-02	10610080/020120/0000029	2407009152			"660049, KRASNOYARSKY KRAI, KRASNOYARSK city, s	stre CHENGD		
13	1		2020-01-02	10012020/020120/0000091	391705872150	IVANOVA ALEXANDRA ANDREEVNA		"238324, KALININGRAD region GURIEVSKY district, P	LA FORE		
14	1		2020-01-02	10012020/020120/0000092	391705872150	IVANOVA ALEXANDRA ANDREEVNA		"238324, KALININGRAD region GURIEVSKY district, P	LA FORE		
15	1		2020-01-02	10013160/020120/0000068	5011021227	"LLC ""KRONOSHPAN""		"140341, MOSCOW region, city of Egorievsk, P. NOVY	"LLC "W		
16	1		2020-01-02	10013160/020120/0000002	5011021227	"LLC ""KRONOSHPAN""		"140341, M.O., CITY EGORIEVSK, P. NOVYY, VLADIE 10	0 "LLC ""D/		
17	1		2020-01-02	10013160/020120/0000002	5011021227	"LLC ""KRONOSHPAN""		"140341, M.O., CITY EGORIEVSK, P. NOVYY, VLADIE 10	0 "LLC "'D/		

Fig. 6. TradeSweep in action - input data (on the right), user query and generated code (on the left)

- To what extent can TradeSweep independently generate valid code proposals? (Section V-C)
- 5) How many rounds of user feedback are needed to produce a valid code proposal? (Section V-C)
- 6) How does the source code library enhance TradeSweep's code generation capabilities? (Section V-D)
- What role does RAG play in improving TradeSweep's code generation? (Section V-D)

A. Preprocessed Dataset Correctness

After preprocessing the three raw datasets (Init) using the code generated by all three baselines (B1 to B3) and TradeSweep, we compared the resulting preprocessed datasets to the Ground Truth (GT), the manually cleaned dataset. As Table IV illustrates, TradeSweep's execution results closely align with GT, demonstrating its effectiveness in preprocessing across the three datasets. Baselines 2 and 3 also produced results similar to TradeSweep, as they benefited from referencing and learning from the code library during code generation. Conversely, the codes generated by Baseline 1 showed considerable deviation from Ground Truth. This discrepancy is due to Baseline 1's reliance on the LLM's independent code generation, which often led to the use of different algorithms and approaches compared to those in the code library.

B. Code Quality

To evaluate the code quality of the functions produced by TradeSweep and the baselines, we utilized Pylint⁶ and Radon⁷ as metrics. Pylint measures errors (e.g., syntax errors),

TABLE IV Evaluation of preprocessed outputs relative to Ground Truth (GT)

		Init	GT	B1	B2	B3	TS
	# of cols	127	26	26			
	# of rows	s 69,134				134	
Teak	NaNs (%)	48.77	7.57	10.94	7.09	7.09	7.09
	incorrect	19.99	0	29.98	9.86	1.03	1.03
	formats (%)						
	typos (%)	9.29	0	25	3.92	3.92	3.92
	# of cols	56	61	61			
	# of rows	145,	217		145	,217	
Grain	NaNs (%)	66.94	62.10	3.76	63.07	66.22	63.57
	incorrect	2.84	0	9.68	0.09	0.08	0.08
	formats (%)						
	typos (%)	18.06	0	14.43	2.37	6.65	3.21
	# of cols	29	23		2	3	
# of rows 3,087,822				3,087	7,822		
Timber	NaNs (%)	3.19	4.55	8.68	4.75	8.24	4.42
	incorrect	87.93	0	64.86	0.63	12.81	0.63
	formats (%)						
	typos (%)	93.19	0	74.69	4.04	6.83	2.76

warnings (e.g., unused import packages), and other aspects for Python functions, assigning a score between 0 (inefficient) to 10 (efficient). Radon calculates *Cyclomatic Complexity*, providing a rank from A (simplest) to F (most complex). For simple data preprocessing tasks, such as "Clean dates to yyyymm-dd format", TradeSweep and the three baselines produced code of comparable quality and complexity, regardless of whether they utilized a code library (see Table V).

However, for more complex tasks, such as comparing two columns and cleaning the target column using

⁶https://www.pylint.org/

⁷https://pypi.org/project/radon/

the most frequent value from the reference column, the use of a code library significantly improves code quality. TradeSweep obtained a higher Pylint score and generated a more complex code; Baseline 1, with a Radon complexity rank of A, produced code that was too simple to fulfill the task. This enhancement is evident in the results, where functions generated with the aid of a code library demonstrated higher quality scores and are more complex compared to Baseline 1, which relied on independent code generation without library references.

As part of TradeSweep, we integrated an assertion check feature to ensure that the generated functions execute correctly before they are presented to the user. If a function encounters execution failures due to exceptions, errors, or other issues, the incorrect function and the corresponding error message are sent back to the LLM for correction. This process is performed in an iterative manner.

Table VI shows the average number of execution failures for TradeSweep compared to the baselines. TradeSweep, by leveraging a code library and data information for the LLM, consistently achieved the lowest number of execution failures across various data preprocessing tasks. In contrast, Baseline 1, which generated functions independently without utilizing a code database, had a higher error rate due to numerous coding errors. Baseline 2 encountered challenges due to the lack of data information; for instance, if a user request was vague (e.g., "Correct misspellings in shipper names" without specifying to apply the function on the "Shipper" column), the LLM erroneously applied the function to a non-existent column, leading to an infinite loop of KeyErrors. Baseline 3 experienced difficulties because the LLM was tasked with interpreting all functions in the code library rather than focusing on the most relevant ones, resulting in a higher likelihood of errors in its generated functions.

C. User Involvement

Table VII records the rate of codes approved by users on the LLM's first attempt. Each value represents the number of codes accepted on the first try, divided by the total number of accepted codes. TradeSweep achieved the highest rate of user-approved codes across all three datasets without needing revisions, exhibiting better performance in generating a function that meets the user's demand on the LLM's first attempt by leveraging knowledge from the code library. In contrast, Baseline 1 often produces unstable results initially, as it generates code independently without the benefit of the code library. This results in frequent misalignment with vague user requests (e.g., "Clean numbers in net weights") and requires more explicit instructions (e.g., "Clean numbers in net weights by removing commas and converting the values to float numbers"). The lowest acceptance rate was noted for Baseline 2 due to the absence of column name information provided to the LLM. Although Baseline 2 used the same algorithm as TradeSweep, the lack of column name information necessitates nearly all preprocessing tasks to undergo revision, highlighting the need for accurate column specification (e.g.,

"Clean numbers in the NetWeight column" rather than "Clean numbers in net weights"). Baseline 3, while also capable of generating correct codes initially, requires slightly more user involvement compared to TradeSweep. In this baseline, although the LLM often selects the correct code for reference, the absence of function descriptions leads to over- or undermodification of reference code, resulting in additional revision requests.

D. Enhancing Code Generation Capabilities

To evaluate the impact of how having a code library can enhance the generation of data preprocessing functions, we compared the overall results of data preprocessed by TradeSweep with Baseline 1, as the primary distinction between the two lies in TradeSweep's use of a code library. After conducting preprocessing on the three datasets, Table VIII illustrates that TradeSweep, by leveraging example code functions from the library, consistently produces suitable and executable code that can preprocess data with much higher correct-value rates (comparable to the Ground Truth data) compared with Baseline 1.

The ability to maintain and expand the code library is crucial for handling more complex and repetitive user requests, as newly generated code functions can be added to the library for future use. The results presented in Table IX suggests that Baseline 3, which provides the entire code library to the LLM, requires significantly more time for code generation compared to TradeSweep. This is due to the increased prompt length, which slows down the process. TradeSweep's use of RAG to retrieve only the top-k most relevant codes from the library reduces prompt length and improves code generation efficiency. In our experiments, we determined that setting the value of k to 3 optimizes performance. Values of $k \ge 4$ result in longer code generation times, while values of $k \leq 2$ increase the likelihood of retrieving irrelevant functions. For instance, if a user requests to "clean the dates," setting k too low may lead to retrieving functions like "clean numbers" or "clean strings", which are less relevant compared to "standardize dates". Thus, k = 3 strikes a balance between retrieval accuracy and prompt length efficiency. In conclusion, Baseline 3, which excluded information retrieval and inputted a lengthy prompt to the LLM, required considerably longer time for generating codes compared to TradeSweep.

E. Case Study

To better understand the differences in code generation performance among various baselines and TradeSweep, we analyzed the initial code proposals generated by each method before any user feedback was provided. All baselines received identical user request inputs for each preprocessing task.

Using the example of cleaning the column "Shipper Country", Fig. 7 shows the code proposal generated by the three baselines and TradeSweep. In the scenario, the user's request was: "Compare values in shipper country with shippers, and clean country names with the same shipper to the most frequent value." The ideal approach would involve using TABLE V

COMPARISON OF CODE QUALITY FOR A SIMPLE TASK (CLEAN DATES TO YYYY-MM-DD FORMAT) VS. COMPLEX TASK (COMPARE TWO COLUMNS AND CLEAN)

		Teak		G	rain	Timber		
		simple task	complex task	simple task	complex task	simple task	complex task	
	# of lines (\downarrow)	7	7	7	6	6	6	
Baseline 1	Pylint (↑)	4.29	3.64	5.71	4.44	5.71	1.67	
	$\begin{array}{ $	Α	А	Α	А	Α		
	# of lines (\downarrow)	8	6	9	6	8	6	
Baseline 2	Pylint (†)	4.17	5.65	5	5.65	5.83	5.65	
Baseline 2# of lines (\downarrow Baseline 2# of lines (\downarrow Baseline 3# of lines (\downarrow	Radon (\downarrow)	A	С	А	С	А	С	
	# of lines (\downarrow)	8	6	10	6	6	6	
Baseline 3	Pylint (↑)	5.83	5.65	5.38	5.65	5.71	5.65	
	Radon (\downarrow)	A	С	А	С	А	С	
	# of lines (\downarrow)	7	6	7	6	8	6	
TradeSweep	Pylint (†)	5.45	5.65	6	5.65	5	5.65	
	Radon (\downarrow)	A	С	A	С	Α	С	

 TABLE VI

 Average number of execution failures in the generated code

	Baseline 1	Baseline 2	Baseline3	TradeSweep
Teak	1.92	$ \infty$	0.17	0.08
Grain	1.72	∞	0.72	0.11
Timber	1.17	∞	0.67	0.08

 TABLE VII

 Acceptance rate of initial code proposals

	Baseline 1	Baseline 2	Baseline3	TradeSweep
Teak	5/12	1/12	9/12	10/12
Grain	5/18	5/18	12/18	12/18
Timber	15/24	5/24	17/24	18/24

the "compare_and_clean" function from the code library and apply it to the "Shipper Country" column while comparing with the "Shipper" column.

1) **TradeSweep**: In TradeSweep, the top three relevant functions were retrieved from the code library and inputted into the LLM. The generated code proposal demonstrated that the LLM effectively learned from the "compare_and_clean" function in the code library. By including data column names, TradeSweep accurately identified the columns intended for the task. The proposed code grouped the data by the "Shipper"

 TABLE VIII

 IMPACT OF CODE LIBRARY IN GENERATING RELEVANT CODE FOR DATA PREPROCESSING TASKS - BASELINE 1 VS. TRADESWEEP

			Teak	Grain	Timber		
Ini	t correct-rate		86.60%	85.94%	9.09%		
GT correct-rate			100% (time-consuming)				
Bacalina 1	correct-rate first-attempt valid rate		73.09%	87.65%	29.58%		
Dasenne 1			41.66%	27.78%	62.5%		
TradeSween	correct-rate		97.19%	97.61%	98.17%		
madesweep	first-attempt valid rate		83.33%	66.67%	75%		

 TABLE IX

 IMPACT OF RAG ON CODE GENERATION - BASELINE 3 VS. TRADESWEEP

	Teak	Grain	Timber
Baseline 3	555.46	369.99	269.34
TradeSweep	167.55	117.40	94.08

column and, within each group, identified the most frequent "Shipper Country" value, modifying other values to match this frequent value.

2) **Baseline 1**: Baseline 1 generated code based on assumptions about the user's request. Although the proposed code successfully grouped the columns and found the most frequent shipper country, it made errors by converting all strings to lowercase and not handling potential NaN values properly. This led to a higher prevalence of NaNs and incorrectly formatted country names. Unlike TradeSweep, Baseline 1 failed to use the most frequent non-NaN value, resulting in a less accurate output.

3) **Baseline 2**: Baseline 2 was provided with candidate codes and the user's request but lacked data information. The code generated attempted to use the "compare_ and_clean" function from the code library, but without knowing the exact column names, it made assumptions. The LLM used column names like "shippers" and "shipper_country" instead of the correct "Shipper" and "Shipper Country," leading to incorrect application of the function. In contrast, TradeSweep accurately applied the function to the intended columns by utilizing the provided data information.

4) **Baseline 3**: In Baseline 3, all functions in the code library were inputted into the LLM without their descriptions. Although the LLM could understand the usage of each function, the excessive length of the prompt, due to including all functions, led to a significant slowdown in code generation—taking 690 seconds. Despite successfully applying the "compare_and_clean" function, the approach of inputting the entire code library resulted in inefficiencies compared to



Fig. 7. Case Study: A snapshot of code proposal generated by the three baselines (B1 - B3) and TradeSweep

TradeSweep's method of using only the top-k relevant functions, which streamlined the process and improved generation speed.

F. Limitations

During the data preprocessing operation, if TradeSweep is unable to identify a relevant function in the code library, it generates a novel code. This limitation causes TradeSweep to behave similarly to SOTA tools, relying solely on the LLM to generate code. Consequently, the user must provide more iterations of feedback and changes for the LLM to produce valid code functions. Another limitation is that the execution outputs of sample data may not provide users with sufficient information to assess the code's performance for certain tasks. For instance, correcting misspellings in company names might result in over- or under-correction, and by observing a limited number of correction examples, users may not confirm whether the code successfully captures all corner cases. In such scenarios, it is helpful to perform an apply-test on the full target data and generate a comparison of names before and after correction. However, applying codes to the full dataset takes much more time than only applying it to sample data.

To address these limitations in the future, several improvements can be made. Expanding the code library with a broader range of pre-defined functions can reduce the likelihood of generating novel code, thereby decreasing the need for user feedback and iterations. Implementing more advanced sampling techniques for the execution outputs can provide users with a more representative subset of the data, helping them better assess the code's performance. Additionally, optimizing the apply-test procedure to efficiently handle larger datasets can expedite the process, enabling quicker feedback and validation for users.

VI. CONCLUSION

The rise of automation and programming support through LLMs has significantly reduced the turnaround time for processing large spreadsheets. Our method, TradeSweep, functions as an LLM-driven data preprocessing agent, retrieving suitable functions from a code library and adapting them to fulfill specific data preprocessing tasks. Our results demonstrate TradeSweep's effectiveness in practical data transformation scenarios. In experiments conducted on three transaction datasets, TradeSweep efficiently performed the requested data preprocessing tasks with minimal user interaction. The ablation study highlights that utilizing a code library and a vector database for information retrieval accelerates the code generation process. Moreover, providing data information and function descriptions to the LLM enhances the accuracy of code generation, reducing the need for user feedback and enabling users with limited programming expertise to achieve effective results.

REFERENCES

- M. Heller, "Large language models and the rise of the AI code generators," InfoWorld, May 23, 2023. https://www.infoworld.com/article/3696970/llmsand-the-rise-of-the-ai-code-generators.html
- [2] P. Ingle, "Top artificial intelligence (AI) tools that can generate code to help programmers (2024)," MarkTechPost, Mar. 14. 2024.https://www.marktechpost.com/2024/03/14/top-artificialintelligence-ai-tools-that-can-generate-code-to-helpprogrammers/
- [3] B. D. Ramel, 06/30/2023, "Top 10 AI 'copilot' tools for visual studio code -," Visual Studio Magazine. https://visualstudiomagazine.com/articles/2023/06/30/vscode-copilots.aspx (accessed Jul. 22, 2024).
- [4] S. Anunaya, "Data preprocessing in data mining: a hands on guide," Analytics Vidhya, Aug. 10, 2021. https://www.analyticsvidhya.com/blog/2021/08/datapreprocessing-in-data-mining-a-hands-on-guide (accessed Jul. 22, 2024).
- [5] A. A. Kandilli, "How is dirty data handled in data analytics?," Medium, Jul. 26, 2022. https://medium.com/@sweephy/how-is-dirty-datahandled-in-data-analytics-1767fb998e37 (accessed Jul. 22, 2024).
- [6] H. Joshi et al. (2024) "Flame: a small language model for spreadsheet formulas", Proceedings of the AAAI Conference on Artificial Intelligence, 38(12), pp. 12995–13003. doi:10.1609/aaai.v38i12.29197.
- [7] J. Li et al. (2024) "Can LLM already serve as a database interface? a big bench for large-scale database grounded text-to-SQLs". In Proceedings of the 37th International Conference on Neural Information Processing Systems (NIPS '23). Curran Associates Inc., Red Hook, NY, USA, Article 1835, 42330–42357.
- [8] P. Lewis et al. 2020. "Retrieval-augmented generation for knowledge-intensive NLP tasks". In Proceedings of the 34th International Conference on Neural Information Processing Systems (NIPS '20). Curran Associates Inc., Red Hook, NY, USA, Article 793, 9459–9474.
- [9] Q. Tang et al., "Self-retrieval: building an information retrieval system with one large language model", arXiv [cs.IR]. 2024.
- [10] C. Fan, M. Chen, X. Wang, J. Wang and B. Huang. "A review on data preprocessing techniques toward efficient and reliable knowledge discovery from building operational Data." Frontiers in Energy Research (2021).
- [11] C. Gopal, "Network visualization and anomaly detection in international timber trade flows". 2021.
- [12] Y. Tian et al. "SpreadsheetLLM: encoding spreadsheets for large language models." arXiv preprint arXiv:2407.09025 (2024).
- [13] J. Li et al. "SkCoder: a sketch-based approach for automatic code generation." 2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)

(2023): 2124-2135.

- [14] N. Jain et al. "Jigsaw: large language models meet program synthesis." 2022 IEEE/ACM 44th International Conference on Software Engineering (ICSE) (2021): 1219-1231.
- [15] K. Zhang, D. Wang, J. Xia, W. Y. Wang, and L. Li. 2024. "ALGO: synthesizing algorithmic programs with LLM-generated oracle verifiers". In Proceedings of the 37th International Conference on Neural Information Processing Systems (NIPS '23). Curran Associates Inc., Red Hook, NY, USA, Article 2389, 54769–54784.
- [16] cognition.ai, "Cognition introducing Devin, the first AI software engineer," cognition.ai. https://www.cognition.ai/blog/introducing-devin
- [17] H. Zhang, Y. Dong, C. Xiao and M. Oyamada. "Large language models as data preprocessors." ArXiv abs/2308.16361 (2023): n. pag.
- [18] D. Qi and J. Wang. "CleanAgent: automating data standardization with LLM-based agents." ArXiv abs/2403.08291 (2024): n. pag.
- [19] S. Lu et al. "ReACC: a retrieval-augmented code completion framework." ArXiv, (2022). Accessed July 25, 2024. /abs/2203.07722.
- [20] Md. R. Parvez, W. U. Ahmad, S. Chakraborty, B. Ray and K. W. Chang. "Retrieval augmented code generation and summarization." ArXiv abs/2108.11601 (2021): n. pag.
- [21] F. Zhang et al. "RepoCoder: Repository-Level Code Completion Through Iterative Retrieval and Generation." Conference on Empirical Methods in Natural Language Processing (2023).
- [22] A. Flaaen et al., "Bill of lading data in international trade research with an application to the COVID-19 pandemic," Finance and Economics Discussion Series, vol. 2021, no. 066, pp. 1–40, Oct. 2021, doi: https://doi.org/10.17016/feds.2021.066.
- "US [23] L. Aratani, imports of 'blood teak' from Myanmar continue despite sanctions," The Guardian, 2023. Available: May 16, https://www.theguardian.com/world/2023/may/16/myanmarteak-wood-import-sanctions
- [24] "Russia smuggling Ukrainian grain to help pay for Putin's war," AP NEWS, Oct. 03, 2022. https://apnews.com/article/russia-ukraine-putin-businesslebanon-syria-87c3b6fea3f4c326003123b21aa78099
- [25] Y. Song et al. 2024. "Enhancing text-to-SQL translation for financial system design". In Proceedings of the 46th International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP '24). Association for Computing Machinery, New York, NY, USA, 252–262. https://doiorg.ezproxy.lib.vt.edu/10.1145/3639477.3639732
- [26] X. Xu, C. Liu, and D. Song. "Sqlnet: generating structured queries from natural language without reinforcement learning." arXiv preprint arXiv:1711.04436 (2017).
- [27] N. Yaghmazadeh, Y. Wang, I. Dillig, and T. Dillig. 2017.

"SQLizer: query synthesis from natural language". Proc. ACM Program. Lang. 1, OOPSLA, Article 63 (October 2017), 26 pages. https://doi.org/10.1145/3133887

- [28] B. Wang et al. "Mac-sql: a multi-agent collaborative framework for text-to-sql." arXiv preprint arXiv:2312.11242 (2024).
- [29] L. Huang et al. "A survey on hallucination in large language models: principles, taxonomy, challenges, and open questions." ArXiv abs/2311.05232 (2023): n. pag.
- [30] J.D. Zamfirescu-Pereira, R. Y. Wong, B. Hartmann, and Q. Yang. 2023. "Why Johnny can't prompt: How non-AI experts try (and fail) to design LLM prompts". In Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems (CHI '23). Association for Computing Machinery, New York, NY, USA, Article 437, 1–21. https://doi.org/10.1145/3544548.3581388
- [31] Z. Xu, S. Jain, and M. Kankanhalli. "Hallucination is inevitable: an innate limitation of large language models." arXiv preprint arXiv:2401.11817 (2024).
- [32] M. Izadi et al. 2024. "Language models for code completion: a practical evaluation". In Proceedings of the IEEE/ACM 46th International Conference on Software Engineering (ICSE '24). Association for Computing Machinery, New York, NY, USA, Article 79, 1–13. https://doi.org/10.1145/3597503.3639138
- [33] A. de Moor, A. van Deursen, and M. Izadi. 2024. "A transformer-based approach for smart invocation of automatic code completion". In Proceedings of the 1st ACM International Conference on AI-Powered Software (AIware 2024). Association for Computing Machinery, New York, NY, USA, 28–37. https://doi.org/10.1145/3664646.3664760
- [34] M. Izadi, R. Gismondi, and G. Gousios. 2022. "Code-Fill: multi-token code completion by jointly learning from structure and naming sequences". In Proceedings of the 44th International Conference on Software Engineering (ICSE '22). Association for Computing Machinery, New York, NY, USA, 401–412. https://doi.org/10.1145/3510003.3510172
- [35] P. Li, Z. Chen, X. Chu, and K. Rong. 2023. "Diff-Prep: differentiable data preprocessing pipeline search for learning over tabular data". Proc. ACM Manag. Data 1, 2, Article 183 (June 2023), 26 pages. https://doi.org/10.1145/3589328
- [36] B. A. Halperin and S. M Lukin. 2024. "Artificial dreams: surreal visual storytelling as inquiry into AI 'hallucination'". In Proceedings of the 2024 ACM Designing Interactive Systems Conference (DIS '24). Association for Computing Machinery, New York, NY, USA, 619–637. https://doi.org/10.1145/3643834.3660685
- [37] D. Nam, A. Macvean, V. Hellendoorn, B. Vasilescu, and B. Myers. 2024. "Using an LLM to help with code understanding". In Proceedings of the IEEE/ACM 46th International Conference on Software Engineering (ICSE '24). Association for Computing Machinery, New York, NY, USA, Article 97, 1–13. https://doi-

org.ezproxy.lib.vt.edu/10.1145/3597503.3639187

- [38] B. Jury, A. Lorusso, J. Leinonen, P. Denny, and A. Luxton-Reilly. 2024. "Evaluating LLMgenerated worked examples in an introductory programming course". In Proceedings of the 26th Australasian Computing Education Conference (ACE '24). Association for Computing Machinery, York. New NY. USA, 77-86. https://doiorg.ezproxy.lib.vt.edu/10.1145/3636243.3636252
- [39] T. Coignion, C. Quinton, and R. Rouvoy. "A performance study of LLM-generated code on Leetcode." In Proceedings of the 28th International Conference on Evaluation and Assessment in Software Engineering, pp. 79-89. 2024.
- [40] S. Roychowdhury. 2024. "Journey of hallucinationminimized generative AI solutions for financial decision makers". In Proceedings of the 17th ACM International Conference on Web Search and Data Mining (WSDM '24). Association for Computing Machinery, New York, NY, USA, 1180–1181. https://doi.org/10.1145/3616855.3635737
- [41] A. Mittal, R. Murthy, V. Kumar, and R. Bhat. 2024. "Towards understanding and mitigating the hallucinations in NLP and speech". In Proceedings of the 7th Joint International Conference on Data Science & Management of Data (11th ACM IKDD CODS and 29th COMAD) (CODS-COMAD '24). Association for Computing Machinery, New York, NY, USA, 489–492. https://doi.org/10.1145/3632410.3633297
- [42] A. Khurana, H. Subramonyam, and P. K Chilana. 2024. "Why and when LLM-based assistants can go wrong: investigating the effectiveness of prompt-based interactions for software help-seeking". In Proceedings of the 29th International Conference on Intelligent User Interfaces (IUI '24). Association for Computing Machinery, New York, NY, USA, 288–303. https://doi.org/10.1145/3640543.3645200
- [43] W. Wang, H. Ning, G. Zhang, L. Liu, and Y. Wang. 2024. "Rocks coding, not development: a human-centric, experimental evaluation of LLMsupported SE tasks". Proc. ACM Softw. Eng. 1, FSE, Article 32 (July 2024), 23 pages. https://doiorg.ezproxy.lib.vt.edu/10.1145/3643758
- [44] R. Nouaji, S. Bitchebe, and O. Balmau. 2024. "Speedy-Loader: efficient pipelining of data preprocessing and machine learning training". In Proceedings of the 4th Workshop on Machine Learning and Systems (EuroMLSys '24). Association for Computing Machinery, New York, NY, USA, 65–72. https://doiorg.ezproxy.lib.vt.edu/10.1145/3642970.3655824
- [45] T. Kim et al. 2024. "FusionFlow: accelerating preprocessing for machine learning data with CPU-GPU cooperation". Proc. VLDB Endow. 17, 4 (December 2023), 863–876. https://doiorg.ezproxy.lib.vt.edu/10.14778/3636218.3636238