

Response Letter

Dear Reviewers,

We express our gratitude to the reviewers for their valuable feedback. In response to the reviewers' comments, we addressed four key areas of improvement: (1) clarifying the setup and default parameter values, (2) enhancing the presentation, (3) including a user study, (4) including additional experiments with more baselines, and (5) including additional experiments varying the quality of the corpus. To address these issues, we have provided a detailed explanation of our responses to the main concerns that were raised by all reviewers. We have also addressed each reviewer's comments individually.

Summary of the Revision Plan

We provide a brief overview of the main changes we made in the revision.

(Overall Point 1) Clarifying the setup and default parameter values. Based on the reviewers' suggestions, we revised our paper in the following ways:

- **Clarifying the setup:** (addressing R1.W2). In the revised paper, we provide a more detailed explanation of the corpus collection process in Section 6.1.3. Specifically, we now explain that the corpora were collected from Kaggle competitions using the Kaggle API. We selected six Kaggle competitions, downloaded all available scripts attached to the competitions, and used them as input corpora. To evaluate the effectiveness of our system, for each competition, we treat each script in the corpus as an input user sketch script and run our algorithm on the remaining scripts as the corpus. We report the average improvement obtained (calculated as explained in Section 6.1.4).

- **Parameterizing our framework:** In the revised paper, we discuss some concrete directions for parameterizing Algorithm 1 based on our experimental results on six real-world datasets (addressing R4.W2, R4.D2). We added the discussion in Section 5.3. In summary, Table 2 lists the possible default parameters based on different corpus properties (such as its size and diversity of scripts within it). In addition, the user-intent threshold is set to $\tau_I = 0.9$ and $\tau_M = 1\%$ as default. We explore how varying the user-intent threshold affects our framework in Section 6.3.2. Lastly, early-checking is set to *on* as a default to ensure the output script is executable. We discuss potential ways to automatically tune these parameters in Section 8.

(Overall Point 2) Enhancing the presentation. To enhance the presentation of the paper, we made the following revisions:

- **Running example:** Following the comments of the reviewers, we have revised our running example to make it clearer and more intuitive (addressing R1.W1, R1.D1, R2.D1, R3.D1). Please refer to the revised introduction for full details. Furthermore, we have explained the purpose of the input script, clarified the representation of Table 1, and elaborated on the rationale behind the added data preparation operation to the modified script (which relies on domain-specific knowledge the user does not possess) to handle outliers. Additionally, we clarified that the modifications to the input script are presented as recommendations, allowing the

user to accept or reject them based on their domain knowledge and modeling objectives.

Furthermore, we clarified in the introduction that to ensure the modifications do not significantly alter the data or impact the performance of the downstream task, the user can specify a threshold to define the extent to which their intent should be preserved. Specifically, as noted in the introduction, we support two measures to evaluate how well the output script preserves the user's intent: (1) a table distance measure, which compares the structured datasets generated by the scripts, and (2) a performance measure that evaluates downstream application outcomes, such as accuracy or fairness of the model.

- **Usability:** (addressing R4.W1, R4.D1). Following the comments of the reviewers, we have included a discussion in the introduction regarding the usability of the proposed system.

We now mention that the proposed system is particularly useful in scenarios where users are working with datasets that are commonly used by other researchers, such as in public policy and medicine. For instance, datasets like Amazon reviews¹ or MIMIC² datasets are widely utilized by researchers worldwide and are often accompanied by open-source data preparation scripts available online (e.g., through platforms such as Papers With Code³ or Kaggle⁴). A researcher new to these open-source datasets can benefit from our approach by quickly identifying standard data preparation steps employed by other researchers. As we demonstrate experimentally, the corpus needed to gain valuable insights can be relatively small, making the system practical for real-life scenarios.

- **Typos:** We thoroughly proofread the paper to eliminate any typos or errors (addressing multiple comments by the reviewers).

(Overall Point 3) A user study. Following the comments of the reviewers, we included a user study to validate that our generated scripts are considered by data scientists more standard compared to the output of the baselines (addressing R1.W1, R2.W1).

To that end, we recruited 34 undergraduate and graduate computer science students who are familiar with data preparation and Python programming. Each participant was given two use cases that include the following:

- An input script.
- Statistics on the prevalence of operations in the corpus (similar to Table 1).
- 5 output scripts generated by LS, and the baselines GPT-3.5, and GPT-4, Auto-Suggest, and Auto-Tables (see details about the new baselines in Overall Point 4 below).

The participants were asked to rank each of the output scripts on a scale of 1 to 5 based on: (a) how standardized the script is w.r.t. the corpus, where standardization represents the commonality of the chosen data preparation steps; and (b) how helpful the script is w.r.t. preserving the user intent, measured by their modeling task. The complete form provided to participants is available at [2].

In summary, LS was rated the most standard and helpful in both cases, among the 5 methods examined. The full results of this study are given and discussed in Section 6.2.2.

¹<https://jmcauley.ucsd.edu/data/amazon/>

²<https://archive.physionet.org/physiobank/database/mimicdb/>

³<https://paperswithcode.com/>

⁴<https://www.kaggle.com/code>

(Overall Point 4) Comparison with related work. We added experiments to compare with two state-of-the-art tabular data transformation recommendation systems, Auto-Suggest [63] and Auto-Tables [43] (addressing R2.W2, R2.W4), which are the academic systems closest to our work. Auto-Suggest [63] uses machine learning models trained on real-world notebooks to predict a single next step for an input table based on its characteristics, while Auto-Tables [43] is the state-of-the-art system that predicts multi-step transformations, among a set of operators. Although neither of the systems are publicly available, we reached out to the authors and obtained their permission to access the source code.

Empirically, we observed that these two methods did not show significant improvement in our script standardization task (Section 6.3.1). This is because they focus on a fixed set of table structural transformations, such as transpose and pivot. However, the real-world scripts that we observed focus on feature engineering and data cleaning tasks (Figure 1a). Our bottom-up approach makes LS robust to the type of data preprocessing tasks.

(Overall Point 5) Quality of the reference scripts and robustness of our approach. We added experiments using low-ranked Kaggle reference scripts to evaluate the effect of quality on script standardness (addressing R1.D3, R2.D5). Our method still managed to improve script standardness, despite the smaller and low-ranked corpora (Section 6.3.3). Although we observed that the improvement on script standardness is relatively low (5%), even lower than using a corpus for a similar dataset (10%), we suggest that if only low-quality corpora are available, it is better to also augment the corpora with scripts processing similar datasets if such scripts are available.

This concludes the summary of the main revisions we made. In the following, you will find detailed responses to the reviewers' comments. Updates made in response to specific comments are highlighted in the revised manuscript using the following colors:

- Changes for multiple reviewers in orange
- Reviewer 1 (R1) in blue
- Reviewer 2 (R2) in magenta
- Reviewer 3 (R3) in purple
- Reviewer 4 (R4) in green

The following is our detailed responses to individual comments from our reviewers and meta reviewer.

Meta-Review

The main WEAKNESSES of the paper seem to be:

- (1) Explanation of the setup and fine-tuning aspects (choice of scripts, fine-tuning of parameters)
- (2) Improvement of presentation (intro, examples, usage of English)
- (3) Missing a detailed discussion of why the proposed method is relevant and practical to real-world data curators, ideally (but not obligatorily) backed up by results of a user study
- (4) Missing a comparison to an alternative way to obtain standardized scripts based on next step recommendation or a clearer argument of why this is not reasonable / possible in the related work section.

Ans. We thank the meta-reviewer and the reviewers for the revision opportunity. As discussed in the summary above and in the responses below, we addressed all required changes and other comments in the revised paper.

Reviewer #1

R1.W1 *Initial examples raise questions as some of the suggested changes are puzzling. So there may be a bit of a disconnect between the quantitative results from evaluation, and the intuitive correctness of the generated scripts.*

Ans. We thank the reviewer for this valuable comment. We have revised our motivating example in the introduction to make it more intuitive, as outlined in Overall Point 2 above. To validate the correctness of our approach, we have included a user study in the revised paper, as mentioned in Overall Point 3. Please refer to our detailed response in Overall Point 3.

R1.W2 *from the evaluation section I am missing a description of how the reference scripts are selected – this is critically important but I don't think the space of possible reference scripts has been explored as a parameter of the problem space.*

Ans. To clarify the setup, please see our response to Overall Point 1. We note that the examined corpora vary in size (as shown in Table 3) and quality. To further assess the impact of corpus quality on performance, we added an experiment varying corpus quality using low-ranked reference scripts from Kaggle. Please refer to our newly added experiment in Overall Point 5.

R1.W3 *As LLM capabilities increase, it would be interesting to see how this solution may incorporate them and thus "survive".*

Ans. We thank the reviewer for this important comment. Indeed, with the rapid improvement of LLMs, especially in code generation, we expect to see improvements in their performance. We believe that LLMs could be integrated into our system to automatically verify that constraints are satisfied without explicitly checking them (i.e., ensuring that the user's intent is preserved to the required extent and that the code is executable). We have added this point to the future work discussion in Section 8. However, as shown by our experiments with recent top-performing LLMs, they are currently far from adequate to modify an input script while correctly preserving the user's intent and embedding domain-specific knowledge extracted from the given corpus. We plan to rerun the experiments with new versions of LLMs in the future to track their improvement in this task.

R1.D1 *example 1 suggests using dropna() instead of imputing. This is a big decision with consequences that an analyst is unlikely to accept. For instance, it reduces the size of the training set, possibly very drastically as this is done across the board. Also a new condition on skinThickness appears which I am struggle to connect with original user's intent, as it was not there at all. it may take care of outliers but again, there may very good reasons why it wasn't there. also normalisation is done across the board and blindly, as a data scientist I would not necessarily endorse this practice.*

Ans. We thank the reviewer for this valuable comment. We have revised our running example to make it more intuitive and clear, as explained in Overall Point 2. Additionally, we have clarified in the introduction that the modifications can be viewed as recommendations for the user, who has the option to accept or reject them. To ensure that the

modifications do not significantly alter the data or impact the performance of the downstream task, the user can specify a threshold to indicate the extent to which their intent should be preserved. Specifically, as highlighted in the introduction, we support two measures to assess how well the output script preserves the user's intent: (1) a table distance measure, which compares the structured datasets generated by the scripts, and (2) a performance measure that evaluates downstream applications, such as model accuracy or fairness.

R1.D2 *following from D1, it seems the approach is to provide suggestions with the assumption that the analysts' choices are insufficient / inadequate. This is not always the case, in Fig. 7 there may be a good reason to create a copy of the df, for example. What I see missing here is explanations for the suggestions given – without those, the analyst is still in the dark on which of the two versions is better.* **Ans.**

We agree with the reviewer that the user's decision is not always inadequate. Our approach aims to offer a way to identify common (i.e., standard) data preparation steps tailored to the user's script, based on a corpus of scripts that the user considers valuable. Intuitively, this allows the user to easily "compare notes". We have revised the introduction to clarify this important point.

Additionally, we now explain in the introduction that our approach provides the user with full transparency regarding what is changed and how the changes affect the user's intent (e.g., data values, model accuracy). This means that adopting the recommended modifications is not an "all-or-nothing" decision (since the user may accept only some of the modifications), and the user is informed of the potential outcomes of adopting the changes before making a decision.

The reviewer raised an interesting point regarding the expandability of the output. We have added a discussion on this in the future work section (Section 8), highlighting it as a potential area for further exploration.

R1.D3 *with ref to W2: surely the quality of the generated script depends on that of the reference corpus? but there seems to be an assumption that those corpora are given, eg in sec 5.1 they are curated but this does not include QA on those scripts... and it seems this is not explored in the experiments* **Ans.** Please refer to our response to R1.W2 above.

Reviewer #2

R2.W1 *I am not convinced that standardization, as defined in this paper, is helpful or practical for data scientists. A user study would be necessary to confirm the initial hypothesis that standardization is helpful.* **Ans.** Following the feedback of the reviewer, we have included a user study in the revised paper. Please refer to Overall Point 3.

R2.W2 *The baselines used in the comparative evaluation could be extended by methods for next-step recommendations. While such approaches do not focus on standardization, depending on the way the recommendations are generated, the resulting scripts may be more standardized than what the current baselines achieve.* **Ans.** Following the reviewer's comments, we included in our revised paper two additional baselines from the literature. Please refer to Overall Point 4 for more details.

R2.W3 *In places, the quality of the English text was sub-par, making it difficult to follow the technical discussion. I recommend carefully checking for spelling and grammatical errors.* **Ans.** We thank the reviewer for this comment. We thoroughly proofread the paper to eliminate any typos or errors.

R2.D1 *In the example used as motivation in the introduction, I do not fully understand the cause of the error (making the model only applicable to a specific country) and, more importantly, how the standardized script mitigates this problem. That is, how does the standardized script avoid the cause of the problem?* **Ans.** Following the reviewer's comment, we modified our motivating example to make it more intuitive and easier to understand. Please refer to our response in Overall Point 2 above.

R2.D2 *Section 1, "Our approach": "computes a similar result" -> similar to what? Similar in terms of what?* **Ans.** We have revised the phrasing in the introduction to clarify that by similarity, we refer to how well the intent of the original script is preserved in the modified script.

R2.D3 *End of Sec. 1: "We showed the system still obtains improvement." -> improvement in terms of what?* **Ans.** We thank the reviewer for this comment. We have clarified that improvement is measured according to our standardization measure. Additionally, we verify that this measure aligns with human judgment through a user study, as mentioned above in Overall Point 3.

R2.D4 *In the experimental setting, when describing the competing methods, please justify why you limit to the used systems. I was expecting some academic solutions as well. The reason is hinted at much later in the discussion. It would be better to state it early on.* **Ans.** In response to the reviewer's comment, we have added two state-of-the-art baselines for next-step recommendations in data preparation to the revised paper. For more details, please see Overall Point 4.

R2.D5 *Overall, the assumptions of a high-quality, clean corpus for the same task and dataset are quite strong. The relaxation performed in Section 6.3.4 only slightly relaxes the assumption on the same dataset to a dataset of the same schema. It would be interesting to see how the method performs as the assumptions are further relaxed.* **Ans.** Following the reviewers' comments, in our revised paper, we included an experiment that varies the corpus quality by using low-ranked reference scripts from Kaggle to evaluate how quality affects performance. Please refer to R1.W2. Our experiments demonstrate that even with a relatively small corpus, we can still extract valuable recommendations (Section 6.3.3). Additionally, in Section 6.3.3, we clarified that the two datasets do not share the same schema. Specifically, we show that even when the corpus contains scripts operating on different datasets with similar information but differing schemas, our system can still provide insightful recommendations to standardize the script.

Nevertheless, we agree with the reviewer that a corpus is still a necessary condition. However, we argue that in fields such as public policy and medicine, researchers worldwide often work with open-source datasets (e.g., Amazon reviews datasets, MIMIC datasets) along with publicly available open-source scripts (e.g., Paper With Code). A researcher new to these open-source datasets would find

our approach helpful in quickly identifying standard data preparation steps used by more experienced researchers. We have clarified this point in the Introduction.

Reviewer #3

R3.D1 *Example 1. In its current form it is difficult for the reader to fully understand the specific purpose behind each command. This issue also applies in Example 2. You need to add more details along with each command. Additionally, Table 1 should be introduced and described at the beginning of Example 1 in order to provide necessary context.* **Ans.** Following the reviewer’s comment, we have revised our running example and clarified the purpose of Table 1. Please refer to our response in Overall Point 2 above.

R3.D2 *Section 6, Paragraph 2: “Both programs compute a “similar” result.” -> somehow also mention the “user intent” notion.* **Ans.** Thank you for this comment. We clarified that similarity refers to the user intent.

R3.D3 *Minor Example 1: “Identify the missing values and impute them handle them..” -> fix syntax Figure 3: Reduce lines weight and/or change colors Figure 4 and 5: Increase size and/or change colors* **Ans.** Thank you for paying detailed attention to our work. We have updated Figure 4, Figure 5, and Figure 6 correspondingly.

Reviewer #4

R4.W1 *Application to real cases of this specific problem/setting can be of limited interest.* **Ans.** Please refer to our discussion in Overall Point 2 regarding the usability of the proposed system.

R4.W2 *Parametrization of the tool seems hard.* **Ans.** Please refer to Overall Point 1 regarding our discussion on parameter tuning.

R4.W3 *The solution changes the intent of the user.* **Ans.** Please refer to our answers for R1.D1 and R1.D2.

R4.D1 *The main problem I see in the work is to which extent standardization is interesting or useful in a real setting. The main point is to generate code that helps to explore a new dataset. However, the tool would be useless to any company doing that, because the corpus would not be available. Section 6.3.4 confirms this. Once the company (or someone else) has analyzed a dataset tens or hundreds of times, either it is already in production or maybe this dataset is not that interesting anymore..* **Ans.** We thank the reviewer for this important comment and agree that the tool is useless without a corpus. However, as demonstrated in Section 6.3.3, the tool remains useful when the corpus contains similar datasets, although the improvement may be reduced. Furthermore, we show that the tool retains its value even when the corpus is small or contains low-quality scripts (Section 6.3.3). As highlighted in Overall Point 2, we have expanded our discussion on the tool’s usability for widely used open-source real-world datasets where multiple scripts are available.

R4.D2 *Allowing to emit data that is not identical to those emitted by the original script sounds also a bit weird, specially if the user has not full control of this. This is also related to the high number of parameters of Algorithm 1, and the concrete consequences of their values. Despite the extensive experiments, I still would not know which values would be more convenient to my dataset if I had to use this tool.* **Ans.** As mentioned in Overall Point 2, the recommended modifications to the user script can be seen as recommendations, which the user

may decide to accept or reject. Furthermore, as mentioned above, to ensure that the modifications do not significantly alter the data or impact the performance of the downstream task, the user can specify a threshold to define the extent to which their intent should be preserved. Specifically, as noted in the introduction, we support two measures to evaluate how well the output script preserves the user’s intent: (1) a table distance measure, which compares the structured datasets generated by the scripts, and (2) a performance measure that evaluates downstream application outcomes, such as model accuracy or fairness. Regarding the parameters of the algorithm, please refer to our response in Overall Point 1 above.

R4.revisions *(1) Provide some summary table with clear hints on how to decide the value of each of the parameters of Algorithm 1, and which are their consequences. (2) You should clarify the meaning of vertical bars (i.e., |) in the formula of Table Jaccard, where it is used twice but seems that with different meaning. (3) Better to also define AST in the text (though it is clear in reference [1]) (4) Section 6.1.2 proposes to introduce some samples of the corpus in the prompt of GPT. Why not using the whole corpus in a specific training of an LLM? Wouldn’t this be possible? (5) “As Example 5 shown” (6) “... ensuring that that: ...”.* **Ans.**

- (1) Please refer to Overall Point 1.
- (2) Thank you for this point. It was a typo, which we fixed.
- (3) We added an explanation of what an AST is in the revised Section 3.
- (4) We agree with the reviewer that the performance of language models can be highly dependent on the problem prompt. We polled 12 graduate students and conducted a GPT prompt survey (including prompts that contain as many scripts in the corpus, within the token limit). We then empirically compared all the submitted prompts and selected the best performing prompt (details in Section 6.1.2).
- (5) Thank you for highlighting the typos. We have corrected them and thoroughly proofread the paper to ensure that no additional errors remain.

Toward Standardized Data Preparation: A Bottom-Up Approach

Eugenie Lai

MIT

Cambridge, Massachusetts, USA

eylai@mit.edu

Brit Youngmann

Technion

Haifa, Israel

brity@technion.ac.il

Yuze Lou

University of Michigan

Ann Arbor, Michigan, USA

yuzelou@umich.edu

Michael Cafarella

MIT

Cambridge, Massachusetts, USA

michjc@csail.mit.edu

Abstract

Data preparation is an important step in any data-related (R2&4) application¹, from academic research to industrial decision-making. Typically, data preparation is not a core contribution of a project — it transforms raw data into a format that supports further innovative work. However, the reality is that data scientists spend much of their time on data preparation. Because data preparation scripts are highly project-specific and often written in general-purpose languages, they are tedious to understand and difficult to verify. As a result, data preparation scripts can be a breeding ground for poor engineering and statistical practices, and even when they are perfectly built, they are difficult to reuse. Ideally, data preparation scripts should serve the project, but otherwise, be as simple and as standard as possible — adhering to the best common practices to process data. We propose a bottom-up script standardization framework that takes a user's data preparation script and transforms it into a more standardized version of itself. Our framework treats the user's input script as a semantic sketch of the user's intent, which can be modified even if the changed script yields a slightly different output. We present an algorithmic framework and developed a prototype system. (R2) We evaluated our approach against state-of-the-art methods, including GPT-4, on six real-world datasets. Our approach improved script standardization by 39.5% while meaningfully preserving the user's intent, while GPT-4 achieved 2.9%.

1 Introduction

Data preparation is a crucial step in various domains, from academic research to data-driven decision making in industry [63, 64]. Custom data preparation programs are often an essential part of data science pipelines, transforming raw data into a usable format for further innovative work. However, data preparation is often overlooked, and analysts rarely detail the process [14]. These “uninteresting” data preparation scripts serve solely to support downstream innovative efforts.

However, data preparation scripts remain extremely difficult to understand, verify, and reuse. Data preparation scripts are highly project-specific and often written in general-purpose programming languages. These attributes make them tedious to

understand and difficult to verify, without substantial effort. Consequently, data preparation scripts can be a breeding ground for poor engineering and statistical practices. For example, *variability* in data preparation was identified as one main reason for the high rate of false positive results in biomedical research [14]. Despite the widespread use of standard datasets, preparation scripts can be constructed to produce desirable experimental results, based on specific data values and characteristics, rendering scientific conclusions not reproducible [27, 59, 62].

In our view, data preparation scripts should be as standardized and boring as possible. The ideal data preparation script does its job and strives for standardization and predictability. The data management community has made many attempts to standardize data preparation, but only in a *top-down* fashion by proposing rules and systems [26, 30, 37, 38, 56, 57]. However, perhaps due to the richness of the project- and data-specific knowledge needed for data preparation, data scientists still resort to general-purpose programming languages [67] such as Python.

Intuitively, a script is *more standard* if it uses data preparation steps that are frequently used in other scripts processing the same or similar datasets. (R1&2) In this work, we adopt a common assumption in crowdsourcing [66]: the majority is presumably correct. Specifically, if a data preparation step appears repeatedly across multiple scripts processing the same dataset, it is likely to be both important and correct. Based on this assumption, we propose a *bottom-up standardization approach*, in which the user inputs a sketch script and gets a modified script as output. This output script aims to perform the same task as the input script, but is more similar to previously written scripts. Crucially, the modified script is allowed to emit data that are not identical to those emitted by the original script. (R1&4) The modifications to the input script can be viewed as recommendations for the user, who has the option to accept or reject them. The following example illustrates how bottom-up standardization facilitates data preparation for data scientists:

EXAMPLE 1. Alex, a data scientist hired to help a medical research team, is writing a data preparation script (Figure 1a) to transform a patient dataset [5]. (R1,2,3) It will be used to train a prediction model that accurately identifies diabetes in young adults aged 18-25. She examines previous data preparation scripts on Papers With Code [19], which were developed for other similar projects. (R1) Table 3 outlines the data preparation steps used in Alex's initial script (s_u), in three example scripts from the corpus (s_1, s_2, s_3), and in

¹Updates made in response to specific comments are highlighted using the following colors: Changes for multiple reviewers in orange; Reviewer 1 (R1) in blue; Reviewer 2 (R2) in magenta; Reviewer 3 (R3) in purple; Reviewer 4 (R4) in green.

```

1. import pandas as pd
2. df = read_csv('diabetes.csv')
3. df = df.fillna(df.median())
4. df = df[df['Age'].between(18,25)]
5. df = pd.get_dummies(df)
6. X = df.drop(['Outcome', axis = 1])
7. y = df['Outcome']

```

(a) Input script.

```

1. import pandas as pd
2. df = read_csv('diabetes.csv')
3. df = df.fillna(df.mean())
4. df = df[df['Age'].between(18,25)]
5. df = df[df['SkinThickness'] < 80]
6. df = pd.get_dummies(df)
7. X = df.drop(['Outcome', axis = 1])
8. y = df['Outcome']

```

(b) Output script.

Figure 1: (R1,2,3) Example input and output scripts.

the modified output script (\hat{s}_u). Alex quickly feels discouraged due to the volume and required domain knowledge. Consequently, she writes a script from scratch. Specifically, she: (1) **Imputes missing values using the median value** (`df.fillna(df.median())`), based on her experience with other projects (Table 1, a_3); (2) Selects only relevant records (`df[df['Age'].between(18,25)]`), based on her modeling objective (Table 1, a_4).

Every step in Alex’s script is seemingly correct, but unfortunately, she fails to consider the common practices to process this dataset. (R1,2,3) **The resulting downstream model is less accurate because outliers in the data were not addressed, leading to distractions that reduce the model’s accuracy. The error, which arises during data preparation due to a lack of domain-specific knowledge that Alex does not possess, is unlikely to be noticed by readers of Alex’s paper or experimental results.** □

Now imagine that Alex uses a *script-standardization* system:

EXAMPLE 2. Alex is writing a data preparation script from scratch, before giving it to a script standardization system. Then she sets a 5% threshold for the maximum % change that the standardization system is allowed to implement. In this case, the value of 5% refers to the maximum change in predictive accuracy allowed in the downstream model. The system responds with a standardized script (Figure 1b) that does the following:

- (1) Recognizes that (R1,2,3) **imputing missing values with average is a more common practice than median imputation for Alex’s dataset** (Table 1, a_3 removed and a_5 added);
- (2) Selects only the relevant patient records, to maintain Alex’s modeling objective (Table 1, a_4 unchanged);
- (3) **Adds additional step to handle outliers using domain-specific knowledge² embedded in other scripts** (`df[df["SkinThickness"] < 80]` – Table 1, a_6 added).

Alex sees the modified output script, and decides to accept all the recommendation changes. Her improved script now uses more standard transformations. She is also pleased that her prediction accuracy has increased slightly due to the changes (although this will not happen in all cases). □

²Triceps skinfold thickness is a measurement used to estimate body fat by assessing the thickness of a fold of skin. A value above 80mm is considered abnormal [44].

Table 1: (R1,2,3) An illustration showing the steps in three historical scripts (s_1, s_2, s_3), the input script s_u , and the output script \hat{s}_u . Green shows added steps, and red shows removed steps.

	Data preparation step	s_u	s_1	s_2	s_3	\hat{s}_u
a_1	import pandas as pd	✓	✓	✓	✓	✓
a_2	df = pd.read_csv('diabetes.csv')	✓	✓	✓	✓	✓
a_3	df = df.fillna(df.median())	✓				
a_4	df = df[df["Age"].between(18,25)]	✓				✓
a_5	df = df.fillna(df.mean())		✓		✓	✓
a_6	df = df[df["SkinThickness"] < 80]		✓	✓	✓	✓
a_7	df = pd.get_dummies(df)	✓	✓	✓	✓	✓

Our Approach: (R2) The goal of our approach is to modify an input script to produce a new script **similar in intent** to the input script, while being more **standard** (i.e., aligning with common data preparation steps from prior scripts). The user begins with a draft program, which our framework treats as a **semantic sketch** of their intent. We adjust the semantics of the new script to balance two objectives: increasing its standardization and maintaining an **agreement with the user’s intent**. A straightforward approach is to leverage powerful code cleaning tools such as Sourcery [7], LLMs such as GPT [51], (R2) **or recommendation systems predicting the next-best data preparation steps such as [43, 63].** However, our experimental results show that while these tools can generate code, they often fail to capture the user’s original intent and struggle to select data preparation steps that are particularly relevant to the dataset, as informed by domain knowledge embedded in the corpus. We argue that to obtain such a system, there are several challenges to overcome: (C1) Defining how to capture script standardization in relation to a given corpus of scripts. (C2) Quantifying how well the user intent is preserved in the modified script. (C3) Navigating the vast search space (which is exponential in the number of data preparation operations present in the corpus) of potential modifications to the input script to find a sequence of changes that produces an executable script, maintains the user’s intent, and achieves the highest degree of standardization.

To address C1, we define the standardness of a script w.r.t. a corpus of scripts. We assume the availability of a corpus of scripts processing the same or similar datasets for different objectives (as done in other frameworks [16, 63]). Our experiments show that even a small corpus yields valuable results. Our optimization objective is to minimize the relative entropy [40] of the data preparation step distribution between the input draft and the script corpus. Relative entropy is a common measure of how one probability distribution (the input script) diverges from a second, expected probability distribution (the scripts in the corpus). We extract these distributions using a novel graph-based script representation, leveraging NLP techniques such as lemmatization [54] and n-grams.

To address C2, we need to determine whether the input and output scripts perform similar tasks by examining their outputs or the results of their downstream applications. (R1) **We use two measures to assess how well the output script preserves the user’s intent and that the modifications to the input script do not significantly alter the data or impact the performance of the downstream task:** (1) a table distance measure, which compares the structured datasets generated by the input and output scripts, and (2) a performance

measure that evaluates downstream application performance, such as model accuracy or fairness.

User intent is treated as a constraint on the similarity of the input script and the modified script: the user provides a threshold indicating how much they are willing to compromise on their intent. The objective is to find a script that stays within this threshold of similarity to the original script, remains executable, and is as standardized as possible (i.e., minimizing the relative entropy).

To address C3, we need an efficient way to find a sequence of modifications to the input script, as an exhaustive search is impractical given the large search space (comprising all combinations of data preparation steps that are present in the corpus), and even greedy approaches present challenges. Our method has two phases: during the offline phase, we represent each script in the corpus as a graph and curate the search space; during the online phase, we employ a search framework with five algorithmic components (including diversity, sampling, and early-checking criteria) to explore enough of the search space while meeting the constraints.

(R4) Usability: The proposed system is particularly useful in scenarios where users are working with datasets that are commonly used by other researchers, such as in public policy and medicine. For instance, datasets such as Amazon reviews³ or MIMIC⁴ are widely used by researchers and are often accompanied by open-source data preparation scripts available online (e.g., on platforms such as Papers With Code⁵). A researcher new to these open-source datasets can benefit from our approach by quickly identifying standard data preparation steps used by other researchers. As we show experimentally, the corpus needed to gain valuable insights can be relatively small, making the system practical for real-life scenarios.

A demonstration of our system usability and its suitability for end-to-end employment was recently shown in [41]. The short paper accompanying the demonstration provides only a brief, high-level description of the system, whereas the present paper provides the theoretical foundations and algorithms underlying the demonstrated system, as well as the experimental study.

Contribution — Our main contributions are as follows:

- We define a new problem of script standardization that helps data consumers write standardized data preparation scripts while preserving their intent (Section 4).
- We present an efficient and effective search framework and apply optimization techniques (Section 5).
- We evaluated our prototype system, LucidScript, against state-of-the-art approaches [6, 7, 15, 51, 55] on six real-world datasets (Section 6). **(R2) Improvement is measured using our standardization measure, and we verified that this measure aligns with human judgment through a user study.** Our approach obtained a 39.5% improvement on script standardness while, **the second best competitor, GPT-4, achieved 2.9% (Table 5).** **(R1&2) We also evaluated our system in scenarios when high-quality on-topic script corpora were not available.** We tested it using a smaller corpus, as well as using a corpus from a different dataset. Our approach still obtained improvements on script standardness.

2 Script Standardization

A data preparation script is a sequence of lines of code that process some dataset for a downstream task, such as training a model, generating visualizations, etc. Let $\mathcal{S} = \{s_1, s_2, \dots, s_n\}$ be a collection of n scripts (referred to as the corpus). Let s_u denote an input user script. We assume that s_u and all scripts in \mathcal{S} process the same dataset D_{IN} . In Section 6, we will empirically show that this assumption can be relaxed. For simplicity, throughout the paper we assume that D_{IN} comprises a single data file. Applying a script s to D_{IN} yields an output dataset D_{OUT}^s . The goal of our system is to modify the input script s_u to yield a new script \hat{s}_u such that they both accept D_{IN} as input, and: (1) Both programs compute a **"similar"** result, **(R3) measured by how well the user intent is preserved;** (2) \hat{s}_u is a more **conventional** program than s_u , according to statistics computed over the corpus \mathcal{S} .

2.1 Measuring Goal One: User Intent

Measurement of the semantic similarity of two scripts is an extensively studied topic, often measured by the output of the operation [11, 48, 68]. In this work, we measure whether two scripts are doing something similar by looking at their outputs, or the outputs of their downstream application. We use these measures to assess how well \hat{s}_u preserves the user's original intent embodied in s_u .

Table Jaccard: An example measure is table distance, which focuses on a structured dataset emitted by the script. This would be appropriate when the user produces databases for visualization, OLAP pipelines, or machine learning training procedures. Jaccard index [46] is commonly used to measure how close the new script's output dataset $D_{OUT}^{\hat{s}_u}$ is to the original user script's output dataset $D_{OUT}^{s_u}$. Formally, let $\Delta_J(D_{OUT}^{s_u}, D_{OUT}^{\hat{s}_u})$ denote the table Jaccard similarity defined as:

$$\Delta_J(D_{OUT}^{s_u}, D_{OUT}^{\hat{s}_u}) = \frac{|D_{OUT}^{s_u} \cap D_{OUT}^{\hat{s}_u}|}{|D_{OUT}^{s_u} \cup D_{OUT}^{\hat{s}_u}|}$$

$\Delta_J(D_{OUT}^{s_u}, D_{OUT}^{\hat{s}_u})$ ranges from 0 to 1, where $\Delta_J = 1$ indicates that the two tables are identical.

EXAMPLE 3. Script \hat{s}_u has an additional step to normalize all strings to lowercase, which is not present in s_u . $|D_{OUT}^{s_u}| = \{\text{'benign'}, \text{'Benign'}, \text{'High Risk'}, \text{'High risk'}, \text{'high risk'}\}$ and $|D_{OUT}^{\hat{s}_u}| = \{\text{'benign'}, \text{'high risk'}\}$. Table jaccard would be $\Delta_J(D_{OUT}^{s_u}, D_{OUT}^{\hat{s}_u}) = \frac{2}{5} = 0.4$. □

Model Performance: Another example measure of user intent is the performance of the downstream data application, which is appropriate when we have a downstream task with an easy-to-measure quality metric, such as accuracy.

Formally, let $\Delta_M(D_{OUT}^{s_u}, D_{OUT}^{\hat{s}_u})$ denote the absolute value of the relative percentage change in model accuracy, which is in $[0\%, 100\%]$. $\Delta_M = 0\%$ means the two datasets are identical.

$$\Delta_M(D_{OUT}^{s_u}, D_{OUT}^{\hat{s}_u}) = \left| \frac{\text{acc}(D_{OUT}^{s_u}) - \text{acc}(D_{OUT}^{\hat{s}_u})}{\text{acc}(D_{OUT}^{s_u})} \right| \times 100\%$$

³<https://jmcauley.ucsd.edu/data/amazon/>

⁴<https://archive.physionet.org/physiobank/database/mimicdb/>

⁵<https://paperswithcode.com/>

EXAMPLE 4. The model accuracy for s_u and its \hat{s}_u is $acc(D_{OUT}^{s_u}) = 0.65$ and $acc(D_{OUT}^{\hat{s}_u}) = 0.67$, respectively. The model performance measure would be $\left| \frac{0.65-0.67}{0.65} \right| \times 100\% = 3.1\%$. \square

Our framework can be generalized to support other user-intent measures, such as earth mover distance or fairness constraints [29], as discussed in Section 8. We chose table Jaccard and model performance as example measures since they are two of the most straightforward metrics to understand.

2.2 Measuring Goal Two: Standardization

Our overall goal is to create new versions of the input script that are as standard and conventional as possible. We use *relative entropy* (also known as the Kullback–Leibler divergence) [40] to measure how standard a script is. Relative entropy is a non-negative statistical distance of how one probability distribution $P(x)$ is different from a reference probability distribution $Q(x)$. We adapt the notation of relative entropy to measure script standardness. Formally, the relative entropy of a script s w.r.t. a corpus \mathcal{S} is defined as:

$$RE(s, \mathcal{S}) = \sum_{x \in \mathcal{X}} P(x) \log\left(\frac{P(x)}{Q(x)}\right)$$

where \mathcal{X} is the sample space shared by both P and Q .

In our problem, \mathcal{X} is the space of all the data preparation steps in script s , while $Q(x)$ is the probability distribution of the steps in corpus \mathcal{S} . A large distance between $P(x)$ and $Q(x)$ would indicate that s uses many data preparation steps that are not commonly used in \mathcal{S} , while a small distance would mean that s does a good job adhering to the standard practices in \mathcal{S} .

Formally defining $P(x)$ and $Q(x)$ is a non-trivial task since scripts need to be in a format that makes it easy to compute these statistics. We describe such a representation in Section 3. We then use $RE(s, \mathcal{S})$ as our optimization objective in Section 4.

Problem Objective. The *script standardization* problem is described by a tuple comprising: (i) A **script corpus** $\mathcal{S} = \{s_1, \dots, s_n\}$; (ii) An **input user script** s_u ; (iii) An **input dataset** D_{IN} ; (iv) A **user-intent parameter** $\Delta(D_{OUT}^{s_u}, D_{OUT}^{\hat{s}_u}) \leq \tau$. The solution is a new, executable script \hat{s}_u that maximizes the standardness measured by $RE(\hat{s}_u, \mathcal{S})$ while ensuring that the output of the resulting script is within τ of the original.

3 Script Representations

We face a multifold challenge — to represent an arbitrary script while also being able to explore the space of legal changes to it. Concretely, the representation should enable us to compute standardization statistics efficiently and numerate the space of legal changes. It should abstract away inessential syntactic differences. Lastly, it should retain sufficient information about the script so that it can be translated back to a script. Inspired by previous work [28, 39, 47], we present a directed acyclic graph (DAG) representation of scripts that satisfies the above requirements for our setting.

We now introduce how we represent the entire search space, including the DAG representation of scripts and transformations that can be applied to DAGs. In Section 4, we will describe how

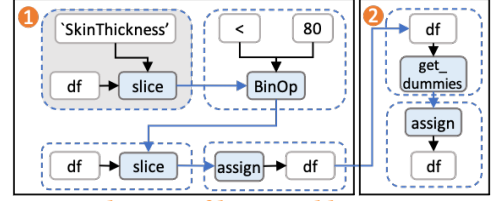


Figure 2: Example DAG of line 5 and line 6 in Figure 1b. The DAG is made of atoms (blue dashed) and edges (blue). Inside an atom, there are AST nodes and edges, where invocation nodes are blue, and data nodes are white. Arrows represent data flows. We color one atom gray to illustrate Example 5.

our script representations are used to compute the probability of observing a particular script s .

DAG Representation: We start by representing a single script using a DAG, in which nodes represent operation invocations, and edges represent data flows. The DAG representation $G_s = (A, E')$ of a script s is based on its Abstract Syntax Trees (ASTs). (R4) An AST is a tree-like data structure used to represent the syntactic structure of a code snippet in a hierarchical and abstract way. Each node corresponds to a construct in the code, such as an operator, a variable, or a control structure, and its children represent the components or sub-expressions of that construct. ASTs focus on the logical relationships and structure of the code, and are widely used in compilers and tools for code analysis, transformation, and optimization. For a formal definition of ASTs for Python, see [1].

A simple script can result in hundreds of AST nodes, and most nodes make sense only in the scope of their context. If we curate the search space directly on the AST nodes without considering their neighbors, the space is huge, and most transformations would be invalid. For example, for line 5 in Figure 1b, replacing the AST node `80` with `'SkinThickness'` would result in an execution error since `<` can only compare numerical values.

Inspired by the natural language processing (NLP) applications [15, 17, 22, 55], where atoms are defined at the word level, we construct atoms at the operation-invocation level. The DAG representation $G_s = (A, E')$ of a script s is made of two components, atoms A and edges E' (Figure 2). An atom a , or atomic unit, in our DAG representation, is an operation invocation. In programming languages, an operation invocation refers to a snippet of code required to call a function, with its arguments. Instead of AST nodes, the use of atoms helps to make the output scripts free of syntax and execution errors and also helps to reduce the search space. Thus, we define each atom $a = (V, E)$ as a collection of AST nodes V and edges E . We simplify AST nodes to two types, data nodes, and invocation nodes (Figure 2). We later perform lemmatization techniques on atoms, which helps us to further control the size of the vocabulary and leverage the semantics of the operator invocations (Section 5).

DEFINITION 1. [Atom and DAG] An *atom*, is the atomic unit of the DAG representation. Using ASTs, we denote an atom by $a = (V, E)$. We can express the DAG $G_s = (A, E')$, with atoms A and edges between atoms E' . Specifically, an atom a consists of one invocation node and its parents that are not invocation nodes. \square

EXAMPLE 5. (All) An atom is an operation invocation. In Figure 2, the right block shows the DAG representation of `df = df[df['SkinThickness'] < 80]`. Specifically, the gray-color atom represents `df['SkinThickness']`, where `df` and `'SkinThickness'`

are AST data nodes and slice is an AST invocation node. The black arrows are the edges E inside the atom, representing data flows. \square

Atoms can be used in 1-gram, n-gram, or in a way learned from the script collection \mathcal{S} . We use atoms in two ways, at the operation-invocation level (1-gram) and the line level (n-gram). In Figure 2, the gray-color atom is an example of 1-gram atoms, while the numbered blocks are examples of n-gram atoms.

From Script to DAG: Some components in the standardization measure of script are now more clear. \mathcal{X} denotes all the data preparation steps in \mathcal{S} . As Figure 2 shows, using the DAG representation, the data preparation steps are decomposed into atoms A and edges E' . Formally, $\mathcal{G} = \{G_{s_1}, \dots, G_{s_n}\}$ is the DAG representation of $\mathcal{S} = \{s_1, \dots, s_n\}$. Denote by \mathcal{V}_A the vocabulary that contains unique atoms $A \in \mathcal{G}$ and by $\mathcal{V}_{E'}$ the vocabulary that contains unique edges $E' \in \mathcal{G}$. We use $\mathcal{V}_{E'}$ to model \mathcal{X} , instead of \mathcal{V}_A , since E' encodes the data flow information, representing the order of the steps.

EXAMPLE 6. (All) Figure 2 shows the DAG representation of $\text{df} = \text{df}[\text{df}[\text{'SkinThickness'}] < 80]$ and $\text{df} = \text{pd.get_dummies}(\text{df})$, which are lines 5-6 in Figure 1b. The data flow between lines 5-6 is represented by the blue edge e' between the numbered blocks, connecting two data nodes. \square

Transformations: Transformations are the actions we perform to change the DAGs. We define two types of transformations, add and delete.⁶ Each transformation takes the following parameters: type, what to change (i.e., one atom and its edges), and where to change (i.e., line number). Next, we formally define the notation of a transformation over G_s .

DEFINITION 2. [Transformation] A transformation is a function of the form $f(\text{type}, a, \{e'_1, \dots, e'_n\}, \text{lineno})$, where $\text{type} \in [\text{add}, \text{delete}]$, $a \in \mathcal{V}_A$, $e'_i \in \mathcal{V}_{E'}$. $\{e'_1, \dots, e'_n\}$ is the set of edges that connect an atom a to its neighbors. $F(G_s) = (f_1, f_2, \dots, f_m)$ is sequence of m transformations that preform on G_s . \square

We discuss how we configure transformations in Section 5.2.

4 Problem Formulation

We have defined the measure of script standardization (Section 2.2) and presented script representations (Definition 1), which help us to explore the space of transformations (Definition 2) that can be applied to the input script. We are now ready to present our standardization objective and the script standardization problem.

DEFINITION 3. [Standardization Objective] The relative entropy of a script s_u w.r.t. a collection of scripts \mathcal{S} .

$$RE(s_u, \mathcal{S}) = \sum_{\mathcal{X}} P(x) \log\left(\frac{P(x)}{Q(x)}\right)$$

where $P(x) = \frac{x}{\sum_{i=1}^{|\mathcal{V}_{E'}|} x_i}$, and $x_i = \sum_{j=1}^{|\mathcal{V}_{E'}|} [e_j = v_i]$, $\forall e_j \in G_{s_u}$.

Similarly, $Q(x) = \frac{x}{\sum_{i=1}^{|\mathcal{V}_{E'}|} x_i}$, and $x_i = \sum_{j=1}^{|\mathcal{V}_{E'}|} [e_j = v_i]$, $\forall e_j \in \mathcal{G}$. \square

$P(x)$ is the probability distribution of the data preparation steps in script s_u , while $Q(x)$ is the probability distribution of the data preparation steps in the script corpus \mathcal{S} (Section 2.2). We use $\mathcal{V}_{E'}$

⁶The edit operation is modeled here as a sequence of delete and add operations.

to derive the model of the data preparation steps \mathcal{X} . We use atoms in Table 1 to illustrate how we compute $P(x)$ and $Q(x)$.

EXAMPLE 7. The DAG of an input script s_u is $G_{s_u} = (A_{s_u}, E'_{s_u})$, where $E'_{s_u} = [(a_0, a_1), (a_1, a_7)]$ are the edges of G_{s_u} . $\mathcal{V}_{E'}$ is the edge count mapping from \mathcal{G} , where $\mathcal{V}_{E'} = \{(a_0, a_1) : 3, (a_1, a_2) : 3, (a_2, a_7) : 2, (a_1, a_7) : 1\}$. For example, the edge $e' = (a_0, a_1)$ appears three times in \mathcal{G} , where $\mathcal{G} = \{G_{s_1}, G_{s_2}, \dots, G_{s_n}\}$ is the DAG representation of the script corpus $\mathcal{S} = \{s_1, s_2, \dots, s_n\}$. \square

We derive the distribution of the data preparation steps in the corpus $Q(x)$ using edge vocabulary $\mathcal{V}_{E'}$ and \mathcal{G} . $Q(x)$ and x are both vectors with a size of $|\mathcal{V}_{E'}|$. We compute x as the number of occurrences for each edge in \mathcal{G} . Then we sum up the occurrences of each edge and divide x by the total occurrences to get $Q(x)$.

EXAMPLE 8. In $Q(x)$, x is derived from $\mathcal{V}_{E'}$, where $x = [3, 3, 2, 1]$ since the number of unique edges $|\mathcal{V}_{E'}|$ is four. The total number of edges in \mathcal{G} is $3 + 3 + 2 + 1 = 9$ so $Q(x) = [\frac{1}{3}, \frac{1}{3}, \frac{2}{9}, \frac{1}{9}]$. \square

Then we compute $P(x)$ from $G_{s_u} = (A_{s_u}, E'_{s_u})$. Similarly, $P(x)$ and x are vectors of the length of the edge vocabulary $|\mathcal{V}_{E'}|$. We compute x as the number of occurrences for each edge in G_{s_u} , where $E' \in \mathcal{V}_{E'}$.

EXAMPLE 9. In $P(x)$, $x = [1, 0, 0, 1]$. The total number of edges in G_{s_u} is $1 + 1 = 2$, so $P(x) = [\frac{1}{2}, 0, 0, \frac{1}{2}]$ and $RE(s_u, \mathcal{S}) = 1.38$. \square

We note that $RE(s_u, \mathcal{S})$ accounts for atoms and edges, but is oblivious to their exact position (i.e., line number) in the script. We explain how we handle this in Section 5.

Revisiting the conceptual definition in Section 2.2, we are now ready to formally define the problem using the script representations in Section 3. The Script Standardization problem aims to find a sequence of transformations F for an input script s_u to make it as standard as possible w.r.t. a collection of scripts \mathcal{S} , (R4) **while ensuring that**: (1) There are no execution errors; (2) No more than k transformations are applied; (3) User intent is preserved. Formally,

DEFINITION 4. [Script Standardization] Given a script s_u and a collection of scripts \mathcal{S} all operating on the data file D_{IN} , a number $k \geq 1$, and a threshold τ we search for a script \hat{s} such that:

- (1) **(sequence length constraint)** \hat{s}_u is obtained by sequentially applying no more than k transformations on s_u .
- (2) **(execution constraint)** \hat{s}_u is executable.
- (3) **(user-intent constraint)** $\Delta(D_{OUT}^{s_u}, D_{OUT}^{\hat{s}_u}) \leq \tau$.
- (4) **(objective function)** $RE(\hat{s}_u, \mathcal{S})$ is minimized. \square

Each transformation f increases or decreases $RE(\hat{s}_u, \mathcal{S})$.

EXAMPLE 10. Continuing with the simple setup in Example 7, the best transformation f is trivial: add a_2 between a_1 and a_7 . This f leads to $P(x) = [\frac{1}{3}, \frac{1}{3}, \frac{1}{3}, 0]$, which gives $RE(\hat{s}_u, \mathcal{S}) = 0.2$. \square

$RE(\hat{s}_u, \mathcal{S})$ is much smaller than the starting $RE(s_u, \mathcal{S})$ in Example 9, indicating that $P(x)$ has come closer to $Q(x)$. This means that the resulting script is more standard w.r.t. the script corpus \mathcal{S} .

To find the best transformations to change s_u , we need to find the sequence that minimizes $RE(\hat{s}_u, \mathcal{S})$. However, the search space can become intractable quickly, as there could be many possible places to apply a transformation in s_u . For example, a transformation to

append a new atom can have multiple possible existing atoms to append to. In addition, possible next-step transformations depend on the current step. As a result, an exhaustive search algorithm is exponential in $(|\mathcal{V}_A| \times |A_{s_u}| + |E'_{s_u}|)$.

Therefore, in the next section, we will present an efficient algorithm for our script standardization problem. This algorithm overcomes the obstacles of the naive exhaustive search algorithm by taking a greedy-like approach while ensuring that a diverse set of executable candidate scripts is taken into consideration.

5 Methodology

5.1 Offline Search Space Curation

In the offline phase, we construct the atom vocabulary \mathcal{V}_A and the edge vocabulary \mathcal{V}_E and obtain the corpus distribution. These components then serve as inputs to the online phase. We take the following steps to construct them from the corpus \mathcal{S} : (1) We parse each script $s_i \in \mathcal{S}$ into its corresponding DAG G_{s_i} , simultaneously building \mathcal{V}_A and \mathcal{V}_E through this process. (2) We compute the probability distribution $Q(x)$ from \mathcal{V}_E , as described in Section 3.

Reducing Vocabulary: A single data preparation step can be represented through various coding operations that convey an equivalent semantic meaning. For instance, when both `df` and `train` are read from the same CSV file, expressions like `df['Age']` and `train['Age']` refer to the same column. We address such semantically equivalent steps by leveraging a rich body of work in NLP [54] and static code lemmatization techniques. For example, we unify the variables created by reading from the same data file to a consistent name across all scripts. Lemmatization significantly reduces the size of the vocabulary and helps the system better understand the semantics of scripts without actually executing them.

The semantic similarity of two scripts is an extensively studied topic in programming languages [11, 48, 68]. LLMs [6, 51] have also been proven to excel at such tasks. More optimizations can be extended without changing other system components.

5.2 Online Phase Search Framework

The next goal is to find a sequence of transformations that standardizes a given script s_u while satisfying the constraints (Definition 4). We model this sequence searching task as a constrained optimization problem with the objective to minimize its relative entropy measured w.r.t. the corpus $RE(s_u, \mathcal{S})$.

Recall in Section 4, a simple greedy approach faces two challenges. First, the optimal sequence may not always score the best during the search process. Discarding all but the best in-progress sequence may eliminate many potentially good candidates. Second, the output script \hat{s}_u must satisfy the constraints, and checking for the constraints only at the end may result in invalid sequences (e.g., execution errors, failure to satisfy user-intent constraints). We propose five optimizations that overcome the limitations and evaluate the impact of these components on our framework (Section 6).

(1) Beam search: We keep K beams rather than a single best, which retains multiple best options during the search.

(2) Diversity measure: We consider transformations that are different enough every time we add one step to the beams, which explores different parts of the search space.

(3) Monotonicity: A transformation sequence cannot go back and change an earlier portion of the script, denoted by the line number. This property ensures that once a script becomes non-executable after a transformation, no further transformation would make the script executable again. This helps reduce the search space.

(4) Checking strategies for constraints: We have three constraints. (1) We use the sequence-length constraint as the stopping criterion to limit the number of modifications applied to the input script. (2) For the execution constraint, we have an early- and late-checking strategy. In early-checking, we remove candidate sequences that lead to non-executable scripts after every transformation is applied. In late-checking, we only check for execution errors after the sequence is fully developed. (3) We check the user-intent constraint at the end. These strategies ensure that our approach always finds valid sequences while keeping runtime in check.

(5) Sampling: When D_{IN} is large, we apply random sampling on the tuples. This helps reduce runtime.

Algorithm 1 is our search framework, which takes as input a script collection \mathcal{S} , a script s_u , beam size K , user-intent threshold τ , and an early-checking switch α and returns the standardized script \hat{s}_u . In line 1, the candidate script set is initiated with the input script. In line 2, the vocabularies and distribution of the data preparation step are computed. In lines 3-8, beam search is applied to find the top K scripts in each transformation step, until the stopping criteria are met. In lines 9-11, constraints are verified and the best result is returned.

In Algorithm 1, `GetSteps()` takes a script s and the search space \mathcal{V}_A , \mathcal{V}_E , and $Q(x)$ as input, and outputs a list of transformations ranked based on the RE score. This step can be divided into three parts: (1) Translate s to its DAG representation $G_s = (A, E')$. (2) Given \mathcal{V}_A , \mathcal{V}_E , enumerate through G_s to find possible next-step transformations. (3) For each possible transformation, compute the RE score and determine where in s to perform each transformation.

Algorithm 1: A meta-level framework

```

Input:  $\mathcal{S}$  (script collection),  $s_u$  (input script),  $seq$  (sequence length),
         $K$  (beam size),  $\tau$  (user-intent threshold),  $\alpha$  (early checking)
Output:  $\hat{s}_u$  (a standardized version of  $s_u$ )
1  $C = \{s_u\}$ ; /* Set of all candidates. */
2  $\mathcal{V}_A, \mathcal{V}_E, Q(x) = \text{CurateSearchSpace}(\mathcal{S})$ ;
3 while not CheckStoppingCriteria( $C, seq$ ) do
4    $C' = \emptyset$ ; /* Set of new candidates. */
5   foreach  $s \in C$  do
6      $\mathcal{F} = \text{GetSteps}(s, \mathcal{V}_A, \mathcal{V}_E, Q(x))$ ; /* A ranked set of
7       next steps based on RE score. */
8      $C' = \text{GetTopKBeams}(C', s, K, \alpha, \mathcal{F})$ ;
9    $C = C'$ 
10  $C = \text{VerifyAllConstraints}(C, \tau)$ ;
11  $\hat{s}_u = \text{GetTopCandidate}(C)$ ;
12 return  $\hat{s}_u$ ;

```

Configuring Transformations: Our approach uses two types of transformations, add and delete, with two types of atom to consider, 1-gram and n-gram (Section 4). To generate all possible transformations, we enumerate the combination of which atom to change and where to apply the change. Configuring *delete transformations* is straightforward. We simply detect all existing atoms in the script and configure a list of delete transformations with the atoms, edges, and positions. Configuring *add transformations* is more elaborate

since the location to add is unknown. For 1-gram atoms, we leverage the edges to find the possible locations to append. For every atom a in the script s , a new atom $a' \in \mathcal{V}_A$ is considered a possible candidate to be appended after a if there exists (a, a') in \mathcal{V}_E . For n -gram atoms, we use the script corpus \mathcal{S} to find possible locations in s . When curating the search space, we retain the relative location of each n -gram atom in \mathcal{S} . When determining where to add a new n -gram atom, we compute their possible positions in s using the relative locations observed in \mathcal{S} .

After a transformation is configured, its impact on the RE score is calculated by marginally updating its vector $P(x)$ instead of performing the actual transformation in DAG. We have yet to describe `GetTopKBeams()`, where we extend each candidate sequence by one transformation. We propose a few search strategies that will be discussed next.

5.3 Extending By One Transformation

Given the list of configured transformations, we have two challenges when deciding which transformation to append to a sequence: (1) The search space is still large due to the exponential number of possible sequences and the large number of unique atoms in \mathcal{S} ; (2) The validity of the execution and user-intent constraints of a candidate sequence is difficult to estimate during search.

Algorithm 2: GetTopKBeams

Input: C (current candidates), s (a script), K (beam size), α (early checking), \mathcal{F} (a ranked set of next steps)
Output: C' (updated candidates)

```

1 valid = True;
2  $C' = C$ ;
3 foreach  $f \in \mathcal{F}$  do
4    $s' = \text{ApplyStep}(s, f)$ ;
5   if  $\alpha$  then valid = CheckIfExecutes( $s'$ );
6   if Score( $s'$ ) < min $s \in C$  Score( $s$ ) or  $|C| \leq K$  then
7     if valid then  $C' = \text{UpdateKBeams}(C', s', K)$ ;
8 return  $C'$ ;
```

We design Algorithm 2 that applies the concepts discussed in Section 5.2 and addresses those challenges with the following. In line 4, transformation f is applied to G_s . In line 5, the transformed script s' is verified for the execution constraint. In lines 6-7, the top K beams are updated by swapping out the script with the highest RE score in C' with s' . The beam search in our search strategy keeps track of multiple candidate sequences with potential best RE scores while keeping the search space tractable. This algorithm also takes a parameter α that controls the early checking for execution constraints. When $\alpha = \text{True}$, early-checking is *on*, which means that the execution constraint is verified every time after a transformation f is applied to G_s . This ensures that the algorithm retains beams that are free of execution errors during the search.

A problem with beam search is that the top-ranked next transformations can be similar. This is because the $RE(s', \mathcal{S})$ score compares the atoms in s' and \mathcal{S} so the high-ranked transformations would suggest changing the same atom. We then design an alternative to Algorithm 2 that incorporates a diversity measure to allow the beams to explore different transformations.

Transformation Diversity Measure: We use K-means clustering [31] to group similar transformations. `ClusterSteps()` takes a

ranked list of transformations, output by `GetSteps()`. `ClusterSteps()` then takes the updated vectors and groups them into M clusters. Within each cluster, we then rank the transformations according to the RE scores and the output \mathcal{F}_M . In `GetDiverseTopKBeams()`, we iterate through each cluster and call Algorithm 2.

Algorithm 3: GetDiverseTopKBeams

Input: C (current candidates), s (a script), K (beam size), α (early checking), \mathcal{F} (a ranked set of next steps)
Output: C' (updated candidates)

```

1  $C' = C$ ;
2  $\mathcal{F}_M = \text{ClusterSteps}(\mathcal{F})$ ;
3 foreach  $\mathcal{F}_{i \in M}$  do
4    $C' = \text{GetTopKBeams}(C', s', \frac{K}{M}, \alpha, \mathcal{F})$ ;
5 return  $C'$ ;
```

(R4) Parameterizing our framework: Algorithm 1 provides the user with a set of parameters to explore and tailor in a flexible way to their use case. Based on our main experiment (Section 6.3.3) and ablation studies (Section 6.4) on six real-world datasets, Table 2 lists the possible default parameters based on different corpus properties. In addition, the user-intent threshold allows the user to control the difference between the data processed by s_u and \hat{s}_u , which is set to $\tau_J = 0.9$ and $\tau_M = 1\%$ by default. We explore how varying the user-intent threshold affects our framework in Section 6.3.2. Lastly, early-checking is set to *on* as default to ensure the output script is executable. We discuss potential ways to automatically tune these parameters in Section 8.

6 Experimental Evaluation

We embody our algorithm in a prototype system, LucidScript (LS), which was written in 3,000 lines of predominantly Python source code. Our current implementation supports straight-line Python scripts. Support for other programming languages and more complex programs can be addressed with more engineering efforts. Our experiment data, code, and artifacts are available [2].

We empirically demonstrate the following claims about our method as embodied in LS:

- (1) An increase in our metrics (relative entropy and user intent) leads to better outcomes for data scientists.
- (2) (All) Our generated scripts are considered to be more standard than scripts generated by alternative approaches.
- (3) LS performs better than competing methods in standardizing data preparation scripts with the user intent. We evaluated this on six real-world datasets with varying properties, such as data file size and corpus size.
- (4) Each component of LS contributes meaningfully to its ability to find a sequence of transformations that satisfies our optimization goal and constraints.
- (5) LS performs well in a reasonable time.
- (6) Bottom-up script standardization can also be applied to identify anomalous data preparation steps.

Table 2: (R4) Parameterization effected by corpus properties.

Corpus properties		Parameters	
Large	Diverse	seq	K
# of scripts > 10	# of uniq. edges > 300	16	3
# of scripts > 10	# of uniq. edges ≤ 300	16	1
# of scripts ≤ 10	# of uniq. edges > 300	8	3
# of scripts ≤ 10	# of uniq. edges ≤ 300	8	1

Table 3: Examined datasets and their DAG statistics.

Statistics	Titanic	House	NLP	Spaceship	Medical	Sales
Scripts	62	49	24	38	47	26
Data files	3	4	3	3	1	6
Data tuples (k)	2.6	4.3	22.7	17.2	0.7	744.3
Data features	25	163	11	29	9	18
Avg # code lines	64	43	19	44	30	39
Uniq. 1-grams	625	659	175	427	281	347
Uniq. n-grams	548	221	48	201	143	114
Uniq. edges	748	632	193	423	220	308

6.1 Experimental Setting

6.1.1 Competing methods. We evaluate each of the following methods based on the performance metrics in Section 6.1.4.

- **Sourcery:** This is a commercial recommendation system to automatically improve code quality and ensure clean code [7].
- **GPT-3.5:** GPT-3.5 [6] is a language model with billions of parameters that enables it to generate text and code.
- **GPT-4:** GPT-4 [51] is an advanced iteration of GPT-3.5.
- **(R2) Auto-Suggest [63]:** This method uses machine learning models trained on real-world notebooks to predict a single next step for an input table based on its characteristics.
- **(R2) Auto-Tables [43]:** This method predicts multi-step transformations, among a set of table-reshaping operators.

The corpus can be accessed by the GPT baselines in two ways. Since GPT models are initially trained with scripts crawled from the internet [6, 51], the Kaggle datasets we used in the experiment should already be included in their training data. Another way for the GPT models to use the specific corpus more explicitly is to incorporate the corpus scripts in the prompt. Due to the limited token length, it is impossible for the prompts to include the entire script corpus. Empirically, we observed that prompts that include more scripts from the corpus achieve better performance. We describe how we select the prompts next.

6.1.2 GPT prompt survey. The performance of language models can be highly dependent on the problem prompt. We ensured that our experiments used a high quality prompt by polling 12 graduate students for prompt ideas, with access to the corresponding script corpus. We then ran all of them and chose the best-performing prompt to compare against our method. As a result, we believe that more typical users of LLMs would obtain worse results than what we report here. The survey and results are available in [2]. The best prompt randomly picks 4 scripts from the corpus and asks the LLMs to enhance the performance of the user script in terms of model accuracy or any other aspect that the LLMs find relevant.

6.1.3 Data collection and analysis. We used the Kaggle API [4] to obtain data preparation scripts in six competitions (Table 3). (R1) Specifically, we downloaded all available scripts attached to the chosen competitions, and used them as the input corpus. To evaluate the effectiveness of our system, for each competition, we treated each script in the corpus as an input user script and ran the system on the remaining scripts as the corpus. We report the average improvement on script standardness (calculated as explained in Section 6.1.4).

- **Titanic:** The *Titanic* competition aims to create a model to predict who survived the Titanic shipwreck.
- **Sales:** The *Predict Future Sale* competition aims to predict total sales for every product and store in the next month.

- **House:** The *House Prices* competition aims to use 79 variables to predict the final price of each home.
- **NLP:** The *NLP with Disaster Tweets* competition aims to distinguish real and fake tweets.
- **Spaceship:** The *Spaceship Titanic* competition aims to predict which passengers were transported from a spaceship.
- **Medical:** The *Pima Indians Diabetes Database* dataset aims to diagnostically predict if a patient has diabetes.

6.1.4 Performance metrics. We evaluated effectiveness using the percentage improvement in relative entropy (Definition 3), defined as % improvement = $\frac{RE(s_u, S) - RE(\hat{s}_u, S)}{RE(s_u, S)}$. We used the user-intent constraint to evaluate how well the user intent is being preserved. Specifically, we varied the Jaccard user intent threshold τ_J between $[0, 1]$. The user intent is preserved more (i.e., given less flexibility to standardize) when τ_J is closer to 1. We varied the absolute percentage difference in the model performance user intent threshold τ_M between $[0\%, 5\%]$. The user intent is preserved more when τ_M is closer to 0%. For each dataset, we iterated through the individual scripts as the user input script s_u and used the rest as the script corpus. The methods output \hat{s}_u . We then calculated % improvement of each \hat{s}_u relative to its original s_u .

6.1.5 Default LucidScript (LS) configuration. We set the LS default configuration as follows: sequence length $seq = 16$, beam size $K = 3$, diversity measure $div = on$, early-checking $\alpha = on$. This configuration obtained the best result. We set the user-intent threshold $\tau_J = 0.9$ and $\tau_M = 1\%$. This is a relatively strict user-intent constraint since, for example, $\tau_M = 1\%$ means that the output script should be as standard as possible while keeping the difference in its model performance within 1%.

6.2 Case Studies & A User Study

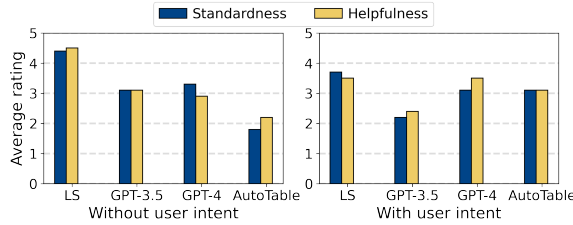
6.2.1 Metric Evaluation. Our first goal is to validate that improvements in our metrics (relative entropy and user intent) actually lead to meaningful outcomes for data scientists. We demonstrate this through the following case study.

Summary: An increase in our standardization metric indicates a more standardized data preparation script while maintaining the original user intent.

Results: We consider the *Titanic* competition. Table 4 shows an input script s_u and two output scripts s_1, s_2 , with their corresponding RE score and user intent measures Δ_J, Δ_M (for table Jaccard similarity and model performance, respectively). In this case, the input script s_u simply loads the dataset. The user intent measures for s_u , as in the input script, are $\Delta_J = 1, \Delta_M = 0\%$. s_1 , one of the potential outputs, involves a common step in which the target variable "Survived" is separated from the features. This data preparation step appears in 50% of the scripts within the corpus. As shown, the corresponding RE score increased by 27%, while remaining within the user intent measures of $\Delta_J = 0.9$ and $\Delta_M = 1\%$. Beyond the steps in s_1 , a more "standard" output script, s_2 , includes two additional common steps to handle missing values in the variables "Age" and "Embark". These steps are present in 45% and 14% of the scripts in the corpus. This modification further increases the RE score by 55% compared to the input script s_u . This example shows that the modified scripts incorporate common and logical data preparation

Table 4: Case study for metrics evaluation, where s_u is the input script, and s_1 and s_2 are the potential outputs.

	Script	RE	Δ_J	Δ_M
s_u	import pandas as pd import numpy as np df = pd.read_csv("/data/titanic/train.csv")	3.02	1	0%
s_1	import pandas as pd import numpy as np df = pd.read_csv("/data/titanic/train.csv") y = df["Survived"] X_train = df.drop("Survived", axis=1)	2.49	0.92	<0.1%
s_2	import pandas as pd import numpy as np df = pd.read_csv("/data/titanic/train.csv") df["Age"].fillna(df["Age"].mean()) df["Embarked"].fillna("S") y = df["Survived"] X_train = df.drop("Survived", axis=1)	1.37	0.90	<0.1%

**Figure 3: User study result.**

steps. It also shows that an increase in the RE score corresponds to the inclusion of such steps while preserving user intent.

6.2.2 User Study. (R1&2) Our next goal is to validate that LS generated scripts that are considered by data scientists more standard compared to the output of the baselines. We recruited 34 undergraduate and graduate computer science students who are familiar with data preparation and Python programming. Each participant was given two use cases, including statistics on the prevalence of data preparation steps in the corpus (similar to Table 1), and 5 scripts output by LS and the baselines. In the without-user-intent case, participants evaluated the methods in a cold start setting, while in the with-user-intent case, participants were given an input script.

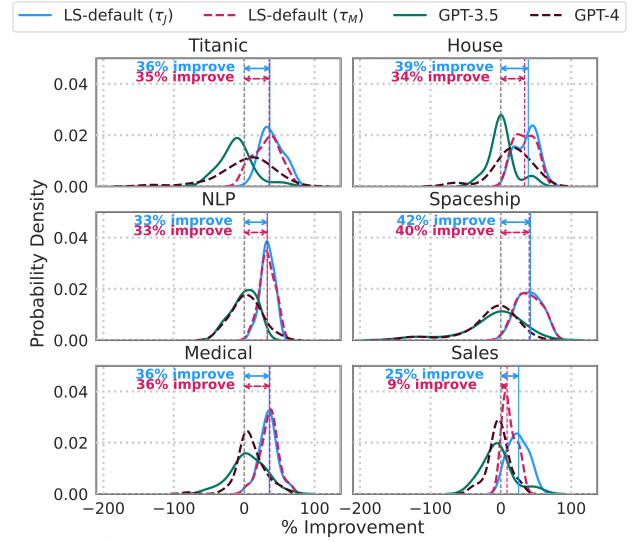
Participants were asked to rank each of the output scripts on a scale of 1 to 5 based on: (a) how standardized the script is w.r.t. the corpus, where standardization represents the commonality of the chosen data preparation steps; and (b) how helpful the script is w.r.t. preserving the user intent, measured by their modeling task. The complete form provided to participants is available at [2].

Summary: LS was rated to be the most standard and helpful.

Results: In the without-user-intent case, it is statistically significant that LS is more standard and helpful than baseline methods⁷. In the with-user-intent case, LS is the most standard and helpful among all 5 methods. Note that we omitted Auto-Suggest from Figure 3 since it had the same results as Auto-Tables for this study.

Table 5: (R1&2) % improvement on Kaggle datasets, with $\tau_J = 0.9$ and $\tau_M = 1\%$. LS is set as default.

Corpus setup	Method	% Improvement			
		min	median	max	mean
Full-size corpus	LS (τ_J)	0.0	33.1	72.3	33.6
	LS (τ_M)	0.0	26.9	70.0	25.8
	GPT-3.5	-57.3	-1.5	58.4	-3.7
	GPT-4	-129.0	0.0	66.1	3.4
	Sourcery	0.0	0.0	0.0	0.0
	Auto-Suggest	0.0	0.0	0.0	0.0
	Auto-Tables	0.0	0.0	0.0	0.0
Small corpus	LS (τ_J)	0.0	18.5	69.3	20.3
	LS (τ_M)	0.0	15.0	69.3	17.1
Different corpus	LS (τ_J)	0.0	11.1	36.4	10.5
	LS (τ_M)	0.0	11.8	36.4	11.2
Low-ranked corpus	LS (τ_J)	0.0	5.4	32.9	7.8
	LS (τ_M)	0.0	5.2	32.9	7.7

**Figure 4: (R3) % improvement distribution. Sourcery, Auto-Suggest, and Auto-Tables are omitted.**

6.3 Evaluating Effectiveness

We aim to answer the following three questions. Is our framework effective in standardizing data preprocessing scripts with the constraints? How does our method compare to the others? And how does our method perform on a variety of datasets?

6.3.1 Script standardization with constraints. We evaluated the % improvement of LucidScript (LS) with the default configuration against the three competing methods on six datasets.

Summary: LS consistently outperforms competing methods by far. None of the state-of-the-art methods can effectively improve script standardization.

Results: In Figure 4, the most optimal result would be a sharp peak on the most right, meaning that all input scripts get an improvement of 100%, while a curve centered around $x = 0$ means that the method does not improve script standardness. A curve extending to the left ($x \leq 0$) means that the method decreases script standardness.

⁷according to a t-test, p -value < 0.05

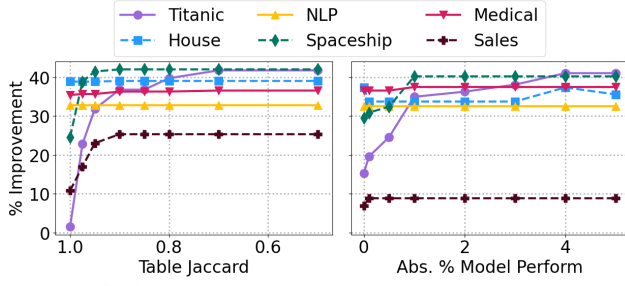


Figure 5: (R3) Median % improvement of LS-default with varied table Jaccard threshold τ_J (left); median % improvement of LS-default with varied model performance threshold τ_M (right).

LS achieves at least 30% improvement in script standardization while preserving user intent, other than Sales (Table 5). GPT models perform similarly with only a small positive average % improvement on Titanic (GPT-4: 7.7%) and House (GPT-4: 6.2%). Moreover, LS guarantees a standardness improvement, whereas GPT models do not produce reliable results, sometimes emitting a script that uses less common transformations than the user’s original (GPT-4: -130%). This is because LS aims to preserve the user intent while standardizing the script using the specific collective knowledge of D_{IN} . Although GPT models are trained on all publicly available scripts, they do not standardize the script with a focus on knowledge specific to D_{IN} .

(R2) We compared LS to two academic solutions, Auto-Suggest [63] and Auto-Tables [43], which are, respectively, the state-of-the-art work for single-step and multi-step table transformation prediction. These methods made a 0% improvement because they focus on a fixed set of table structural transformations, such as transpose and pivot. However, the real-world scripts that we observed focus on feature engineering and data cleaning tasks (Figure 1a). Our bottom-up approach makes LS robust to the type of data preprocessing tasks. We also compared with Sourcery, a commercial product for code cleaning, but observed that it consistently made no difference on all measures, as it focuses on syntax standardization.

6.3.2 Evaluating with user intent. We experimented with two user-intent measures, table Jaccard and model performance. We evaluated how performance changes with several user intent thresholds, ranging from more lenient to more strict. Since competing methods do not have this property, we evaluated only LS.

Summary: LS can adhere to a wide range of user intent while improving script standardness in the six datasets.

Results: We varied table Jaccard τ_J between $[0, 1]$ (Figure 5). $\tau_J = 1$ means the D_{OUT}^{su} and D_{OUT}^s are identical, which is the most strict user intent constraint, where the intent should be preserved in its entirety. We found that as this constraint becomes more relaxed (i.e., τ_J becomes smaller), LS can make the input scripts more standardized. In all six data sets, once $\tau_J < 0.7$, the % improvement plateaus, so we cap the figure at $\tau_J = 0.5$. We attribute this to the default setting that has a maximum sequence length of 16, which limits the maximum number of transformations that can be applied to s_u .

We varied model performance τ_M between $[0\%, 5\%]$ (Figure 5). $\tau_M = 0\%$ means the $acc(D_{OUT}^{su})$ and $acc(D_{OUT}^s)$ are identical. In

general, on all six datasets, as the constraint becomes more relaxed, LS can make more improvements to script standardization. We observed a non-monotonic trend in House, which could be attributed to the complex interaction between $acc(D_{OUT})$ and D_{OUT} . Although $acc(D_{OUT})$ is influenced by D_{OUT} , the correlation is complex. Some small changes in the content of D_{OUT} may yield a significant improvement or decrease in $acc(D_{OUT})$.

We would expect that when the system has more freedom to change the script — that is, the constraint on the user’s original intent is relaxed — we would see highly standardized output scripts. In contrast, when the system is limited to producing an output that is extremely fidelity to the user’s original script, it has little flexibility, and we would expect to see output scripts that cannot be standardized very much. That is in fact what we see in Figures 5. As we relax the constraints on user intent, we see that LS indeed usually produces output scripts that are more standardized.

6.3.3 Script standardization using various corpus settings. We consider three corpus scenarios to examine the robustness (Table 5).

Summary: Overall, LS managed to find common transformations that improve script standardness, given a variety of corpora.

Results: We experimented with the following three scenarios.

Using a small corpus. We consider the scenario where only a limited number of scripts are available. We sampled 10 scripts from each dataset. Overall, LS maintained at least a 13% improvement in all datasets for both user-intent constraints, except Sales. We observed the most performance reduction on Titanic, House, and Medical. This is because their original corpus size is larger (Table 3) compared to the other datasets, providing a richer search space and potential opportunities to improve standardization.

Using a corpus from a different dataset that has a similar schema.

We consider the scenario where a corpus for the identical dataset is not available, so the user resorts to a corpus on a similar dataset. We experimented with improving Spaceship scripts using Titanic corpus. LS still achieves an 11% improvement. This is because the competition setup and D_{IN} for Spaceship and Titanic share many similarities: they both predict survival, and there are identical column names in D_{IN} . This shows that LS can be used successfully by applying one script corpus to a different target, although — unsurprisingly — the quality improvement is reduced.

(R1&2) Using a low-ranked corpus. We consider the scenario where the corpus consists solely of low-ranked scripts. Specifically, the corpus uses the bottom 30% of the scripts, determined by their count of votes received on Kaggle. Despite being challenged by a smaller, low-ranked corpus, LS still managed to find common transformations that improve script standardness by 5%. Our takeaway here is that when only low-quality scripts are available, it is more effective to use scripts that process a similar dataset, as this approach leads to better system performance (as shown in Table 5).

6.4 Ablation Studies

We isolate each piece of our algorithm and evaluate their contribution to the overall effectiveness of LS.

6.4.1 Varied sequence lengths. Sequence length is the maximum number of transformations in a sequence. We use it as the stopping criterion. We set the sequence length seq to $\{2, 4, 8, 16\}$.

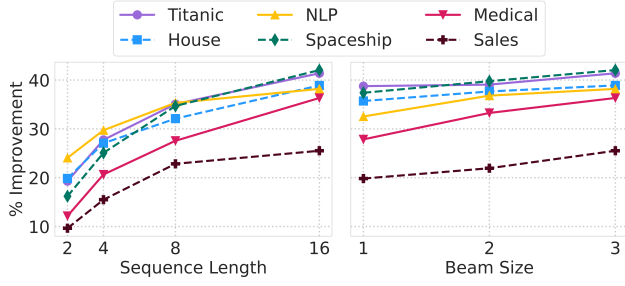


Figure 6: (R3) Median % improvement of LS-default with varied sequence lengths (left) and varied beam sizes (right).

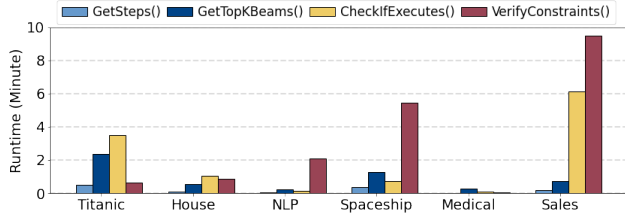


Figure 7: Median runtime breakdown on $seq = 16$.

Summary: LS makes scripts more standardized as the maximum sequence length increases.

Results: As Figure 6 (R4) shows, on the six datasets, the median % improvement increases as the maximum number of sequence length. As seq increases, LS improves standardization at a faster rate at the start, and then the improvement slowly starts to plateau from $seq = 8$ to $seq = 16$. This could be attributed to the size of the script corpus and the diversity of atoms and edges (Table 3). We observed the best performance at $seq = 16$ on Titanic and Spaceship, which are two datasets with the most unique atoms and edges, with a % improvement of 41.4% and 42% respectively.

6.4.2 Varied beam sizes. The number of beams K in $\{1, 2, 3\}$ is how many in-progress sequences LS retains.

Summary: LS performs better as K increases.

Results: Figure 6 shows that in all six data sets, the beam search enables LS to increase the % improvement steadily as K increases. We observe that the beam search effect is the strongest on Medical, achieving an increase of 8.5% from $K = 1$ to $K = 3$, and the weakest on Titanic, obtaining an increase of 2.6%. On average, the beam search increases the % improvement by 5% across the six datasets.

6.5 Evaluating Efficiency

We isolate each component and discuss factors that could affect search performance and end-to-end latency.

On average, LS takes several minutes to generate a standardized version of a given input script (on average, one script takes 7.2 minutes end-to-end). Although this process may take a few minutes, it involves generating an entire script instead of just a single-step recommendation. This time frame is typical and comparable to other systems that generate full pipelines hands-free [49, 50, 53, 58, 65]. Future work will focus on improving runtime by reducing the search space, possibly by grouping semantically similar operations. See Section 8 for further discussion.

Summary: In general, the latency of the search steps (GetSteps(), GetTopKBeams()) increases as unique atoms and edges increase.

```

1 import pandas as pd
2 df = pd.read_csv('diabetes.csv')
3 df = df.fillna(df.median())
4 df = df[df["Pregnancies"] < 10]
5 df = df[df["Age"].between(18,25)]
6 df = pd.get_dummies(df)
7 df['Outcome_copy'] = df['Outcome']
8 update = df.sample(20).index
9 df.loc[update, 'Outcome_dup'] = 0
10 X = df.drop(['Outcome'], axis=1)
11 y_train = df['Outcome']

1 import pandas as pd
2 df = pd.read_csv('diabetes.csv')
3 df = df.fillna(df.median())
4 df = df[df["Pregnancies"] < 10]
5 df = df[df["Age"].between(18,25)]
6 df = pd.get_dummies(df)
7 df['Outcome']
8 update = df.sample(20).index
9 df.loc[update, 'Outcome_dup'] = 0
10 X = df.drop(['Outcome'], axis=1)
11 y_train = df['Outcome']

```

Figure 8: Example input s_u (left) and output scripts \hat{s}_u (right). s_u has created a column correlated to the target with random noises (line 7-9). Green shows the added steps, while red shows the steps recognized to be removed from s_u .

The size of the data file D_{IN} also significantly affects the latency of the checking steps (CheckIfExecutes(), VerifyConstraints()).

Results: As shown in Figure 7, LS has the lowest latency on Medical and NLP, which are the two datasets with the lowest number of atoms and edges. The small number of atoms and edges implies a more narrow search space, which makes LS faster to find the potential next steps and converge to a standardized script. The type of ML models the users choose mainly has an impact on the latency of VerifyConstraint() since that is where our framework checks for the user-intent constraint, which includes model performance. This is the main reason why NLP has a longer latency than Medical. The size of the data file affects the latency from the perspective of constraint verification. The original latency is 20x slower on Sales compared to the other datasets before sampling is applied. This is because Sales has 744k tuples in D_{IN} , which is 30x more than NLP, the second largest D_{IN} . LS checks for the constraints by running the scripts that inevitably load D_{IN} . The larger D_{IN} , the more constraint checking contributes to the latency. We can sufficiently reduce the latency caused by large datasets with sampling.

Overall, the latency is influenced by both the algorithmic components of LS and the characteristics of the dataset. The median latency can be less than a minute per script in Medical and can take up to 17.5 minutes on Sales. We also note that the current latency depends on this specific version of the system prototype, which could be improved with more engineering efforts. For example, we can use parallelism, running each beam on a separate thread. In general, data preparation is a time-consuming task that is challenging even for experienced data scientists [20, 21]. In contrast, we offer a unique and effective solution that can significantly standardize data preparation scripts that require minimum effort from data consumers regardless of their backgrounds.

6.6 Case Study: Target Leakage Detection

Script standardization is the process of simultaneously removing out-of-the-ordinary steps from the input script and adding the missing common practices to the input script. Target leakage is another example of an out-of-the-ordinary step that should be replaced with one that is more conventional and correct. As a result, our system can potentially identify it (Figure 8).

6.6.1 Study setup. We performed a qualitative case study.

Target leakage. As one of the common data preparation mistakes in practice, target leakage occurs when the target variable or information derived from it is included in the feature set. This leads to

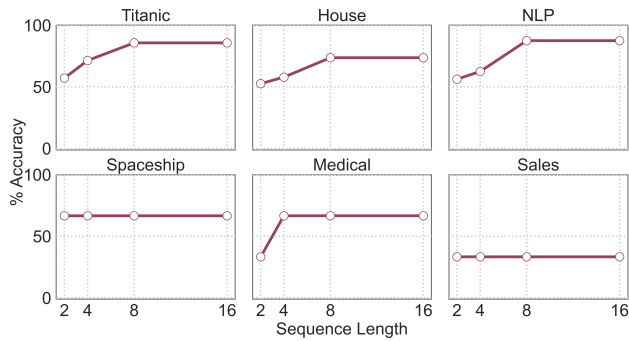


Figure 9: Accuracy of LS-default for target leakage detection.

overly optimistic performance estimates during model evaluation, as the model has access to information that would not be available in a real-world scenario. For each dataset, we sample 10% of the real-world scripts and use GPT-4 [51] to programmatically inject target leakage code snippets (ground truth). Prompts used are in [2].

Performance metric. We evaluate the effectiveness using the accuracy of detecting the ground truth snippets. Specifically, an output script is correct when it satisfies the constraints (Definition 4), and the ground truth snippet is labeled as to be removed.

6.6.2 Study results. We vary sequence lengths to study how many steps would LS take to detect the ground truth snippets. Figure 9 shows the effectiveness of our approach in target leakage detection, where over 66% of ground truth snippets are discovered within 8 steps for all datasets except Sales.

7 Related Work

One-Step Recommendations: Automating data preparation pipelines can be achieved by recommending the next best data preparation step [8, 32, 33, 35, 63]. These solutions utilize AI and crowd-sourcing [33], web data [8], available data science notebooks [63], and learning models [63] to generate recommendations and predict user intent. Close to our work, the authors of [63] proposed Auto-Suggest [63], which “learns” data preparation steps from data scientist workflows and recommends the next step by predicting the next operator and its parameters. Our work is distinct in that it seeks to *standardize* an input script as much as possible while preserving user intent, rather than simplifying the task of analysts by predicting their next steps.

Multi-Step Recommendations: Researchers have recently explored multi-step automation for data preparation pipelines. Our work is the most similar to (R2) Auto-Tables [43] and Learn2Clean [12, 13]. Auto-Tables uses machine learning models to predict a sequence of table structural transformations for a given input table. Learn2Clean [12, 13] employs reinforcement learning to explore all possible strategies to process a given dataset to maximize the performance of a specific machine learning model. Although multi-step recommendations can be an application of our work, fundamentally, we solve a different problem. Instead of recommending a sequence of data preparation steps for a given user intent in [12, 13, 64], we take the user’s input script merely as a sketch and tweak the semantics within a threshold to improve script standardness.

Program Synthesis: Program synthesis is a technique that seeks to automatically generate code based on user intent, which can be

specified by various forms of constraints such as I/O examples. Studies have demonstrated that several components of data preparation can be automated using programming by example, including spreadsheet automation [18], code generation in R [23], and SQL query generation [61]. Closer to our work, AutoPandas [10] takes input/output examples that demonstrate the desired transformation and uses a neural-backed synthesis engine to generate Python programs. Recently, the use of large language models for program synthesis has become increasingly popular. OpenAI Codex [17], a GPT [15] language model fine-tuned on public Python code on GitHub, shows its capabilities in code writing by generating programs using docstrings. It can also be used to power many downstream program synthesis applications, such as GitHub Copilot [3], which suggests code and functions in real time to save programmers’ time on repetitive code. The key difference between our problem and this line of work is that in program synthesis, the semantics and functionality of the program should not deviate from user intent, while we use user intent as guidance and allow it to be refined. This flexibility enables analysts to explore and integrate the collective knowledge embedded in the script collection. We evaluated the GPT models [6, 51] on the script standardization problem, without putting constraints on the intent of the user.

Code style transfer Learning-based approaches often analyze common coding styles found in codebases by converting code into representations like one-hot token embeddings or parse trees. They then derive formatting conventions and generate explicit [9, 45] or abstract [52] rules. Another recent system, DUETCS [16], aims to facilitate automatic code style transfer. However, these methods typically require the user to pre-define the coding style or focus solely on formatting issues. In contrast, our system automatically extracts abstract styles from the corpus. The goal is to preserve user intent by modifying the script to align with the most prevalent (i.e., standard) style in the corpus.

Automated Machine Learning (AutoML): This line of work focuses on automated processes that use machine learning algorithms and techniques to build, optimize, and select predictive models for a given dataset. AutoML is used to automate the process of training and tuning machine learning models, allowing users to quickly and easily generate more accurate models with less effort and time [34, 60]. Tools such as Auto-sklearn [24, 25], Auto-WEKA [60], and DeepLine [34] focus on automating machine learning model building. These systems tackle the challenge of automatic selection and optimization of machine learning models and their hyperparameters, and combine multiple types of machine learning model into a single model [34]. We are situated upstream in the machine learning model-building pipeline, helping analysts prepare their datasets before feeding them into machine learning models.

8 Conclusion and Limitations

We introduce a new problem and propose an efficient framework to standardize data preparation scripts with minimal effort. When interpreting the generated script, several limitations should be considered. First, its quality may be influenced by factors such as the quality of the input script and the user intent measure and threshold. Second, the generated script may not be fully ready for production,

so human validation is necessary to ensure that it performs as intended and that no unintended biases are introduced [29].

Our framework currently supports scripts written in Python. In principle, it can be generalized to other programming languages. However, this would require extending the script-to-DAG and DAG-to-script procedures, as well as the lemmatization techniques used. We note that other systems on data preparation recommendations have also focused on Python [63, 64].

Standardizing data preparation scripts is challenging due to the diverse ways in which the same operations can be expressed across different frameworks (e.g., Pandas, PySpark, Modin, DuckDB). To accurately measure the similarity between various operations, it is crucial to use a high-level abstraction or knowledge base that ignores specific details such as the frameworks or their versions. Our current implementation uses simple NLP solutions [54]. Although we recognize existing work in this area and suggest that LLMs may be effective, our system does not yet incorporate these advanced approaches. This highlights a key opportunity for extending our framework. However, despite using LLM-based solutions as baselines, which is expected to perform well in this context, our framework demonstrated superior performance. (R1) *With the rapid improvement of LLMs, especially in code generation, we expect to see improvements in their performance. We plan to integrate LLMs into our system to automatically verify that constraints are satisfied without explicitly checking them (i.e., ensuring that the code is executable). We also plan to rerun the experiments with new versions of LLMs in the future to track their improvement.*

Another limitation is that the user must provide a script corpus. We note that similar assumptions are made by other frameworks [16, 63]. In addition, the dataset may be updated, or different datasets that have the same schema become available (e.g., MIMIC datasets [36, 42]). As we showed, even a small corpus or one containing scripts that process similar datasets can still yield valuable results. To address situations where obtaining a script corpus is challenging, one potential solution is to generate scripts using LLMs. Additionally, scripts authored by domain experts could be weighted differently (e.g., using the vote counts of Kaggle scripts).

The user also sets a threshold that indicates how much they are willing to compromise on their intent in exchange for a more standardized script. Our experiments show that the modified script not only stays within this threshold, but also, in many cases, improves model performance. (R4) *Along with the user-intent threshold, our framework provides a set of parameters for the user to customize. While the parameters offer flexibility, they may require effort to tune. A possible extension to this work is an algorithm that optimizes configurations, such as exploring user intent thresholds and returning the Pareto curve.*

Our framework currently supports two methods for evaluating how well user intent is preserved: table Jaccard and model performance. Future work will focus on incorporating additional metrics, such as assessing the semantic similarity between scripts, comparing their bags of operations, or evaluating the fairness of the model [29]. (R1&4) *Lastly, an interesting direction for further exploration that we intend to pursue is providing explanations for why a transformation is recommended to the input script. The explanation would inform the user about the frequency of this operation in the corpus, its impact on the user intent, and the rationale behind*

it, such as scaling the features to improve the performance of a learning model.

References

- [1] ASTs – abstract syntax trees. <https://docs.python.org/3/library/ast.html>. Accessed: 2023-01-30.
- [2] Bottom-up standardization for data preparation: full technical report, code, and results for review. <https://github.com/ey-l/bottom-up-script-standardization>.
- [3] GitHub Copilot your ai pair programmer. <https://github.com/features/copilot>. Accessed: 2023-01-30.
- [4] Kaggle API. <https://github.com/Kaggle/kaggle-api>. Accessed: 2023-01-30.
- [5] Kaggle pima indians diabetes database. <https://www.kaggle.com/datasets/ucml/pima-indians-diabetes-database>. Accessed: 2023-01-30.
- [6] OpenAI introducing chatgpt. <https://openai.com/blog/chatgpt>. Accessed: 2023-01-30.
- [7] Sourcery automatically improve python code quality. <https://sourcery.ai/>. Accessed: 2023-01-30.
- [8] Z. Abedjan, J. Morcos, I. F. Ilyas, M. Ouzzani, P. Papotti, and M. Stonebraker. Dataxformer: A robust transformation discovery system. In *2016 IEEE 32nd International Conference on Data Engineering (ICDE)*, pages 1134–1145, 2016.
- [9] M. Allamanis, E. T. Barr, C. Bird, and C. Sutton. Learning natural coding conventions. In *Proceedings of the 22nd acm sigsoft international symposium on foundations of software engineering*, pages 281–293, 2014.
- [10] R. Bavishi, C. Lemieux, R. Fox, K. Sen, and I. Stoica. Autopandas: Neural-backed generators for program synthesis. *Proc. ACM Program. Lang.*, 3(OOPSLA), oct 2019.
- [11] I. D. Baxter, A. Yahin, L. Moura, M. Sant’Anna, and L. Bier. Clone detection using abstract syntax trees. In *Proceedings. International Conference on Software Maintenance (Cat. No. 98CB36272)*, pages 368–377. IEEE, 1998.
- [12] L. Berti-Equille. Learn2clean: Optimizing the sequence of tasks for web data preparation. In *The World Wide Web Conference, WWW ’19*, page 2580–2586, New York, NY, USA, 2019. Association for Computing Machinery.
- [13] L. Berti-Equille. Reinforcement learning for data preparation with active reward learning. In S. El Yacoubi, F. Bagnoli, and G. Pacini, editors, *Internet Science*, pages 121–132, Cham, 2019. Springer International Publishing.
- [14] R. Botvinik-Nezer, F. Holzmeister, C. Camerer, et al. Variability in the analysis of a single neuroimaging dataset by many teams. pages 84–88, 2020.
- [15] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei. Language models are few-shot learners, 2020.
- [16] B. Chen and Z. Abedjan. Duets: Code style transfer through generation and retrieval. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*, pages 2362–2373. IEEE, 2023.
- [17] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. d. O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, A. Ray, R. Puri, G. Krueger, M. Petrov, H. Khlaaf, G. Sastry, P. Mishkin, B. Chan, S. Gray, N. Ryder, M. Pavlov, A. Power, L. Kaiser, M. Bavarian, C. Winter, P. Tillet, F. P. Such, D. Cummings, M. Plappert, F. Chantzis, E. Barnes, A. Herbert-Voss, W. H. Guss, A. Nichol, A. Paino, N. Tezak, J. Tang, I. Babuschkin, S. Balaji, S. Jain, W. Saunders, C. Hesse, A. N. Carr, J. Leike, J. Achiam, V. Misra, E. Morikawa, A. Radford, M. Knight, M. Brundage, M. Murati, K. Mayer, P. Welinder, B. McGrew, D. Amodei, S. McCandlish, I. Sutskever, and W. Zaremba. Evaluating large language models trained on code, 2021.
- [18] X. Chen, P. Maniatis, R. Singh, C. Sutton, H. Dai, M. Lin, and D. Zhou. Spreadsheetcoder: Formula prediction from semi-structured context, 2021.
- [19] P. W. Code. Papers with code. <https://paperswithcode.com/>. 2015. Accessed: April 12, 2023.
- [20] T. Dasu, T. Johnson, S. Muthukrishnan, and V. Shkapenyuk. Mining database structure; or, how to build a data quality browser. In *Proceedings of the 2002 ACM SIGMOD International Conference on Management of Data, SIGMOD ’02*, page 240–251, New York, NY, USA, 2002. Association for Computing Machinery.
- [21] D. Deng, R. C. Fernandez, Z. Abedjan, S. Wang, M. Stonebraker, S. Madden, M. Ouzzani, and N. Tang. The data civilizer system.
- [22] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2018.
- [23] Y. Feng, R. Martins, O. Bastani, and I. Dillig. Program synthesis using conflict-driven learning, 2017.
- [24] M. Feurer, K. Eggenberger, S. Falkner, M. Lindauer, and F. Hutter. Auto-sklearn 2.0: Hands-free autml via meta-learning, 2020.
- [25] M. Feurer, A. Klein, K. Eggenberger, J. Springenberg, M. Blum, and F. Hutter. Efficient and robust automated machine learning. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015.
- [26] H. Galhardas, D. Florescu, D. Shasha, and E. Simon. Ajax: An extensible data cleaning tool. *SIGMOD Rec.*, 29(2):590, may 2000.
- [27] A. Gelman and E. Loken. The garden of forking paths: Why multiple comparisons can be a problem, even when there is no. *Fishing Expedition” or “p-hacking” and the Research Hypothesis was Posited ahead of Time.” Technical report, Department of Statistics, Columbia University*, 2013.
- [28] S. Graftberger, J. Stoyanovich, and S. Schelter. Lightweight inspection of data preprocessing in native machine learning pipelines. In *Conference on Innovative Data Systems Research (CIDR)*, 2021.
- [29] S. Guha, F. A. Khan, J. Stoyanovich, and S. Schelter. Automated data cleaning can hurt fairness in machine learning-based decision making. *IEEE Transactions on Knowledge and Data Engineering*, 2024.
- [30] W. R. Harris and S. Gulwani. Spreadsheet table transformations from examples. *ACM SIGPLAN Notices*, 46(6):317–328, 2011.
- [31] J. A. Hartigan and M. A. Wong. Algorithm as 136: A k-means clustering algorithm. *Journal of the royal statistical society. series c (applied statistics)*, 28(1):100–108, 1979.
- [32] Y. He, K. Ganjam, K. Lee, Y. Wang, V. Narasayya, S. Chaudhuri, X. Chu, and Y. Zheng. Transform-data-by-example (tde): Extensible data transformation in excel. In *Proceedings of the 2018 International Conference on Management of Data, SIGMOD ’18*, page 1785–1788, New York, NY, USA, 2018. Association for Computing Machinery.
- [33] J. Heer, J. M. Hellerstein, and S. Kandel. Predictive interaction for data transformation. In *CIDR*, 2015.
- [34] Y. Heffetz, R. Vainshtein, G. Katz, and L. Rokach. Deepline: Autml tool for pipelines generation using deep reinforcement learning and hierarchical actions filtering. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD ’20*, page 2103–2113, New York, NY, USA, 2020. Association for Computing Machinery.
- [35] Z. Jin, Y. He, and S. Chaudhuri. Auto-transform: Learning-to-transform by patterns. *Proc. VLDB Endow.*, 13(12):2368–2381, jul 2020.
- [36] A. E. Johnson, T. J. Pollard, L. Shen, L.-w. H. Lehman, M. Feng, M. Ghassemi, B. Moody, P. Szolovits, L. Anthony Celi, and R. G. Mark. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3(1):1–9, 2016.
- [37] S. Kandel, A. Paepcke, J. Hellerstein, and J. Heer. Wrangler: Interactive visual specification of data transformation scripts. In *Proceedings of the sigchi conference on human factors in computing systems*, pages 3363–3372, 2011.
- [38] S. Kandel, R. Parikh, A. Paepcke, J. M. Hellerstein, and J. Heer. Profiler: Integrated statistical analysis and visualization for data quality assessment. In *Proceedings of the International Working Conference on Advanced Visual Interfaces, AVI ’12*, page 547–554, New York, NY, USA, 2012. Association for Computing Machinery.
- [39] K. Kennedy. Use-definition chains with applications. *Computer Languages*, 3(3):163–179, 1978.
- [40] S. Kullback and R. A. Leibler. On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86, 1951.
- [41] E. Lai, Y. Lou, B. Youngmann, and M. Cafarella. Lucidscript: Bottom-up standardization for data preparation. *Proceedings of the VLDB Endowment (PVLDB)*, 17(12):4317–4320, 2024.
- [42] J. Lee, D. J. Scott, M. Villarreal, G. D. Clifford, M. Saeed, and R. G. Mark. Open-access mimic-ii database for intensive care research. In *2011 Annual International Conference of the IEEE Engineering in Medicine and Biology Society*, pages 8315–8318, 2011.
- [43] P. Li, Y. He, C. Yan, Y. Wang, and S. Chaudhuri. Auto-tables: Relationalize tables without using examples. *ACM SIGMOD Record*, 53(1):76–85, 2024.
- [44] W. Li, H. Yin, Y. Chen, Q. Liu, Y. Wang, D. Qiu, H. Ma, and Q. Geng. Associations between adult triceps skinfold thickness and all-cause, cardiovascular and cerebrovascular mortality in nhanes 1999–2010: a retrospective national study. *Frontiers in cardiovascular medicine*, 9:858994, 2022.
- [45] V. Markovtsev, W. Long, H. Mougard, K. Slavov, and E. Bulchev. Style-analyzer: fixing code style inconsistencies with interpretable unsupervised algorithms. In *2019 IEEE/ACM 16th International Conference on Mining Software Repositories (MSR)*, pages 468–478. IEEE, 2019.
- [46] A. H. Murphy. The finley affair: A signal event in the history of forecast verification. *Weather and Forecasting*, 11(1):3 – 20, 1996.
- [47] M. H. Namaki, A. Floratou, F. Psallidas, S. Krishnan, A. Agrawal, Y. Wu, Y. Zhu, and M. Weimer. Vamsa: Automated provenance tracking in data science scripts. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 1542–1551, 2020.
- [48] I. Neamtii, J. S. Foster, and M. Hicks. Understanding source code evolution using abstract syntax tree matching. In *Proceedings of the 2005 international workshop on Mining software repositories*, pages 1–5, 2005.
- [49] E. Oikarinen, K. Puolamäki, S. Khoshrou, and M. Pechenizkiy. Supervised human-guided data exploration. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 85–101. Springer, 2019.
- [50] B. Omidvar-Tehrani, A. Personnaz, and S. Amer-Yahia. Guided text-based item exploration. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*, pages 3410–3420, 2022.
- [51] OpenAI. Gpt-4 technical report, 2023.
- [52] T. Parr and J. Vinju. Towards a universal code formatter through machine learning. In *Proceedings of the 2016 ACM SIGPLAN International Conference on Software Language Engineering*, pages 137–151, 2016.
- [53] A. Personnaz, B. Youngmann, and S. Amer-Yahia. Eda4sum: Guided exploration of data summaries. *Proceedings of the VLDB Endowment (PVLDB)*, 15(12):3590–3593, 2022.

- 2022.
- [54] J. Plisson, N. Lavrac, D. Mladenic, et al. A rule based approach to word lemmatization. In *Proceedings of IS*, volume 3, pages 83–86, 2004.
 - [55] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever. Language models are unsupervised multitask learners. 2019.
 - [56] V. Raman and J. M. Hellerstein. Potter's wheel: An interactive data cleaning system. In *Vldb*, volume 1, pages 381–390, 2001.
 - [57] C. Scaffidi, B. Myers, and M. Shaw. Intelligently creating and recommending reusable reformatting rules. In *Proceedings of the 14th International Conference on Intelligent User Interfaces*, IUI '09, page 297–306, New York, NY, USA, 2009. Association for Computing Machinery.
 - [58] M. Seleznova, B. Omidvar-Tehrani, S. Amer-Yahia, and E. Simon. Guided exploration of user groups. *Proceedings of the VLDB Endowment (PVLDB)*, 13(9):1469–1482, 2020.
 - [59] J. P. Simmons, L. D. Nelson, and U. Simonsohn. False-positive psychology: Undisclosed flexibility in data collection and analysis allows presenting anything as significant. *Psychological science*, 22(11):1359–1366, 2011.
 - [60] C. Thornton, F. Hutter, H. H. Hoos, and K. Leyton-Brown. Auto-weka: Combined selection and hyperparameter optimization of classification algorithms, 2012.
 - [61] C. Wang, A. Cheung, and R. Bodik. Synthesizing highly expressive sql queries from input-output examples. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI 2017, page 452–466, New York, NY, USA, 2017. Association for Computing Machinery.
 - [62] J. M. Wicherts, C. L. Veldkamp, H. E. Augusteijn, M. Bakker, R. Van Aert, and M. A. Van Assen. Degrees of freedom in planning, running, analyzing, and reporting psychological studies: A checklist to avoid p-hacking. *Frontiers in psychology*, page 1832, 2016.
 - [63] C. Yan and Y. He. Auto-suggest: Learning-to-recommend data preparation steps using data science notebooks. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, SIGMOD '20, page 1539–1554, New York, NY, USA, 2020. Association for Computing Machinery.
 - [64] J. Yang, Y. He, and S. Chaudhuri. Auto-pipeline: Synthesizing complex data pipelines by-target using reinforcement learning and search, 2021.
 - [65] B. Youngmann, S. Amer-Yahia, and A. Personnaz. Guided exploration of data summaries. *arXiv preprint arXiv:2205.13956*, 2022.
 - [66] M.-C. Yuen, I. King, and K.-S. Leung. A survey of crowdsourcing systems. In *2011 IEEE third international conference on privacy, security, risk and trust and 2011 IEEE third international conference on social computing*, pages 766–773. IEEE, 2011.
 - [67] A. X. Zhang, M. Muller, and D. Wang. How do data science workers collaborate? roles, workflows, and tools. *Proc. ACM Hum.-Comput. Interact.*, 4(CSCW1), may 2020.
 - [68] J. Zhang, X. Wang, H. Zhang, H. Sun, K. Wang, and X. Liu. A novel neural source code representation based on abstract syntax tree. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*, pages 783–794. IEEE, 2019.