

1 **Checklist**

2 1. For all authors...

- 3 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
4 contributions and scope? [Yes]
- 5 (b) Did you describe the limitations of your work? [Yes] Section C in the supplementary
6 materials
- 7 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 8 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
9 them? [Yes]

10 2. If you are including theoretical results...

- 11 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 12 (b) Did you include complete proofs of all theoretical results? [N/A]

13 3. If you ran experiments (e.g. for benchmarks)...

- 14 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
15 mental results (either in the supplemental material or as a URL)? [Yes] Code and data
16 are released at <https://github.com/shirley-wu/daco>
- 17 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
18 were chosen)? [Yes] Section 5 in the main content, and Section D in the supplementary
19 materials
- 20 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
21 ments multiple times)? [No] Did not perform experiments of multiple random seeds
22 due to resource constraints
- 23 (d) Did you include the total amount of compute and the type of resources used (e.g., type
24 of GPUs, internal cluster, or cloud provider)? [Yes] Section D in the supplementary
25 materials

26 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- 27 (a) If your work uses existing assets, did you cite the creators? [Yes] Section 3 in the main
28 content
- 29 (b) Did you mention the license of the assets? [Yes] Spider [1] releases their dataset using
30 Apache-2.0 license
- 31 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
32 Dataset is released at <https://github.com/shirley-wu/daco>
- 33 (d) Did you discuss whether and how consent was obtained from people whose data you're
34 using/curating? [Yes] The annotators are trained and we have obtained their consent
- 35 (e) Did you discuss whether the data you are using/curating contains personally identifiable
36 information or offensive content? [Yes] The data is public data and does not contain
37 personally identifiable information or offensive content

38 5. If you used crowdsourcing or conducted research with human subjects...

- 39 (a) Did you include the full text of instructions given to participants and screenshots, if
40 applicable? [No] Annotations are done through internal annotators and instructions are
41 potentially confidential information
- 42 (b) Did you describe any potential participant risks, with links to Institutional Review
43 Board (IRB) approvals, if applicable? [N/A]
- 44 (c) Did you include the estimated hourly wage paid to participants and the total amount
45 spent on participant compensation? [No] Annotations are done through internal annota-
46 tors and costs are potentially confidential information

47 **A Appendix**

48 Website is at <https://shirley-wu.github.io/daco/index.html>. Data and code are released
49 at <https://github.com/shirley-wu/daco>. Croissant metadata record is at <https://github.com/shirley-wu/daco/blob/main/data/croissant.json>. We license our resources under
50 Apache-2.0 license.
51

52 We thereby state that we bear all responsibility in case of violation of rights, etc., and confirmation of
53 the data license.

54 **B Dataset Documentation**

55 Below are dataset documentation following the framework from datasheets for datasets:

56 **Motivation:**

- 57 • *For what purpose was the dataset created?* - For the novel task of data analysis as explained
58 in the main content.
- 59 • *Who created the dataset and on behalf of which entity?* - This dataset is created during a
60 collaboration of ByteDance AI Lab and University of California, Los Angeles.
- 61 • *Who funded the creation of the dataset?* - ByteDance

62 **Composition:**

- 63 • *What do the instances that comprise the dataset represent?* - Each instance contains a
64 database of tabular data, a question, a reasoning process including code snippets, and a final
65 answer. Everything is in text format, except the database is stored as `pd.DataFrame`
- 66 • *How many instances are there in total?* - As detailed in Table 1 in the main content.
- 67 • *Does the dataset contain all possible instances or is it a sample of instances from a larger*
68 *set?* - The dataset is not a sample from a larger set.
- 69 • *What data does each instance consist of?* - Raw data.
- 70 • *Is there a label or target associated with each instance?* - Yes, as explained in Section 3 in
71 the main content.
- 72 • *Is any information missing from individual instances?* - No
- 73 • *Are relationships between individual instances made explicit?* - N/A
- 74 • *Are there recommended data splits?* - Yes. We split the dataset randomly, and encourage
75 people to follow this split for reproductivity. We also curate human annotations only for the
76 test set.
- 77 • **Are there any errors, sources of noise, or redundancies in the dataset?** - Yes. The input
78 questions and answer annotations are generated by ChatGPT, which will inevitably contain
79 errors. We try to manage the affect by manually filtering the questions, and by curating a
80 test set of human refined answer annotations.
- 81 • *Is the dataset self-contained, or does it link to or otherwise rely on external resources?* -
82 Self-contained.
- 83 • *Does the dataset contain data that might be considered confidential?* - No.
- 84 • *Does the dataset contain data that, if viewed directly, might be offensive, insulting, threaten-*
85 *ing, or might otherwise cause anxiety?* - No.
- 86 • *Does the dataset identify any subpopulations?* - No.
- 87 • *Is it possible to identify individuals, either directly or indirectly from the dataset?* - No.
- 88 • *Does the dataset contain data that might be considered sensitive in any way?* - No.

89 **Collection process:** as described in Section 3 in the main content.

90 **Preprocessing/cleaning/labeling:** we release the raw text data and do not perform any preprocess-
91 ing/cleaning/labeling of the texts.

92 **Uses:**

- 93 • *Has the dataset been used for any tasks already?* - The data analysis task, as in the main
94 content
- 95 • *Is there a repository that links to any or all papers or systems that use the dataset?* -
96 <https://github.com/shirley-wu/daco>
- 97 • *What (other) tasks could the dataset be used for?* - As in the main content
- 98 • *Is there anything about the composition of the dataset or the way it was collected and*
99 *preprocessed/cleaned/labeled that might impact future uses?* - N/A
- 100 • *Are there tasks for which the dataset should not be used?* - N/A

101 **Distribution:**

- 102 • *Will the dataset be distributed to third parties outside of the entity on behalf of which the*
103 *dataset was created?* - Yes
- 104 • *How will the dataset will be distributed?* - <https://github.com/shirley-wu/daco>
- 105 • *When will the dataset be distributed?* - Already released
- 106 • *Will the dataset be distributed under a copyright or other intellectual property (IP) license,*
107 *and/or under applicable terms of use (ToU)?* - Yes, Apache-2.0 license
- 108 • *Have any third parties imposed IP-based or other restrictions on the data associated with*
109 *the instances?* - No
- 110 • *Do any export controls or other regulatory restrictions apply to the dataset or to individual*
111 *instances?* - No

112 **Maintenance:**

- 113 • *Who will be supporting/hosting/maintaining the dataset?* - The dataset is not planned to be
114 a dynamic dataset, but the authors will keep maintaining the github repo
- 115 • *How can the owner/curator/manager of the dataset be contacted?* - Github or email
116 xueqing.wu@cs.ucla.edu
- 117 • *Is there an erratum?* - No
- 118 • *Will the dataset be updated?* - No unless to correct errors
- 119 • *Will older versions of the dataset continue to be supported/hosted/maintained?* - N/A
- 120 • *If others want to extend/augment/build on/contribute to the dataset, is there a mechanism*
121 *for them to do so?* - Yes, please feel free to do that as long as citing our work and following
122 the license

123 C Limitations

124 While we put forth to construct the first of its kind dataset, DACO, for the comprehensive data
125 analysis task, the challenging nature of the data analysis process (which often requires certain
126 domain expertise) itself presents two major limitations of this work: **(1) It is expensive to build
127 expert-annotated dataset.** In our work, the large-scale annotations are automatically generated by
128 GPT-4, and their quality cannot always be well guaranteed. Although we curate a test set refined
129 by humans, those answers are initially generated by GPT-4, which may introduce biases to human
130 annotators during the refinement process for the final answer. Annotations curated by experts from

131 scratch might have higher quality, but they are indeed quite costly and can create other sort of
 132 alignment problems. **(2) It is nontrivial to evaluate model generations.** Evaluating the quality of
 133 the analyses is by itself challenging and requires data science expertise. For automatic evaluation, we
 134 use ChatGPT to rank the helpfulness of model generations, which can partially but not perfectly align
 135 with human preference. Additionally, ChatGPT cannot always robustly evaluate the correctness of
 136 model generations. We use an off-the-shelf NLI model to evaluate the entailment probability between
 137 human-refined ground truths and the model generations, which can partially reflect the correctness of
 138 model generations. However, the entailment probability prediction can sometime propagate errors
 139 which lead to false positives or negatives. We make efforts to alleviate such an issue by additionally
 140 collecting human evaluations, which are supposed to better reflect the answer quality, despite that
 141 humans can occasionally exhibit subjective evaluation patterns. Notice that our annotators do not
 142 fully check the correctness of the generated answers, where we task them to focus more on the
 143 helpfulness metrics defined in this work.

144 D Implementation Details

145 For **zero-shot API-based systems** including ChatGPT and GPT-4, we evaluate two settings, directly
 146 reading the table content, and using code generation. For the former setting, we linearize the table
 147 content into text representation as model input. Due to token limit, we feed the first 20 rows as input,
 148 which covers the full content of 93% tables. For the code generation setting, we employ the pipeline
 149 described in Figure 2(b) in the main content. When the generated code causes a syntax or runtime
 150 error, we re-sample the model until the generated code can be executed. We allow up to 5 resamplings
 151 for each turn. We use the `gpt-3.5-turbo-16k-0613` API for ChatGPT and `gpt-4-32k` API for
 152 GPT-4. We limit the number of total coding turns maximally at 9. For annotation generation where
 153 GPT-4 self-correction is allowed, we limit the number of self-correction within 2 for each turn and 4
 154 for the whole session.

155 For **finetuned models** including SFT, RLHF and fine-grained RLHF, we use CodeGeeX2-6B [2] as
 156 the base model. We first train the SFT model using GPT-4 annotations, and then train our RLHF
 157 models on top of the SFT model. When training R_{a+c} and R_r , we initialize the model from the SFT
 158 model. When training our fine-grained RLHF model, we initialize the value model V from R_{a+c} ,
 159 and initialize the policy model π from the SFT model. In inference, we use nucleus decoding with
 160 $p = 0.9$ and temperature = 1.0. Similarly, we allow up to 5 resamplings when the generated code
 161 causes an error. The SFT model is trained with 8 A100 GPU for about 4 hours. The RLHF models
 162 are trained with 8 A100 GPU for about 18 hours. Detailed hyper-parameters are in Table 1. The
 163 only hyper-parameter we tune is λ for fine-grained RLHF. We experiment with 0.8, 0.9 and 1.0 and
 164 discover that 1.0 works the best.

	SFT	RL
learning rate	1e-5	2e-6
gradient accumulation	4	4
total steps	600	200
λ	-	1.0
γ	-	1.0

Table 1: Hyperparameters.

165 E Qualitative Examples

166 We show final answers generated by SFT and RLHF in Figure 1. RLHF better focuses on user query,
 167 while SFT tends to display generic statistics that are less relevant to user query.

168 We show examples of code generations in Figure 2. We also report their reward scores from
 169 contribution RM and regularization RM.

As a game developer, I want to balance the stats of different Pokemon to ensure fair gameplay.

SFT	RLHF
<p>Findings</p> <ol style="list-style-type: none">1. There are 8 attributes (height, weight, hp, attack, defense, sp atk, sp def, speed) in the dataset that have strong positive or negative correlations, indicating that certain attributes are related.2. The pairs of attributes with strong positive or negative correlations are as follows:<ul style="list-style-type: none">- height and weight (correlation: 0.65)- hp and attack (correlation: 0.45)- hp and defense (correlation: 0.44)- hp and speed (correlation: 0.34)3. The pairs of attributes with strong negative correlation include:<ul style="list-style-type: none">- sp atk and defense (correlation: -0.6)- sp def and speed (correlation: -0.5) <p>Suggestions</p> <ol style="list-style-type: none">1. For balance purposes, you could consider assigning a higher weight to taller Pokemon and a higher percentage of the total weight for Pokemon, as height and weight are strongly correlated with each other.2. To ensure that the gameplay remains diverse and fair, you should also assign appropriate attack, defense, sp atk, sp def, and speed values to the Pokemon, taking into consideration their correlations with the other attributes.3. Since the Pokemon attributes are strongly correlated, consider implementing a more complex system that takes advantage of this relationship, such as a network or genetic algorithm-based approach to balance the attributes of diverse Pokemon species.	<p>Findings</p> <ol style="list-style-type: none">1. The majority of Pokemon are Water (38.13%), followed by Normal (24.72%) and Grass (13.82%).2. Fire, Psychic, and Steel types are relatively less frequent in the dataset.3. Types with high positive correlations with main stats (HP, attack, defense, and speed) include Water, Normal, and Grass types, while types with high negative correlations include Fire, Psychic, and Steel types. <p>Suggestions</p> <ol style="list-style-type: none">1. Consider balancing the game around the most common types (Water, Normal, and Grass) to make the gameplay more diverse and interesting for different types of strategies.2. Since types with high positive correlations have higher HPs, attacks, and speeds, it may be beneficial to focus on the Pokemon with those stats to balance the game accordingly.3. Monitor the performance of high-positive correlation types (Water, Normal, and Grass) to ensure their representation in the game doesn't lead to an imbalance in terms of stats or gameplay strategy.4. Monitor the performance of high-negative correlation types (Fire, Psychic, and Steel) to ensure their representation in the game isn't detrimental due to a potential imbalance in game mechanics.

Figure 1: Case study.

170 F GPT Prompts

171 Here we show the prompts we use for ChatGPT and GPT-4. Prompt for query generation is in Table 2.
172 Prompt for helpfulness annotation collection is Table 3. Prompt for helpfulness evaluation is Table 4.

173 References

- 174 [1] Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene
175 Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. Spider: A large-scale
176 human-labeled dataset for complex and cross-domain semantic parsing and text-to-SQL task. In
177 *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*,
178 pages 3911–3921, Brussels, Belgium, October–November 2018. Association for Computational
179 Linguistics.
- 180 [2] Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Zihan Wang, Lei Shen,
181 Andi Wang, Yang Li, Teng Su, Zhilin Yang, and Jie Tang. Codegeex: A pre-trained model for
182 code generation with multilingual evaluations on humaneval-x. *CoRR*, abs/2303.17568, 2023.

I will examine the relationships from the correlation matrix, including the most attended concerts and the conductor's experience, as well as the orchestras with the highest official ratings and their conductors.

```
# Get the top 3 most attended concerts
top_attended = merged_df.nlargest(3, "Attendance")

# Get conductor names and age with most attended concerts
top_conductors = top_attended[["Name", "Age", "Year_of_Work", "Attendance" ].sort_values("Attendance",
ascending=False)

print("Top 3 Conductors with Most Attended Concerts:")
print(top_conductors)
print()

# Get top 3 orchestras with the highest official ratings
top_rated_orchestras = merged_df.nlargest(3, "Official_ratings_(millions)")

# Get conductor names and orchestras with the highest official ratings
top_rated_conductors = top_rated_orchestras[["Orchestra", "Name", "Age", "Year_of_Work",
"Official_ratings_(millions)"]].sort_values("Official_ratings_(millions)", ascending=False)

print("Top 3 Orchestras with Highest Official Ratings and Their Conductors:")
print(top_rated_conductors)
print()
```

Contribution RM: 5.05 Regularization RM: 5.15 Two RMs combined: 5.10

(a) A good case that receives high scores from both contribution RM and regularization RM.

Next, I will join the tables 'Order_Deliveries', 'Actual_Orders', 'Products', and 'Customers' to create a single DataFrame with the necessary information to analyze the most popular products across delivery routes.

```
# Join tables
orders_deliveries = db["Order_Deliveries"]
actual_orders = db["Actual_Orders"]
products = db["Products"]
customers = db["Customers"]

# Merge tables
merged_df = pd.merge(orders_deliveries, actual_orders, on="actual_order_id")
merged_df = pd.merge(merged_df, products, left_on="actual_order_id", right_on="product_id")
merged_df = pd.merge(merged_df, customers, left_on="driver_employee_id", right_on="customer_id")

# Display a sample of 5 rows
print(merged_df.head())
```

Contribution RM: -2.63 Regularization RM: 6.59 Two RMs combined: 1.98

(b) A bad case that receives low score from contribution RM and high score from regularization RM.

Based on the analysis, provide a conclusion summarizing the findings and suggesting possible reasons for the high resettlement rates in the top 5 countries.

```
print("Conclusion:")
print("The average resettlement rate in the top 5 countries is 593.8%, which is quite high compared to the
average resettlement rate of 553.8% in the entire dataset.")
print("This is due to the high resettlement rates in the top 5 countries, which may be due to the following
reasons:")
print("1. The top 5 countries have a higher population, which may indicate a higher demand for migration.")
print("2. The top 5 countries have specialized in Industry or Commerce, which might require a high number of
workers to be migrated.")
print("3. The top 5 countries have a significant history of migration, which may indicate a diverse population
and a need for migration.")
```

Contribution RM: 8.35 Regularization RM: -5.87 Two RMs combined: 1.24

(c) A reward hacking case that receives high score from contribution RM and low score from regularization RM.

Figure 2: **Qualitative examples** of code generations, and their scores assigned by reward models.

I have a database of [database title]. I am a stakeholder and I am analyzing the database to make a decision. Who am I and what decision might it be? List 10 possibilities in a numbered list.

Each point should introduce who I am and briefly explain my intention in this format: As a/the [who I am], I want to [explain my intention]

Examples:

Based on the extracurricular activities database:

1. As the dean of student affairs, I want to decide on extracurricular activities to promote or cut
2. As the department head, I want to decide on faculty advisor assignments
3. As the school administrator, I want to review and revise faculty activity engagement

Based on a diabetes database:

1. As a healthcare policy maker, I want to decide on healthcare resource allocation
2. As a NIH official, I want to decide on medical research funding
3. As a health insurance actuary, I want to improve health insurance pricing strategy
4. As a health provider, I want to decide on patient care and treatment

Based on an allergy database:

1. As a catering manager, I want to plan meal options
2. As the school principal, I want to plan allergy awareness programs
3. As an administrator in the Student Affairs or Housing department, I want to decide on housing assignments
4. As the school administrator, I want to improve campus emergency preparedness
5. As the school principal, I want to develop policies for allergy accommodations

Based on a Home Equity Line of Credit (HELOC) product database, you can:

1. As the credit risk manager, I want to modify the credit underwriting policy

The database is as follows:

Database `[title]` has [x] tables. Table names are: [aaa], [bbb], [ccc]

Table `[caption]` has [x] rows and [y] columns. Column are:
`[column name]`, example values: [value 1], [value 2], [value 3], [value 4], [value 5]

...

Table 2: Prompt for **query collection**.

I have a database of [database title]. As a [stakeholder role], I want to [describe intention].

Given below two findings/conclusions, which one is more helpful to my analysis?

- * [answer bullet point 1]
- * [answer bullet point 2]

Your response should be in the following format:

- * Reasoning: <explain your reasoning here>
- * Answer: <repeat the more helpful finding here>

Table 3: Prompt for **helpfulness annotation collection**.

I have a database of [database title]. As a [stakeholder role], I want to [describe intention].

I have hired two data analysts to perform the analysis, and they gave me two different reports (listed below). Each report consists of two lists, one for findings and one for suggestions. Which one is more helpful to my analysis? When evaluating helpfulness, you should consider the following three rubrics in decreasing priority: (1) relevance to my analysis goal; (2) insightfulness; and (3) diversity of perspectives, especially for suggestions.

Your response should be in the following format. Note: <answer> should be either Report-1 or Report-2

* Answer: <answer>

* Reasoning: <explain your reasoning here>

The reports are as follows:

Report-1

[report 1]

Report-2

[report 2]

Table 4: Prompt for **helpfulness evaluation**.