CodeXData: Do Code Generating Language Models Understand Data?

Anonymous ACL submission

Abstract

Large language models (LLMs) are effective at code generation. Some code tasks, such as data wrangling or analysis, can be data-dependent. We introduce two novel taxonomies to characterize (1) the extent to which a code generation task depends on data and (2)the effect of data redaction. We curate two new datasets for Python code generation from natural language for data-centric tasks. We evaluate these datasets by varying configurations over our taxonomies and find that code generation performance varies based on the task class, data redaction, and prompting strategy. This is the first empirical measurement of the impact of data in the NL-to-code setting using LLMs for data-centric tasks.

1 Introduction

005

007

012

017

021

033

037

Large language models (LLMs) like OpenAI's Codex model (Chen et al., 2021) have demonstrated impressive results on coding tasks, including code snippet generation and competitive coding (Li et al., 2022a). Within this space, some code generation tasks are dependent on a concrete dataset. Data wrangling and tabular data analysis, often found in computational notebooks and in spreadsheets, are two such examples.

It is difficult to complete a task such as "*extract the street name from each address*" without knowing the data the task is targeting–can we expect code generating models to work without data? To answer this question we study whether augmenting prompts with data can unlock latent capabilities within LLMs. "*Extract the street name from each address*" is an under-specified query, but when paired with address data, Codex can synthesise the correct regular expression.

Although prompting with data has potential to improve code generation, there are practical considerations that limit data availability. Resources such as bandwidth, memory, and time all impose constraints; a client can be limited to a subset of a large dataset, or the prompt can be limited to a small set of tokens. Privacy concerns can restrict the types of data that can be used, such as personally identifiable information, and where or how the data can be transmitted. The data source itself may impose constraints, such as data obtained from ever-changing streams where only a snapshot can be provided to a model.

041

042

043

044

045

047

049

051

054

057

058

060

061

062

063

064

065

066

067

068

069

070

071

072

074

075

076

078

079

Code generation using LLMs is therefore a combination of the data required by a given task, and the data available for the task.

We introduce two new taxonomies for code generation from natural language to capture the nuances of these data requirements.

The first concerns the extent to which interpreting a task depends on the data presented to the task. A task can be *data independent*, where the query itself provides sufficient information for the model, such as "*return the unique values from column Location*". A task can be *data dependent*, such as extracting the street name from an address (described previously), where the query lacks data-specific information required by the model. Finally, a task can be *externally dependent*, where the query requires knowledge outside the task (and a particular API if a certain API is considered) such as "*create a new column that verifies that Year is an election year in the USA*".

Second, we define a taxonomy for redacting task data. Given our focus on queries over tabular data, we define our taxonomy in terms of the column headers and data rows of a table. Following data redaction, a task has access (1) to either the original column headers or simply an anonymous schema (generated headers such as "Coll"), and (2) to none, some, or all of the table rows.

To explore the impact of our taxonomies on code generation, we curate two new datasets. Each datapoint in our datasets consists of a textual query, a data input (column-major-flat table), and

- 120 121 122
- 122 123

124

125

126 127

12

128

queries based on questions arising on the Stack-Overflow and Mr. Excel² forums.
Using these datasets and taxonomies we evaluate CODEXDATA, our algorithm to turn natural language into code for querying tabular data.
CODEXDATA queries Codex with a prompt that combines both task description and task data to generate Python code that uses the Pandas API as shown in Figure 1. We chose Codex because there is significant evidence for its performance from prior research and from practical usage in Github

an expected correct output (extra columns). Our

first dataset is sourced from questions on Stack-

Overflow¹ requesting help to solve a problem in

spreadsheets. These tasks are not specific to the

languages such as Python or SQL that are predom-

inantly used to train language models. To better

understand performance on type-specific atomic

tasks or the impact of noise in the data, we created

a second synthetic dataset with curated types and

Copilot (Friedman, 2021). We select Pandas because Codex was trained on a significant amount of Python, and Pandas is likely the most common table manipulation library for Python. As our study is a first step towards understand-

As our study is a first step towards understanding the relationship of code generation with data in LLMs, we scoped our exploration to enable a focused analysis, and to mitigate confounders. Expanding to multi-table inputs or hierarchical tables requires further research into how to encode the data, and we do not have any guarantees that our conclusions will hold in other settings.

In summary, our central contribution is the first empirical study that measures the impact of data (single column-major-flat table) in the NL-to-code (NL-to-Pandas) setting using LLMs (Codex) for data-centric tasks. To support this, we make the following additional contributions:

- We define two novel data-centric task taxonomies. The first characterises the data required to complete a task. The second characterises the data redacted for a given task.
- We curate two new datasets to explore the impact of our taxonomies on code generation. The first data set is a collection of real spreadsheet problems from StackOverflow. The second dataset is a synthetic bench-

mark designed to evaluate performance on the different types of data, and the impact of noise.

129

130

131

132

133

134

135

136

137

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

157

158

159

160

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

• We study how data redaction affects the performance (pass@k) for different task classes. On data-dependent tasks, redaction halves pass@k, whereas it makes little difference to data independent tasks (Figure 2). Overall, removing the data leads to a drop of 20-45% (depending on the dataset), but it takes just one example row, with our prompting method, to restore performance to within 5% of the full data performance (Figure 3).

2 Background and Related work

Many systems based on LLMs, such as Codex (Chen et al., 2021) or PaLM (Chowdhery et al., 2022), tackle code generation. These pre-trained models can be enhanced with the goal of improving the performance of code generation in specific domains (e.g. for Python or SQL: Jigsaw (Jain et al., 2022), CodeT (Chen et al., 2022), CodexDB (Trummer, 2022), CCTest (Li et al., 2022b)). There is a body of work focusing on exploring fundamental questions related to code generation with LLMs, such as the impact of training on a large number of diverse tasks (Chan et al., 2022), the length generalization (Anil et al., 2022), and the compositional generalization (Shi et al., 2022). It has been shown that LLMs can also accomplish tasks for which they were not explicitly trained. Narayan et al. (2022) show how to convert data tasks to text generation. Additional examples of such tasks that LLMs can do include reverse engineering (Pearce et al., 2022), API generation (Hadi et al., 2022), and generation of development tools (Bareiss et al., 2022).

We explore the question of: *how does the inclusion of data impact code generation for datacentric tasks?* We focus on tasks where we can use the Python Pandas API on column-major-flat tables to produce new additional data columns.

This focus differentiates us from prior work related to data-centric code generation. Related prior work has focused on SQL tasks (Rajkumar et al., 2022) from the SPIDER benchmark (Yu et al., 2018) or in-place data transformations that do not yield a new dataframe column (Narayan et al., 2022). Prior datasets for Python generations (e.g. APPS (Hendrycks et al., 2021) or HumanEval (Chen et al., 2021)) are not data-centric,

¹https://stackoverflow.com/questions/tagged/excelformula

²https://www.mrexcel.com/board/forums/excelquestions.10/



Figure 1: This figure illustrates the components of the CODEXDATA system that transforms the user input table and query into a list of valid completions. The data is redacted to anonymize private information. A data description is extracted from the input table. The resulting redacted data, data description and query are used to construct a prompt which is fed to a code synthesis LLM such as Codex, generating multiple possible completions. The outputs of these completions are then validated and the first k valid completions (along with the outputs) are returned to the user.

while the tasks from prior Pandas generations (e.g. Jain et al. (2022) or Zan et al. (2022)) do not have sufficient and comparable complexity, or the amount of data that we have used in our analysis (See Section 4.1 for details).

179

185

186

190

191

194

195

Furthermore, we also explore the impact of both noise and data redaction on data-centric tasks by introducing novel datasets. Other related work has either left out consideration of these important issues, or similarly to Narayan et al. (2022), they have only briefly explored the role of noisy table attributes in LLM performance for entity matching. Narayan et al. (2022) also identify the challenges of data privacy in LLM-based data wrangling as future research, but they do not experimentally explore the impact of privacy-related decisions on performance.

3 CodeXData: Technical Approach

197Algorithm 1 presents our inference algorithm,198with a high-level overview shown in Figure 1. The199algorithm takes as input a user query Q as text,200user input table T as a Pandas dataframe, and the201target cardinality k of distinct completions to gen-202erate. To ensure termination within a reasonable

Algorithm 1 CodeXData Inference Algorithm

Input: Explicit: user query Q, input table T, cardinality k. Implicit: completion limit k_{max} (with $k \le k_{max}$).

Output: Pair of lists (C, O), with $|C| = |O| \le k$, of unique completions and their corresponding outputs.

1, K)	procedure CODEADATA(Q ,	1:
▷ redact the input table	$R \leftarrow \text{Redact}(T)$	2:
▷ extract data description	$D \leftarrow \text{Descript}(T)$	3:
▷ prompt creation	$P \leftarrow \text{Prompt}(Q, R, D)$	4:
▷ initialize budget, caches	$B, C, O \leftarrow k_{max}, [], []$	5:
do	while $B > 0 \land C < k$	6:
▷ sample completion	$c \leftarrow \text{CODEX}(P)$	7:
▷ decrement budget	$B \leftarrow B - 1$	8:
\triangleright execute against T	$o \leftarrow \text{EXEC}(c, T)$	9:
$\notin C$) then	if VALIDATE $(o) \land (c$	10:
\triangleright append completion to C	$C \leftarrow C + [c]$	11:
\triangleright append output to O	$O \leftarrow O + [o]$	12:
	return (C, O)	13:

time, we set a limit k_{max} on the number of calls to CODEX ($k_{max} = 8k$ by default).

As shown in Figure 1, the query Q can be "create a new column with the first names" and T can be DataFrame({"Names":["Mary M", "Emma G"]}). k can be set to 10.

At a high-level the algorithm: *redacts* the input table (line 2) using a REDACT; *extracts* a data description of the table using a DESCRIPT (line 203

3); combines the query Q, redacted data R, and 212 data description D to create a prompt P (line 4); queries CODEX repeatedly using this prompt until the target completions are reached or we exceed the budget of calls (lines 6-12). Each completion c(line 7) is executed on the input table (line 9) using a EXEC procedure, and if the completion is new and its output o satisfies a VALIDATE procedure, the two are accumulated in C and O. The lists of completions and outputs are returned to the user (line 13). We now describe the key components in detail.

213

214

215

216

217

218

219

222

233

234

235

241

242

244

245

246

247

249

253

254

3.1 Data Redactor Procedure: REDACT

As discussed in Section 1, the need for and extent of data redaction depends on privacy and resource constraints. Users can control redaction via procedure REDACT, which creates a redacted data Rfrom the original input table T.

We explore variations of R to examine the trade-off between data redaction and model performance. For our running example, the redacted data in no data (+ anonymous headers) and no data (+ actual headers) case would be $R = \text{DataFrame}(\{\text{"Anon":[]}\}) \text{ and } R =$ DataFrame({"Names":[]}) resp. The user can also provide a subset of rows as R =DataFrame({"Names":["Mary M"]}) or the full dataframe with no redaction.

3.2 Data Descriptor Procedure: DESCRIPT

A descriptor procedure DESCRIPT extracts a data description D from the input table T. The description can consist of: number of elements in a column, column type, basic predicates satisfied by all elements of a column (e.g. all lowercase) or a regular expression that captures a pattern of the elements in the column (Padhi et al., 2018).

3.3 Prompt Creation Procedure: PROMPT

The procedure PROMPT creates a textual prompt by combining the user's query, redacted data, and data description. We consider three strategies for creating this prompt in addition to the baseline prompt. The baseline prompt simply concatenates the redacted data R, and the textual query Q.

Instruction. Instruction prompts (Mishra et al., 2021) bias the model by appending instructions to the baseline prompt, emphasizing that the comple-257 tion should have a particular property, e.g. "write an executable, type correct pandas solution". (See Table 1 in Appendix D.)

Example-usage. These prompts add concrete pandas usage examples in a few-shot style. We prepend to the baseline prompt examples that emphasize a property. (See Table 2 in Appendix D.)

261

262

263

265

266

267

268

269

270

271

272

273

274

275

276

277

278

281

282

289

290

291

292

293

294

295

297

299

300

301

302

303

304

305

Data-description. Data-description prompts augments the baseline prompts with the data description D. (See Table 3 in Appendix D.)

3.4 Completion Procedure: CODEX

The procedure CODEX queries Codex's Da Vinci engine, passing the prompt P, a temperature derived from k, and stop sequences found to allow Codex to generate at least one solution while not usually using the entire token budget. The completion then undergoes clean-up, such as truncating at the first comment after executable code, and removing whitespace and comments. Further details can be found in Appendix E.5.

3.5 Execution Procedure: EXEC

The procedure EXEC constructs an executable program by (1) rewriting the completion c to expose the likely answer, obtaining c'; and (2) appending c' to the input table T. It executes the program in a sandbox to obtain the final output o. Further details can be found in Appendix E.6.

3.6 Validating Completions: VALIDATE

We validate the extracted output (e.g., for expected output type) using an output validator VALIDATE. The completions that pass the validator are called valid completions and are returned to the user.

4 **Experiments**

We design our experiments to study how CodeX-Data depends on the data and task class. In particular, we want to understand how the performance changes with (a) the level of data redac $tion^3$, (b) the data description present in the query, (c) the noise level of the data, across the three task classes. The task classes are of increasing difficulty: (1) data-independent (IND), (2) datadependent (DEP), (3) external-dependent (EXT). We also evaluate how different prompting strategies can bridge the gap between the data required and the data available.

4.1 Our Datasets

One of our main contributions is the set of Pandas datasets we curate to answer our research

³Full-data access, no-data access, subset-data access (in all our experiments we use the first row only as subset).

357

358

questions. We curate two datasets, SOFSET and TYPESET, which consist of real world problems 307 across diverse domains. In addition to the three 308 task classes mentioned before, we classify our dataset tasks based on data types (strings, numbers names, dates, units, addresses and mixed types) 311 that are popular for data-centric tasks (Table 4 in 312 the appendix). Our datasets are larger and likely 313 more diverse than either of two existing Pandas 314 datasets, JIGSAW (Jain et al., 2022) and CERT 315 (Zan et al., 2022). JIGSAW has 79 unique tasks 316 with a median of 7 rows in the data and CERT has 317 100 unique tasks with a median of 3 rows. On the 318 other hand, SOFSET and TYPESET have 200 and 319 64 unique tasks, with a median of 10 and 20 rows, respectively.

322 SOFSET. This is a collection of 200 real world tasks from StackOverflow that have the tag "Ex-323 celFormulas" where the user has asked about 324 spreadsheets tasks. We either use the data the 325 326 user has added to their query, or if it is ambiguous, we create new data that matches the query. 327 We also add extra data to make the data have at 328 least 10 rows. We manually summarise the ver-330 bose ask into a concise query based on the title of the post and annotate with the desired output. We remove post identifiers for anonymization. Across 332 all task classes, SOFSET contains complex real world problems with various corner cases, especially the ones in DEP and EXT task classes.

TYPESET. To evaluate performance across the three task classes, we curate the TYPESET dataset with type-specific tasks. This dataset is a collection of 64 tasks (queries plus pandas solutions) mostly of type dates and addresses. These type-specific columns are sourced from Sherlock (Hulsebos et al., 2019), AutoType (Yan and He, 2018) or Wikidata (Vrandecic and Krötzsch, 2014). We manually annotate the queries along with Pandas solutions and outputs.

336

337

338

339

340

341

342

343

346**TYPESETNOISY.** We create a noisy subset of the347TYPESET with 131 tasks to evaluate how perfor-348mance changes in the presence of noise in the349data. We intentionally altered approximately 20%350of the values in the input columns to investigate351the model behaviour in scenarios inspired by real352world scenarios. We induce the following noise353scenarios: corruption, by inserting a space, dash,354or a new line at a random position in the input355strings; missing, by replacing values in an input356column with an empty string; and mixformat, by

changing column formats, specifically we mix ints and floats or combine different date/address formats.

4.2 Evaluation Metrics

We report correctness based on whether the generated code produces the expected output. Consider a single datapoint consisting of a query Q, an input table T, and an expected correct output o. The probability that at least one of k inferred outputs is correct is called pass@k (Chen et al., 2021). More formally, pass@k is the probability that $o \in O$ if we sample (C, O) from CODEXDATA(Q, T, k).

To measure this probability empirically for each datapoint, we compute up to 2k valid programs by sampling (C, O) from CODEXDATA(Q, T, 2k). We count the number s of occurrences of o in O, and hence compute an estimate of pass@k as $1 - \binom{2k-s}{k} / \binom{2k}{k}$ (Chen et al., 2021). (By computing 2k completions the estimate has lower variance than by simply computing k completions.) Each pass@k we report on a whole dataset is the average of pass@k over all its datapoints.

Additionally, we calculate a partial form of the metric written as pass@k(X%). It is defined on datapoints and datasets in the same way as pass@k except that correctness of an output column is relaxed from exact to partial matching: at least X% of the inferred output matches the expected output. Hence, pass@k(100%) is the same as pass@k. This partial metric is important in our setting as partially correct output columns may still be valuable to the user, although we need to be mindful when designing the user interaction.

For detailed error analysis, we report a set of evaluation metrics considering code quality and errors. We conduct a sensitivity analysis of pass@k and report observed standard deviations for each task class with and without data redaction. Due to space limitations, we report these in the appendix (see Figure 13 in Appendix F).

4.3 Empirical Insights

We evaluate performance across tasks in the two datasets: SOFSET and TYPESET for different data redactions and prompting strategies. For the data redaction experiments, we use anonymized version of column names for the no-data (anonymous columns) redaction and proper column names for the other three redactions. The prompting experiments are done with full-data access unless specified otherwise. All evaluation results are



Figure 2: Impact on pass@k with data redaction for SOFSET and TYPESET. The results shown are grouped by the different task classes: data-independent (IND), data-dependent (DEP) and external-dependent (EXT). In each group, the different bars represent different levels of redaction: (a) full data, (b) a subset of the data (only one row) (c) no data and proper column names and (d) no data and anonymous column names. Redaction has negligible effect on IND tasks, and a large impact on DEP and EXT tasks.

averages of a single CODEXDATA run per data-407 point, computing 2k valid completions to estimate 408 pass@k or pass@k(X%). In the appendix, we re-409 port stability across several runs and detailed error 410 analyses including invalid completions.

411

412

413

414

415

416

417

Performance varies with task complexity. The pass@k for tasks in TYPESET is better than for SOFSET tasks across all experiments (Figures 2, 3, and 4). This is expected as SOFSET problems are more complex, especially for task classes DEP and EXT.

418 The level of data redaction has a high impact on **pass**@k. Figure 2 shows how pass@k (for k =419 (1, 5, 10) varies with the task and data redaction 420 for SOFSET and TYPESET. We observe that for 421 both datasets, irrespective of the task class, perfor-422 mance drops when the data redaction is stricter-423 the bars from (a) to (d) are in increasing degree 424 of redaction. The performance with no redaction 425 is almost always the best, and it degrades consid-426 erably when no data is available. However, the 427 performance degradation from no redaction (a) to 428 partial redaction (only first row (b)) is relatively 429 small in most cases, which implies that sampling 430 rows is almost as effective as providing all the 431 data. The difference between (c) no data + proper 432 column names and (d) no data + anonymous col-433 umn names is also relatively small which indicates 434 that having the proper column names is not as ben-435 eficial as having a sample of the data. 436

Performance on data-dependent and externaldependent tasks is highly affected by the level of data redaction. The different groups in Figure 2 show data redaction's high impact on performance across data-dependent task classes for both the datasets. The tasks under DEP and EXT classes are more complex, data-dependent and multi-step (e.g. doing complex operations after extracting the year from a column). We observe that the difference in performance between no and partial redaction (a and b) versus full redaction (c and d) is large for DEP and EXT classes. The performance difference is even more for TYPESET than SOF-SET since the latter has more complex tasks.

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

The impact of varying data redaction on dataindependent tasks is negligible. As expected, we do not observe a noticeable change in performance for tasks that are data-independent (IND), but the presence of data does provide extra context that is helpful. The IND tasks in TYPESET perform the best when only a subset of data is available (only first row) which further reinforces that these tasks do not depend on data.

Data description prompts counter data redaction. Figure 3 shows the interplay of data descriptions in the prompt under different data redactions. For each pair of same colored lines, the dashed line represents the result when a data description (in this case type information) is added to the original prompt. The data description prompt increases



Figure 3: Interplay between the varying amount of data redaction and how data description in the prompt helps with performance. For each redaction, we evaluate the performance by adding a data description to the prompt in the form of type information. Type information is especially helpful when we only have access to anonymous columns and just a subset of the data. The extracted data descriptions do not directly depict private information.



Figure 4: Impact of different prompting strategies on SOFSET and TYPESET compared to the baseline. We present a subset of prompting strategies that perform better than the baseline, detailed results in Figure 11 and Figure 12. Data description prompts (with type information and regular expression patterns that match the column elements) outperform the baseline for both the datasets.

performance in all cases but is particularly valuable in the case where all data is removed and columns anonymized. Surprisingly, adding a data description also has a positive impact in the case where a subset of the data is present. Thus we can use data description prompts to counter the effect of data redaction. The x-axis in Figure 3 represents how pass@10 varies with percentage of examples passed. Interestingly, we see a sharper decrease in pass@10(X%) in Figure 3 for SOFSET compared to TYPESET due to presence of complex corner case examples in the latter.

467

468

469

470

471

472

473 474

475

476

477

478

479 **Careful prompt engineering helps.** We explore

different prompting strategies (subsection 3.3) and complete results are given in Appendix (Figure 11, Figure 12). Figure 4 presents the prompting strategies that perform better than baseline. We observe that carefully extracting the data description (in the form of regular expressions, type information and information about the number of elements in column) helps for both the datasets. Interestingly, instructing the model to be concise, executable and generate iterative solutions also consistently improves performance across the two datasets. The best performing prompting strategy for SOFSET (regular expressions + type informa480

481

482

483

484

485

486

487

488

489

490

491

492



Figure 5: Impact on pass@10(X%) of the different noisy scenarios such as corrupt data values, missing values and mixed formats in TYPESETNOISY compared to the original data in TYPESET (baseline).

tion) has a pass@10 of 0.56 compared to the baseline pass@10 of 0.53. TYPESET has a relative performance improvement of 4%. We see that the performance gain is not huge. Prompting can also sometimes hurt as seen in case of regex(subset) in Figure 12c. This prompt adds too much data description (one regex per element) which likely ends up confusing the model. So prompting strategies are beneficial but only when curated carefully.

493

494

495

497

498

499

501

502

503

505

507

509

510

Performance varies considerably per data-type. The performance between data-types differs in pass@k, length and number of retrieved completions as well as number of executable completions. Tasks with string data-type for example have an average pass@k of 0.83 compared to 0.48 for tasks with dates data-type. Overall tasks with units data-type are most successful, address data-type least successful. (See Table 13)

Noisy data impacts performance. We evaluate 511 the impact on performance with noise additions in 512 TYPESETNOISY. Corresponding to the approxi-513 mately 20% noise in the data we see a drop in 514 pass@10(75%) and higher (Figure 5). Additional 515 difference compared to baseline is attributed to 516 runtime errors caused by the noise (e.g., an in-517 crease in IndexErrors seen in Table 12 in the ap-518 pendix).

Generation efficiency depends on data redaction and task class. Figure 6 shows the performance of pass@10 versus total number of completions from CODEX (including the invalid ones)
to generate 20 valid completions, for different
task classes and levels of data redactions across
datasets. Irrespective of the dataset and task class,



Figure 6: Interplay between pass@10 and total number of retrieved completions from CODEX for the different datasets (shapes), forms of data redaction (color) and task categories (labels). **Lower right**: Low pass@10 and many completions area: most results for runs with highly redacted data and data-dependent tasks lie here. **Upper left**: High pass@10 and few completions: most results for runs with no redaction and/or dataindependent tasks lie here. **Lower left/middle**: Low pass@10 and few completions: most results for runs for external-dependent tasks lie here.

access to data greatly reduces the number of generations needed; for instance, with full-data we need fewer than 50 completions irrespective of the task class. We need more completions for datadependent tasks versus data-independent tasks, and in general sampling more completions results in poorer pass@10. 527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

5 Conclusion and Future Work

Our work highlights the importance of data for code generation on data-centric tasks. We have empirically shown that the presence of data influences both the system's capability to generate good solutions and its efficiency. We have shown that optimizing prompts can improve performance, and several extensions are possible. Handling a broader problem space (e.g., multitable inputs, hierarchical table inputs, new table outputs) raises interesting challenges in how to best represent the data. We hypothesise that the gap between generations with no-data and fulldata cannot be fully eliminated, and thus it would be interesting to investigate how to integrate with approaches to privacy preservation.

References

550

553

554

555

559

563

564

568

569

570

571

573

576

577

578

579

583

584

586

588

590

591

592

594

595

596

597

599

602

605

- Philip E Agre et al. 1997. Lessons learned in trying to reform AI. Social science, technical systems, and cooperative work: Beyond the Great Divide, 131.
- Cem Anil, Yuhuai Wu, Anders Johan Andreassen, Aitor Lewkowycz, Vedant Misra, Vinay Venkatesh Ramasesh, Ambrose Slone, Guy Gur-Ari, Ethan Dyer, and Behnam Neyshabur. 2022. Exploring length generalization in large language models. *ArXiv*, abs/2207.04901.
- Patrick Bareiss, Beatriz Souza, Marcelo d'Amorim, and Michael Pradel. 2022. Code generation tools (almost) for free? a study of few-shot, pre-trained language models on code. *ArXiv*, abs/2206.01335.
- Emily M Bender and Alexander Koller. 2020. Climbing towards nlu: On meaning, form, and understanding in the age of data. In *Proceedings of the 58th annual meeting of the association for computational linguistics*, pages 5185–5198.
- Jun Shern Chan, Michael Martin Pieler, Jonathan Jao, J'er'emy Scheurer, and Ethan Perez. 2022. Few-shot adaptation works with unpredictable data. *ArXiv*, abs/2208.01009.
- Bei Chen, Fengji Zhang, A. Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. 2022. Codet: Code generation with generated tests. *ArXiv*, abs/2207.10397.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating large language models trained on code.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek B Rao, Parker Barnes, Yi Tay, Noam M. Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Benton C. Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari,

Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier García, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Oliveira Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Díaz, Orhan Firat, Michele Catasta, Jason Wei, Kathleen S. Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. Palm: Scaling language modeling with pathways. *ArXiv*, abs/2204.02311. 607

608

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

Michael Droettboom, Roman Yurchak, Hood Chatham, Dexter Chua, Gyeongjae Choi, Marc Abramowitz, casatir, Jan Max Meyer, Jason Stafford, Madhur Tandon, Michael Greminger, Grimmer Kang, Chris Trevino, Wei Ouyang, Joe Marshall, Adam Seering, Nicolas Ollinger, Ondřej Staněk, Sergio, Teon L Brooks, Jay Harris, Alexey Ignatiev, Seungmin Kim, Paul m. p. P., jcaesar, Carol Willing, Cyrille Bogaert, Dorian Pula, Frithjof, and Michael Jurasovic. 2022. Pyodide: A Python distribution for WebAssembly (0.19.0).

Nat Friedman. 2021. [link].

- Mohammad Abdul Hadi, Imam Nur Bani Yusuf, Ferdian Thung, Kien Gia Luong, Lingxiao Jiang, Fatemeh Hendijani Fard, and David Lo. 2022. On the effectiveness of pretrained models for api learning. 2022 IEEE/ACM 30th International Conference on Program Comprehension (ICPC), pages 309–320.
- Patrick M. Haluptzok, Matthew Bowers, and Adam Tauman Kalai. 2022. Language models can teach themselves to program better. *ArXiv*, abs/2207.14502.
- Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Xiaodong Song, and Jacob Steinhardt. 2021. Measuring coding challenge competence with apps. *ArXiv*, abs/2105.09938.
- Madelon Hulsebos, Kevin Hu, Michiel Bakker, Emanuel Zgraggen, Arvind Satyanarayan, Tim Kraska, Çağatay Demiralp, and César Hidalgo. 2019. Sherlock: A deep learning approach to semantic data type detection. In *Proceedings of the* 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. ACM.
- Naman Jain, Skanda Vaidyanath, Arun Iyer, Nagarajan Natarajan, Suresh Parthasarathy, Sriram Rajamani, and Rahul Sharma. 2022. Jigsaw: Large language models meet program synthesis. In *International Conference on Software Engineering (ICSE)*.
- Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal

747

748

719

Lago, Thomas Hubert, Peter Choy, Cyprien de Masson d'Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. 2022a. Competition-level code generation with alphacode.

- Zongjie Li, Chaozheng Wang, Zhibo Liu, Haoxuan Wang, Shuai Wang, and Cuiyun Gao. 2022b. Cctest: Testing and repairing code completion systems. *ArXiv*, abs/2208.08289.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, Yejin Choi, and Hannaneh Hajishirzi. 2021. Reframing instructional prompts to gptk's language. *arXiv preprint arXiv:2109.07830*.

679

681

683

691

705

706

710

712

713

714

715

718

- Avanika Narayan, Ines Chami, Laurel Orr, and Christopher R'e. 2022. Can foundation models wrangle your data? *ArXiv*, abs/2205.09911.
- Saswat Padhi, Prateek Jain, Daniel Perelman, Oleksandr Polozov, Sumit Gulwani, and Todd Millstein. 2018. Flashprofile: a framework for synthesizing data profiles. *Proceedings of the ACM on Programming Languages*, 2(OOPSLA):1–28.
- Hammond A. Pearce, B. Tan, Prashanth Krishnamurthy, Farshad Khorrami, Ramesh Karri, and Brendan Dolan-Gavitt. 2022. Pop quiz! can a large language model help with reverse engineering? *ArXiv*, abs/2202.01142.
- Nitarshan Rajkumar, Raymond Li, and Dzmitry Bahdanau. 2022. Evaluating the text-to-sql capabilities of large language models. *ArXiv*, abs/2204.00498.
- Advait Sarkar. 2022. Is explainable AI a race against model complexity? *arXiv preprint arXiv:2205.10119*.
- Advait Sarkar, Mateja Jamnik, Alan F Blackwell, and Martin Spott. 2015. Interactive visual machine learning in spreadsheets. In 2015 IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC), pages 159–163. IEEE.
- Kensen Shi, Joey Hong, Manzil Zaheer, Pengcheng Yin, and Charles Sutton. 2022. Compositional generalization and decomposition in neural program synthesis. *ArXiv*, abs/2204.03758.
- Immanuel Trummer. 2022. Codexdb: Generating code for processing sql queries using gpt-3 codex. *ArXiv*, abs/2204.08941.
- Denny Vrandecic and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. *Commun. ACM*, 57(10):78–85.
- Albert Webson and Ellie Pavlick. 2022. Do promptbased models really understand the meaning of their prompts? In *Proceedings of the 2022 Conference of the North American Chapter of the Association for*

Computational Linguistics: Human Language Technologies, pages 2300–2344, Seattle, United States. Association for Computational Linguistics.

- Jack Williams, Carina Negreanu, Andrew D Gordon, and Advait Sarkar. 2020. Understanding and inferring units in spreadsheets. In 2020 IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC), pages 1–9. IEEE Computer Society.
- Cong Yan and Yeye He. 2018. Synthesizing typedetection logic for rich semantic data types using open-source code. In *Proceedings of the 2018 International Conference on Management of Data*, SIG-MOD '18, page 35–50, New York, NY, USA. Association for Computing Machinery.
- Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir R. Radev. 2018. Spider: A largescale human-labeled dataset for complex and crossdomain semantic parsing and text-to-sql task. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018, pages 3911–3921. Association for Computational Linguistics.
- Daoguang Zan, Bei Chen, Dejian Yang, Zeqi Lin, Minsu Kim, Bei Guan, Yongji Wang, Weizhu Chen, and Jian-Guang Lou. 2022. Cert: Continual pretraining on sketches for library-oriented code generation. *ArXiv*, abs/2206.06888.

751

752

758

759

763

764

766

768

770

771

774

775

776

777

778

779

783

784

785

789

790

791

794

795

796

A Ethics Statement

There are broad ethical impacts resulting from the creation of AI models that attempt to generate code solutions from natural language descriptions and these are discussed in detail in previous papers including Codex (Chen et al., 2021), AlphaCode (Li et al., 2022a), and PaLM (Chowdhery et al., 2022). These impacts include overreliance, misalignment between what the user expressed and what they intended, potential for bias and under/over representation in the model results, economic impacts, the potential for privacy and security risks, and even environmental considerations. All of these considerations also apply to the work described here. Our focus is to highlight how the presence of data improves the performance of these models but it is important to note that the quality of the data used in the prompt will impact whether the resulting generation exhibits bias, exposes private data, etc. We explore the overall impact of providing data as part of the prompt but do not conduct a more focused analysis of determining how bias in the prompt data might influence the resulting code generation, a task we leave for future work.

We are wary of using the word 'understand' in this paper. It has been correctly argued that language models do not really 'understand' language in the sense of connecting language's syntactic content with the semantics of the physical world (Bender and Koller, 2020; Webson and Pavlick, 2022). There have long been critics of the use of such terms in AI research (Agre et al., 1997). Nonetheless, large language models have shown themselves in certain situations to be capable of the syntactic manipulation of language which in humans we take to be commonsense evidence of understanding. This is the less contentious manner in which we use the word. Thus our intention in using the word 'understand' is not to claim that models can connect data with real-world concepts, but rather that the model can manipulate language about data in a useful manner, where 'useful' is defined by our quantitative benchmarks. These benchmarks aim to be reflective of a qualitative notion of utility, but as with all quantitative benchmarks the process of operationalising a qualitative notion inevitably requires some reductionism.

This paper does not directly contribute to a tool built on the assumed capabilities of language models to understand data, but nonetheless, it is motivated by their potential applications in such tools. These tools may be deployed in many data applications such as databases, spreadsheets, and business intelligence applications. Depending on the audience of the tool, various interaction design concerns arise. Explainability of the model is a key consideration, and the tool should offer decision support to evaluate mispredictions and potential next steps (Sarkar, 2022). Previous research of non-experts using inference driven tools for data manipulation has shown the importance of tool design in the critical appreciation of the model and its limitations, and in the potential cost of errors (Williams et al., 2020; Sarkar et al., 2015). As an exploratory paper without a concrete application, we do not encounter these issues, but the project has nonetheless been reviewed by our institution's ethics board, which checks for compatibility with Microsoft's Responsible AI standard.⁴

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

There is the question of the sources of data and of consent to use the data in the manner exhibited in this paper. We have reviewed each of the datasets we have included in this paper to ensure that our use is compatible with the intent of the authors and publishers. Our datasets have also been reviewed by our institution's ethics board to review that this is an ethical use.

⁴https://www.microsoft.com/en-us/ai/ responsible-ai

B Datasheets for Datasets

B.1 MOTIVATION

827

829

830

832

833

834

836

840

841

845

848

850

851

852

855

857

859

861

864

870

872

873

874

For what purpose was the dataset created? Was there a specific task in mind? Was there a specific gap that needed to be filled? Please provide a description.

Answer: We created the new dataset to be able to evaluate how well LLMs perform on datacentric tasks when they are given various degrees of access to data. LLMs have been evaluated on Python generations in prior work - on puzzles, e.g. (Haluptzok et al., 2022), on programming contests, e.g. (Li et al., 2022a), on Python problems, e.g. (Chen et al., 2021) or (Hendrycks et al., 2021), and even Pandas tasks, e.g. (Zan et al., 2022) or (Jain et al., 2022). We are contributing this dataset as it contains larger dataframes and more diverse tasks than prior work, allowing us to draw meaningful conclusions for our research questions.

Who created this dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)?

Answer: to be added after ARR reviews to not break anonymity

What support was needed to make this dataset? (e.g.who funded the creation of the dataset? If there is an associated grant, provide the name of the grantor and the grant name and number, or if it was supported by a company or government agency, give those details.)

Answer: to be added after ARR reviews to not break anonymity

B.2 COMPOSITION

What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)? Are there multiple types of instances (e.g., movies, users, and ratings; people and interactions between them; nodes and edges)? Please provide a description.

Answer: Queries and data from Stackoverflow or synthetic data.

How many instances are there in total (of each type, if appropriate)?

Answer: We have collected a benchmark of 780 datapoints.

Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set? If the dataset is a sample, then what is the larger set? Is the sample representative of the larger set (e.g., geographic coverage)? If so, please describe how this representativeness was validated/verified. If it is not representative of the larger set, please describe why not (e.g., to cover a more diverse range of instances, because instances were withheld or unavailable).

Answer: Part of the dataset (SOFGoldSet) is a sample from all the queries in Stackoverflow under #ExcelFormula. The dataset was sampled such that we prioritize highly rated questions/answers, so we believe they are representative for real problems users have. The other dataset (GoldType-Set) was synthetically created with queries that are common for each type (e.g. extracting the year from a date) - as the dataset has fine-grained types (e.g. names and addresses) we have a small set of possible queries available.

What data does each instance consist of? "Raw" data (e.g., unprocessed text or images) or features? In either case, please provide a description.

Answer: Each datapoint comprises of a query (e.g. "count the number of clothing items"), data (e.g. a table of products and sales) and hand-annotated metadata (e.g. column names, type of query, type of data access).

Is there a label or target associated with each instance? If so, please provide a description. *Answer*: We provide the expected output (values in the newly generated column) and for GoldType-Set we also provide a Pandas solution.

Is any information missing from individual instances? If so, please provide a description, explaining why this information is missing (e.g., because it was unavailable). This does not include intentionally removed information, but might include, e.g., redacted text.

Answer:No

Are relationships between individual instances made explicit (e.g., users' movie ratings, social network links)? If so, please describe how these relationships are made explicit. *Answer*:No

Are there recommended data splits (e.g., training, development/validation, testing)? If so, please provide a description of these splits, explaining the rationale behind them.

Answer: This dataset is hand-curated for evaluation only.

Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide

979

a description.

928

929

930

931

932

933

934

936

937

938

941

943

945

949

951

952

953

956

957

960

961

962

963

964

965

967

968

969

971

972

973

975

976

977

978

Answer: Part of the dataset has been augmented (the data has been corrupted) to help us explore the impact of naturally occuring noise in users' work. We added missing data, corrupted data values, surplus columns and mixed formats.

Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g., websites, tweets, other datasets)? If it links to or relies on external resources, a) are there guarantees that they will exist, and remain constant, over time; b) are there official archival versions of the complete dataset (i.e., including the external resources as they existed at the time the dataset was created); c) are there any restrictions (e.g., licenses, fees) associated with any of the external resources that might apply to a future user? Please provide descriptions of all external resources and any restrictions associated with them, as well as links or other access points, as appropriate.

Answer: Self-contained

Does the dataset contain data that might be considered confidential (e.g., data that is protected by legal privilege or by doctor-patient confidentiality, data that includes the content of individuals' non-public communications)? If so, please provide a description.

Answer: No

Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety? If so, please describe why.

Answer: No

Does the dataset relate to people? If not, you may skip the remaining questions in this section. *Answer*: No

Does the dataset identify any subpopulations (e.g., by age, gender)? If so, please describe how these subpopulations are identified and provide a description of their respective distributions within the dataset.

Answer: No

Is it possible to identify individuals (i.e., one or more natural persons), either directly or indirectly (i.e., in combination with other data) from the dataset? If so, please describe how. *Answer*: No

Does the dataset contain data that might be considered sensitive in any way (e.g., data that reveals racial or ethnic origins, sexual orientations, religious beliefs, political opinions or union memberships, or locations; financial or health data; biometric or genetic data; forms of government identification, such as social security numbers; criminal history)? If so, please provide a description. *Answer*: No

B.3 COLLECTION

How was the data associated with each instance acquired? Was the data directly observable (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly inferred/derived from other data (e.g., part-ofspeech tags, model-based guesses for age or language)? If data was reported by subjects or indirectly inferred/derived from other data, was the data validated/verified? If so, please describe how.

Answer: Directly observable

Over what timeframe was the data collected? Does this timeframe match the creation timeframe of the data associated with the instances (e.g., recent crawl of old news articles)? If not, please describe the timeframe in which the data associated with the instances was created. Finally, list when the dataset was first published.

Answer: The StackOverflow posts date back to 2006 and we curated the dataset throughout 2022.

What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or sensor, manual human curation, software program, software API)? How were these mechanisms or procedures validated?

Answer: We manually sourced the datapoints and manually annotated with metadata.

What was the resource cost of collecting the data? (e.g. what were the required computational resources, and the associated financial costs, and energy consumption - estimate the carbon footprint. See Strubell *et al.*(?) for approaches in this area.)

Answer: The costs were negligible (no computation cost) and the annotations were sourced among the research group.

If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic, probabilistic with specific sampling probabilities)?

Answer: Deterministic - we sampled based on highest rated posts in StackOverflow after filtering to the tag ExcelFormula and checking manu-

ally that the posts were genuine tasks (not related to Excel behaviour).

Who was involved in the data collection process (e.g., students, crowdworkers, contractors) and how were they compensated (e.g., how much were crowdworkers paid)?

Answer: Researchers in the team (no extra compensation).

Were any ethical review processes conducted (e.g., by an institutional review board)? If so, please provide a description of these review processes, including the outcomes, as well as a link or other access point to any supporting documentation.

Answer: Yes, the dataset passed our institution's on-boarding and publishing reviews

Does the dataset relate to people? If not, you may skip the remainder of the questions in this section.

Answer: No

1029

1030

1031

1032

1033

1034

1035

1037

1038

1040

1041

1042

1043

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1072

1073

1074

1076

1077

1078

1079

Did you collect the data from the individuals in question directly, or obtain it via third parties or other sources (e.g., websites)? *Answer*: Third party (StackOverflow)

Were the individuals in question notified about the data collection? If so, please describe (or show with screenshots or other information) how notice was provided, and provide a link or other access point to, or otherwise reproduce, the exact language of the notification itself.

Answer: No, all data was public and we didn't infringe any licensing

Did the individuals in question consent to the collection and use of their data? If so, please describe (or show with screenshots or other information) how consent was requested and provided, and provide a link or other access point to, or otherwise reproduce, the exact language to which the individuals consented.

Answer: Not applicable

If consent was obtained, were the consenting individuals provided with a mechanism to revoke their consent in the future or for certain uses? If so, please provide a description, as well as a link or other access point to the mechanism (if appropriate)

Answer: Not applicable

Has an analysis of the potential impact of the dataset and its use on data subjects (e.g., a data protection impact analysis)been conducted? If so, please provide a description of this analysis, including the outcomes, as well as a link or other access point to any supporting documentation. *Answer*: No

B.4 PRE-PROCESSING / CLEANING / LABELING

Was any preprocessing/cleaning/labeling of the data done(e.g.,discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)? If so, please provide a description. If not, you may skip the remainder of the questions in this section.

Answer: We manually created the queries (either from scratch or based on the original Stack-Overflow description), we augmented the data (we added more rows, corner cases, augmentations), we extracted metadata (e.g. questions type, type of column referencing).

Was the "raw" data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support unanticipated future uses)? If so, please provide a link or other access point to the "raw" data.

Answer: No.

Is the software used to preprocess/clean/label the instances available? If so, please provide a link or other access point.

Answer: We will provide access to the software we built for data processing, but the majority of the work carried involved manual annotations.

B.5 USES

Has the dataset been used for any tasks already? If so, please provide a description.

Answer: The current work is the first use of the dataset.

Is there a repository that links to any or all papers or systems that use the dataset? If so, please provide a link or other access point. *Answer*: Not applicable

What (other) tasks could the dataset be used for?

Answer: The dataset could be used the evaluate code generation for other languages (e.g. R or SQL) on data-centric tasks, query generation (from the code).

Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses? For example, is there anything that a future user might need to know to avoid uses

14

1216

1217

1218

that could result in unfair treatment of individu-1130 als or groups (e.g., stereotyping, quality of service 1131 issues) or other undesirable harms (e.g., financial 1132 harms, legal risks) If so, please provide a descrip-1133 tion. Is there anything a future user could do to 1134 mitigate these undesirable harms? 1135 1136

Answer: No.

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

Are there tasks for which the dataset should not be used? If so, please provide a description. Answer: No.

B.6 DISTRIBUTION

Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization) on behalf of which the dataset was created? If so, please provide a description.

Answer: Yes, it will be distributed widely.

How will the dataset will be distributed (e.g., tarball on website, API, GitHub)? Does the dataset have a digital object identifier (DOI)? Answer: Github

When will the dataset be distributed? Answer: At the end of the anonymity period.

Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)? If so, please describe this license and/or ToU, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or ToU, as well as any fees associated with these restrictions.

Answer: Yes, TBC

Have any third parties imposed IP-based or other restrictions on the data associated with If so, please describe these rethe instances? strictions, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these restrictions.

Answer: No.

Do any export controls or other regulatory restrictions apply to the dataset or to individual instances? If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any supporting documentation.

Answer: No.

B.7 MAINTENANCE

Who is supporting/hosting/maintaining the dataset?

Answer: The authors, but we will also accept PRs from the community if they would like to extend the dataset and they pass our quality checks.

How can the owner/curator/manager of the dataset be contacted (e.g., email address)? Answer: via Github

Is there an erratum? If so, please provide a link or other access point.

Answer: No.

Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete in-If so, please describe how often, by stances)? whom, and how updates will be communicated to users (e.g., mailing list, GitHub)?

Answer: Yes, we will consider any requests and if any mistakes are pointed out we will update regularly (monthly to start with).

If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances (e.g., were individuals in question told that their data would be retained for a fixed period of time and then deleted)? If so, please describe these limits and explain how they will be enforced.

Answer: Not applicable.

Will older versions of the dataset continue to be supported/hosted/maintained? If so, please describe how. If not, please describe how its obsolescence will be communicated to users. Answer: No.

If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so? If so, please provide a description. Will these contributions be validated/verified? If so, please describe how. If not, why not? Is there a process for communicating/distributing these contributions to other users? If so, please provide a description. Answer: Yes, via PRs.

Examples of problems for different С 1219 scenarios 1220

Description	Prompt Example		
	import pandas as pd		
Data Indopendent (IND)	df = pd.DataFrame()		
Data-Independent (IND)	df["String"] = ["Anna"]		
	# create a new column with the String column capitalized		
	import pandas as pd		
	df = pd.DataFrame()		
Data-Dependent (DEP)	df["Name"] = ["Anna"]		
	df["DoB"] = ["20/10/1989"]		
	# extract the year from the date		
	import pandas as pd		
External knowledge (EVT)	df = pd.DataFrame()		
External knowledge (EXT)	df["State"] = ["WA"]		
	# expand the state name.		
	import pandas as pd		
Anonymous header and no data	df = pd.DataFrame()		
Anonymous header and no data	df["Col1"] = []		
	# Create a column with the first number		
	import pandas as pd		
Proper header	df = pd.DataFrame()		
i toper neader	df["Address"] = []		
	# Create a column with the first number from the address		
	import pandas as pd		
Proper header and data sample	df = pd.DataFrame()		
Toper neader and data sample	df["Address"] = ["6088 184A ST"]		
	# Create a column with the first number		
	import pandas as pd		
Full data	df = pd.DataFrame()		
	df["Address"] = ["6088 184A ST", "4234 32A ST"]		
	# Create a column with the first number		

Figure 7: Prompt examples for the different taxonomies. Top three rows cover task taxonomies. For class IND, the query contains all information required for the completion: the input column name and the task. For class DEP, the query does not specify the date format or column. For class EXT, the model needs external world knowledge about US states. Bottom four rows cover data availability.

Description	Prompt & Completion Example			
	import pandas as pd			
	df = pd.DataFrame()			
Completion for a prompt	df["Col1"] = []			
with no data in the	# Create a column with the first number			
dataframe	df["Col2"] = [1]			
	import pandas as pd			
	df = pd.DataFrame()			
Completion for a prompt	df["Col1"] = ["6088 184A ST"]			
with only the first row of	# Create a column with the first number			
the dataframe	df['Col2'] = df['Col1'].str.split().str[0]			
	import pandas as pd			
	df = pd.DataFrame()			
Completion for a premet	df["Address"] = []			
with only a direct refer	# Create a column with the first number from the address			
with only a direct refer-	<pre>df['Number'] = df['Address'].str.extract(</pre>			
ence to the columns	'(d+)', expand = False)			

Figure 8: Prompt examples and completions results for type errors in the TypeGoldSet using different data availability.

Description	Prompt & Completion Example			
	import pandas as pd			
	df = pd.DataFrame()			
Completion for a prompt	df["Col1"] = []			
with no data in the	# Create a new column with the time stamp			
dataframa	<pre>df['New_Col'] = pd.datetime.now()</pre>			
uatarrame	print(df)			
	import pandas as pd			
	df = pd.DataFrame()			
Completion for a prompt	df["Col1"] = ["Fri, 06 Feb 1987 06:00:37 GMT"]			
with only the first row of	# Create a new column with the time stamp			
the dataframe	<pre>df['Col2'] = pd.to_datetime(df.Col1)</pre>			
	import pandas as pd			
	df = pd.DataFrame()			
Completion for a prompt	df["Date"] = []			
with only a direct refer	# Create a new column with the time stamp from the dates			
ance to the columns	<pre>df['Time Stamp'] = pd.to_datetime(df['Date']).apply(</pre>			
	lambda x: x.timestamp())			

Figure 9: Prompt examples and completions results for value errors in the TypeGoldSet using different data availability.

Description	Prompt & Completion Example			
	import pandas as pd			
	df = pd.DataFrame()			
Original	df["Col1"] = [41, 163,]			
	df["Col2"] = [83, 60,]			
	# Create a new column with the sum of the numeric columns			
	df['Col3'] = df['Col1'] + df['Col2']			
	import pandas as pd			
	df = pd.DataFrame()			
	df["Col1"] = [41, "163",]			
Mixed types and numeric	df["Col2"] = ["83", 59.9,]			
delta	# Create a new column with the sum of the numeric columns			
	<pre>df['sum'] = df['Col1'].astype(float) +</pre>			
	<pre>df['Col2'].astype(float)</pre>			
	import pandas as pd			
	df = pd.DataFrame()			
	df["Col1"] = [41, "",]			
Empty strings	df["Col2"] = ["", "",]			
Empty sumgs	# Create a new column with the sum of the numeric columns			
	<pre>df['Sum'] = df.apply(pd.to_numeric,</pre>			
	errors='coerce').sum(axis=1)			
	print(df)			
	import pandas as pd			
	df = pd.DataFrame()			
	<pre>df["String"] = ["mailto: andrew_wilson@gmx.net",</pre>			
Original	<pre> "mailto: emily_marie_lopez@gmx.net"]</pre>			
Oliginal	# Create a new column with the email address extracted from the string column			
	<pre>df['Email'] = df['String'].str.extract(</pre>			
	'([a-zA-Z0-9+-]+@[a-zA-Z0-9-]+\\.[a-zA-Z0-9]+)',			
	expand=True)			
	import pandas as pd			
	df = pd.DataFrame()			
	<pre>df["String"] = ["andrew_wilson@gmx.net",</pre>			
Text corruption	<pre> "-mailto: emily_marie_lopez@gmx.net"]</pre>			
	# Create a new column with the email address extracted from the string column			
	<pre>df['Email'] = df['String'].str.extract(</pre>			
	r'(?<=mailto:) (\\S+)')			

Figure 10: Prompt examples and completions results for corruptions in TypeGoldSet

1261

1262

1264

1265

1266

1267

1268

1221

D **Prompt Experiments**

This section includes concrete prompt examples for the different prompting schemes. Refer to tables Table 1 for instruction prompts, Table 2 for example-usage prompts, Table 3 for data description prompts.

Е **Implementation details**

Statistics about our datasets **E.1**

E.2 Task classes

We distinguish three task classes:

(IND) data-independent

For these tasks there is a solution that can be applied to the data without knowing the content or format. Hence, if there is a Pandas or default Python function that can be applied with default parameters the task is of class (IND). Examples are simple date operations, like to extract days or hours, with default formats that can be inferred by the pandas function "to_date" are of class A.

(DEP) data-dependent

For these tasks the model needs to know the data content and/or format to find the solutions. For example, if we would need a special parameter, filtering or any kind of parameterization derived from the content and the format of the data the task it is of class (DEP). Examples are data operations on dates with non-default format, mixed formats or corrupted values. In these cases we need either a special parameter (for example errors=coerce) or other transformations derived from the data.

(EXT) external-dependent

The model can not find a solution of the task only knowing the data content and the format. Its need additional open world knowledge. This knowledge needs to be beyond the knowledge about the pandas and python syntax, APIs and libraries. Examples are operations in names where we need knowledge about what a surname, middle name or first name is. Or operations on addresses where we might need to know what a zip code or a state is.

E.3 Anonymous columns

We say a column header (or column name) is anonymous when it does not give any information about the content or type of the columns. For example the column name "Names" indicates that the column contains names or the column name "Dates" indicates that the column contains date

values. On the other hand a column name like "Data" or "Info" does give no information about the content or type of the column.

1269

1270

1271

1272

1273

1274

1275

1276

1277

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1290

1291

1292

1294

1295

1296

1297

1298

1299

1300

1301

1302

1303

1304

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

E.4 Adding noise

For our investigations on how much use the mode is making of the actual data in the dataframes and to check the quality of the solution on non-perfect data, we added additional noise. This additional noise is:

Corrupting values:

With a probability of 20% a value in an input column will have a corrupted value. For string like objects this corruptions will be an insertion of a space, a dash or a new line in a random position in the corresponding string.

Mixing formats:

For numerical values, we mix integer and float, for date values we mix different date formats and for names and addresses with permute the tokens.

Missing values:

With a probability of 20% a value in an input column will be replaced with a empty string which represents a missing value in all our settings.

E.5 Obtaining completions

Parallelization. For efficiency, we request multiple completions from Codex per iteration in Alg.1. To try to minimize both inference time and the load on OpenAI's servers, we adapt the batch size to an estimate of the probability that the next completion is valid. The batch size used in each iteration is $n = \min(\lceil r/p \rceil, B, L)$, where r = k - |C|is the number of valid completions still to obtain, B is the remaining completion budget, and L is a parallelization limit enforced by the Codex API. The probability estimate p is updated after each iteration by counting the number of valid and invalid completions in that iteration's batch.

Since pass@k is calculated only from valid completions, it is not influenced by either parallelization or batch size adaptation. We additionally report the average "pool" size (valid and invalid completions) to measure the cost of retrieving valid completions using the above approach in all our experiments.

Temperature. We use the same k-dependent temperature t as described in Chen et al. (2021); i.e. for k = 1, t = 0.2; for $1 < k \le 5, t = 0.4$, for $5 < k \le 50, t = 0.6$; otherwise t = 0.8.

Prompt Name	Prompt Example (Query Q)			
instruct-concise	write a concise, short and idiomatic pandas solution for the following query: create a new column with the last names.			
instruct-execute	write an executable and type correct pandas solution for the following query: create a new column with the last names.			
instruct-function	write a pandas function that create a new column with the last names.			
instruct-iterative	create a new column with the last names. <i>Prioritize the iterative programming style</i> with use of for loops, while loops and iterative data structures			
instruct-lambda	create a new column with the last names. <i>Prioritize the functional programming style</i> with use of lambdas, comprehensions, and generators			
instruct-domain	create a new column with the last names. <i>Prioritize the use of pandas operators that deal with string manipulation.</i>			
assert-concise create a new column with the last names. <i>The following example demonstrate</i> to write a concise, short and idiomatic pandas solution.				
assert-execute create a new column with the last names. <i>The following example demonstrates how to write an executable and type correct pandas solution.</i>				
assert-function	The following example demonstrates how to write a pandas function that create a new column with the last names.			
assert-iterative create a new column with the last names. <i>The following example demonstrates he prioritize the iterative programming style with use of for loops, while loops and iterative data structures.</i>				
assert-lambda	create a new column with the last names. <i>The following example demonstrates how to prioritize the functional programming style with use of lambdas, comprehensions, and generators.</i>			
assert-domain	create a new column with the last names. <i>The following example prioritizes the use of pandas operators that deal with string manipulation.</i>			
ask-concise	how to create a new column with the last names? how to write a concise, short and idiomatic pandas solution?			
ask-execute	<i>how to</i> create a new column with the last names? <i>how to write an executable and type correct pandas solution</i> ?			
ask-function	how to write a pandas function that create a new column with the last names?			
ask-iterative	how to create a new column with the last names? how to prioritize the iterative programming style with use of for loops, while loops and iterative data structures?			
ask-lambda	how to create a new column with the last names? how to prioritize the functional programming style with use of lambdas, comprehensions, and generators?			
ask-domain	how to create a new column with the last names? how to prioritize the use of pandas operators that deal with string manipulation?			

Table 1: Instruction prompt experiments include command-style prompts, assert-style prompts and ask-style prompts illustrated on the query *create a new column with the last names*. The different prompts vary along two dimensions: the style in which the model is queried and the property we want the completions to emphasize. The phrasing style we explore are instruct, assert and ask. Examples of properties include executable, type correctness and using operators from the problem domain.

1317Stop sequences. The most effective stop sequence1318we found that allows Codex to generate at least1319one solution while not usually using the entire to-1320ken budget is a blank line followed by a line com-1321ment; i.e. $\n\n\#$. Further, to keep Codex from1322generating what appears to be the rest of a forum1323post after a code snippet, we also use the stop se-1324quence </code>.

Completion cleanup. Having forum posts appar-1325 ently in Codex's training data means some com-1326 pletions would raise SyntaxError exceptions 1327 when executed due to formatting artifacts, and 1328 therefore be invalid. Instead, to make the most of 1329 the completion budget, we replace formatting arti-1330 facts. In particular, we replace HTML escape se-1331 quences such as < and " with Python 1332 operators and delimiters. 1333

Cleanup additionally removes unnecessary whitespace, blank lines and comments, and truncates completions at n# when it appears after executable code.

1334

1335

1336

1337

1338

E.6 Executing completions

Rewriting. Completions returned by Codex do not 1339 clearly indicate which variables or expressions are 1340 intended to be the answer to a query. This must be 1341 inferred from the shape of the code. We found that 1342 an effective way to identify and expose the likely answer is to search backwards to find the last unin-1344 dented (i.e. top-level) statement that has one of a 1345 few forms, and rewrite the completion so that its 1346 last statement is an assignment to a fresh identifier 1347 varout. The statement forms and rewrites are 1348

Prompt Name	Prompt Example
	import pandas as pd
usage-function	<pre>def startswith_at(txt):</pre>
	return txt.startswith('@')
	df = pd.DataFrame("Names":["Charles Moore", "Anna Green"])
	# create a new column with the last names.
	import pandas as pd
	df = pd.DataFrame("a":[1, 2, 3])
	# loop over the rows of a df using pd.iloc, create a function to iterate over the df.
usage_iterative	for i in df.index:
usage nerative	<pre>val = df[col].iloc[i]</pre>
	<pre>df = pd.DataFrame("Names":["Charles Moore", "Anna Green"]) # create a new column with the last names.</pre>
	import pandas as pd
	df = pd.DataFrame("a":[1, 2, 3])
	# group df on column b and keep half of the elements at random
usage-lambda	<pre>dfout = df.groupby('b').apply(lambda x:x.sample(frac=0.5))</pre>
	df = pd.DataFrame("Names":["Charles Moore", "Anna Green"])
	# create a new column with the last names.

Table 2: Example-usage prompt examples illustrated on the query *create a new column with the last names*. Example-usage prompts append example code snippets to the query that demonstrate how to solve the pandas problem in various ways.

• var = expr: append the statement var_{out} = var to the completion

1349

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

1361

1362

1363

1364

1365

- $var[expr_i] = expr$: append the statement $var_{out} = var$ to the completion
- print (*expr*, ...): replace this statement and the rest of the completion with $var_{out} = expr$
- *expr*: replace this statement and the rest of the completion with $var_{out} = expr$

Rewriting also inserts import statements for common libraries (e.g. import numpy as np).

The rewritten completion is appended to the code that defines the input dataframe to create a completed program. The completed program and the output variable name var_{out} are sent to a sandbox for execution.

Sandboxing. Because of security risks inherent 1366 in running LLM-generated code, we run com-1367 pleted programs in a sandbox. Our sandbox is a 1368 JavaScript web service that runs Python programs in Pyodide (Droettboom et al., 2022), a Python 1370 distribution for WebAssembly. While Python pro-1371 grams running in Pyodide have access to the host's 1372 network resources, they at least are isolated from 1373 other host resources including its filesystem, of-1374 fering some level of protection from malicious or 1375 accidentally harmful completions. 1376

After running the code, the sandbox returns the	
value of <i>var_{out}</i> .	

1377

1378

1379

1388

1389

1390

E.7 Evaluation

For a completion to be considered a correct solu-1380 tion in the calculation of pass@k, its actual output 1381 must match the expected output. Matching can-1382 not be the same as equality and still conform to 1383 a reasonable notion of correctness; for example, 1384 the natural breakdown of a solution might gener-1385 ate intermediate columns in the actual output that 1386 are not in the expected output. 1387

The actual output is allowed to vary from the expected output in the following ways and still match the expected output:

- Extra columns 1391
- Different column order 1392
- Different column headers
 1393
- Number expected; actual is a number within small relative error (default 0.01)
 1394
- Number expected; actual is a string that parses as a number within small relative error
 1396
 1397
 1398
- Boolean expected; actual is number 0 or 1
- Boolean expected; actual is a string that represents a truth value 1400

Prompt name	Prompt template				
	df = pd.DataFrame("Position":[3, 14, 25])				
column info	Given position number n, create a new column with the n th letter of the alphabet.				
column-mio	There is 1 column in the input data. The column Position has 10 entries, 10 of				
	which are unique.				
	df = pd.DataFrame("Position":[3, 14, 25])				
	Given position number n, create a new column with the n th letter of the alphabet.				
	Returns				
type-info					
	pd.Dataframe with new columns. Data columns are as follows:				
	Position integers (as int)				
	Letter strings (as str)				
	df = pd.DataFrame("Position":[3, 14, 25])				
regex-info (exact)	Given position number n, create a new column with the n th letter of the alphabet.				
	10 of 10 elements in column Position match the regex format $[0-9]$ + where				
	examples of these elements include 8 and 23.				
	df = pd.DataFrame("Body Fat":["4%", "14%"])				
	create a new column that writes Body Builder if Body Fat is less than 5%, else if it is				
regex-info (partial)	less than 13% writes Athletic, else writes No Data Yet.				
.8	9 of 10 (90%) elements in column Body Fat match the regex format $[0-9]+\%$				
	where an example of these elements is (40%) . 1 of 10 (10%) elements in column				
	Body Fat match the regex format NA where an example of these elements is NA.				
	di = pd.DataFrame("Position":[3, 14, 25])				
	Given position number n, create a new column with the n th letter of the alphabet.				
	Returns				
regex-type-info					
	pd.Dataframe with new columns. Data columns are as follows:				
	Position integers (as int) that match the regex $[0-9]+$				
	Letter strings (as str)				
	di = pd.DataFrame("Position":[3, 14, 25])				
	Given position number n, create a new column with the n ^{an} letter of the alphabet.				
	There is 1 column in the input data. The column Position has 10 unique entries.				
regex-type-col-info	Keturns				
	pd.Dataframe with new columns. Data columns are as follows:				
	Position integers (as int) that match the regex $[0-9]$ +				
	Letter strings (as str)				
nroportion (taxtual)	u = pu. Datarialle ("POSILIOII"; [5, 14, 25])				
properties (textual)	The solumn Desition has the following properties: All elements contain disits				
	The countin rostion has the following properties: All elements contain digits.				
(f	$\alpha I = p\alpha. \forall atarrame("Position": [3, 14, 25])$				
properties (format)	Given position number n, create a new column with the n letter of the alphabet.				
	The column Position satisfies ContainsDigits()				

Table 3: Data-description prompts illustrated on the queries *Given the position number n, create a new column with the nth letter of the alphabet.* Data-description prompts aim to extract task-specific structural information about the data input provided by the user. We automatically compute information about the number of elements in the columns, the type of each of the column elements (one of string, integer, boolean) and properties satisfied by the elements. In addition, we compute a regular expression that captures a pattern of the elements in the column.

• String expected; actual is a string that differs only in case

Allowed string truth value representations, allowed relative error, and whether string matching is case-sensitive are (optionally) overridden per data point as appropriate.

1402

1403

1404

1405

1406

1407

F Detailed evaluation summaries of different scenarios

1408

1409

1410

F.1 Stability analysis

To estimate the stability of the CODEXDATA al-1411 gorithm on pass@k we ran each task in TYPESET 1412 and SOFSET 10 times and estimated the pass@5. 1413 In Fig.13 we show the corresponding means and 1414 standard deviations for SOFSET (top) and TYPE-1415 SET (bottom). Due to the effort to run CODEXfor 1416 each example we were only able to calculate these 1417 numbers for k=5. 1418

Data type	Example of task	TYPESET	SOFSET
strings	sub-string extracts: df["col1"].str[4:10]	16	116
numbers	divisions e.g. df["col1"] / df["col2"]	14	34
names	surname in upper-case df["Givenname"].str.upper()[0]	12	13
dates	extract month e.g. pd.DatetimeIndex(df['date']).month	14	32
units	conversions	40	1
addresses	extract the state df["col1"].str[:-2]	18	4
mixed types	combine columns df["col1"] + df["col2"].astype(str)	14	0

Table 4: Number of tasks per data type in TYPESET and SOFSET with an example.

F.2 Detailed error analyses 1419





TYPESET Instruction Prompts (Commands)

Figure 11: This figure presents prompt results for the instruction type prompts for SOFSET and TYPESET datasets. We evaluate the different ways of phrasing the instruction: as a command, as an assertion and as a question. We observe that instruction prompts that are phrased as a command perform slightly better than instruction prompts phrased as an assertion or question. A speculation can be made about the impact of training data: not many natural language comments are typically phrased as questions.



Figure 12: This figure presents prompt results for the example usage and data description prompting strategy for SOFSET and TYPESET datasets. We observe that example usage prompts always perform worse or same as the baseline. Amongst the data description prompts, augmenting type information leads to better performance in both the datasets.

Notation				
Metric	Description	Туре		
pass@k	see metrics	unit test		
pass@k(%)	see metrics			
NumUniqueValid	Number of unique completions that are valid (have no run-			
	time errors and have the correct output type).			
NumSuccess	Number of completions that produce the correct output	unit test		
	type and values (pass unit test).			
NumError	Number of runtime errors (completions with no valid out-	runtime		
	put produced).			
NumSyntaxErrors	Number of syntax errors (completion generated is no valid	runtime		
-	python code).			
NumTypeErrors	Number of type errors in the generated completions (oper-	runtime		
	ations/function not allowed for type).			
NumValueErrors	Number of value errors (value caused error in operation/-	runtime		
	function).			
NumIndexErrors	Number of error involving the dataframe index.	runtime		
NumAttributeErrors	Number of attribute errors (referred object attribute does	runtime		
	not exists for example).			
NumNameErrors	Number of name errors (function/variable does not exists).	runtime		
NumKevErrors	Number of key errors (column does not exist).	runtime		
NumRegexErrors	Number of regular expression errors.			
NumAssertionErrors	Number of assertion errors in the completions.			
NumCorrectOutputType	Number of completion that produce the correct output type	semantic		
	(a new column).			
NumTypeMismatch	Number of mismatches of the output type of the comple-	semantic		
~	tion and the ground truth output. Please note that the mis-			
	match values are only computed for completions that are			
	valid/execute.			
NumValueMismatch	Number of mismatches of the column values from the gen-	semantic		
	erated by the completion and the ground truth.			
DiffColNumber	Number of completions with different output column num-	semantic		
	ber than ground truth output.			
InputColumnsError	Number of completions that do not use the input columns	semantic		
F F -	from the dataframe.			
ExtraColumnsError	Number of completions that use columns that are not rele-	semantic		
	vant for the task.			
NumCompletions	Number of completions retrieved	stats		
CompletionLenght	String length of generated completions.	stats		
comprehendingin	Sumg lengui or generated completions.	Stats		

Table 5: Short description of the used metrics. For each experiment we calculate the above metrics and report them as average over all completions and all examples per evaluation. Please see the Metrics section for further details.

	full-data	no-data prop-col	subset-data	no-data anon-col
Pass@K	0.673	0.583	0.6016	0.5687
Pass@K75pct	0.8273	0.7533	0.7512	0.7217
Pass@K70pct	0.8519	0.7761	0.7874	0.7709
Pass@K50pct	0.933	0.8584	0.8642	0.8259
Pass@K30pct	0.9524	0.906	0.9069	0.8625
Pass@K25pct	0.9531	0.9123	0.9152	0.8662
NumUniqueValid	0.8086	0.6397	0.7529	0.6604
NumSuccess	0.3135	0.2292	0.2952	0.2242
NumError	0.1833	0.3454	0.2384	0.3274
NumSyntaxErrors	0.0072	0.0159	0.0074	0.0167
NumTypeErrors	0.0214	0.0306	0.0252	0.0327
NumValueErrors	0.0405	0.0655	0.0714	0.0632
NumIndexErrors	0.007	0.0109	0.007	0.0088
NumAttributeErrors	0.0098	0.0146	0.0106	0.0209
NumNameErrors	0.0045	0.0094	0.0062	0.0087
NumKeyErrors	0.0081	0.018	0.0068	0.0715
NumRegexErrors	0.0	0.0	0.0	0.0
NumAssertionErrors	0.0	0.0	0.0	0.0
NumCorrectOutputType	0.8086	0.6397	0.7529	0.6604
NumTypeMismatch	0.0081	0.0149	0.0087	0.0122
NumValueMismatch	0.4951	0.4105	0.4577	0.4362
DiffColNumber	0.0264	0.0313	0.0201	0.0301
DiffTableNumber	0.0264	0.0313	0.0201	0.0301
CompletionLenght	9.5091	9.5446	9.5564	9.6114
NumCompletions	28.7045	40.0606	31.25	40.1053

Full-data vs no-data experiments SOFSET

Table 6: Detailed results for SOFSET for task class (IND) and different levels of data redaction. Metrics shown are average numbers across the sampled completions.

	full-data	ata no-data prop-col subset-data n		no-data anon-col	
Pass@K	0.302	0.1213	0.1384	0.1354	
Pass@K75pct	0.4031	0.1739	0.2294	0.1564	
Pass@K70pct	0.4352	0.1859	0.2793	0.1901	
Pass@K50pct	0.6094	0.376	0.4767	0.3741	
Pass@K30pct	0.7121	0.4835	0.605	0.4546	
Pass@K25pct	0.7121	0.4835	0.605	0.4546	
NumUniqueValid	0.6071	0.2997	0.5895	0.2926	
NumSuccess	0.0828	0.0296	0.0497	0.0243	
NumError	0.3761	0.6936	0.398	0.7017	
NumSyntaxErrors	0.0112	0.0207	0.0105	0.0209	
NumTypeErrors	0.172	0.3263	0.1473	0.259	
NumValueErrors	0.0653	0.123	0.1023	0.0928	
NumIndexErrors	0.0018	0.0063	0.0049	0.0083	
NumAttributeErrors	0.0398	0.0876	0.0333	0.0865	
NumNameErrors	0.0029	0.0061	0.0036	0.0034	
NumKeyErrors	0.0074	0.0162	0.0052	0.1619	
NumRegexErrors	0.0	0.0	0.0	0.0	
NumAssertionErrors	0.0	0.0003	0.0	0.0	
NumCorrectOutputType	0.6071	0.2997	0.5895	0.2926	
NumTypeMismatch	0.0168	0.0068	0.0126	0.0057	
NumValueMismatch	0.5244	0.2701	0.5398 0.26		
DiffColNumber	0.0328	0.0413	0.0419	0.0261	
DiffTableNumber	0.0328	0.0413	0.0419	0.0261	
CompletionLenght	9.6305	9.8837	9.6271	9.9495	
NumCompletions	41.4474	80.3421	42.3421	88.3243	

Full-data vs no-data experiments SOFSET

Table 7: Detailed results for SOFSET for task class (**DEP**) and different levels of data redaction. Metrics shown are average numbers across the sampled completions.

	full-data	ata no-data prop-col subset-data		no-data anon-col
Pass@K	0.2078	0.0962	0.1369	0.1124
Pass@K75pct	0.4314	0.2181	0.2546	0.2097
Pass@K70pct	0.4641	0.2835	0.2869	0.2258
Pass@K50pct	0.6079	0.3806	0.5297	0.3277
Pass@K30pct	0.7273	0.5732	0.6814	0.5658
Pass@K25pct	0.7273	0.5732	0.6814	0.5658
NumUniqueValid	0.6669	0.4943	0.6079	0.4372
NumSuccess	0.0559	0.011	0.0314	0.0244
NumError	0.3193	0.4991	0.3826	0.5566
NumSyntaxErrors	0.0118	0.0156	0.0126	0.0174
NumTypeErrors	0.0492	0.1167	0.0433	0.1196
NumValueErrors	0.1019	0.1487	0.1449	0.1448
NumIndexErrors	0.0025	0.013	0.0058	0.0211
NumAttributeErrors	0.0248	0.0834	0.0283	0.0688
NumNameErrors	0.0126	0.0145	0.0114	0.0127
NumKeyErrors	0.0073	0.0106	0.0064	0.0937
NumRegexErrors	0.0	0.0	0.0	0.0
NumAssertionErrors	0.0	0.0	0.0	0.0
NumCorrectOutputType	0.6669	0.4943	0.6079	0.4372
NumTypeMismatch	0.0138	0.0066	0.0095	0.0061
NumValueMismatch	0.611	0.4833	0.4833 0.5765	
DiffColNumber	0.0487	0.077	0.0621	0.0548
DiffTableNumber	0.0487	0.077	0.0621	0.0548
CompletionLenght	9.5894	9.8389	9.6326	9.8971
NumCompletions	36.4516	55.2581	41.3548	65.2258

Full-data vs no-data experiments SOFSET

Table 8: Detailed results for SOFSET for task class (**EXT**) and different levels of data redaction. Metrics shown are average numbers across the sampled completions.

	full-data	no-data prop-col	subset-data	no-data anon-col
Pass@K	0.55	0.4	0.5881	0.5663
Pass@K75pct	0.7741	0.6263	0.7013	0.6915
Pass@K70pct	0.7741	0.6263	0.7013	0.6915
Pass@K50pct	0.7741	0.6263	0.7013	0.6923
Pass@K30pct	0.7741	0.6263	0.7263	0.6923
Pass@K25pct	0.7741	0.6263	0.7263	0.6923
NumUniqueValid	0.8481	0.6374	0.793	0.6897
NumSuccess	0.4008	0.2669	0.3914	0.3257
NumError	0.1399	0.3498	0.2006	0.2987
NumSyntaxErrors	0.0025	0.0064	0.0083	0.0104
NumTypeErrors	0.0455	0.1061	0.069	0.0583
NumValueErrors	0.0173	0.0277	0.0146	0.1008
NumIndexErrors	0.0	0.0004	0.0	0.0
NumAttributeErrors	0.0	0.0197	0.0006	0.0097
NumNameErrors	0.0	0.0047	047 0.0 0.00	
NumKeyErrors	0.0042	0.002	0.0	0.0034
NumRegexErrors	0.0	0.0	0.0	0.0
NumAssertionErrors	0.0	0.0	0.0	0.0
NumCorrectOutputType	0.8481	0.6374	0.793	0.6897
NumTypeMismatch	0.012	0.0128	0.0063	0.0116
NumValueMismatch	0.4474	0.3705	0.4016	0.364
DiffColNumber	0.0087	0.0169	0.0043	0.0103
DiffTableNumber	0.0087	0.0169	0.0043	0.0103
InputColumnsError	0.0961	0.0539 0.0721 0.0		0.0222
ExtraColumnsError	0.0	0.0	0.0	0.0
CompletionLenght	9.5944	9.7442	9.6465	9.6858
CompletionSize	25.15	45.35	28.15	41.8

Full-data vs no-data experiments TYPESET

Table 9: Detailed results for TYPESET for task class (IND) and different levels of data redaction. Metrics shown are average numbers across the sampled completions.

	full-data	no-data prop-col	subset-data	no-data anon-col
Pass@K	0.9243	0.3454	0.8071	0.3333
Pass@K75pct	0.9243	0.3454	0.8071	0.4876
Pass@K70pct	0.9243	0.3618	0.8071	0.4876
Pass@K50pct	0.9243	0.3748	0.8071	0.5293
Pass@K30pct	0.9243	0.375	0.8071	0.5293
Pass@K25pct	0.9243	0.375	0.8071	0.5293
NumUniqueValid	0.7122	0.325	0.6844	0.5172
NumSuccess	0.4548	0.1282	0.386	0.1739
NumError	0.2822	0.674	0.3156	0.4828
NumSyntaxErrors	0.0	0.0159	0.0	0.0
NumTypeErrors	0.1031	0.1846	0.0854	0.2164
NumValueErrors	0.0567	0.1867	0.0875	0.0808
NumIndexErrors	0.0	0.0213	0.0	0.0129
NumAttributeErrors	0.0138	0.1038	0.0202	0.0851
NumNameErrors	0.0	0.0	0.0	0.0023
NumKeyErrors	0.0	0.0	0.0	0.0181
NumRegexErrors	0.0	0.0	0.0	0.0
NumAssertionErrors	0.0	0.0	0.0	0.0
NumCorrectOutputType	0.7122	0.325	0.6844	0.5172
NumTypeMismatch	0.0056	0.001	0.0	0.0
NumValueMismatch	0.2574	0.1968	0.2984	0.3432
DiffColNumber	0.0054	0.013	0.006	0.006
DiffTableNumber	0.0054	0.013	0.006	0.006
InputColumnsError	0.0889	0.0448	0.0464	0.1982
ExtraColumnsError	0.0	0.0	0.0	0.0
CompletionLenght	9.5041	9.8964	9.7128	9.8413
CompletionSize	30.5	68.75	35.25	47.25

Full-data vs no-data experiments TYPESET

Table 10: Detailed results for TYPESET for task class (**DEP**) and different levels of data redaction. Metrics shown are average numbers across the sampled completions.

	full-data	no-data prop-col	subset-data	no-data anon-col
Pass@K	0.854	0.4407	0.7824	0.2383
Pass@K75pct	0.8947	0.5174	0.8131	0.2677
Pass@K70pct	0.8947	0.5174	0.8131	0.2677
Pass@K50pct	0.8975	0.5174	0.8131	0.2677
Pass@K30pct	0.8976	0.5174	0.8131	0.2677
Pass@K25pct	0.8976	0.5174	0.8131	0.2677
NumUniqueValid	0.7303	0.4097	0.6981	0.292
NumSuccess	0.3118	0.1299	0.2794	0.0728
NumError	0.2639	0.5855	0.2994	0.7016
NumSyntaxErrors	0.0087	0.0186	0.0049	0.0053
NumTypeErrors	0.0807	0.2443	0.1063	0.131
NumValueErrors	0.0677	0.0889	0.0493	0.0846
NumIndexErrors	0.0036	0.0009	0.0004	0.0023
NumAttributeErrors	0.0057	0.0245	0.0063	0.0145
NumNameErrors	0.0521	0.1623	0.0821	0.0083
NumKeyErrors	0.0027	0.0095	0.0042	0.0057
NumRegexErrors	0.0	0.0	0.0	0.0
NumAssertionErrors	0.0	0.0002	0.0	0.0
NumCorrectOutputType	0.7303	0.4097	0.6981	0.292
NumTypeMismatch	0.0058	0.0048	0.0025	0.0064
NumValueMismatch	0.4186	0.2798	0.4187	0.2192
DiffColNumber	0.0174	0.0182	0.0089	0.0103
DiffTableNumber	0.0174	0.0182	0.0089	0.0103
InputColumnsError	0.3928	0.1325	0.3077	0.0901
ExtraColumnsError	0.0	0.0	0.0	0.0
CompletionLenght	9.4444	9.8718	9.5836	10.1252
CompletionSize	29.0	72.6111	32.5556	95.3784

Full-data vs no-data experiments TYPESET

Table 11: Detailed results for TYPESET for task class (EXT) and different levels of data redaction. Metrics shown are average numbers across the sampled completions.

Noise level experiment							
	missing	mixformat	corruption	baseline			
Pass@K	0.2189	0.4613	0.2557	0.6423			
Pass@K75pct	0.477	0.5435	0.6777	0.7017			
Pass@K70pct	0.5868	0.5435	0.6777	0.7017			
Pass@K50pct	0.6563	0.5436	0.7044	0.7127			
Pass@K30pct	0.724	0.5436	0.7044	0.7189			
Pass@K25pct	0.7713	0.5436	0.7044	0.7303			
NumUniqueValid	0.7905	0.8084	0.7503	0.8112			
NumSuccess	0.0552	0.1788	0.0987	0.3083			
NumError	0.2037	0.1864	0.2476	0.1833			
NumSyntaxErrors	0.0038	0.0056	0.0081	0.0095			
NumTypeErrors	0.0369	0.038	0.0122	0.0358			
NumValueErrors	0.0674	0.0467	0.0869	0.0678			
NumIndexErrors	0.033	0.0021	0.0034	0.0004			
NumAttributeErrors	0.0034	0.0433	0.0049	0.0115			
NumNameErrors	0.0023	0.0022	0.0	0.0			
NumKeyErrors	0.0021	0.0028	0.0	0.0025			
NumRegexErrors	0.0	0.0	0.0	0.0			
NumAssertionErrors	0.0	0.0	0.0	0.0			
NumCorrectOutputType	0.7905	0.8084	0.7503	0.8112			
NumTypeMismatch	0.0058	0.0052	0.0021	0.0055			
NumValueMismatch	0.9283	0.7826	0.8667	0.6284			
DiffColNumber	0.0207	0.0311	0.0289	0.0214			
DiffTableNumber	0.0207	0.0311	0.0289	0.0214			
InputColumnsError	0.0716	0.0276	0.0864	0.1536			
ExtraColumnsError	0.0	0.0	0.0	0.0			
CompletionLenght	9.6438	9.5977	9.7633	9.5553			
NumCompletions	27.9167	27.1875	37.2143	26.9318			

Table 12: Code completions results on TYPESETNOISY for different levels of noise. Metrics shown are average numbers across the sampled completions.

	renormance per data-type on Trieber						
	units	strings	numeric	names	mixed	dates	address
Pass@K	0.9652	0.8333	0.5461	0.4809	0.5	0.4892	0.3205
Pass@K75pct	0.9798	0.9722	0.7128	0.7181	0.5587	0.5799	0.6193
Pass@K70pct	0.9798	0.9722	0.713	0.7181	0.5587	0.58	0.6646
Pass@K50pct	0.9798	0.9722	0.7938	0.7266	0.5587	0.5877	0.6928
Pass@K30pct	0.9798	0.9722	0.7938	0.7303	0.5587	0.5877	0.7354
Pass@K25pct	0.9798	0.9722	0.7938	0.7303	0.5587	0.62	0.7543
NumUniqueValid	0.7004	0.9084	0.543	0.8675	0.6703	0.6406	0.8546
NumSuccess	0.2792	0.5707	0.246	0.2807	0.2835	0.1515	0.1791
NumError	0.2963	0.0879	0.4507	0.1311	0.3192	0.3559	0.1365
NumSyntaxErrors	0.0032	0.0	0.0017	0.0028	0.0037	0.0035	0.0136
NumTypeErrors	0.1273	0.0	0.209	0.0028	0.1575	0.1229	0.0
NumValueErrors	0.0264	0.0239	0.1493	0.021	0.0276	0.0704	0.0638
NumIndexErrors	0.0018	0.0023	0.0	0.0087	0.0	0.0061	0.0074
NumAttributeErrors	0.0053	0.0051	0.0011	0.0016	0.0048	0.0401	0.0016
NumNameErrors	0.0938	0.0	0.0	0.0	0.003	0.0013	0.0023
NumKeyErrors	0.0	0.0024	0.0	0.0014	0.0	0.0004	0.0036
NumRegexErrors	0.0	0.0	0.0	0.0	0.0	0.0	0.0
NumAssertionErrors	0.0	0.0	0.0	0.0	0.0	0.0	0.0
NumCorrectOutputType	0.7004	0.9084	0.543	0.8675	0.6703	0.6406	0.8546
NumTypeMismatch	0.0033	0.0037	0.0062	0.0014	0.0105	0.0035	0.0088
NumValueMismatch	0.4212	0.3377	0.297	0.5868	0.3868	0.4891	0.6755
DiffColNumber	0.0097	0.0138	0.0043	0.0355	0.0047	0.0071	0.0124
DiffTableNumber	0.0097	0.0138	0.0043	0.0355	0.0047	0.0071	0.0124
InputColumnsError	0.6608	0.0	0.0137	0.0067	0.3417	0.0088	0.0313
ExtraColumnsError	0.0	0.0	0.0	0.0	0.0	0.0	0.0
CompletionLenght	9.3607	9.7139	9.8027	9.6785	9.6382	9.5844	9.6848
NumCompletions	29.25	22.1389	47.25	23.1071	35.0769	40.9032	24.4872

Performance per data-type on TYPESET

Table 13: Code completions results on TYPESET by data-types. Metrics shown are averaged across the sampled completions.



Figure 13: Mean pass@5 and standard deviation for data independent tasks full-data; data independent tasks no-data; data dependent tasks full-data, data dependent tasks no-data; external knowledge tasks full-data; external knowledge tasks no-data.