# Source2Synth: Synthetic Data Generation and Curation Grounded in Real Data Sources

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

023 024

025

Paper under double-blind review

## Abstract

Large Language Models still struggle in challenging scenarios that leverage structured data, complex reasoning, or tool usage. In this paper, we propose *Source2Synth*: a new method that can be used for teaching LLMs new skills without relying on costly human annotations. *Source2Synth* takes as input a custom data source and produces synthetic data points with intermediate reasoning steps grounded in real-world sources. *Source2Synth* improves the dataset quality by discarding low-quality generations based on their answerability. We demonstrate the generality of this approach by applying it to two challenging domains: we test reasoning abilities in multi-hop question answering (MHQA), and tool usage in tabular question answering (TQA). Our method improves performance by 25.51% for TQA on WikiSQL and 22.57% for MHQA on HotpotQA compared to the fine-tuned baselines.

1 INTRODUCTION

Large Language Models (LLMs) (Devlin et al., 2019; Chowdhery et al., 2022; Brown et al., 2020; Vaswani et al., 2017) have risen to popularity due to their remarkable ability to digest and generate human-like text (Radford et al., 2018). However, they still struggle with more complex tasks such as multi-step reasoning, tool use and manipulating or processing structured data. For many of these tasks there exists source data - for example, existing structured data on the web -, but little information of how to use these data to solve a task. In principle, one can achieve performance improvements during fine-tuning by enriching the data with human annotations collected for specific tasks. However, this is an expensive and time-consuming process (Touvron et al., 2023) subject to human-errors and bias.

In this paper, we propose *Source2Synth*, a general approach to produce *synthetic data grounded in external real-world sources*. Basing the data generation process on real-world sources steers the examples to be more realistic, diverse, and factually correct. We showcase our method on two challenging tasks: multi-hop questions based on sources from the web and tabular question answering using SQL as a tool. In both cases, models trained following *Source2Synth*'s pipeline achieve improved performance without relying on human annotations, resulting in a scalable data generation method for complex tasks, and present increased abilities at tackling corner cases.

Source2Synth consists of three stages: Dataset Generation, Dataset Curation, and Model Finetuning,
 see Figure 1. At the Dataset Generation stage, we start by selecting a data source (such as tables on
 the web or Wikipedia articles) to ground our synthetic generation in realistic information. Then, our
 method selects a seed topic to trigger the generation and condition all its components - for example a
 specific entity in a Wikipedia article or a factual statement about a table. Given the seed topic, the
 method then produces the full example: the instruction (e.g., question), the reasoning chain to arrive
 at the answer (e.g., the steps of multi-hop question answering, or tool use) and the answer itself.

- At the *Data Curation* stage, the constructed synthetic dataset is split into two slices: the first half is used to fine-tune the LLM. We use this intermediate model to curate the second half of the synthetic dataset via an imputation and a filtering step by rejection sampling. For imputation, we blank some parts of the given example in order to get a more natural and cohesive entry. For filtering, we reject examples that cannot produce the correct answer in k = 3 tries. This provides a higher quality curated dataset for the second fine-tuning stage, resulting in a final better performing model on a given task.
  - To demonstrate the generality of our approach, we apply it to two different domains:

- Answering *tabular-based questions* by learning how to use SQL as a tool;
- Answering *multi-hop questions* by performing multi-step reasoning and information extraction.

To summarize, our key contributions are:

- We introduce a new method for generating synthetic examples aligned with the target task, given a real-world data source as context.
- We introduce a curation method based on filtering and imputation which yields higher quality data and improved task performance.

## 2 RELATED WORK

068 069

054

056

059

060 061

062

063

064

065 066 067

Synthetic Data Generation using LLMs A number of works propose different strategies to generate 071 synthetic datasets leveraging pre-trained language models. Some of these works rely on knowledge-072 probing by first providing a prompt and letting the model either generate the continuation of a prefix 073 or predict missing words in a close-style template (Schick & Schütze, 2020; Schick & Schütze, 2021; Petroni et al., 2019; Jiang et al., 2019). Other works introduce a variety of ways to improve the quality 074 of synthetic data by using model-based or human filtering (Schick & Schütze, 2021; Liu et al., 2022; 075 Li et al., 2024; Thoppilan et al., 2022). Our method however does not rely on human annotations, and 076 we improve the quality of the synthetic data by leveraging the LLM itself. Furthermore, our selection 077 of the seed topic is automated and we use real data as a starting point. We note that some recent work also leverages real-world data for specific cases, such as a corpus from the web to construct 079 high-quality synthetic data (Nguyen et al., 2024) or open-source code snippets to generate diverse instruction data for code generation (Wei et al., 2024; Dubey et al., 2024). In our case, we proposes 081 a general framework which can be applied across tasks and we do not require a back-translation 082 approach or an initial finetuning to come up with the seed.

See Liu et al. (2024) for a thorough overview of synthetic data research and references therein.

084 Teaching LLMs to Use Tools Enabling LLMs to use different tools can extend their abilities towards 085 manipulating structured data, retrieving information from external sources, or interacting with APIs. 086 Even though the goal of our work is not specifically to teach models to use tools, but to develop a 087 general synthetic data generation approach, we consider this to be a by-product. As an example, we 880 demonstrate how our method can be applied so that LLMs use SQL. Various works augment LLMs with general tools or API calls (Parisi et al., 2022; Schick et al., 2023; Tang et al., 2023), while some propose to interweave intermediate reasoning steps with API calls (Gao et al., 2023; Cai et al., 2024; Paranjape et al., 2023) which improves performance on more complex tasks. Finally, handling unseen 091 tools at test time has also been tackled (Paranjape et al., 2023; Mekala et al., 2024). See Mialon et al. 092 (2023) and Qin et al. (2023) for an in-depth review of augmented tool-use. 093

Teaching LLMs to use SQL The above approaches usually tool usage is restricted to inputs that are
 strings or numbers. However, using structured data (like tables and graphs) during post-training can
 be useful to enhance the LLM's capabilities in complex tasks. A particular tool of interest is SQL
 since it enables aggregating information from tabular data. There exist a variety of benchmarks that
 have been proposed to assess LLMs abilities to generate SQL as well as their performance on tabular based question answering leveraging SQL tasks (Li et al., 2023a; Zhong et al., 2017). Alternatively,
 handling tabular data directly by LLMs has also been tried (Herzig et al., 2020; Gemmell & Dalton,
 2023), and tabular question answering benchmarks have been proposed (Pasupat & Liang, 2015).

101

## 102

## 3 Method

104 105

Source2Synth produces high-quality synthetic examples grounded in external real-world data sources,
 which can be fed to the LLM as step-by-step examples at fine-tuning. Source2Synth is composed of three stages: Dataset Generation, Dataset Curation, and Model fine-tuning.



Figure 1: **Overall** *Source2Synth* **Method.** In the *Dataset Generation* step we first choose a data source to build our dataset from. For each example we select a seed topic to condition the generation on, and use the data source and seed together to construct the example. The resulting synthetic dataset is sliced in two: slice 0 is used to fine-tune an intermediate version of the LLM (*LLMSynth*), and we use *LLMSynth* to curate slice 1 through filtering and/or imputation during the *Dataset Curation* step. The resulting curated dataset is of higher quality and aligned with the user's design. At the *Model Finetuning* stage, the final LLM (*LLMCurated*) is trained on the curated synthetic dataset, which can then be used to provide good performance on the task of interest.

134 135 136

137

## 3.1 DATASET GENERATION

Data source selection The generation process begins by selecting a data source. This can be an already existing dataset re-purposed for a given task, a collection of existing data points that we would like to leverage to construct a new dataset, or structured information (e.g. graphs, tables). There is no need for human annotations on the entries, as *Source2Synth* will enrich it with extra instructions.

Seed In order to create a given example of our new synthetic dataset, we first generate a *seed* topic as the initial trigger for the generation process, which is chosen conditioned on a randomly selected portion of the source data. The seed inspires the creation of the entry and dictates how the source data will be used. In addition, the randomness of the seed ensures variety in the generated data.

Dataset construction In order to tackle complex tasks, LLMs can leverage a step-by-step approach (Wei et al., 2022) that divides reasoning into smaller sub-tasks plus instructions on how to merge back each step into the final one. In *Source2Synth*, we leverage the seed to build synthetic data step-by-step, decomposing into such intermediate steps in order to arrive at an answer for a given question. This reasoning chain can then be used as supervision by providing it as the target in the synthetically generated training examples.

152 153

# 3.2 DATASET CURATION

During curation, *LLMSynth* is then used to improve the quality of the second slice of the dataset using imputation plus a filtering step. We observe that tasks that require compositionality during the construction of the synthetic entry benefit from imputation since it allows the LLM to reformulate into something of higher likelihood under the model's distribution. For example, in MHQA, the merging of two sub questions might produce a multi-hop question that sounds artificial. Following the effort in aligning LLMs with human preferences for a more natural use, we observe by looking at the perplexity of imputation that partially reconstructing *Q* leads to more human-like questions.. After these steps, we obtain the final curated dataset (shown in purple in Figure 1). **Data filtering** During filtering, *LLMSynth* is used to predict the output of the given synthetic example using k = 3 tries. If the output cannot be predicted at least once, it is assumed the example is low quality and is not included in the final curated dataset.

Data Imputation We also consider an imputation process, which involves blanking parts of the augmented data points and using the LLM to fill in the blanks, to replace those fields. This is to provide cleaner data which is less unnatural.

## 170 3.3 MODEL FINE-TUNING

At this stage, we fine-tune on the curated synthetic dataset, initializing from a base or instructiontuned version of the LLM. We use our dataset for supervised training of both the reasoning chain and the final answer. The resulting model *LLMCurated* is then ready to perform the desired task.

174 175 176

177

171

172

173

## 4 SOURCE2SYNTH'S APPLICATIONS

The general pipeline described above can be used to produce examples for the task at hand and to teach LLMs new skills. To demonstrate the impact of *Source2Synth*, we apply it to two challenging tasks where LLMs struggle which are also areas of great interest for the community: multi-hop question answering and tabular question answering.

182 183 184

## 4.1 MULTI-HOP QUESTION ANSWERING

In multi-hop question answering (MHQA), we generate a dataset of multi-hop question-answer pairs,
 enriched with the reasoning chain that is used to answer the question. The chain consists of question
 decomposition into sub questions with answers, plus the entity that links them.

See Figure 2 for an overview of the procedure and Figure 3 - Right for an example response from the model that underwent the *Source2Synth*'s pipeline.

190 191

192

## 4.1.1 DATASET GENERATION

**Data source selection** For multi-hop question answering, we pick English Wikipedia (Wikipedia contributors, 2004) as the data source, since it contains articles in natural language as well as additional meta-information like links to related articles. The data generation process starts by randomly selecting an initial article, denoted as  $D_1$ , among all available Wikipedia articles. For each  $D_1$  we collect  $n \ge 2$  related articles.

**Seed** An MHQA seed topic corresponds to an entity E retrieved from  $D_1$ . The seed in MHQA doubles also as the "hop" in the multi-hop question Q that we aim to generate, since E links the n = 2 subquestions that compose Q. For example, in Figure 2, we sample "The Moon" article at random, denoted by  $D_1$ , and the corresponding entity, denoted by E, is "Apollo 11" (displayed in blue). Then, we pick "Neil Armstrong" as  $D_2$  from the pool of related articles, since it contains a paragraph where the entity "Apollo 11" is included.

204 Dataset construction We prompt an instruction-tuned language model to generate two questions: a 205 question  $Q_1$  based on  $D_1$  and whose answer is the selected entity E, and a second question  $Q_2$  based 206 on  $D_2$  such that its main topic is E. See Figures 18 and 19 for the exact prompts. For example, 207 in Figure 2,  $Q_1$  = "What was the spaceflight that first landed humans on the Moon?", the hop is E = "Apollo 11" and  $Q_2 =$  "Who was the commander of Apollo 11?". We then prompt the LLM to 208 merge the two questions, in order to generate the final two-hop question Q by using the entity as a 209 conceptual link (hop). The exact prompt is given in Figure 17. For MHQA, we generated a total of 210 1250 synthetic questions starting from a collection of 50 Wikipedia articles (randomly selected). 211

212

214

## 213 4.1.2 DATASET CURATION

In the MHQA experiments, the curation step removed around 13% of the questions originally generated.



Figure 2: Source2Synth synthetic data generation process for multi-hop question answering. 237 The method first randomly picks one article  $D_1$ , in this case with title "The Moon". At the Seed stage, 238 an entity E is selected from  $D_1$ 's pool of entities, "Apollo 11". Then, documents are sampled from 239 the related documents pool of  $D_1$  such that E is present, and  $D_2$ , "Neil Armstrong", is selected. A 240 question  $Q_1$  is then generated from  $D_1$  with the constraint that the answer  $A_1$  is the entity itself. A 241 second question  $Q_2$  is then generated from  $D_2$  with the constraint that its main topic is the entity. We 242 then prompt an LLM to merge the two questions based on the link/entity they have in common to 243 produce the final question, reasoning chain and answer that comprise the training example. 244

246

**Data filtering** We check if the predicted answer matches the answer in the synthetically generated example, and if after k = 3 tries tries the LLM has not supplied the correct answer we filter out the entry entirely. See Figure 3 - Left for an example of model inference.

**Data Imputation** For MHQA, we blank  $Q_1$  and provide the LLM with Q,  $Q_2$ , E, and the relative doc sample  $D_1$  as context when asking it to reconstruct  $Q_1$ . The new candidate  $Q'_1$  for  $Q_1$  is then assessed: if A' (the answer to the new multi-hop question Q' resulting from piecing together  $Q'_1$ and  $Q_2$ ) matches A (the original answer to Q) then we keep the example. We find that asking the model to reconstruct parts of the multi-hop question in-context results in a more natural and cohesive question, thus removing some of the unnaturalness of the text that can occur from automatically generated and merged examples (see Appendix **??** for more details).

- 257
- 258 259

## 4.2 TABULAR QUESTION ANSWERING

260 261 262

263

In Tabular question answering (TQA) we generate a question-answer dataset where each question is based on a table from the data source. Generated training examples are hence enriched with annotations built from automatically-generated interesting facts retrieved from the table.

Data source selection In the TQA case, we use 4k unlabeled tables from the WikiSQL training dataset as sources (Zhong et al., 2017).

Seed We then prompt an instruction-tuned language model to generate a statement based on the table.
 This statement corresponds to our seed topic for the generation and is a pertinent interesting fact or set of observations in natural language that can be derived from the table. The prompt is given in Figure 13.



Figure 3: Left: Example Source2Synth Response on MHQA (closed book inference). We show the model's response (reasoning steps and answer) to a multi-hop input question (yellow). The colours highlight the generation of the corresponding augmented entries: the decomposition into sub questions  $Q_1$  and  $Q_2$  in green, the seed  $A_1$  in blue, and the final answer  $A_2$  in red.

**Right: Example** *Source2Synth* **Response on TQA.** We show the model's response (SQL call and final answer) to the tabular input question (yellow). The coloured parts highlight the generation of the corresponding augmented entries: SQL in green, and the final answer in red.



Figure 4: Source2Synth synthetic data generation process for Tabular question answering. The
method first generates the seed, which is a fact based on the table (shown in blue). Given the seed
and table, an *SQL query* is then generated (in green) as well as its translation into natural language
(the question *Q*). Then the SQL is executed on the table to obtain the answer A.

# 324 4.2.1 DATASET CONSTRUCTION

We next generate an SQL-statement by zero-shot prompting the LLM: we provide the table and the seed (factual statement) as context, see Figure 14 for the exact prompt. Given the produced SQL statement, it is then executed using the Python library *sqlite3*<sup>1</sup> to obtain an SQL answer formatted as a table. If the generated statement is invalid, we discard it and re-generate. We generate a total of 10k SQL statements based on the source tables. Checking the statements for validity (i.e. refusing non-executable SQL statements) brings the dataset size to 8k (per slice).

332333 4.2.2 DATASET CURATION

333 334 335

340 341

342 343

344

345

346 347

348

367

368

369 370

372

373

374

375

376 377 In Tabular QA, the curation process consists only of the filtering step. After curation, we keep 2160 (27%) of the original examples in slice 1.

**Data filtering** We check if the predicted answer of *LLMSynth* fine-tuned on slice 0 matches the answer in the synthetically generated example, and if after k = 3 tries the model has not supplied the correct answer we filter out the entry entirely. See Figure 3 - Right for an example of model inference.

## 5 EXPERIMENTAL SETUP

We test our method on two domains: *tabular question answering* and *multi-hop question answering*. For each, we use *Source2Synth* to generate and curate a high quality dataset suitable for fine-tuning, and compare our method to a number of baselines.

5.1 MULTI-HOP QA SETUP

Data To evaluate *Source2Synth* on MHQA, we evaluate it on HotpotQA (Yang et al., 2018): a
 benchmark based on Wikipedia containing 113,000 examples of multi-hop QA pairs, split in train, test, and validation sets.

A comparison question entails comparing the same concept between n objects (e.g. "Who is the 353 tallest student in class?"), while a bridge question builds on a logical and/or causal link and requires 354 deriving statements to get to the answer (e.g. "What is the height of the student that topped the 355 entry exam?" - this requires first identifying the student that topped the exam). The hop length is the 356 number of comparison objects for comparison questions or the number of links for bridge questions. 357 In our case, we chose n = 2 to be consistent with HotpotQA. The test set consists of 7,405 entries, 358 split evenly between bridge and comparison questions. We only generate synthetic data for bridge 359 questions, since they pose a bigger challenge to current LLMs and to counterbalance this disparity, 360 we include 500 comparison questions from HotpotQA 's training dataset in our fine-tuning dataset.

Metrics We measure the performance using soft exact match (soft-EM) as the metric. Soft-EM is 1 if the generated output contains the golden answer and 0 otherwise.

Model In MHQA experiments we chose Llama-2 70B-Chat and we fine-tune *Source2Synth* and various other baseline methods initializing from this model. *Source2Synth* is trained with 1250 synthetic examples, unless noted otherwise, in addition to the 500 HotpotQA examples above.

**Baselines** We compare *Source2Synth* to the following baselines:

- Instruction-tuned LLM: using LLama 2 70B-Chat for the task in a zero-shot manner.
- Fine-tuned LLM (HotpotQA only): fine-tuning from the base model on 500 HPQA examples from the training split.
- *LLMSynth* (Synthetic dataset only): training our model with 1250 synthetic examples from Slice 0 (see Figure 1), *without* the data curation step.
- *LLMSynth* (Synthetic and HotpotQA ): training with the uncurated synthetic data in addition to the 500 HPQA examples.

<sup>&</sup>lt;sup>1</sup>https://www.sqlite.org

378 For all the models listed, we tested them using two prompting methods: a zero-shot and a three-shot 379 CoT prompt, see the Appendix E for details. 380

381 5.2 TABULAR QA SETUP 382

**Data** We conduct evaluations with the WikiSOL (Zhong et al., 2017)'s validation split. WikiSOL 384 consists of a corpus of 80,654 hand-annotated examples of natural language questions, SQL queries, and SQL tables created from 24,241 tables extracted from Wikipedia. The validation split contains 385 386 7,857 examples after removing non-executable SQL tables, see Appendix B for more details.

387 **Metrics** We measure performance using the exact match (EM) and the soft-EM metrics. The EM 388 metric equals 1 if the golden answer is equal to the generated answer and 0 otherwise.

389 Model For TQA, we use the Starchat-beta language model Li et al. (2023b) from Huggingface as the 390 initial language model (batch size 32, 100 steps, lr 0.0001, linear warm-up). The Starchat model is an 391 instruction-tuned LLM with 16 billion parameters trained to act as a helpful coding assistant. This 392 model is a fine-tuned version of StarCoder Li et al. (2023b), a LLM which was pre-trained and then 393 fine-tuned on a large code corpus, which contains SQL statements, and successively fine-tuned on 394 35B Python tokens.

**Baselines** We compare the performance of our *Source2Synth* method against a variety of baselines. The baselines consist of prompting the Starchat-beta instruction-tuned language model as follows: 397

- Zero-shot Table QA: prompt with the task instruction, the table and the question in a zero-shot fashion. See Figure 9 for the prompt.
- One-Shot No Context QA: prompt with the task instruction and a one-shot example containing a question and answer, together with the actual question for the model to answer. See Figure 10 for the prompt.
- One-Shot Table QA: prompt that includes the table for both the one-shot example and the question to be answered. We use one-shot due to LLM context length and the typically large size of the tables. See Figure 11 for the prompt.
- One-shot Table+SQL QA: the prompt includes an example containing the table and question, and an instruction suggesting that the model can leverage an SQL tool. We then execute the predicted SQL to obtain the answer. See Figure 12 for the prompt.
  - *LLMSynth*: Fine-tune the model with synthetic data *without* applying the data curation step.
- RESULTS 6 412
  - 6.1 MULTI-HOP QUESTION ANSWERING

Table 1: Evaluation of Source2Synth on Multi-hop question answering. The models shown are 416 fine-tuned with 500 entries from HotpotQA ('HotpotQA ") and/or 1250 entries from the Source2Synth 417 Synthetic Dataset ("Synthetic Dataset"). Using Source2Synth curated synthetic data in combination 418 with HotpotQA (last row) works best. 419

420			
421	Method	0-shot	3-shot CoT prompt
422	Instruction-tuned LLM (LLama 2 70B-Chat)	40.45%	44.13%
423	fine-tuned LLM (HotpotQA only)	53.22%	58.40%
424	LLMSynth (Synthetic dataset only)	52.31%	56.70%
425	LLMSynth (Synthetic and HotpotQA)	57.46%	62.73%
426	LLMCurated (Synthetic and HotpotQA)	65.23%	66.05%

42 427

399

400

401

402

403

404

405

406

407

408

409

410 411

413

414 415

**Overall performance of** Source2Synth on MHQA We report the experimental results in Table 1. We 428 include the baselines of the vanilla instruction-tuned LLM (0-shot and 3-shot, please see Prompt 16), a 429 fine-tuned LLM using only the HPQA 500 examples from the train split (second row), and LLMSynth 430 which only uses the uncurated synthetic data for fine-tuning (third row). All fine-tuned methods 431 outperform the instruction-tuned model (first row). Using only synthetic data or only HotpotQA Table 2: Analysis of MHQA bridge and comparison questions with respect to level of difficulty.
We evaluate models on the full train dataset (where questions are labelled with easy, medium and hard). *Source2Synth* outperforms both the baseline and the model fine-tuned on HotpotQA, yielding an LLM capable of handling hard questions of both types.

	Bridge		Comparison			
Model	Hard	Medium	Easy	Hard	Medium	Easy
Llama2-70B-Chat	14.5%	27.2%	30.1%	66.6%	71.3%	73.2%
Fine-tuned LLM (HotpotQA only)	20.1%	29.8%	34.3%	74.5%	78.3%	82.1%
LLMCurated-1250	31.3%	35.6%	39.7%	83.1%	85.7%	87.8%

data for fine-tuning demonstrates worse performance than when combined, whether the synthetic data is curated, as in *LLMCurated* (fifth row) or not, as in *LLMSynth* (fourth row). Once we use the full *Source2Synth* pipeline to obtain the curated synthetic dataset for fine-tuning we see further performance improvements *LLMCurated* (fifth row) over not curating the data (fourth row).



Figure 5: Synthetic Data scaling performance. We show how the performance of *Source2Synth* changes with respect to MHQA data mix size, both before and after curation. During the curation step, the following percentages of samples were removed: 7% for 500, 8% for 750, 11% for 1250. *LLMSynth* (before curation) performs worse than *LLMCurated* (after curation) despite having more samples – but both approaches improve with more data.

Analysis of performance on different question types and levels of difficulty We study the capabilities of our model by analysing the performance of LLM-Curated-1250 with particular focus on the
type and difficulty of the questions – namely hard/medium/easy bridge and comparison questions. We
compare the performance of the base model, the model fine-tuned on HotpotQA , and *Source2Synth*according to the difficulty level, as provided by the HotpotQA train dataset. We also subdivide the
results according to the type of question (bridge vs. comparison). Results are given in Table 2.

We observe that *Source2Synth* performs better across all types of questions and difficulties, with an average overall gain of 12.4% on the base LLM and a 7.5% gain compared to the LLM fine-tuned on HotpotQA . In particular, by applying our method, the resulting model is able to achieve +16.8% and +16.5% on hard bridge and comparison questions respectively, when comparing to the baseline.
Furthermore, it is interesting to see substantial improvement on comparison-type questions, despite not explicitly targeting those during synthetic generation. Hard questions pose a greater challenge to the reasoning abilities of LLMs and these results introduce *Source2Synth* as a possible method for further improvement.

486 Table 3: Tabular question answering. The models are fine-tuned using Source2Synth curated 487 synthetic data only. Performance comparison on the WikiSQL evaluation dataset.

488			
489	Method	Exact Match	Soft-EM
490	One-Shot No Context QA (Starchat-beta LLM)	0.25%	16.22%
491	Zero-shot Table QA (Starchat-beta LLM)	1.83%	20.07%
492	One-Shot Table QA (Starchat-beta LLM)	2.03%	31.06%
494	Une-snot Table+SQL QA (Starchat-beta LLM) UMSwath (Synthetic dataset only)	12.30%	34.13% 34.21%
495	<i>LLMCurated</i> (Synthetic dataset only)	34.50%	42.80%

496

497 498

499

500

501

502

503

504

505

506

507

Scaling performance We also report scaling performance in Figure 5. We study how performance evolves when adding more synthetic data in the fine-tuning data mix - that already includes 500 samples from the HPQA train split. We perform the analysis on *LLMSynth* and *LLMCurated* to show the impact of the curation technique. In both cases and in all data mixes, we see that applying the Source2Synth pipeline results in a stronger model on the task. For the LLMSynth model fine-tuned on uncurated samples we see that providing more synthetic examples leads to a steady improvement in performance across all data sizes, for both zero-shot and three-shot prompting variants. LLMCurated follows a similar trend, but consistently outperforms the uncurated version of the model, for all training set sizes. Overall, we observe that using our synthetic data generation pipeline to construct more data brings further performance gains in the task.

#### 508 6.2 TABULAR QUESTION ANSWERING 509

510 We report the experimental results for Tabular question answering in Table 3. Firstly, we see that 511 providing no context about the table when prompting the instruction-tuned StarChat language model 512 has very poor performance (first row), with an EM metric of 0.25%. This is expected, since questions 513 in WikiSQL require information contained in the table to answer, while the model does not have any 514 other information except for the general knowledge stored in its parameters. However, even if we 515 pass the table as part of the prompt, the performance does not improve much due to its difficulties 516 to digest structured data. For example, passing in a zero-shot fashion (second row) only has an EM 517 metric of 1.83%. While passing an example of table usage in a one-shot fashion (third row) improves the soft-EM metric, the EM metric is still very low (2.03%). Hence, this is still very challenging for 518 the model. Thirdly, the performance increases once we provide a one-shot example containing the 519 relevant table and SQL query (fourth row), with an EM of 12.3%. The ability to use the SQL tool 520 improves performance markedly. 521

522 We obtain a significant increase in performance when we fine-tune the StarChat model using the 523 Source2Synth curated data (last row), with an EM of 34.5%. Our full method performs significantly better than fine-tuning the StarChat language model using synthetic data without curation, LLMSynth 524 (second to last row) which has an EM of 23.86%, although that still outperforms the other baselines 525 by a large margin as well, indicating the utility of our Source2Synth synthetic data generation scheme. 526

527 528

529

#### 7 LIMITATIONS

530 In this paper, our applications use a single seed or table per query to derive questions. However, 531 Source2Synth can be extended to more complex scenarios e.g. multiple hops or queries that require 532 multi-table tool-use. This can be done by looping the dataset generation steps and feeding the result 533 of the previous step as input to the next one. Our method could also be improved with more clever 534 sampling techniques. We consider this to be an interesting avenue of future research.

536

#### CONCLUSION 8

- 537 538
- In this paper, we introduce *Source2Synth*, a new method for generating and curating high-quality synthetic data grounded in real data sources. We demonstrate its utility on two tasks that pose

significant challenges for LLMs: multi-hop reasoning and tabular question answering with SQL. Our
 work could also be beneficial in other low-data regimes and on other tasks and in diverse fields.

References

543

544

546

547

548

549

550

- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- Tianle Cai, Xuezhi Wang, Tengyu Ma, Xinyun Chen, and Denny Zhou. Large language models as tool makers, 2024. URL https://arxiv.org/abs/2305.17126.

551 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam 552 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam 553 Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James 554 Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, 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. 558 Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon 559 Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, 561 Slav Petrov, and Noah Fiedel. Palm: Scaling language modeling with pathways, 2022. URL https://arxiv.org/abs/2204.02311.

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
   bidirectional transformers for language understanding, 2019. URL https://arxiv.org/abs/
   1810.04805.
- 567 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha 568 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston 569 Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, 570 Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris 571 McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton 572 Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David 573 Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, 574 Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip 575 Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme 576 Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, 577 Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, 578 Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, 579 Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianvu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz 582 Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence 583 Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas 584 Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, 585 Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan 588 Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia 592 Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla,

594 Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek 595 Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, 596 Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent 597 Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, 598 Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe 600 Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya 601 Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex 602 Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei 603 Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew 604 Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley 605 Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin 606 Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, 607 Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt 608 Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide 610 Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, 611 Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily 612 Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix 613 Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank 614 Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, 615 Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid 616 Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen 617 Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-618 Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste 619 Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, 620 Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik 621 Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly 622 Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, 623 Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, 624 Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria 625 Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, 626 Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle 627 Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, 628 Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, 629 Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, 630 Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia 631 Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, 632 Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, 633 Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan 634 Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara 635 Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh 636 Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, 637 Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, 638 Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan 639 Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, 640 Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe 641 Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, 642 Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, 644 Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, 645 Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, 646 Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd 647 of models, 2024. URL https://arxiv.org/abs/2407.21783.

- Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. Pal: Program-aided language models, 2023. URL https://arxiv.org/abs/ 2211.10435.
- Carlos Gemmell and Jeffrey Dalton. Generate, transform, answer: Question specific tool synthesis
   for tabular data, 2023. URL https://arxiv.org/abs/2303.10138.
- Jonathan Herzig, Pawel Krzysztof Nowak, Thomas Müller, Francesco Piccinno, and Julian Eisenschlos. TaPas: Weakly supervised table parsing via pre-training. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4320–4333, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.398. URL https://aclanthology.org/2020.acl-main.398.
- Zhengbao Jiang, Frank F. Xu, J. Araki, and Graham Neubig. How can we know what language models know? *Transactions of the Association for Computational Linguistics*, 8:423–438, 2019. URL https://api.semanticscholar.org/CorpusID:208513249.
- Jinyang Li, Binyuan Hui, Ge Qu, Binhua Li, Jiaxi Yang, Bowen Li, Bailin Wang, Bowen Qin, Ruiying
  Geng, Nan Huo, Xuanhe Zhou, Chenhao Ma, Guoliang Li, Kevin C. C. Chang, Fei Huang, Reynold
  Cheng, and Yongbin Li. Can Ilm already serve as a database interface? a big bench for large-scale
  database grounded text-to-sqls, 2023a.
- 668 Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, 669 Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, Terry Yue 670 Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João Monteiro, 671 Oleh Shliazhko, Nicolas Gontier, Nicholas Meade, Armel Zebaze, Ming-Ho Yee, Logesh Kumar 672 Umapathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo Wang, Rudra Murthy, Jason 673 Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey, Zhihan Zhang, 674 Nour Fahmy, Urvashi Bhattacharyya, Wenhao Yu, Swayam Singh, Sasha Luccioni, Paulo Villegas, Maxim Kunakov, Fedor Zhdanov, Manuel Romero, Tony Lee, Nadav Timor, Jennifer Ding, Claire 675 Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank Mishra, Alex Gu, Jennifer Robinson, 676 Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contractor, Siva Reddy, Daniel Fried, 677 Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Arjun 678 Guha, Leandro von Werra, and Harm de Vries. Starcoder: may the source be with you!, 2023b. 679 URL https://arxiv.org/abs/2305.06161. 680
- Kian Li, Ping Yu, Chunting Zhou, Timo Schick, Omer Levy, Luke Zettlemoyer, Jason Weston, and Mike Lewis. Self-alignment with instruction backtranslation, 2024. URL https://arxiv.org/ abs/2308.06259.
- Alisa Liu, Swabha Swayamdipta, Noah A. Smith, and Yejin Choi. WANLI: Worker and AI collaboration for natural language inference dataset creation. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 6826–6847, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-emnlp.508. URL https://aclanthology.org/2022.findings-emnlp.508.
- Ruibo Liu, Jerry Wei, Fangyu Liu, Chenglei Si, Yanzhe Zhang, Jinmeng Rao, Steven Zheng, Daiyi
   Peng, Diyi Yang, Denny Zhou, and Andrew M. Dai. Best practices and lessons learned on synthetic
   data, 2024. URL https://arxiv.org/abs/2404.07503.
- Dheeraj Mekala, Jason Weston, Jack Lanchantin, Roberta Raileanu, Maria Lomeli, Jingbo Shang, and Jane Dwivedi-Yu. Toolverifier: Generalization to new tools via self-verification, 2024. URL https://arxiv.org/abs/2402.14158.
- Grégoire Mialon, Roberto Dessi, Maria Lomeli, Christoforos Nalmpantis, Ramakanth Pasunuru,
   Roberta Raileanu, Baptiste Roziere, Timo Schick, Jane Dwivedi-Yu, Asli Celikyilmaz, Edouard
   Grave, Yann LeCun, and Thomas Scialom. Augmented language models: a survey. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL https://openreview.net/forum?
   id=jh7wH2AzKK. Survey Certification.

702 Thao Nguyen, Jeffrey Li, Sewoong Oh, Ludwig Schmidt, Jason Weston, Luke Zettlemoyer, and Xian 703 Li. Better alignment with instruction back-and-forth translation, 2024. URL https://arxiv. 704 org/abs/2408.04614. 705

- Bhargavi Paranjape, Scott Lundberg, Sameer Singh, Hannaneh Hajishirzi, Luke Zettlemoyer, and 706 Marco Tulio Ribeiro. Art: Automatic multi-step reasoning and tool-use for large language models, 2023. URL https://arxiv.org/abs/2303.09014. 708
- 709 Aaron Parisi, Yao Zhao, and Noah Fiedel. Talm: Tool augmented language models, 2022. URL 710 https://arxiv.org/abs/2205.12255.
- 711 Panupong Pasupat and Percy Liang. Compositional semantic parsing on semi-structured tables. 712 In Chengqing Zong and Michael Strube (eds.), Proceedings of the 53rd Annual Meeting of 713 the Association for Computational Linguistics and the 7th International Joint Conference on 714 Natural Language Processing (Volume 1: Long Papers), pp. 1470–1480, Beijing, China, July 715 2015. Association for Computational Linguistics. doi: 10.3115/v1/P15-1142. URL https: 716 //aclanthology.org/P15-1142. 717
- Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H. Miller, 718 and Sebastian Riedel. Language models as knowledge bases? In Conference on Empirical Methods 719 in Natural Language Processing, 2019. URL https://api.semanticscholar.org/CorpusID: 720 202539551. 721
- 722 Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen, Ning Ding, Ganqu Cui, Zheni Zeng, Yufei Huang, 723 Chaojun Xiao, Chi Han, Yi Ren Fung, Yusheng Su, Huadong Wang, Cheng Qian, Runchu Tian, 724 Kunlun Zhu, Shihao Liang, Xingyu Shen, Bokai Xu, Zhen Zhang, Yining Ye, Bowen Li, Ziwei 725 Tang, Jing Yi, Yuzhang Zhu, Zhenning Dai, Lan Yan, Xin Cong, Yaxi Lu, Weilin Zhao, Yuxiang 726 Huang, Junxi Yan, Xu Han, Xian Sun, Dahai Li, Jason Phang, Cheng Yang, Tongshuang Wu, Heng Ji, Zhiyuan Liu, and Maosong Sun. Tool learning with foundation models, 2023. URL 727 https://arxiv.org/abs/2304.08354. 728
- 729 Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language under-730 standing with unsupervised learning. 2018. 731
- 732 Timo Schick and Hinrich Schütze. Few-shot text generation with natural language instructions. In Conference on Empirical Methods in Natural Language Processing, 2020. URL https: 733 //api.semanticscholar.org/CorpusID:238260199. 734
- 735 Timo Schick and Hinrich Schütze. Generating datasets with pretrained language models. ArXiv, 736 abs/2104.07540, 2021. URL https://api.semanticscholar.org/CorpusID:233241169. 737
- Timo Schick and Hinrich Schütze. Exploiting cloze questions for few shot text classification and 739 natural language inference, 2021. URL https://arxiv.org/abs/2001.07676.

738

- 740 Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke 741 Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach 742 themselves to use tools. In Thirty-seventh Conference on Neural Information Processing Systems, 743 2023. URL https://openreview.net/forum?id=Yacmpz84TH. 744
- 745 Qiaoyu Tang, Ziliang Deng, Hongyu Lin, Xianpei Han, Qiao Liang, Boxi Cao, and Le Sun. Toolalpaca: Generalized tool learning for language models with 3000 simulated cases. ArXiv, 746 abs/2306.05301, 2023. URL https://api.semanticscholar.org/CorpusID:259108190. 747

748 Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze 749 Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, YaGuang Li, Hongrae Lee, Huaixiu Steven 750 Zheng, Amin Ghafouri, Marcelo Menegali, Yanping Huang, Maxim Krikun, Dmitry Lepikhin, 751 James Qin, Dehao Chen, Yuanzhong Xu, Zhifeng Chen, Adam Roberts, Maarten Bosma, Vincent Zhao, Yanqi Zhou, Chung-Ching Chang, Igor Krivokon, Will Rusch, Marc Pickett, Pranesh 752 Srinivasan, Laichee Man, Kathleen Meier-Hellstern, Meredith Ringel Morris, Tulsee Doshi, 753 Renelito Delos Santos, Toju Duke, Johnny Soraker, Ben Zevenbergen, Vinodkumar Prabhakaran, 754 Mark Diaz, Ben Hutchinson, Kristen Olson, Alejandra Molina, Erin Hoffman-John, Josh Lee, 755 Lora Aroyo, Ravi Rajakumar, Alena Butryna, Matthew Lamm, Viktoriya Kuzmina, Joe Fenton,

756 757 758	Aaron Cohen, Rachel Bernstein, Ray Kurzweil, Blaise Aguera-Arcas, Claire Cui, Marian Croak, Ed Chi, and Quoc Le. Lamda: Language models for dialog applications, 2022. URL https://arxiv.org/abs/2201.08239.
759 760 761 762 763 764 765 766 766 767 768 769 770 771	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023. URL https://arxiv.org/abs/2307.09288.
772 773 774 775	Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In <i>Neural Information Processing Systems</i> , 2017. URL https://api.semanticscholar.org/CorpusID:13756489.
776 777 778	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, Quoc Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. <i>CoRR</i> , abs/2201.11903, 2022. URL https://arxiv.org/abs/2201.11903.
779 780	Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. Magicoder: Empowering code generation with oss-instruct, 2024. URL https://arxiv.org/abs/2312.02120.
781 782 783 784	Wikipedia contributors. Plagiarism — Wikipedia, the free encyclopedia, 2004. URL https: //en.wikipedia.org/w/index.php?title=Plagiarism&oldid=5139350. [Online; accessed 22-July-2004].
785 786 787 788	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In <i>Conference on Empirical Methods in Natural Language Processing</i> , 2018. URL https://api.semanticscholar.org/CorpusID:52822214.
789 790 791	Victor Zhong, Caiming Xiong, and Richard Socher. Seq2sql: Generating structured queries from natural language using reinforcement learning. <i>CoRR</i> , abs/1709.00103, 2017.
792	
793	
794	
795	
790	
708	
790	
800	
801	
802	
803	
804	
805	
806	
807	
808	

#### MORE RESULTS ON PROMPT ENGINEERING A

Prompt Type	Model Accuracy (soft-EM)
0-shot	40.45%
Role	22.34%
1-shot	26.65%
Few-shots (5-shots)	21.83%
Role (1-shot)	28.29%

Table 4: MHQA prompts sweep. Overview of the model's accuracy across different prompt strategies. Role "You are a QA-robot. Answer the following question:". Model used: Llama-2-70B-Chat, Dataset: HotpotQA test.

#### В SQL NON-EXECUTABLE CODE FILTERING

We discard incorrect SQL statements - i.e. whose execution with sqlite3<sup>2</sup> leads to an error. Discarded proportion: out of 50 tables, we generate 800 seed statements and the number of valid (executable) SQL statements was 658.

#### С ILLUSTRATED EXAMPLE: ADAPTING TO A NEW TASK

There are (at max.) three components in *Source2Synth*'s pipeline that need to be changed in order to adapt it for a new task: prompts, seed, and data source. If the new task builds on document- or table-use, some parts of the two applications showcased in the main body of the paper can be reused. We proceed illustrating how to adapt prompts, seed, and data source in case of the following new task: generating code (Python functions) to compute statistics on spreadsheets. Since this task builds on tables, we integrate TQA in it by leveraging its seed generation system. We believe that this is a meaningful and diverse example as it showcases a shorter tasks that does not produce a QA dataset and that leverages structured data. Steps for adaptations: 

- 1. Select the data source The data source is the collection of spreadsheets of interest to the user. This can be a custom dataset, or a public one (like WikiSQL).;
- 2. Select the seed It is important to condition the seed generation so that it reflects the goal of the user. For example since we would like to compute statistics on the tables in the dataset, it is important to define which statistics are of interest to the user (max / min value, average, median, etc...);
- 3. Changing the seed generation prompt so that it conditions the generation of the statement to focus on the interests of the user. For example, in practice we would update the prompt in Fig.13 as follows - please see Fig. 6;
- 4. Using the seed to generate code Adapt the data generation prompt to output Python functions biased on the statement produced - please see Fig. 7.

<sup>2</sup>https://www.sqlite.org

870

871 872

873 874

875

876

877

878

879

882 883

884 885 886

887 888

889

890

891

892

893

894

895

896 897

899

900

901

902

903

904

905

906

907 908

909

910

911 912 913

914 915

916 917

## Generating a seed in the new task.

Please generate an interesting statement about this table. The statement is a fact about one of the columns in the following table. The statement should include one of these metrics: {statistics\_of\_interest\_list}

{table}

As a result of this, an interesting statement about the table and the metrics is:

## Figure 6: Prompt adaptation for seed generation (new task)

## Generating code in the new task.

Please generate a Python function based on the table and statement below. The function must solve the task described by the statement. The table is an input variable to the function.  $\{table\}$ 

Statement: {seed}

As a result of this, a Python function that solves the problem described by the statement is:

Figure 7: Prompt adaptation for output generation (new task)

#### D ILLUSTRATED EXAMPLE: ADAPTING TO A NEW DOMAIN

Similarly to across tasks, there are three components to adapt for a new domain application: prompts, seed, and data source. We illustrate an example of QA task in a medical domain. Steps for adaptations:

- 1. Select the data source The data source is a collection of medical documents (either private or public like PubMed <sup>3</sup>);
- 2. Generating the OA dataset Since the task is the same as MHQA (but just applied to a different domain - medical), once the data source is adapted we can leverage the same pipeline to generate the seed and QA pairs. Please see Fig. 2 the seed and construction stages for a visual reminder of the pipeline, and 18, 19 for the prompts;

## Comparing Q pre- and post- imputation

### Before imputation:

Q: "What pet did the poet and father of mathematician Ada Lovelace had when he was a student at Trinity out of resentment for rules forbidding pet dogs like his beloved Boatswain?"  $Q_1$ : "What pet did the poet Lord Byron had when he was a student at Trinity out of resentment for rules forbidding pet dogs like his beloved Boatswain?  $Q_2$ : "Who is the father of mathematician Ada Lovelace?" E: "Lord Byron" A: "A bear  $D_1$ : "Lord Byron also kept a tame bear while he was a student at Trinity out of resentment for rules forbidding pet dogs like his beloved Boatswain." After imputation: Q': "What pet did the poet and father of mathematician Ada Lovelace had when he was a student at Trinity?  $Q'_1$ : "What pet did the poet Lord Byron had when he was a student at Trinity?" Figure 8: Comparing Q pre- and post- imputation

#### Ε **PROMPTS USED IN OUR EXPERIMENTS**

<sup>&</sup>lt;sup>3</sup>https://huggingface.co/datasets/ncbi/pubmed,https://pubmed.ncbi.nlm.nih.gov/about/

## Zero-shot Table QA prompt.

Answer the following question using the table below. You may leverage an SQL tool.

Q: {question}

 $\{table\}$ 

## Figure 9: Zero-shot Table QA prompt for the TQA task.

## One-Shot No context QA prompt.

Example Q: What was the last year where this team was part of the US A-league?
A: 2004
Now do the same for the following question.
Q: {question}

## Figure 10: One-Shot No context QA prompt for the TQA task.

Season	Team	League_apps	Goals
923	Swindon Town	55	3
922	Swindon Town	14	4
921	Swindon Town	24	11
920	Swindon Town	26	16
919	Swindon Town	20	10
914	Swindon Town	23	12
913	Swindon Town	24	18
912	Swindon Town	12	9
911	Swindon Town	20	16
910	Swindon Town	30	19
909	Swindon Town	33	19
908	Swindon Town	34	28
907	Swindon Town	30	17
Q: How A: 97	y many league ap	pearances were	there between 1907 and 1909 (inclusive)?
low do	the same for the	e following tabl	e and question.

Figure 11: One-shot Table QA prompt for the TQA task.



## One-shot Table+SQL QA prompt.

-- Example --Answer the following question using the table below. You may leverage an SQL tool. The table is stored in a variable 'sql\_table' and has the following schema: Season | Team League\_apps | Goals 1923 |Swindon Town 55 3 1922 |Swindon Town | 14 4 Q: How many league appearances were there between 1907 and 1909 (inclusive)? SQL: SELECT SUM(League\_apps) FROM sql\_table WHERE Season BETWEEN 1907 AND 1909 | Result result | 97 Now do the same for the following table and question. {table}

Q: {question}

Figure 12: One-shot Table+SQL QA prompt for the TQA task.

## Generating a seed in TQA.

Please generate an interesting statement about this table. The statement is a fact about one of the columns in the following table. {table}

An interesting statement as a result of this is:

Figure 13: Prompt used to induce a pertinent and interesting seed topic in TQA. This is done zero-shot.

## Generating meaningful SQL in TQA.

Please generate SQL statements for the following table:

 $\{table\}$ 

Seed:  $\{seed\}$ 

An interesting SQL statement as a result of this is

Figure 14: Prompt used to induce a meaningful SQL statement given the table and seed for the TQA task. This is done zero-shot.

## Generating a question in TQA.

I want to convert an SQL statement into a question. Here is the original table: {table} SQL: {SQL}

What is the question that this SQL statement would be the answer to?

Figure 15: Prompt used to induce a meaningful question using the table and generated SQL query for the TQA task. This is done zero-shot.

1026	
1027	
1028	
1029	Three-shot CoT prompt used at evaluation time on MHQA.
1030	Answer the following multi-hop question 'Q' by decomposing it into 'Q1' and 'Q2' and
1031	solving them step-by-step. Learn from the following 3 examples. As shown in the following
1032	example:
1033	Example #1
1034	'Q' = 'Who was the commander of the spaceflight that first landed humans on the Moon?'
1035	1 Splitting 'Q' into 'Q1' and 'Q2':
1036	'Q1': 'What was the spaceflight that first landed humans on the Moon?';
1037	'Q2': 'Who was the commander of [A1]?';
1038	2. Answering Q1:
1039	The answer 'A1' to 'Q1' : 'What was the spaceflight that first landed humans on the Moon?'
1040	is 'Apollo 11'. 'A1' = 'Apollo 11'
1041	3. Substituting A1 to Q2:
1042	'Q2': 'Who was the commander of Apollo 11?',
1043	4 Answers Q2:
1045	The answer 'A2' to Q2': 'Who was the commander of Apollo 11?' is 'Neil Armstrong'.
1046	(A2) = (A) = (Neil Armstrong)
1047	Example #2
1048	'Q' = 'What' is the main ingredient in the flagship product of Ferrero?'
1049	1. Splitting $(\Omega)$ into $(\Omega)$ and $(\Omega)$ .
1050	'Q1': 'What is the flagship product of Ferrero?'
1051	'Q2': 'What is the main ingredient in $[A1]$ ?'
1052	2. Answering Q1:
1053	The answer 'A1' to 'Q1': 'What is the flagship product of Ferrero?' is Nutella':'A1' = Nutella'
1054	3 Substituting A1 to $\Omega^2$ :
1055	'Q2': 'What is the main ingredient in Nutella?',
1056	4. Assessed OD
1057	4. Answers Q2: The answer 'A2' to Q2': 'What is the main ingredient in Nutella?'.
1058	A2' = A' = Hazelnuts
1059	Example #3
1060	Example #5
1061	(Q') = (Who) was the Roman Emperor when Jesus was born?'
1062	1. Splitting 'Q' into 'Q1' and 'Q2': ' $(O1)$ ' 'When was Jesus born? '
1063	'Q2': 'Who was the Roman Emperor in [A1]?'
1064	$2 - A $ remains $\alpha = 0.1$ .
1065	The answer 'A1' to 'Q1': 'When was Jesus born?' is 1 BCE. 'A1' = 1 BCE
1067	
1069	3. Substituting A1 to Q2: 'O2' · 'Who was the Roman Emperor in 1 BCE?'
1069	g- · ··································
1070	4. Answers Q2: The answer $(A2)$ to $O2'$ . (When use the Remar Error error in 1 PCE?)
1071	'A2' = 'A' = 'Caesar Augustus'
1072	
1073	You MUST apply this structure when asked to answer a multi-hop question 'Q'. Now answer the multi-hop question 'O' as shown in the examples above
1074	Q: {question}
1075	
1076	
1077	Figure 16: <i>Three-shot CoT prompt</i> used at evaluation time in MHQA.

1080	
1001	
1082	Prompt used to marge O1 and O2 in MHOA
1084	rompt used to intege Qr and Q2 in writeA.
1085	Merge 'Q1' and 'Q2' into a single multi-hop bridge question 'Q'.
1086	Learn nom the following 5 examples. As shown in the following example.
1087	Example #1
1088	$(\Omega^{1})$ "What was the spaceflight that first landed humans on the Moon?"
1089	'Q2': "Who was the commander of Apollo 11?"
1090	
1091	Solution: 1 Answer O1: 'A1' is "Apollo 11"
1092	2. If 'A1' is in 'Q2' print(A1); 'A1' = Apollo 11 is in 'Q2' so I print "Apollo 11"
1093	3. Since you found 'A1' in 'Q2', rewrite 'Q2' so that you delete 'A1' and substitute 'Q1'
1094	Rewriting Q2. Original 'Q2': "Who was the commander of Apollo 11?". Since 'A1' is in
1095	'Q2', I delete it and write 'Q1' there. Rewritten 'Q2': "Who was the commander of the
1096	spaceflight that first landed humans on the Moon?"
1097	The single multi-hop question is therefore the rewritten 'O2'.
1098	$^{\circ}Q^{2} = ^{\circ}Q^{\circ} = "Who was the commander of the spaceflight that first landed humans on the$
1099	Moon?"
1100	Example #2
1101	
1102	'Q1': What is the flagship product of Ferrero?
1103	Solution:
1104	1. Answer Q1; 'A1' is "Nutella"
1105	2. If 'A1' is in 'Q2' print(A1); 'A1' = "Nutella" is in 'Q2' so I print "Nutella" 3. Since you found 'A1' in 'Q2' rewrite 'Q2' so that you dolote 'A1' and substitute 'Q1'
1106	5. Since you found A1 in Q2, rewrite Q2 so that you delete A1 and substitute Q1 there;
1107	Rewriting Q2. Original 'Q2': "What is the main ingredient in Nutella?".
1108	Since 'A1' is in 'Q2', I delete it and write 'Q1' there. Bowritten ' $(O2')$ : "What is the main ingredient in the flagship product of Forroro?"
1109	Rewritten Q2. What is the main ingredient in the hagship product of Perfeto:
1110	The single multi-hop question is therefore the rewritten 'Q2'. 'Q2' = 'Q' = "What is the main ingredient in the flagship product of Ferrero?"
1111	the main ingredient in the nagonip product of reffero.
1112	Example #3
1113	'O1': "When was Jesus born?"
1114	'Q2': "Who was the Roman Emperor in 1 BCE?"
1115	Colution
1116	1. Answer Q1; 'A1' is "1 BCE"
1117	2. If 'A1' is in 'Q2' print(A1); 'A1' = 1 BCE is in 'Q2' so I print "1 BCE"
01118	3. Since you found 'A1' in 'Q2', rewrite 'Q2' so that you delete 'A1' and substitute 'Q1' there:
1119	Rewriting Q2. Original 'Q2': "Who was the Roman Emperor in 1 BCE?". Since 'A1' is in
1120	'Q2', I delete it and write 'Q1' there. Rewritten 'Q2': "Who was the Roman Emperor when
1121	Jesus was born?"
1122	The single multi-hop question is therefore the rewritten 'Q2'.
1123	(Q2) = (Q) = "Who was the Roman Emperor when Jesus was born?"
1124	
1125	You MUST apply this structure when asked to merge 'Q1' and 'Q2'.
1120	Now merge $(\hat{Q}\hat{1})$ and $(\hat{Q}\hat{2})$ into a single multi-hop bridge question $(\hat{Q}\hat{2})$ .
1127	$\mathbb{Q}_2^{\prime}$ : {question1}
1128	\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\
1129	
1130	Figure 17: Prompt used to merge O1 and O2 in MHOA.
1131	

## Generating Q1 in MHQA.

Identify one entity in the following text. Come up with a question so that the answer to this question is the entity chosen earlier. The question must be based on the following text. Write your results as 'Question:' and then the question and 'Entity:' and then the entity.

Text: {document\_one}

Figure 18: Prompt used to generate  $Q_1$ .  $Q_1$  is generated such that its answer A1 = E where E is the entity retrieved.

## Generating Q2 in MHQA.

Come up with a question based on the following text that contains the word: {entity}

Text: {document\_two}

Figure 19: Prompt used to generate  $Q_2$ .  $Q_2$  is generated such that its main topicis E where E is the entity retrieved.