

---

# Grounding Code Generation with Input-Output Specifications

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Large language models (LLMs) have demonstrated significant potential in code  
2 generation. However, the code generated by these models occasionally deviates  
3 from the user’s intended outcome, resulting in executable but incorrect code. To  
4 mitigate this issue, we propose GIFT4CODE, a novel approach for the instruction  
5 fine-tuning of LLMs specifically tailored for code generation. Our method leverages  
6 synthetic data produced by the LLM itself and utilizes execution-derived feedback  
7 as a key learning signal. This feedback, in the form of program input-output  
8 specifications, is provided to the LLM to facilitate fine-tuning. We evaluated our  
9 approach on two challenging data science benchmarks, ARCADE and DS-1000.  
10 Our results suggest that the method enhances the LLM’s alignment with user  
11 intentions, reducing the incidence of executable but incorrect outputs.

## 12 1 Introduction

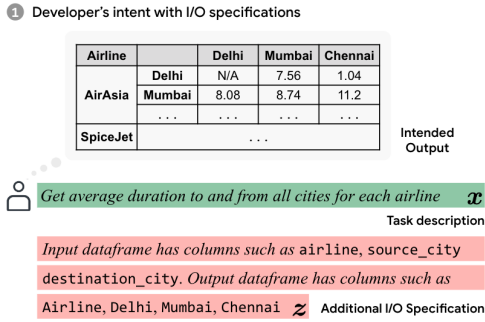
13 Large language models (LLMs) trained on code have demonstrated tremendous success as AI pair  
14 programmers in assisting developers writing code (Chen et al., 2021a; Austin et al., 2021; Li et al.,  
15 2023; Chowdhery et al., 2022; Li et al., 2022; Nijkamp et al., 2022; Fried et al., 2022; Li et al.,  
16 2023). Developers often interact with code LLMs using succinct natural language (NL) intents (*e.g.*  $x$   
17 in Fig. 1) to describe their tasks (Barke et al., 2022; Ross et al., 2023). However, NL intents are  
18 often ambiguous (Yin et al., 2022b). This ambiguity can be problematic in complex tasks, such as  
19 manipulating Pandas DataFrames or PyTorch Tensors (Lai et al., 2022).

20 Auxiliary input-output (I/O) specifications, ranging from concrete I/O examples to high-level NL  
21 summaries (*e.g.* **red text** in Fig. 1), offer a natural way to reduce this ambiguity (Gulwani et al.,  
22 2015; Balog et al., 2016; Jain et al., 2022; Yin et al., 2022a). Prior to the emergence of LLMs,  
23 auxiliary specifications served as essential problem descriptions in program synthesis (Gulwani, 2016;  
24 Devlin et al., 2017; Shi et al., 2020). Real-world synthesis systems like FlashFill are testimony to  
25 the adoption and effectiveness of I/O specifications (Gulwani, 2011; Gulwani et al., 2012). In this  
26 work, we consider the problem of LLM-based code generation when the LLM has access to both a  
27 natural-language intent and an additional I/O specification.

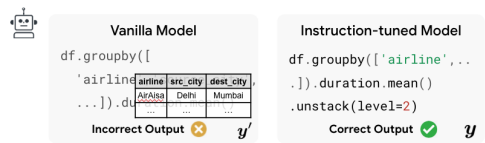
28 However, code LLMs often fall short on following intents with additional complex semantic con-  
29 straints like I/O specifications out-of-the-box, leading to plausible solutions that fail to satisfy the  
30 constraints (*e.g.*  $y'$ , Fig. 1). Such a lack of *alignment* between the user’s intent and the model’s  
31 predictions (Chen et al., 2021a) could pose unnecessary burden on developers who are then required  
32 to fix the generated code (Bird et al., 2023). Therefore, we posit that addressing this misalignment by  
33 *grounding* the code generated by LLMs to the provided specifications is of paramount importance.

34 Instruction fine-tuning has emerged as an effective strategy to tackle the issue of misalignment (Wei  
35 et al., 2021; Sanh et al., 2021; Chung et al., 2022). Classical approaches for instruction tuning typically

## Code generation using Intents with I/O Specifications



## 2 Predictions from different code LLMs



## Instruct-tuning with synthetic Intents and Code

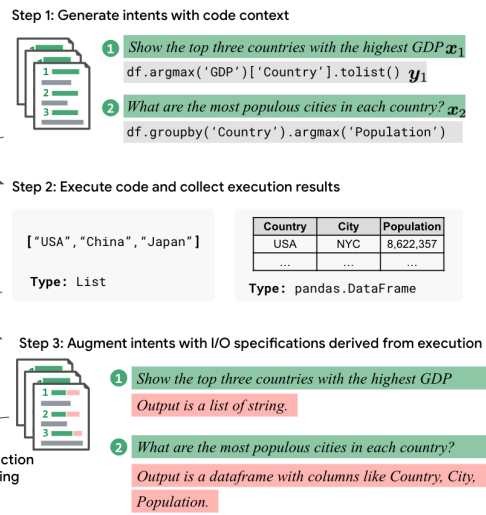


Figure 1: **Left:** Illustration of how developers prompt code LLMs with NL intents and I/O specifications to generate code with complex outputs (pandas Dataframes). Vanilla code LLMs fail to understand extra I/O specifications. **Right:** Our proposed instruction tuning approach uses synthetic intents and code solutions, where intents are augmented with I/O specifications derived from program execution results. Models trained on the synthetic data could better follow a developer’s intent.

36 require a substantial amount of parallel labeled data of NL intents and gold model responses. The  
 37 process of gathering such data is labor-intensive and time-consuming. Recent studies have suggested  
 38 that generating synthetic instruction-following data using the LLM itself is a promising approach to  
 39 improve alignment, with empirical success on natural language text generation tasks (Wang et al.,  
 40 2022a; Honovich et al., 2022a; Taori et al., 2023; Peng et al., 2023, *inter alia*).

41 In this paper we build upon the recent success of instruction tuning using synthetic data and fine-tune  
 42 code LLMs to follow NL intents with additional I/O specifications. Unlike existing approaches,  
 43 our key insight is to leverage program execution for synthetic data generation. First, in contrast to  
 44 other open-ended text generation tasks where assessing the quality of target responses is challenging,  
 45 the quality of synthetic code generation data can be easily improved using heuristics such as code  
 46 executability (Yin et al., 2022c). Moreover, from the program execution states one could derive  
 47 precise and *aligned* I/O specifications that can be included in the intents to supervise a model to  
 48 follow those extra semantic constraints (Fig. 1, *Right*). In other words, when fine-tuned on such  
 49 synthetic data, a model learns to *ground* NL task descriptions to program execution states expressed  
 50 as I/O specifications (Berant et al., 2013).

51 We apply our grounded instruction fine-tuning for code (GIFT4CODE) method to two challenging  
 52 natural language to code generation applications: synthesizing complex pandas programs in com-  
 53 putational notebooks (ARCADE, Yin et al. (2022b)) and answering data science questions on Stack  
 54 Overflow (DS-1000, Lai et al. (2022)). First, we demonstrate the value of leveraging program  
 55 execution information by showing that strong code LLMs can already be significantly improved by  
 56 up to 10% absolute on ARCADE after fine-tuning on intents and executability-filtered code solutions  
 57 *without* including any I/O specifications in synthetic data. Then, to further align model predictions to  
 58 various types of user-provided I/O specifications, we derive those specifications at different levels of  
 59 abstraction from code execution results. This ranges from concrete input/output examples to succinct  
 60 natural language summaries of target variables (specifications in Fig. 1). By fine-tuning on parallel  
 61 data of intents with I/O constraints and their target code solutions, the model is better at following a  
 62 developer’s intents while producing code that is more likely to execute to the desired outcome.

## 63 2 GIFT4CODE: Learning to Follow Intents with I/O Specifications

64 In this section we elaborate on GIFT4CODE, our proposed approach to fine-tune code LLMs to better  
 65 follow developers’ natural language intents along with I/O specifications, using synthetic parallel

66 data. Fig. 1(Right) illustrates an overview of GIFT4CODE. We first synthesize a collection of intents  
67 with code solutions via few-shot prompting (§2.1), and then execute model-predicted code to derive  
68 I/O specifications from execution results (§2.2). Finally, we fine-tune the code LLM to predict code  
69 solutions given intents inlined with I/O specifications (§2.3).

## 70 2.1 Generating Synthetic Intents and Code Solutions

71 **Programmatic Contexts** We initialize a program state given some programmatic context  
72 and generate a series of contextualized NL-to-code problems for that context. As an  
73 example, the synthetic problems in Fig. 1 (Right) could have the contextual code `df =`  
74 `pd.read_csv("world_statistics.csv")`, which initializes the DataFrame variable `df`, sub-  
75 sequently used in the generated synthetic examples. The fact that our problems are contextualized  
76 sets our approach apart from existing instruction-tuning methods for text generation models (Wang  
77 et al., 2022a; Honovich et al., 2022a), where synthetic examples do not depend on any particular  
78 contexts. In our case, we mine those programmatic contexts from real-world code repositories, such  
79 as tabular datasets (e.g., `.csv`) used in data science notebooks on Github (§3).

80 **Creating Initial NL Intents** Given a programmatic context  $c$ , we few-shot prompt an LLM to  
81 create a sequence of natural language intents  $\{x_i\}$  (e.g.  $x_1, x_2$  in Fig. 1(Right)). A problem  $x_i$  that  
82 appears later in the sequence might depend on the earlier ones  $\{x_{<i}\}$  (Nijkamp et al., 2022; Yin et al.,  
83 2022b). To generate NL intents, we use a “generalist” LLM instead of the code LLM that we aim  
84 to improve, since predicting intents conditioned on some context is similar to other text generation  
85 tasks, which could be better handled by a LM trained on general-purpose text data (Zelikman et al.,  
86 2022). The “generalist” LLM is a state-of-the-art general-purpose large language model. It achieves  
87 competitive results with GPT-4 on a variety of NL reasoning tasks.<sup>1</sup> Empirically, we observe that the  
88 problems generated by this LLM encompass a wide range of tasks relevant to the given programmatic  
89 context.

90 **Predicting Code Solutions** After generating an intent  $x$ , we then prompt the code LLM to get a  
91 code solution  $y$  for  $x$  (e.g.  $y_1$  in Fig. 1(Right)). Specifically, a prompt to the LLM is the concate-  
92 nation of the programmatic context  $c$  and the intent  $x$ , with additional few-shot demonstrations of  
93  $\{c', x', y'\}$ . Since many NL intents can be ambiguous and there could exist multiple alternative  
94 solutions (e.g. without additional I/O specifications, the intent in green in Fig. 1(Left) could be an-  
95 swered using tables with different layouts; see more in Yin et al. (2022b)), we therefore draw multiple  
96 candidate code solutions  $\{y\}$  for each intent. Intuitively,  $\{y\}$  could have a variety of alternative  
97 solutions for  $x$ , each leading to different execution results. This equips the model with the capacity  
98 to predict code for the same task but with different user-provided I/O specifications.

## 99 2.2 Code Execution and Inference of I/O Specifications

100 Given the set of synthetic problems  $\{x, y\}$  generated by few-shot prompting, we execute the  
101 code for each problem (step 2, Fig. 1(Right)) and derive I/O specifications from the execution results  
102 as additional semantic constraints to be included in the intents (step 3, Fig. 1(Right)).

103 Specifically, for each candidate solution  $y$  of an intent, we first execute its original programmatic  
104 context  $c$ , followed by executing  $y$ . We trace the execution to collect the set of input and output  
105 variables in  $y$ , denoted as  $\{v\}$ , which are used to derive I/O specifications (details below). Executing  
106 code with arbitrary programmatic contexts collected from the wild is highly non-trivial due to issues  
107 such as library dependency. However, the use of synthetic data alleviates the need for a complex  
108 environment setup.

109 Given the set of input and output variables extracted from execution results, we formulate an I/O  
110 specification, denoted as  $z$ , which serves as additional information to augment a developer’s intent,  
111 thereby providing a more comprehensive problem description. The level of detail and the style  
112 of these I/O specifications can vary based on the complexity of the problem and the developer’s  
113 preferences. In this work, we investigate three distinct types of I/O specifications, each characterized  
114 by its own linguistic style and level of abstraction, as illustrated in Tab. 1.

<sup>1</sup>Details are publicly available in Anonymous (2023b). The model is now publicly available as an API, but was only privately accessible at the time of submission. Anonymized for double-blind review.

| Spec. Type   | Description                                | Example I/O Specification  |
|--------------|--|--|
| TypeDesc     | Variable type name                         | Generate a variable with name <code>df</code> and type <code>pandas.DataFrame</code>   |
| I/O Examples | Concrete I/O examples                      | Output variable <code>df</code> :<br><pre> Bangalore(float) Chennai(float) Delhi(float) Hyderabad(float) Kolkata(float) Hyderabad(float) Kolkata(float) ...</pre> <pre> ----- ----- ----- ----- ----- ----- </pre> <pre>  nan   1.04   8.08   3.62   7.56   7.56   8.32  </pre> <pre>  1.18   nan   11.96   6.80   6.31   8.75  </pre> <pre>  8.46   11.10   nan   9.19   9.52   10.32  </pre> |
| I/O Summary  | LLM-generated NL summaries of I/O examples | Given the user intent and the code, the salient columns (at most given 3) in the input dataframe are <code>airline</code> , <code>source_city</code> , <code>destination_city</code> . The output dataframe has columns (at most given 3) such as <code>Delhi</code> , <code>Mumbai</code> , <code>Chennai</code> .  |

Table 1: Types of I/O specifications proposed in this work at different levels of abstraction. Example specifications are for the intent in Fig. 1(Left). Only the output specifications for I/O Examples are shown for brevity.

115 First, as a simple baseline, we utilize the variable type (TypeDesc, Tab. 1) as the I/O specification.  
 116 Next, we incorporate the concrete values of the input/output variables into the specification, which  
 117 we refer to as I/O Examples. This is reminiscent of classical program synthesis using I/O exam-  
 118 ples (Gulwani et al., 2012; Alur et al., 2013; Balog et al., 2016). However, in our scenario, these  
 119 I/O examples are used in conjunction with natural language (NL) intents to define the problem, in  
 120 line with Jain et al. (2022). Given that the majority of the problems in our synthetic dataset involve  
 121 complex Python objects such as pandas DataFrames, we simplify the I/O specification to include  
 122 only partial variable states (e.g. by excluding some rows and columns in large DataFrames).

123 In our effort to generate a more natural variety of I/O specifications that closely resemble the style of  
 124 specifications in developers’ NL intents, we employ an LLM to summarize the values of input/output  
 125 variables  $\{v\}$  into a succinct natural language description  $z$  (I/O Summary). Intuitively, the NL I/O  
 126 summary includes salient information in the variables that can best clarify the original intent (e.g. the  
 127 subset of columns in a DataFrame that are most relevant to solve a problem, as in Tab. 1, Bottom).

128 Specifically, we few-shot prompt the generalist LLM to generate  $z$ , using information from its  
 129 programmatic context  $c$ , the intent  $x$ , the code solution  $y$ , as well as I/O variables  $\{v\}$ , i.e.  $z \sim$   
 130  $P_{LLM}(\cdot | c, x, y, \{v\})$ . We then update the intent  $x$  by appending  $z$  to it. The few-shot exemplars  
 131 used for prompting cover example I/O summaries for various types of Python objects, such as nested  
 132 container types (e.g. nested dicts), along with more complex objects like pandas DataFrames and  
 133 `pytorch` or `tensorflow` Tensors.

### 134 2.3 Fine-tuning Code LLMs to Follow Intents with I/O Specifications

135 Our approach, GIFT4CODE, aims to fine-tune code LLMs to generate code that adheres closely to  
 136 the desired intents which are supplemented by I/O specifications. In our synthetic training data, each  
 137 example  $\langle c, x, y \rangle$  consists of a programmatic context  $c$ , an intent  $x$  augmented with I/O specifications,  
 138 and the corresponding code solution  $y$ . During fine-tuning, the code LLM learns to generate code that  
 139 not only satisfies the provided intents but also respects the specified I/O constraints, while leveraging  
 140 any relevant information in the programmatic contexts. In other words, we optimize  $P_{LLM}(y | c, x)$ .  
 141 It is worth noting that the code LLM that undergoes this optimization is different from the “generalist”  
 142 LLM employed to generate the NL intents and I/O specification  $z$ .

## 143 3 Experiments

144 The core research question explored in this section is whether GIFT4CODE enhances the LLM’s  
 145 ability to follow developers’ NL intents with *complex* I/O specifications. While common code  
 146 generation benchmarks like HumanEval and MBPP Chen et al. (2021a); Austin et al. (2021) feature  
 147 simple algorithmic tasks (e.g., sorting) utilizing basic Python data types (e.g., lists), thus allowing for  
 148 the use of concrete I/O examples as specifications, they lack the diverse and complex I/O specifications

149 that we aim to explore. For more open-ended tasks such as data science programming, the output  
150 data type is more complex and diverse (e.g., Pandas DataFrames, PyTorch tensors). Hence, we apply  
151 our method to two different data science code generation applications.

152 **ARCADE** (Yin et al., 2022b) is a benchmark of natural language to code generation in interactive  
153 data science notebooks. Each evaluation notebook consists of a series of interrelated NL-to-code  
154 problems in data wrangling (e.g. “*Min-max normalize numeric columns*”) and exploratory data  
155 analysis (e.g. intents in Fig. 1) using the pandas library. ARCADE features succinct NL intents to  
156 reflect the style of ephemeral queries from developers when prompting LLMs for code completion.  
157 More than 50% of the dataset’s problems are under-specified, which means that additional I/O  
158 specifications could provide extra clarification. To construct programmatic contexts for synthetic  
159 training data generation, we scraped 7,500 CSV files that are used in public Jupyter notebooks. Each  
160 context contains a DataFrame import statement, for example, `df = pd.read_csv(·)`, followed by  
161 an NL description of the DataFrame to help the LLM understand its content. We generated 6 intents  
162 for each programmatic context and sampled 5 candidate code solutions for each intent. Roughly 60%  
163 of the code samples were executable. After filtering based on executability and API diversity (§2.1),  
164 we obtained around 20K synthetic training examples.

165 **DS-1000** (Lai et al., 2022) is a benchmark of data science problems sourced from Stack Overflow  
166 (SO). Compared to ARCADE, problems in DS-1000 feature a wider variety of I/O types, such  
167 as `numpy/scipy Arrays` and `pytorch/tensorflow Tensors`, making it particularly appealing to  
168 evaluate our instruction tuning approach aimed at generating code following I/O specifications.  
169 However, in contrast to ARCADE which features succinct NL intents, DS-1000 follows the typical  
170 style of detailed problem descriptions found in SO posts. These elaborate descriptions often include  
171 additional information such as task background and descriptions of unsuccessful attempts, providing  
172 a more complex intent structure, with an average length of 140 words. Given that such elaborate  
173 intents may not reflect the style of developers’ prompts to code LLMs, we do not focus on generating  
174 intents with similar styles. Instead, we held-out 500 problems in DS-1000 and use their annotated  
175 intents as training data, while evaluating on the remaining problems.<sup>2</sup>

### 176 3.1 Setup

177 **Base Code LLM** We use a strong decoder-only code language model with 62B parameters. The  
178 model was first pre-trained on a collection of 1.3T tokens of web documents and github code data,  
179 and was then fine-tuned on a disjoint set of 64B Python code tokens together with 10B tokens from  
180 Python Jupyter notebooks (Anonymous, 2023a).<sup>3</sup>

181 **Learning Methods** We evaluated the performance of both the baseline and instruction-tuned  
182 models across a range of data formats, as shown in Tab. 2. For each I/O specification type, we  
183 augmented the intents and few-shot exemplars with specifications of the corresponding type. Similarly,  
184 at test time, we augmented the intents with the same type of I/O specifications. The baseline models  
185 are tested under both zero-shot and few-shot prompting. For the latter, we manually created exemplars  
186 for all types of specifications. These exemplars were prepended to the prompt when querying the  
187 LLM for code generation during inference.

188 **Simulate Noisy I/O Specifications at Test Time** At testing time, the generation of I/O Summary  
189 underwent a minor modification from the process detailed in §2.2. We remove the concrete input/out-  
190 put variable states  $\{v\}$  to produce noisy I/O summaries, simulating scenarios where users might give  
191 noisy I/O specifications (Devlin et al., 2017). While the “generalist” LLM uses the code solution to  
192 generate noisy I/O summaries, we remark that the code LLM, which we aim to evaluate, does not  
193 have access to the ground truth solution. In other words, the “generalist” LLM acts merely as a “data  
194 labeler” to create I/O summaries in prompts in order to construct the evaluation dataset. It is also a  
195 common practice in program synthesis to derive specifications from ground truth solutions, which  
196 then serve as the sole input to the model during its evaluation (Balog et al., 2016).

---

<sup>2</sup>We only use the annotated intents, while the code solutions and I/O specifications are still predicted by the LLM. We ensure the training and evaluation problems are disjoint and from different SO posts.

<sup>3</sup>Model details are available in Anonymous (2023a), but withheld from this submission for review.

| Methods                           | ARCADE        |              |                |              | DS-1000       |
|-----------------------------------|---------------|--------------|----------------|--------------|---------------|
|                                   | <i>pass@5</i> |              | <i>pass@20</i> |              | <i>pass@1</i> |
|                                   | No Context    | Full Context | No Context     | Full Context |               |
| <b>Zero-shot Prompting</b>        |               |              |                |              |               |
| Code LLM (no spec.)               | 12.45         | 24.67        | 19.85          | 37.47        | 22.62         |
| <b>Few-shot Prompting</b>         |               |              |                |              |               |
| Code LLM (no spec.)               | 15.96         | 30.98        | 26.35          | 42.30        | 23.92         |
| + TypeDesc                        | 16.58         | 29.68        | 29.68          | 42.30        | 25.90         |
| + I/O Examples                    | 19.85         | 32.47        | 30.79          | 43.23        | <b>26.41</b>  |
| + I/O Summary                     | <b>23.75</b>  | <b>37.11</b> | <b>34.50</b>   | <b>46.75</b> | 26.25         |
| <b>Synthetic Data Fine-tuning</b> |               |              |                |              |               |
| Code LLM (no spec.)               | 20.78         | 34.33        | 33.40          | 46.94        | 24.56         |
| + TypeDesc                        | 21.52         | 36.73        | 33.58          | 48.61        | 27.35         |
| + I/O Examples                    | 25.23         | 42.30        | 38.03          | 53.99        | 28.66         |
| + I/O Summary                     | <b>28.01</b>  | <b>43.79</b> | <b>43.04</b>   | <b>55.47</b> | <b>29.34</b>  |
| StarCoder 15B                     | 11.75         | 22.38        | 17.24          | 32.52        | 26.52         |
| WizardCoder 15B                   | 12.45         | 24.04        | 18.58          | 34.30        | 27.35         |

Table 2:  $pass@k$  on ARCADE and DS-1000. For each type of I/O specification in Tab. 1 (e.g. +I/O Summary), intents are augmented with I/O specifications of that type (e.g. intents inline with I/O summary) in fine-tuning data or few-shot exemplars. At test time, input intents use the same type of I/O specifications.

197 **Metrics** We adopted the  $pass@k$  metrics as defined in Chen et al. (2021a); Austin et al. (2021),  
198 which is calculated as the fraction of problems with at least one correct sample given  $k$  samples.  
199 Following Yin et al. (2022a), we drew 50 samples to calculate  $pass@5$  and  $pass@20$  to reduce the  
200 variance in ARCADE. Similar to Lai et al. (2022), we drew 40 samples to calculate  $pass@1$  on  
201 DS-1000. Consistent with the original works’ settings, the sampling temperature was set to 0.8 for  
202 ARCADE and to 0.2 for DS-1000 respectively.

### 203 3.2 Main Results

204 Tab. 2 presents the  $pass@k$  results on ARCADE and DS-1000. We evaluate both few-shot prompting  
205 and fine-tuning with synthetic data. Specifically, for ARCADE we evaluate on two versions of the  
206 dataset. First, we consider the original version where an intent is prefixed by prior notebook cells as  
207 its programmatic context (**Full Context**), as well as a **No Context** ablation to simulate the scenario  
208 where users query a code LLM using an intent without any context. This no-context setting is more  
209 challenging, where the zero-shot performance of the base code LLM is nearly halved. The standard  
210 errors in all cells of the table are less than 0.5%, and are excluded for clarity in presentation.

211 In our few-shot prompting experiments, we observe that  $pass@k$  generally improves with more  
212 detailed I/O specifications. Interestingly, on ARCADE, the improvements from prompting using  
213 I/O specifications compared to the baseline where no I/O specifications were used (**no spec**), are  
214 more notable in the more challenging no-context scenario (e.g.  $15.96 \mapsto 23.75$  v.s.  $30.98 \mapsto 37.11$   
215 for +I/O Examples). This trend suggests that additional specifications could provide more valuable  
216 clarifications when adequate programmatic contexts are lacking.

217 Next, we fine-tune the base code LLM using our synthetic parallel data using different types of I/O  
218 specifications. Interestingly, without using any I/O specifications in the synthetic intents, on ARCADE  
219 the model already registers significant improvements compared to the zero- and few-shot settings.  
220 The model-predicted code solutions are filtered using executability heuristics, which helps improve  
221 the quality of the synthetic data, and a model fine-tuned on such data could generally be better at  
222 following users’ intents, even without I/O specifications. Moreover, by fine-tuning the model to  
223 follow intents with additional I/O specifications, we observe significantly better results. We also  
224 remark that instruction fine-tuning using natural language I/O summaries (+I/O Summary) yields  
225 the best results on both datasets. Intuitively, those I/O summaries could encode salient information in  
226 target input and output variables through natural language descriptions, which could make it easier  
227 for the model to capture patterns in the data as compared to other more elaborate versions such as  
228 using concrete I/O examples.

229 We also evaluated Starcoder (Li et al., 2023) and its instruction tuned variant WizardCoder (Luo et al.,  
230 2023) on ARCADE and DS-1000. The result shows that GIFT4CODE is a more effective instruction  
231 tuning method in the data science domain. This is especially observed by the fact that GIFT4CODE  
232 offers much more relative improvement to the base model than the gains WizardCoder boasts over  
233 StarCoder. Overall, our results demonstrate that GIFT4CODE significantly improves the performance  
234 of code LLMs in following intents with I/O specifications at varying level of abstraction.

## 235 4 Related Work

236 **Execution Guided Code Generation** One area of study primarily focuses on utilizing execution  
237 as I/O examples, facilitating the synthesis of programs that align with the intended behavior. Gulwani  
238 (2016) involves synthesizing intended programs in an underlying domain-specific language (DSL)  
239 from example based specifications. This method has been further explored and adapted to different  
240 applications in subsequent studies (Devlin et al., 2017; Chen et al., 2018; Bunel et al., 2018). Another  
241 strand of research (Chen et al., 2021b; Wang et al., 2018; Ellis et al., 2019) leverages intermediate  
242 execution results to guide the search of programs. More recently, there have been attempts to utilize  
243 program execution results to verify and select code samples predicted by LLMs, either during auto-  
244 regressive decoding to prune search space (Zhang et al., 2023), or by few-shot prompting (Chen et al.,  
245 2023) and post-hoc reranking (Shi et al., 2022; Ni et al., 2023).

246 **Instruction Fine-tuning** Instruction fine-tuning is a widely adopted approach to address the mis-  
247 alignment issue in LLM-generated content. LLMs such as FLAN (Wei et al., 2021), which excel  
248 at understanding and executing instructions from prompts, are trained on labeled training data. Re-  
249 inforcement learning with human feedback (RLHF) aims to mitigate the amount of labeling effort  
250 using model-based reward (Ouyang et al., 2022). Other works also confirmed the effectiveness of  
251 using instructional data in the fine-tuning stage (Mishra et al., 2021; Sanh et al., 2021; Chung et al.,  
252 2022; Wang et al., 2022b). To lower labeling cost, several recent works explored the possibility of  
253 automatic instruction generation (Ye et al., 2022; Zhou et al., 2022; Honovich et al., 2022b). In  
254 particular, SELF-INSTRUCT (Wang et al., 2022a) demonstrated that LLMs can be further improved  
255 by utilizing its own generation of instruction data. Our work differs from this line by considering  
256 execution-based specifications. Additionally, recent works attempted to distill instruction following  
257 data from more capable LLMs that have already been instruction-tuned (Honovich et al., 2022a; Taori  
258 et al., 2023; Chiang et al., 2023; Peng et al., 2023). In contrast, GIFT4CODE generates synthetic data  
259 from vanilla LLMs that have not gone through instruction-tuning.

260 **Synthetic Data from LLMs** Besides generating data for instruction following, a number of recent  
261 studies have also harnessed general-purpose LLMs to generate realistic synthetic data in areas where  
262 labeled data limited, such as language understanding and clinical research (Rosenbaum et al., 2022a;  
263 Tang et al., 2023; Borisov et al., 2022; Liu et al., 2022; Rosenbaum et al., 2022b; Josifoski et al.,  
264 2023). To improve the quality of synthetic data extracted from LLMs, such approaches usually apply  
265 a rejection sampling procedure and filter predictions based on domain-specific heuristics such as  
266 logical consistency (Bhagavatula et al., 2022; Yin et al., 2022c). GIFT4CODE is in spirit of this line  
267 in that it leverages program execution feedback to filter code predictions (Xu et al., 2020).

## 268 5 Conclusion

269 We have presented GIFT4CODE, a framework for instruction fine-tuning large language models of  
270 code in which the training is guided by execution based specifications. Empirically, we demonstrated  
271 how our approach enhances the quality of generated code, substantially improving accuracy on two  
272 challenging data science benchmarks, ARCADE and DS-1000.

## 273 References

- 274 Rajeev Alur, Rastislav Bodík, Garvit Juniwal, Milo M. K. Martin, Mukund Raghothaman, Sanjit A.  
275 Seshia, Rishabh Singh, Armando Solar-Lezama, Emina Torlak, and Abhishek Udupa. Syntax-  
276 guided synthesis. *2013 Formal Methods in Computer-Aided Design*, pp. 1–8, 2013.
- 277 Anonymous. Anonymous code language model. 2023a.

- 278 Anonymous. Anonymous language model. 2023b.
- 279 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,  
280 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language  
281 models. *arXiv preprint arXiv:2108.07732*, 2021.
- 282 Matej Balog, Alexander L. Gaunt, Marc Brockschmidt, Sebastian Nowozin, and Daniel Tarlow.  
283 Deepcoder: Learning to write programs. *ArXiv*, abs/1611.01989, 2016.
- 284 Shraddha Barke, Michael B James, and Nadia Polikarpova. Grounded copilot: How programmers  
285 interact with code-generating models. *arXiv preprint arXiv:2206.15000*, 2022.
- 286 J. Berant, A. Chou, R. Frostig, and P. Liang. Semantic parsing on Freebase from question-answer  
287 pairs. In *Empirical Methods in Natural Language Processing (EMNLP)*, 2013.
- 288 Chandra Bhagavatula, Jena D. Hwang, Doug Downey, Ronan Le Bras, Ximing Lu, Keisuke Sakaguchi,  
289 Swabha Swayamdipta, Peter West, and Yejin Choi. I2d2: Inductive knowledge distillation with  
290 neurologic and self-imitation. *ArXiv*, abs/2212.09246, 2022.
- 291 Christian Bird, Denae Ford, Thomas Zimmermann, Nicole Forsgren, Eirini Kalliamvakou, Travis  
292 Lowdermilk, and Idan Gazit. Taking flight with copilot: Early insights and opportunities of  
293 ai-powered pair-programming tools. *Queue*, 20(6):35–57, jan 2023. ISSN 1542-7730. doi:  
294 10.1145/3582083. URL <https://doi.org/10.1145/3582083>.
- 295 Vadim Borisov, Kathrin Sessler, Tobias Leemann, Martin Pawelczyk, and Gjergji Kasneci. Language  
296 models are realistic tabular data generators. *ArXiv*, abs/2210.06280, 2022.
- 297 Rudy Bunel, Matthew J. Hausknecht, Jacob Devlin, Rishabh Singh, and Pushmeet Kohli. Leveraging  
298 grammar and reinforcement learning for neural program synthesis. *ArXiv*, abs/1805.04276, 2018.
- 299 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde, Jared Kaplan, Harrison  
300 Edwards, Yura Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger,  
301 Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick  
302 Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter,  
303 Philippe Tillet, Felipe Petroski Such, David W. Cummings, Matthias Plappert, Fotios Chantzis,  
304 Elizabeth Barnes, Ariel Herbert-Voss, William H. Guss, Alex Nichol, Igor Babuschkin, S. Arun  
305 Balaji, Shantanu Jain, Andrew Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa,  
306 Alec Radford, Matthew M. Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder,  
307 Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating  
308 large language models trained on code. *ArXiv*, abs/2107.03374, 2021a.
- 309 Xinyun Chen, Chang Liu, and Dawn Xiaodong Song. Execution-guided neural program synthesis. In  
310 *International Conference on Learning Representations*, 2018.
- 311 Xinyun Chen, Dawn Xiaodong Song, and Yuandong Tian. Latent execution for neural program  
312 synthesis beyond domain-specific languages. In *Neural Information Processing Systems*, 2021b.
- 313 Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models to  
314 self-debug. *ArXiv*, abs/2304.05128, 2023.
- 315 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng,  
316 Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna:  
317 An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.
- 319 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam  
320 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh,  
321 Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam M.  
322 Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Benton C. Hutchinson, Reiner Pope,  
323 James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm  
324 Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier García, Vedant Misra,  
325 Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret  
326 Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick,



- 327 Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica  
328 Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan  
329 Saeta, Mark Díaz, Orhan Firat, Michele Catasta, Jason Wei, Kathleen S. Meier-Hellstern, Douglas  
330 Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. Palm: Scaling language modeling with pathways.  
331 *ArXiv*, abs/2204.02311, 2022.
- 332 Hyung Won Chung, Le Hou, S. Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang,  
333 Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac  
334 Suzgun, Xinyun Chen, Aakanksha Chowdhery, Dasha Valter, Sharan Narang, Gaurav Mishra,  
335 Adams Wei Yu, Vincent Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed Huai  
336 hsin Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei.  
337 Scaling instruction-finetuned language models. *ArXiv*, abs/2210.11416, 2022.
- 338 Jacob Devlin, Jonathan Uesato, Surya Bhupatiraju, Rishabh Singh, Abdel rahman Mohamed, and  
339 Pushmeet Kohli. Robustfill: Neural program learning under noisy i/o. *ArXiv*, abs/1703.07469,  
340 2017.
- 341 Kevin Ellis, Maxwell Nye, Yewen Pu, Felix Sosa, Joshua B. Tenenbaum, and Armando Solar-Lezama.  
342 Write, execute, assess: Program synthesis with a repl. In *Neural Information Processing Systems*,  
343 2019.
- 344 Daniel Fried, Armen Aghajanyan, Jessy Lin, Sida Wang, Eric Wallace, Freda Shi, Ruiqi Zhong,  
345 Wen-tau Yih, Luke Zettlemoyer, and Mike Lewis. InCoder: A generative model for code infilling  
346 and synthesis. *arXiv preprint arXiv:2204.05999*, 2022.
- 347 Sumit Gulwani. Automating string processing in spreadsheets using input-output examples. In  
348 *ACM-SIGACT Symposium on Principles of Programming Languages*, 2011. URL <https://api.semanticscholar.org/CorpusID:886323>.  
349
- 350 Sumit Gulwani. Programming by examples - and its applications in data wrangling. In *Dependable*  
351 *Software Systems Engineering*, 2016.
- 352 Sumit Gulwani, William R. Harris, and Rishabh Singh. Spreadsheet data manipulation using examples.  
353 *Commun. ACM*, 55:97–105, 2012.
- 354 Sumit Gulwani, José Hernández-Orallo, Emanuel Kitzelmann, Stephen H Muggleton, Ute Schmid,  
355 and Benjamin Zorn. Inductive programming meets the real world. *Communications of the ACM*,  
356 58(11):90–99, 2015.
- 357 Or Honovich, Thomas Scialom, Omer Levy, and Timo Schick. Unnatural instructions: Tuning  
358 language models with (almost) no human labor. *ArXiv*, abs/2212.09689, 2022a.
- 359 Or Honovich, Uri Shaham, Samuel R. Bowman, and Omer Levy. Instruction induction: From few  
360 examples to natural language task descriptions. *ArXiv*, abs/2205.10782, 2022b.
- 361 Naman Jain, Skanda Vaidyanath, Arun Iyer, Nagarajan Natarajan, Suresh Parthasarathy, Sriram  
362 Rajamani, and Rahul Sharma. Jigsaw: Large language models meet program synthesis. In  
363 *Proceedings of the 44th International Conference on Software Engineering*, pp. 1219–1231, 2022.
- 364 Martin Josifoski, Marija Sakota, Maxime Peyrard, and Robert West. Exploiting asymmetry for  
365 synthetic training data generation: Synthie and the case of information extraction. *ArXiv*,  
366 abs/2303.04132, 2023.
- 367 Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Scott Wen  
368 tau Yih, Daniel Fried, Sida Wang, and Tao Yu. Ds-1000: A natural and reliable benchmark for  
369 data science code generation. *ArXiv*, abs/2211.11501, 2022.
- 370 Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou,  
371 Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, Terry Yue  
372 Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João Monteiro,  
373 Oleh Shliakhko, Nicolas Gontier, Nicholas Meade, Armel Zebaze, Ming-Ho Yee, Logesh Kumar  
374 Umaphathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo Wang, Rudra Murthy, Jason  
375 Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey, Zhihan Zhang,

376 Nourhan Fahmy, Urvashi Bhattacharyya, W. Yu, Swayam Singh, Sasha Luccioni, Paulo Villegas,  
377 Maxim Kunakov, Fedor Zhdanov, Manuel Romero, Tony Lee, Nadav Timor, Jennifer Ding, Claire  
378 Schlesinger, Hailey Schoelkopf, Jana Ebert, Tri Dao, Mayank Mishra, Alexander Gu, Jennifer  
379 Robinson, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contractor, Siva Reddy, Daniel  
380 Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, Sean M. Hughes, Thomas  
381 Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder: may the source be with  
382 you! 2023.

383 Yujia Li, David H. Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom,  
384 Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien  
385 de, Masson d’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven  
386 Gowal, Alexey, Cherepanov, James Molloy, Daniel Jaymin Mankowitz, Esme Sutherland Robson,  
387 Pushmeet Kohli, Nando de, Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-level  
388 code generation with alphacode. *Science*, 378:1092 – 1097, 2022.

389 Qi Liu, Zihuiwen Ye, Tao Yu, Phil Blunsom, and Linfeng Song. Augmenting multi-turn text-to-sql  
390 datasets with self-play. In *Conference on Empirical Methods in Natural Language Processing*,  
391 2022.

392 Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma,  
393 Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with evol-  
394 instruct. *ArXiv*, abs/2306.08568, 2023. URL [https://api.semanticscholar.org/CorpusID:  
395 259164815](https://api.semanticscholar.org/CorpusID:259164815).

396 Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. Cross-task generalization  
397 via natural language crowdsourcing instructions. In *Annual Meeting of the Association for  
398 Computational Linguistics*, 2021.

399 Ansong Ni, Srini Iyer, Dragomir R. Radev, Ves Stoyanov, Wen tau Yih, Sida I. Wang, and Xi Victoria  
400 Lin. Lever: Learning to verify language-to-code generation with execution. *ArXiv*, abs/2302.08468,  
401 2023.

402 Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Haiquan Wang, Yingbo Zhou, Silvio Savarese,  
403 and Caiming Xiong. Codegen: An open large language model for code with multi-turn program  
404 synthesis. 2022.

405 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong  
406 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton,  
407 Luke E. Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Francis Christiano, Jan  
408 Leike, and Ryan J. Lowe. Training language models to follow instructions with human feedback.  
409 *ArXiv*, abs/2203.02155, 2022.

410 Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with  
411 gpt-4. *ArXiv*, abs/2304.03277, 2023.

412 Andrew Rosenbaum, Saleh Soltan, Wael Hamza, Amir Saffari, Macro Damonte, and Isabel Groves.  
413 Clasp: Few-shot cross-lingual data augmentation for semantic parsing. In *AAACL*, 2022a.

414 Andrew Rosenbaum, Saleh Soltan, Wael Hamza, Yannick Versley, and Markus Boese. Linguist:  
415 Language model instruction tuning to generate annotated utterances for intent classification and  
416 slot tagging. In *International Conference on Computational Linguistics*, 2022b.

417 Steven I. Ross, Fernando Martinez, Stephanie Houde, Michael J. Muller, and Justin D. Weisz. The  
418 programmer’s assistant: Conversational interaction with a large language model for software  
419 development. *Proceedings of the 28th International Conference on Intelligent User Interfaces*,  
420 2023.

421 Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai,  
422 Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, Manan Dey, M Saiful Bari, Canwen  
423 Xu, Urmish Thakker, Shanya Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal V.  
424 Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica,  
425 Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj,  
426 Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Févry, Jason Alan Fries, Ryan Teehan,

- 427 Stella Rose Biderman, Leo Gao, Tali Bers, Thomas Wolf, and Alexander M. Rush. Multitask  
428 prompted training enables zero-shot task generalization. *ArXiv*, abs/2110.08207, 2021.
- 429 Freda Shi, Daniel Fried, Marjan Ghazvininejad, Luke Zettlemoyer, and Sida I. Wang. Natural  
430 language to code translation with execution. *ArXiv*, abs/2204.11454, 2022.
- 431 Kensen Shi, David Bieber, and Rishabh Singh. Tf-coder: Program synthesis for tensor manipulations.  
432 *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 44:1 – 36, 2020. URL  
433 <https://api.semanticscholar.org/CorpusID:214605958>.
- 434 Ruixiang Tang, Xiaotian Han, Xiaoqian Jiang, and Xia Hu. Does synthetic data generation of llms  
435 help clinical text mining? *ArXiv*, abs/2303.04360, 2023.
- 436 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy  
437 Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model.  
438 [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca), 2023.
- 439 Chenglong Wang, Kedar Tatwawadi, Marc Brockschmidt, Po-Sen Huang, Yi Mao, Oleksandr Polo-  
440 zov, and Rishabh Singh. Robust text-to-sql generation with execution-guided decoding. *arXiv:  
441 Computation and Language*, 2018.
- 442 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and  
443 Hannaneh Hajishirzi. Self-instruct: Aligning language model with self generated instructions.  
444 *ArXiv*, abs/2212.10560, 2022a.
- 445 Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Anjana  
446 Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, Eshaan Pathak,  
447 Giannis Karamanolakis, Haizhi Gary Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby  
448 Kuznia, Krma Doshi, Maitreya Patel, Kuntal Kumar Pal, M. Moradshahi, Mihir Parmar, Mirali  
449 Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang Karia,  
450 Shailaja Keyur Sampat, Savan Doshi, Siddharth Deepak Mishra, Sujan Reddy, Sumanta Patro,  
451 Tanay Dixit, Xudong Shen, Chitta Baral, Yejin Choi, Noah A. Smith, Hanna Hajishirzi, and Daniel  
452 Khashabi. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp  
453 tasks. In *Conference on Empirical Methods in Natural Language Processing*, 2022b.
- 454 Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du,  
455 Andrew M. Dai, and Quoc V. Le. Finetuned language models are zero-shot learners. *ArXiv*,  
456 abs/2109.01652, 2021.
- 457 Silei Xu, Sina J. Semnani, Giovanni Campagna, and Monica S. Lam. Autoqa: From databases to  
458 q&a semantic parsers with only synthetic training data. *ArXiv*, abs/2010.04806, 2020.
- 459 Seonghyeon Ye, Doyoung Kim, Joel Jang, Joongbo Shin, and Minjoon Seo. Guess the instruction!  
460 flipped learning makes language models stronger zero-shot learners. *ArXiv*, abs/2210.02969, 2022.
- 461 Pengcheng Yin, Wen-Ding Li, Kefan Xiao, A. Eashaan Rao, Yeming Wen, Kensen Shi, Joshua  
462 Howland, Paige Bailey, Michele Catasta, Henryk Michalewski, Alex Polozov, and Charles Sutton.  
463 Natural language to code generation in interactive data science notebooks. *ArXiv*, abs/2212.09248,  
464 2022a.
- 465 Pengcheng Yin, Wen-Ding Li, Kefan Xiao, Abhishek Rao, Yeming Wen, Kensen Shi, Joshua Howland,  
466 Paige Bailey, Michele Catasta, Henryk Michalewski, Alex Polozov, and Charles Sutton. Natural  
467 language to code generation in interactive data science notebooks. 2022b.
- 468 Pengcheng Yin, John Frederick Wieting, Avirup Sil, and Graham Neubig. On the ingredients of an  
469 effective zero-shot semantic parser. In *Annual Conference of the Association for Computational  
470 Linguistics (ACL)*, Dublin, Ireland, May 2022c. URL <https://arxiv.org/abs/2110.08381>.
- 471 E. Zelikman, Qian Huang, Gabriel Poesia, Noah D. Goodman, and Nick Haber. Parsel: A (de-  
472 )compositional framework for algorithmic reasoning with language models. 2022.
- 473 Shun Zhang, Zhenfang Chen, Yikang Shen, Mingyu Ding, Joshua B. Tenenbaum, and Chuang Gan.  
474 Planning with large language models for code generation. *ArXiv*, abs/2303.05510, 2023.

475 Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and  
476 Jimmy Ba. Large language models are human-level prompt engineers. *ArXiv*, abs/2211.01910,  
477 2022.