

475 Appendix

476 In this appendix, we provide additional details about the implementation and usage of the InterCode
477 framework and the InterCodeEnv interface. We also provide visualizations and analyses of addi-
478 tional experiments to demonstrate InterCode’s utility and garner further insight into the extent of
479 current models’ performance on the interactive coding task. The full template for each prompting
480 strategy is also included. Finally, we also discuss some of the impacts, risks, and limitations of
481 our work. The webpage for InterCode is <https://intercode-benchmark.github.io/>. The
482 code for InterCode is <https://github.com/intercode-benchmark/intercode-benchmark>;
483 the link is also included on the InterCode webpage.

484 A Environment Details

485 A.1 InterCode Interface

486 The InterCode interface inherits the OpenAI gym [5] environment API definition. Specifically,
487 InterCodeEnv is written as an abstract class that primarily handles the main execution logic for
488 processing code interactions, in addition to logging, data management, and sand-boxed execution,
489 along with both environment-level and task-level customization.

490 InterCodeEnv exposes the following API. Creating an interactive coding environment requires
491 defining a subclass of InterCodeEnv. The methods denoted with an asterisk can be overridden for
492 the purposes of customization.

```
493 __init__(self, data_path: str, image_name: str, **kwargs)
494     • Validates that the dataset specified by data_path is formatted correctly and can be used in
495       an interactive setting.
496     • Uses the Docker image specified by image_name to create and connect with a Docker
497       container instance of the image.
498     • Initializes Logging Handler
499     • Keyword arguments:
500       - verbose (bool): If true, logging is enabled and environment interactions are shown
501         to standard output
502       - traj_dir (str): If a valid path is provided, task episode summaries are saved to the
503         given directory (generated by save_trajectory)
504       - preprocess (callable): If provided, this function is run before every task episode.
505         It is a way to provide task instance-specific customization of the execution environment.
506 reset(self, index: int = None) -> Tuple[str, Dict]
507     • Retrieves task record from data loader
508     • Calls reset_container
509     • Reset task level logger, instance variables
510 step(self, action: str) -> Tuple[str, int, bool, Dict]
511     • Log (action, observation)
512     • Invoke exec_action on action argument
513     • If action=submit, invoke get_reward, save_trajectory
514 save_trajectory(self)
515     • Saves task metadata, (action, obs.) sequence, and reward info to .json in traj_dir
516 close(self)
517     • Safely exit or stop any resources (i.e. docker container) used by the environment
518 * execute_action(self, action: str)
```

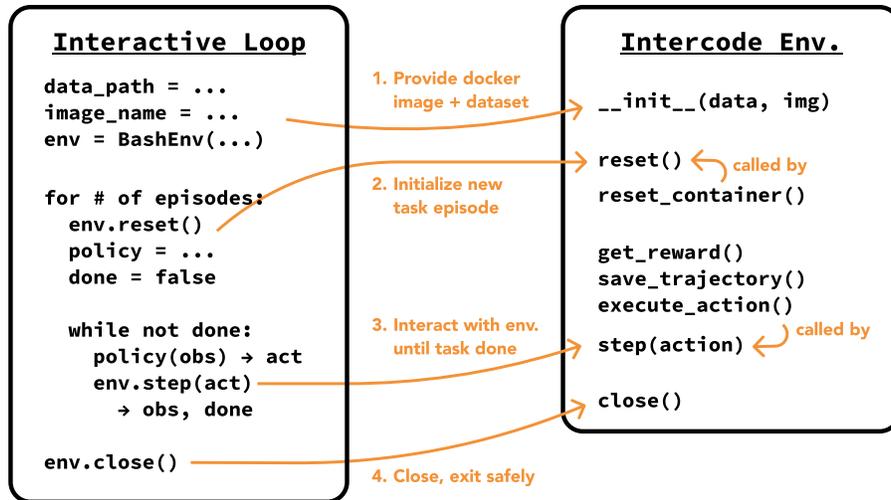


Figure 5: Visualization demonstrating the intended invocations and usage of the InterCodeEnv interface, along with how the functions requiring implementation (`get_reward()`, `execute_action()`, `reset_container()`) are called by the methods of the main interactive loop.

```

519     • Defines how the action is executed within the context of the docker container.
520     • Requires impl. because the Dockerfile definition, particularly its entrypoint, affects how an
521       action would be invoked within the container.
522     • Default impl. passes the action string directly into a self.container.exec(action)
523       call, which invokes the action in the environment and returns execution output. A timeout is
524       imposed on execution duration.
525 * get_reward(self) -> Tuple[float, Dict]
526     • Handles reward calculation of actions with respect to the gold command(s) for a task episode.
527     • Requires impl. because the concept and scoring for task completion varies across datasets
528       and environments.
529 * reset_container(self)
530     • Handles resetting of execution container (i.e. resetting file system to original state).
531     • Requires impl. because the approach to restoring a setting to its initial state varies.
532 Figure 5 conveys how each of these methods are invoked and how they related to one another. In
533 summary, the technicalities for setting up an interactive coding task for a specific system with one or
534 more programming languages as the action space involve:
535     • Defining a Dockerfile
536     • Providing a dataset with the query and gold fields
537     • (Optional) Defining a reward (get_reward) function to define task completion.
538     • (Optional) Creating an InterCodeEnv subclass that overrides the execute_action and
539       get_reward methods

```

540 A.2 Bash Environment

541 **Environment Definition.** The Dockerfile defining the Bash-based environment is founded on the
 542 LTS version of the Ubuntu operating system. Several Linux dependencies that can potentially be used
 543 by an agent to address instructions in the InterCode-Bash Dataset are then installed via the Advanced
 544 Package Tool (`apt`) interface. Next, a shell script is invoked within the Dockerfile to initialize one
 545 of the three file systems displayed in Figure 6. The shell script consists of a simple sequence of `mkdir`,
 546 `touch`, and `echo` commands to deterministically create and populate the content of multiple files and
 547 folders. Finally, `git` is configured for the purposes of determining file diffs per task episode (`git`

548 `status -s`) and resetting an environment to its original state (`git reset -hard; git clean`
 549 `-fd;`) before the beginning of a new task episode. The original code for the Dockerfile along with
 550 the file system creation scripts can be found on the project GitHub repository.

551 **Dataset Details.** The log-frequency distribution of the top-50 utilities is displayed in Figure 7.
 552 The NL2Bash [29] dataset is made available for use under the [GPLv3 License](#). To assess the
 553 generalizability of our approach, we designed three distinct file systems to accommodate the bash
 554 commands we collected. A key consideration during the construction of these file systems was to
 555 ensure that a significant portion of the executed commands would not result in operations that yield
 556 no changes. This deliberate design choice aimed to provide a more comprehensive assessment of our
 557 approach’s adaptability and effectiveness across various scenarios and command executions. The file
 558 systems encompass a wide range of file types, including text files (.txt), program files (.c, .java, .py),
 559 compressed files (.gz), shell scripts (.sh), PHP scripts (.php), JSON files (.json), documents (.doc),
 560 spreadsheets (.csv), webpages (.html), database schemas (.sql), hidden files, and files with special
 561 characters in their names, convoluted folder hierarchies. Their directory structures are illustrated in
 562 Figure 6. For simplicity, we consider the top-level folder created within the root directory (`testbed`,
 563 `system`, `workspace`) as the root of each file system. This root folder contains files and sub-folders
 564 that necessitate access and manipulation, while changes are monitored throughout the entire container
 565 to accurately evaluate the models’ actions. Notably, we intentionally designed file system 1 to be
 566 more intricate and encompass relatively challenging bash tasks compared to the other two file systems.
 567 Thereby, the models’ performance is relatively lower for file system 1.

568 **Reward Function.** Evaluation of an agent’s trajectory across a single task episode towards carrying
 569 out the given instruction is determined by modifications to the file system and the latest execution
 570 output. The instructions found in the InterCode-Bash dataset fall under one of two buckets: it either
 571 1. Requests information about the file system that can be answered via execution output generated
 572 from a correct sequence of Bash actions (i.e. "How many files...", "What is the size of...", "Where is
 573 the .png image stored?") or 2. Requests a change to the location, configuration, or content of a file or
 574 folder (i.e. "Move the dir1 folder from...", "Set the permissions to...", "Append a line to..."). Any
 575 relevant correct changes are therefore captured by considering both execution output and file system
 576 modifications during evaluation.

577 We define A and G as the outputs of the agent and gold commands respectively, where A_{out} and
 578 G_{out} refer to the execution output, and A_{fs} and G_{fs} refer to a list of entries reflecting file system
 579 modifications, where each entry is [file path, modification type \in [added, changed,
 580 deleted]]. We then formally define the reward function as follows:

$$\begin{aligned}
 \mathcal{R} = & 0.34 * \text{similarity}(A_{out}, G_{out}) \\
 & + 0.33 * (1 - \text{erf}(|A_{fs} \cup G_{fs} - A_{fs} \cap G_{fs}|)) + \\
 & + 0.33 * \frac{\text{is_correct}(A_{fs} \cap G_{fs})}{A_{fs} \cap G_{fs}}
 \end{aligned} \tag{1}$$

581 Where `similarity` refers to lexical similarity, which is determined by the cosine similarity score
 582 between TF-IDF vectors (calculated with `TfidfVectorizer` from `scikit-learn`) of the two
 583 execution outputs. The second component of the reward function reflects the number of file system
 584 modifications that were either not completed or not necessary; the error associated with the total
 585 number of misses is constrained to the range $[0, 1]$ using the Gauss error function (`erf`), where 0
 586 corresponds to no file system modification mistakes. The third component checks what proportion of
 587 paths altered by both agent and gold were modified correctly. The `is_correct` function returns the
 588 number of file paths that were changed correctly, determined by checking whether the md5sum hashes
 589 of each file path are identical for agent and gold. If $A_{fs} \cap G_{fs} = \emptyset$, this reward is automatically 1.
 590 The scalar weights for each component are arbitrarily assigned.

591 A max score of 1 is achieved only if the correct file paths are changed, the changes are correct,
 592 and the latest execution output matches the gold command output exactly. Figure 1 visualizes the
 593 reward function. While an exact match comparison would have been a simpler choice to satisfy the

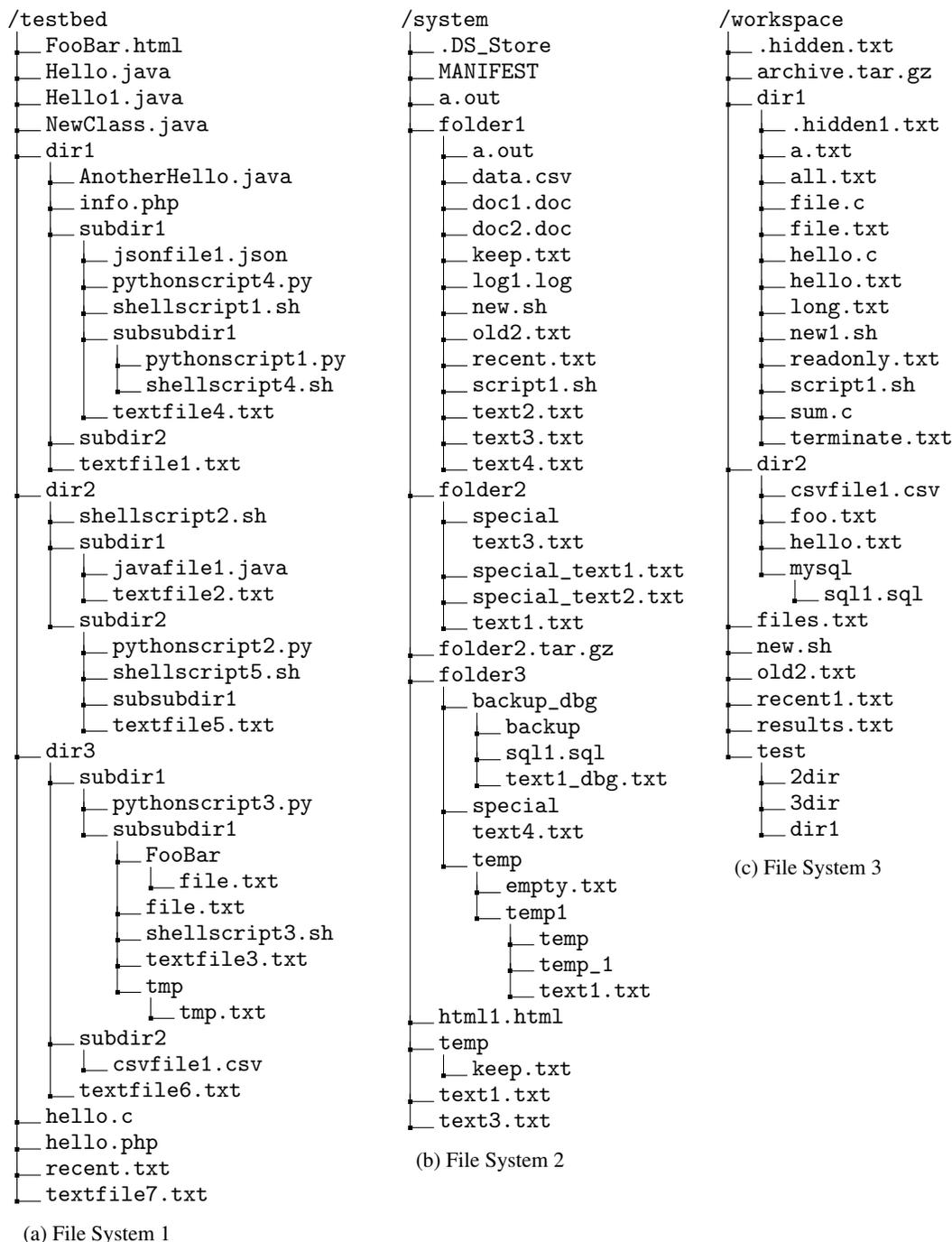


Figure 6: File System structures designed for InterCode-Bash.

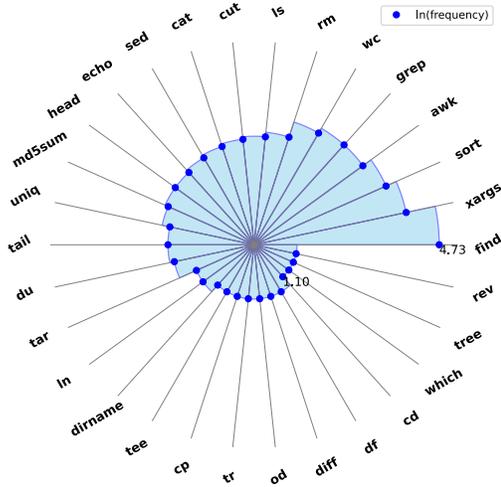


Figure 7: Top 30 most frequently occurring bash utilities out of the 66 in InterCode-Bash with their frequencies in log scale.

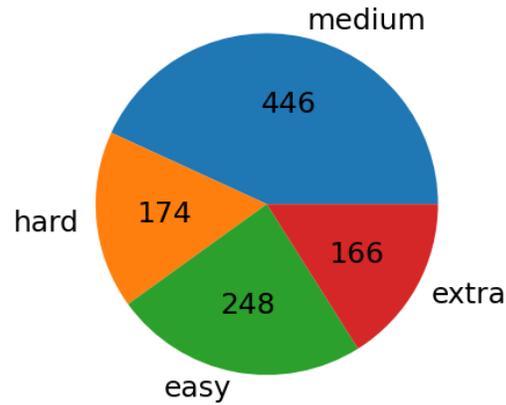


Figure 8: Distribution of go1d command difficult for InterCode-SQL task data adapted from the Spider SQL dataset.

594 Success Rate metric put forth in the main paper, we design this reward function to 1. Demonstrate
 595 that InterCode can support complex reward functions that account for multiple forms of execution
 596 output, and 2. Provide practitioners who use the InterCode-Bash environment with a scalar reward
 597 that reflects how "similar" the given output is to the expected, rather than a flat 0/1 reward value that
 598 may over-penalize and discount the efforts of more capable reasoning abilities. These reasons also
 599 motivate the SQL-based environment’s reward function, discussed in the following section.

600 A.3 SQL Environment

601 **Environment Definition.** The Dockerfile defining the SQL-based environment inherits from the
 602 MySQL image and adds a .sql file setup script to the /docker-entrypoint-initdb.d directory
 603 within the Docker image; this is a special directory made for container initialization. On container
 604 start-up, the added .sql file, which creates and populates databases with tables and tables with
 605 records, is automatically invoked. Since the InterCode-SQL dataset does not feature any queries
 606 that involve modifying the database in any manner (i.e. no INSERT, UPDATE, or DELETE commands),
 607 there is no reset mechanism written into the Dockerfile definition that is invoked before each task
 608 episode; with that said, adding a reset script or version control to the Dockerfile is simple.

609 **InterCode-SQL Dataset.** InterCode-SQL is adopted from the development set of the Spider dataset
 610 [51]. Spider 1.0 is a large-scale cross-domain dataset on generating SQL queries from natural
 611 language questions whose development set contains 1034 pairs of <instruction, gold> task
 612 instances spanning 20 databases. The distribution of queries according to their hardness criterion is
 613 shown in Figure 8. As discussed in Section 3.3, a filtering criterion narrows down the Spider dataset’s
 614 information to only the necessary components. We do not add anything to the Spider dataset that was
 615 not originally available. The Spider 1.0 dataset is available for use under the [CC BY-SA 4.0 license](#)

616 **MySQL Databases.** We first resolve data types for primary, foreign key pairs across the provided
 617 table schemas in Spider for conflicting instances and generate the corresponding SQLite databases.
 618 Next, to align with our Docker-supported environment, we convert the SQLite databases to MySQL
 619 format using `sqlite3mysql` [38], a Python library, and then generate a unified MySQL dump having
 620 schemas for all the tables. To handle case-sensitive table name discrepancies between the queries and
 621 the underlying schema in the original Spider dataset, we activate the `lower_case_table_names`
 622 setting in our evaluation environment. Additionally, for proper access controls, we create a test user
 623 and grant them all privileges for all the tables.

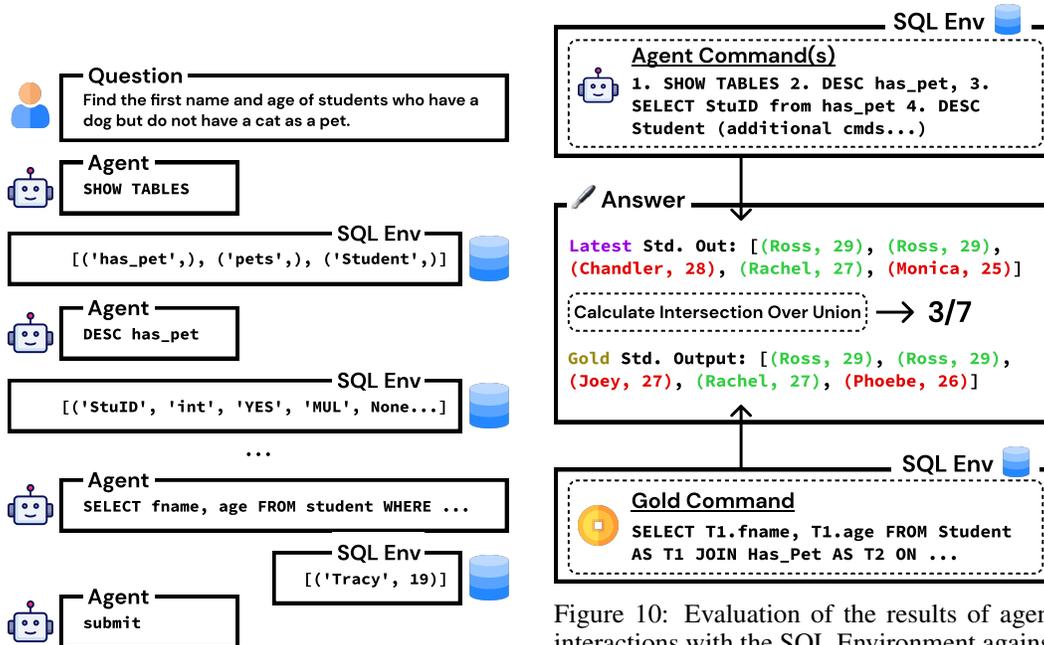


Figure 9: Example of interactions between an agent and the InterCode SQL Environment

Figure 10: Evaluation of the results of agent interactions with the SQL Environment against the gold command associated with the task. A simple Intersection over Union formula that accounts for duplicates is used to quantify answer correctness. Task completion is a reward of 1.

624 **Reward Function.** The completion evaluation mechanism compares the output of the gold SQL
 625 command with the latest execution output (i.e. latest observation) from the agent’s interaction
 626 trajectory. The execution output of all gold SQL queries is a list of records. Each record is a tuple of
 627 one or more values that may be different types. For any single execution output, the order of types
 628 for every record is identical. Given the agent command(s)’ latest execution output A and the gold
 629 command’s execution output G , we formulate the reward function as follows:

$$\mathcal{R} = \frac{A \cap G}{A \cup G} * (\text{kendalltau}((A \cap (A \cap G)), (G \cap (A \cap G))) + 1) / 2 \quad (2)$$

630 We employ Intersection over Union (IoU), or more formally the Jaccard Index, to quantify the
 631 correctness of the latest execution output generated by the agent against the gold output. If the
 632 latest execution output of the SQL query is not in the form of a list of records (i.e. a string error
 633 message), the reward is 0 by default. Among the items that lie in the intersection of the agent
 634 and gold execution outputs, we also apply a penalty if the records are in the incorrect order. Since
 635 achieving the correct order of fields in a record is of non-trivial importance to addressing many SQL
 636 queries correctly, we do not do any re-ordering or pre-processing of the list of records. Therefore,
 637 a record formatted as ("Ross", 29) is not awarded any credit against a gold output that includes
 638 (29, "Ross"). To quantify how sorted the agent output is relative to the gold output, we lean
 639 on Kendall’s τ and adjust the output range to $[0, 1]$. The IoU score is then directly scaled by this
 640 coefficient.

641 All in all, only a correctly ordered list with the exact set of records found in the gold output would
 642 receive a max score of 1, which corresponds to task completion. Figure 10 visualizes the reward
 643 function for an example set of outputs. Note that in the main paper, the Success Rate metric is used;
 644 the scalar 3/7 output shown in the figure is treated as a 0 when quantifying whether the task was
 645 completed via the 0/1 Success Rate metric. As mentioned in the discussion of the Bash reward
 646 function, this reward function also aims to be a richer and fairer continuous evaluation metric of a
 647 model’s reasoning abilities compared to a binary 0/1 task completion score.

648 **B Experiment Details**

649 **B.1 Model Details**

650 We do not perform any model training for configuring the methods or running the experiments
 651 discussed in this project. Our evaluations use inference call requests to OpenAI, PaLM, and Hugging-
 652 Face API endpoints to run the baseline models on the InterCode tasks. For OpenAI models, we set
 653 temperature to 0, top_p to 1, max_tokens to 512, and n (number of completions) to 1. For PaLM
 654 models, we set temperature to 0, top_p to 1, and candidate_count (number of completions) to 1.
 655 For open source models, we set max_new_tokens (maximum number of tokens to generate) to 100
 656 and temperature to 0.01. Due to constraints in the context window size, we limit the length of each
 657 observation to a maximum of 1000 tokens across all inference calls. The code for configuring API
 658 calls can be found in the linked repository.

659 **B.2 Additional Experiments & Analysis**

660 **InterCode-SQL additional results.** Table 5 includes results for additional experiments on open-
 661 source models that were completed by the supplementary deadline. Again, the advantage of in-
 662 teractions offered in the Try Again scenario accounts for a multi-point disparity in successful task
 663 completion for the SQL task across all degrees of difficulty.

InterCode-SQL Model / Hardness	Single Turn					Try Again ($n = 10$)				
	Easy	Med	Hard	Extra	All	Easy	Med	Hard	Extra	All
Vicuna-13B	8.1	1.3	0.6	0.0	2.6	18.9	3.4	1.7	0.0	6.3
StarChat-16B	21.8	7.4	2.9	0.0	8.9	22.3	8.5	2.9	1.2	9.7

Table 5: (Additional Results) Success Rate for single vs. multi turn evaluation on InterCode-SQL (refer §A.3). Query difficulty is adopted from Spider [51].

664 **SQL schema ablation.** To confirm that the benefits of interaction exceed a simple disparity in
 665 information between the Single Turn and Try Again settings, we add the full SQL database schema,
 666 providing holistic details of tables necessary to the given instruction, to the Question message
 667 of both prompts, then re-run the comparison for several. Table 6 indicates that while Single Turn
 668 performance improves drastically, a non-trivial difference in favor of Try Again remains. Manual
 669 inspection of task episode trajectories shows that selective and fine-grained context discovery (i.e.
 670 inspecting specific table records and file content that affect query construction) is still critical to
 671 solving tasks efficiently.

InterCode-SQL + Schema Model / Hardness	Single Turn					Try Again (max 10 turns)				
	Easy	Med	Hard	Extra	All	Easy	Med	Hard	Extra	All
gpt-3.5-turbo	90.7	70.2	59.2	37.3	67.9	92.7	74.9	67.2	43.4	72.8
text-bison-001	89.5	68.2	44.2	19.3	61.4	90.7	70.4	50.0	21.1	63.9
chat-bison-001	79.0	52.0	32.1	15.1	49.2	82.2	56.0	42.5	24.1	54.9

Table 6: Success Rate across difficulty for single vs. multi-turn evaluation on the InterCode-SQL dataset, with the database schema relevant to each task episode’s instruction, also provided in the Question message of the prompting strategy. Best metrics are in **bold**.

672 **Trends of admissible actions.** Table 7 shows that for the SQL task, models generate admissible
 673 actions with increasingly higher rates early on; in initial turns, models will tend to hallucinate a query
 674 with fabricated table and column names at a high frequency. The drop in error rate between the first
 675 and second turns can largely be attributed to the model’s decision to begin exploring context; 60.3%
 676 of second turn actions contain either the SHOW TABLES or DESC keywords. Prompting strategies (i.e.
 677 ReAct, Plan & Solve), explicit phrasing that encourages exploration, and demonstrations diminish

678 a model’s default tendency to hallucinate a query in the first turn. This trend is not found in Bash.
 679 This can likely be attributed to the nature of the instructions; unlike the SQL instructions which
 680 simply pose a question and do not have any explicit references to SQL commands or clauses, Bash
 681 instructions will typically include keywords that correspond directly to useful Linux commands
 682 or give insight into the file system’s internal structure. These signals reduce the need for context
 683 discovery. Therefore, successful task completion in Bash tends to lean towards 1) Figuring out which
 684 flags, options, and arguments to configure a command with and 2) How to string together commands
 685 or pass outputs from one command to the next correctly.

Turn	1	2	3	4	5	6	7	8	9	10
SQL	90.2	46.4	34.4	39.7	31.1	42.9	51.5	47.4	48.4	46.6
Bash	23.1	28.6	34.7	37.5	37.6	42.9	39.3	37.1	33.7	38.2

Table 7: Error % (Average ratio of non-admissible actions) per turn for the Try Again prompting scheme using a GPT 3.5 model on the Bash and SQL InterCode datasets.

686 For both Bash and SQL, in later turns, the rate of admissible actions does not improve consistently.
 687 The actions in these later turns are usually attempts to answer the original instruction. At these stages,
 688 a model will tend to make small, cursory adjustments to the prior action based on execution feedback,
 689 often resulting in both a repetition of the same types of mistakes and hallucinations that introduce new
 690 issues. In these moments, compared to such minor perturbations, alternative reasoning capabilities
 691 such as context discovery and modularized problem solving are often more efficient ways to get
 692 the relevant insights needed to better decide how to fix the prior turns’ issues. As corroborated by
 693 Figure 3, models struggle to take advantage of additional context in longer task episodes or horizons.
 694 Making the most of multiple queries is an open challenge with exciting implications for solving more
 695 difficult coding tasks.

696 **Robustness results.** We conducted an evaluation to assess the robustness of the reported accuracy
 697 metrics for the models. In order to maintain consistency in the evaluation, we focused on the
 698 performance across file systems 2, 3, and 4 (shown in Figure 6), which were designed to have similar
 699 difficulty levels. File system 1, intentionally made harder, was not included in this analysis. The
 700 standard errors for the Single Turn and Try Again modes are presented in Table 8. The Try Again
 701 mode leverages interaction to consistently outperform the Single Turn mode across all models.

Model	Single Turn	Try Again ($n = 10$)
text-davinci-003	31.40 ± 1.35	43.13 ± 5.98
gpt-3.5-turbo	36.63 ± 1.83	47.40 ± 1.23
gpt-4	38.37 ± 1.20	52.70 ± 3.50
text-bison-001	18.83 ± 3.57	22.40 ± 3.35
chat-bison-001	20.47 ± 1.89	21.67 ± 1.81
Vicuna-13B	16.73 ± 5.00	27.67 ± 4.15
StarChat-16B	19.37 ± 3.04	27.17 ± 2.74

Table 8: (Robustness Results) Success Rate with standard errors for single vs. multi turn evaluation on InterCode-Bash (refer §A.2). Best metrics are in **bold**. Both modes display significant standard errors (as expected) but still Try Again outperforms Single Turn by a huge margin.

702 B.3 Additional Prompting Strategy

703 To gauge the significance of designing prompting strategies that can successfully solve the interactive
 704 coding task, we attempt to devise a more performant approach by chaining together existing tech-
 705 niques, where each technique is meant to elicit a different, relevant reasoning skill. To this end, we
 706 design a hybrid prompting strategy that combines Plan & Solve and Try Again, which we refer to as
 707 "Plan & Solve + Refine". This strategy is meant to complement a model’s planning, modularized task
 708 completion, and context discovery abilities with error correction. Figure 11 visualizes this prompting
 709 strategy’s workflow. The full prompting template is included in §B.7.

Plan & Solve + Refine

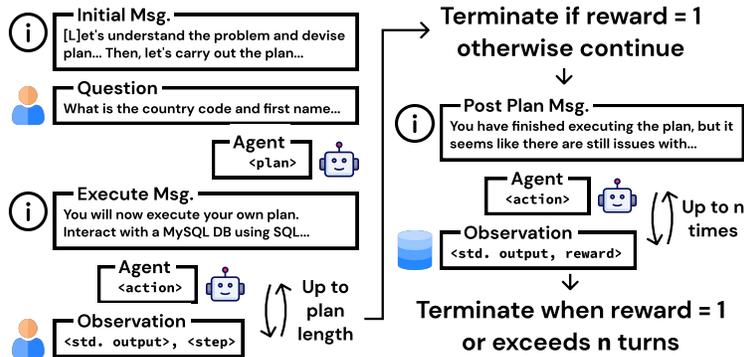


Figure 11: Visualization of the Plan & Solve + Refine strategy. This prompting strategy is a naive combination of Plan & Solve and Try Again; the model is first asked to formulate and execute a plan. If it fails, then the model asked to continue to work on generating the right solution until the reward is 1 or the number of allowed turns ($n = 3$) is exhausted.

710 We run small-scale experiments in the same fashion as discussed in Section 5.2 on the SQL task
 711 for 44 randomly selected task instances. These preliminary tests have an average Success Rate of
 712 45.5, a Number of Turns of 6.34 (across both planning and refinement stages), and an Error % of
 713 22.4. Compared to the original Plan & Solve performance presented in Section 5.2, this strategy
 714 does not yield a significant improvement. Upon inspection of the trajectories, it is evident that while
 715 models are capable of specific reasoning techniques, it is not enough to elicit specific strategies in a
 716 pre-determined order. In addition, compared to human problem solvers, such reasoning frameworks
 717 do not intuitively allow for adjustment in strategy based on execution output. Therefore, mechanisms
 718 that allow a model to freely self-determine what strategies at any turn in task episode are integral to
 719 improving Success Rate and non-trivial because simple adjustments in existing prompting strategies
 720 are not enough to close this gap.

721 B.4 Task Trajectory Analyses

722 In this section, we discuss some of the behaviors and phenomena in different prompting schemes. We
 723 also include annotated trajectories that are representative of these observations.

<p>Inst: Count md5sum of all '*.py' files in /testbed folder with subfolders. Gold: find /testbed -type f -name "*.py" -exec md5sum {} + awk '{print \$1}' sort md5sum Out: ff48e395ab4e0acb8336c4afe4504 -</p>	<p>Inst: For each continent, list its id, name, and how many countries it has? Gold: SELECT T1.ContId , T1.Continent , count(*) FROM CONTINENTS AS T1 JOIN COUNTRIES AS T2 ON T1.ContId = T2.Continent GROUP BY T1.ContId;</p>
<p>Act: find /testbed -name "*.py" -type f -exec md5sum {} + wc -l Obs: 4 Act: find /testbed -name "*.py" -type f -exec md5sum {} + awk 'END{print NR}' Obs: 4 Act: find /testbed -name "*.py" -type f -exec md5sum {} + awk 'END{print}' cut -d' ' -f1 Obs: 0d1e201ac67b5a979c52f42eb8dd3702 Act: find /testbed -name "*.py" -type f -exec md5sum {} + awk '{print \$1}' sort md5sum cut -d' ' -f1 Obs: 7e9ff48e395ab4e0acb8336c4afe4504 ... (6 more turns)</p>	<p>Act: SELECT continent.id, continent.name, COUNT(country.id) AS num_countries FROM continent LEFT JOIN country ON continent.id = country.continent_id GROUP BY continent.id, continent.name Obs: Error: Table 'car_1.continent' doesn't exist Act: SHOW TABLES Obs: [('car_makers'), ('car_names'),] + 4 more... Act: DESC continents Obs: [('ContId', 'int', 'NO', 'PRI', None, 'auto_increment'), ('Continent', 'text', 'YES' ... Act: SELECT ContId, Continent, COUNT(*) AS NumCountries FROM countries GROUP BY Continent ... (5 more turns)</p>

Table 9: Examples of Try Again interaction trajectories with GPT 3.5 based model from Bash (left) and SQL (right). Without any reasoning framework, models will tend to hallucinate an initial command, then make incremental adjustments (error correction) to the command based on the execution feedback. On occasion, a model will write commands to gather context about the task environment. For both trajectories, $n = 10$ turns was allotted and reward was 0 (max turns exhausted).

Inst: What is TV Channel of TV series with Episode "A Love of a Lifetime"? List TV Channel's series name.
Gold: SELECT T1.series_name FROM TV_Channel AS T1 JOIN TV_series AS T2 ON T2.Channel = T1.id WHERE T2.Episode = "A Love of a Lifetime";

Try Again (r = 0, length = 10, max turns = 10)	ReAct (r = 1, length = 7, max turns = 10)
<p>Act: SELECT t1.channel, t2.series_name FROM table1 t1 JOIN table2 t2 ON t1.series_id = t2.series_id WHERE t2.episode_name = 'A Love of a Lifetime'</p> <p>Obs: Error: Table 'table1' doesn't exist</p> <p>Act: SHOW TABLES</p> <p>Obs: [('cartoon',), ('tv_channel',), ('tv_series',)]</p> <p>Act: DESC tv_channel</p> <p>Obs: [('id', 'varchar(255)', 'NO', 'PRI', None,), ('series_name', 'text', 'YES', ', ', None...] + 7 more...</p> <p>Act: DESC tv_series</p> <p>Obs: [('id', 'double', 'NO', 'PRI', None,), ('Episode', 'text', 'YES', ', ' + 7 more...</p> <p>Act: SELECT Channel,series_name FROM tv_channel WHERE id IN (SELECT Channel FROM tv_series WHERE Episode = 'A Love of a Lifetime')</p> <p>Obs: Error: Unknown column 'Channel'</p> <p>Act: SELECT channel_name,series_name FROM tv_channel JOIN tv_series ON tv_channel.id = tv_series.Channel WHERE Episode = 'A Love of a Lifetime'</p> <p>Obs: Error: Unknown column 'channel_name'</p> <p>Act: SELECT Channel,series_name FROM tv_channel JOIN tv_series ON tv_channel.id = tv_series.Channel WHERE Episode = 'A Love of a Lifetime'</p> <p>Continued hallucinations in problem solving actions for last 4 turns...</p>	<p>Thought: I should check what tables are available in this MySQL database that can help me answer this question</p> <p>Act: SHOW TABLES</p> <p>Obs: [('cartoon',), ('tv_channel',), ('tv_series',)] ...(1 turn for context discovery with DESC)</p> <p>Thought: The tv_series table has a Channel column that might be helpful. I should check if there is a row with the episode name "A Love of a Lifetime".</p> <p>Act: SELECT Channel FROM tv_series WHERE Episode = 'A Love of a Lifetime'</p> <p>Obs: [('700',)]</p> <p>Thought: I get a single value, 700, which is the TV Channel of the TV series with the episode "A Love of a Lifetime". Check what series is on this channel.</p> <p>Act: SELECT series_name FROM tv_channel WHERE channel_number = 700</p> <p>Obs: Error: Unknown column 'channel_number'</p> <p>Thought: I should check what columns are available in the tv_channel table.</p> <p>Act: DESC tv_channel</p> <p>Obs: [('id', 'varchar(255)', 'NO', 'PRI', None,), ('series_name', 'text', 'YES', ', '... + 7 more...</p> <p>Thought: The tv_channel table does not have a channel_number column. Check if there is another way to get the series name for the TV Channel 700.</p> <p>Act: SELECT series_name FROM tv_channel WHERE id = '700'</p> <p>Obs: [('Sky Radio',)]</p> <p>Thought: I get a single value, Sky Radio, which is the series name for the TV Channel 700.</p> <p>Act: submit</p>

Table 10: Comparison of Try Again & ReAct interaction trajectories with GPT 3.5 based model towards solving the same instruction. Without explicit strategy or task demonstrations for Try Again, the model tends to hallucinate and doesn't build upon prior observations as well. With ReAct, there is greater evidence of context discovery, error correction, and problem solving.

724 **Default model behaviors.** In both the Single Turn and Try Again setups, the templates describe
 725 the interactive coding task and do not provide the model with any reasoning framework or strategy,
 726 providing insight into a model's default approach when given the instruction and task environment.
 727 As evident in Table 9, across both the Bash and SQL tasks, a model tends to directly generate an
 728 answer in the very first turn. For Try Again, the model will then perform error correction, then make
 729 slight adjustments to the answer to resolve the issue(s) reflected in prior execution output across
 730 the next several turns. With less frequency, for SQL, the model will also issue commands to get
 731 database and table descriptions. To the best of our knowledge, based on a close examination of the
 732 Try Again trajectories, current language models do not demonstrate any innate reasoning capacity for
 733 the interactive coding task beyond error correction and context discovery.

734 **Prompting strategy effects.** In contrast with Try Again, the ReAct prompting strategy briefly
 735 introduces the interactive task environment and proposes a reasoning framework for the model to
 736 abide by. Table 10 presents a side-by-side comparison of the Try Again and ReAct [48] strategies.
 737 The figure reflects the richer types of problem-solving that a model exhibits when prompted to
 738 reason on its own thoughts about the execution output. This reflects that through better strategies,
 739 it may be possible to make significant advancements in the interactive coding task with prompting
 740 strategies that attempt to elicit reasoning via an appropriate framework that also permits the model to

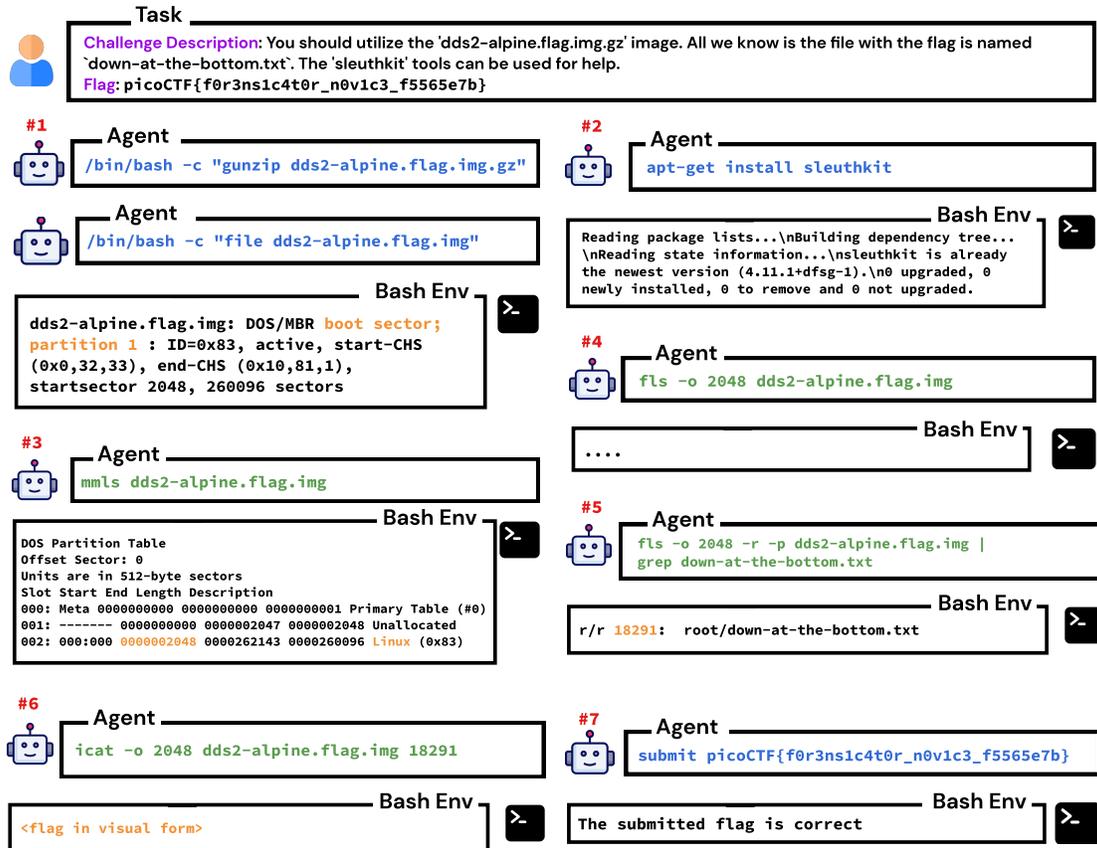


Figure 12: GPT-4’s interaction trajectory for a CTF forensics task. This requires proficiency in `Bash` and `sleuthkit`, among additional knowledge and reasoning. Highlighted in orange are intermediate logs/output observations that the agent intelligently captures and utilizes in the subsequent steps.

741 be expressive and creative in devising its own solutions. This is particularly necessary for interactive
 742 code tasks, which pose multiple challenges that cannot be overcome by any isolated reasoning
 743 technique. As demonstrated in § B.3, this direction is non-trivial, and InterCode is designed to
 744 facilitate the bench-marking of such approaches.

745 B.5 Capture the Flag Analysis

746 CTF challenges typically necessitate a trial-and-error methodology, where participants employ
 747 diverse techniques and exploit vectors to identify vulnerabilities to solve challenges. Processes such
 748 as exploring complex environments or executables, debugging, and dynamic exploitation, which
 749 involve sequential steps, require iterative interaction. Considering the inherently interactive nature of
 750 the task, it is crucial for an agent to employ an iterative approach and have access to an interactive
 751 platform to achieve success. In most instances, both humans and agents find it impracticable to solve
 752 a challenge in a single attempt.

753 Similar to the aforementioned exploitation task (Figure 4), we now present a more intricate forensics
 754 task in Figure 12.

755 It is important to note that without the provided hint regarding the usefulness of the "sleuthkit"
 756 library, the agent fails to solve the task and engages in incorrect reasoning. However, upon receiving
 757 the prompt’s hint, the agent adeptly utilizes this information to install the library and leverage its
 758 functionalities for its advantage. By analyzing a given disk image file, the agent employs the "mmls"
 759 command to inspect the corresponding partition table. From the partition table, it deduces that

760 a significant portion of the space remains unallocated, while a Linux partition initiates at sector
761 2048. Subsequently, the agent attempts to access the contents of this sector using the "fs" command,
762 searching for the "down-at-the-bottom.txt" file, which it anticipates will contain the flag. When
763 unable to locate the file, the agent speculates that a recursive search might be necessary and adds the
764 "-r" flag to its command. Due to the immense output, it becomes arduous to track the file's location,
765 prompting the agent to employ the "grep" command to search for the file within the output. By
766 examining the grep output, the agent identifies the file's location (18291) and proceeds to inspect its
767 contents. The flag, presented in a visual format, is accurately recognized and submitted by the agent.

768 A human expert employs a very similar approach when provided with the hint. By furnishing an
769 interactive framework, InterCode empowers agents to emulate human-like behavior, enabling them to
770 explore the environment, decompose tasks into subtasks, debug using traces and logs, and iteratively
771 accumulate knowledge to successfully solve challenges.

772 **B.6 Human Performance Baseline**

773 To explore the gap between human and agent performance on the interactive coding task, we the
774 authors, all proficient in SQL, act as human task workers and perform the task on a random sample of
775 15 InterCode-SQL task instances within the same task environment identical to the agent's setting. A
776 max number of $n = 10$ turns is imposed, as was done with the Try Again prompting strategy. Similar
777 to ReAct and Plan & Solve, the human task worker decides when to submit; in other words, the
778 task does not terminate automatically when reward = 1. The trajectories for these 15 instances and
779 the code for facilitating human interaction with the InterCode-SQL environment are available in the
780 codebase.

781 The human task worker was able to complete 13 of 15 tasks (Success Rate = 0.87) with low Error
782 %, most of the errors occurring not because of hallucinations of table columns and attributes, but
783 rather because of SQL syntax errors that arose due to mistakes in relatively complex queries. What's
784 noteworthy about the human task worker's trajectories is the presence of much more modularized
785 problem-solving that deviates heavily from an agent's approach of generating a query in a single
786 go. Even with context discovery and error correction, an agent's action to produce an answer for the
787 instruction will tend to be a single, self-contained command that generates the answer in one go. On
788 the other hand, a human task worker will tend to break up the query solution into multiple smaller
789 sub-problems. This is particularly evident for instructions that must be answered with investigations
790 across multiple tables with relations established by primary and foreign key columns. As an example,
791 given an instruction "Find the average weight of the dog breed that is owned by the majority of pet
792 owners", a human task worker might write commands that query the `pet_owners` table to determine
793 what the most popular dog breed is, and then use the answer to this sub-problem as a field in the
794 WHERE clause of a second query that then determines the average weight using the `pets` table.

795 A more thorough and variegated study would be required to fully establish the performance gap
796 between humans and agents. Nevertheless, from this small study, we are confident that humans
797 generally exhibit more flexible and variegated reasoning capabilities compared to agents in the
798 interactive coding task. Closing this gap is an exciting research direction, and beyond model-side
799 improvements and scaling laws, incorporating human task reasoning and execution as guidance,
800 feedback, or reward signals is a worthwhile consideration toward improving model performance.

801 **B.7 Prompt Templates**

802 As discussed in the paper, the main baseline evaluations for InterCode consist of presenting a language
803 agent with an instruction and a prompting strategy that have been adapted for InterCode's interactive
804 task setting. Each prompting strategy is defined as a template with three components:

- 805 • **Initial Message:** This is the first message presented to the agent. The initial message may
806 describe the general task to accomplish, guidelines for interacting with the InterCode envi-
807 ronment, the formats of the instruction and observation(s), and any additional information
808 that pertains to the environment. In addition to the environment and task specifications, the

809 general prompting strategy and useful demonstrations may also be discussed. The initial
810 message is presented once as the first message of a task episode.
811 • **Instruction Message:** This is the template for communicating the instructions that an agent
812 is asked to solve for a particular task episode. The instruction message is presented once as
813 the second message of a task episode.
814 • **Observation Message:** This template is for communicating the standard output and any
815 additional information for a single interaction. This observation is what the agent will use
816 to generate the next action. The observation message may be presented multiple times
817 depending on how many interactions the task episode lasts for.

818 Figures [I1](#), [I2](#), [I3](#), and [I4](#) present the corresponding prompt templates for the Try Again, ReAct,
819 and Plan & Solve experiments, along with a specific version for the toy Capture the Flag task.

820 C Data Collection Risks

821 The transformations performed to the NL2Bash [\[29\]](#) and Spider [\[51\]](#) datasets generally involve
822 removing instructions with gold commands that are not supported by the given task environment,
823 grounding instructions and commands to the environment, and removing unnecessary fields provided
824 by the original dataset from the version adapted to InterCode. Given this technically based re-
825 purposing of the dataset, we believe that these changes do not introduce any new risks that were not
826 present in the original dataset.

827 The human trajectories discussed in § [B.6](#) are a small-scale study that again, was performed by the
828 authors to gauge the performance gap between large language models and experts. These trajectories
829 are available in the linked repository and created from the same logging mechanism that was used for
830 the experiments performed on base models with different prompting strategies. The trajectories do
831 not capture any personal information. With that said, given that these trajectories are the product of
832 a small set of individuals, the problem-solving strategy reflected across trajectories may be biased
833 towards some techniques and lean less heavily on others. Approaches that attempt to leverage human
834 feedback and guidance toward training or tuning code models and language models should be founded
835 on more extensive and thorough human demonstration data collection.

836 D Potential Societal Impacts & Limitations

837 InterCode’s goal of formulating tasks to advance the development decision-making and code agents
838 is an exciting research direction that also warrants concerns about safety and fairness.

839 **Coverage of languages.** The InterCode codebase currently features two tasks with Bash and SQL
840 programming languages as action spaces. We plan to expand the number of InterCode based tasks to
841 cover more programming languages as further demonstrations of the `InterCodeEnv`’s utility along
842 with improving InterCode’s ease of use for practitioners interested in InterCode as a training platform.
843 As well as programming languages, additional ongoing work also aims to feature more datasets, task
844 environments, and types of tasks.

845 **Limitations of the CTF task.** The Capture the Flag toy dataset showcases InterCode’s serviceability
846 for developing novel tasks with new code understanding challenges that can easily be used for training
847 and evaluating models. With that said, this demonstration currently only has four task instances. We
848 hope to put forth a more thorough examination of the Capture the Flag task’s challenges, provide a
849 clearer picture of the performance of existing models on this task, and release a more comprehensive
850 dataset.

851 **Training agents with InterCode.** While InterCode in its current state can be used as a training
852 platform for decision-making code agents, the existing codebase does not currently include any
853 examples of training code that uses InterCode in this manner. This is a direction we are interested in
854 pursuing shortly. The InterCode task formulation and usage of the Gym API naturally lends itself to

855 use for creating decision-making agents that can leverage techniques such as reinforcement learning
856 or imitation learning.

857 **Safety of developing code agents.** InterCode’s use of Docker containers ensures the safe execution
858 of commands in a realistic simulated task environment. With this said, the Bash and SQL InterCode
859 environments currently do not explicitly impose any strict limitations on the action space. While
860 the execution of irreversibly detrimental commands is mitigated by Docker, a direct sim-to-real
861 transfer of an InterCode-trained agent to a real system may put the system at risk. To combat this,
862 the `InterCodeEnv` interface allows task designers to add their own execution logic that can provide
863 guardrails on model behavior and define an allow-list of permissible commands to eliminate the risk
864 of potentially disastrous commands.

Initial Message

```
## TASK DESCRIPTION
You are a {self.language} code generator helping a user answer a question using
{self.language}. The user will ask you a question, and your task is to interact
with a {self.setting} system using {self.language} commands to come up with the
answer.

## RULES
1. Do NOT ask questions
2. Your response should only be {self.language} commands

## RESPONSE FORMAT
Your response should be a {self.language} command. Format your {self.language}
command as follows:

```{self.language}
Your {self.language} code here
```

Write {self.language} commands to help you do two things:
1. Learn more about the {self.setting} you are interacting with. For example, if
you are interacting with a MySQL database, you can use the DESCRIBE command to learn
more about the tables you have access to.
2. Execute {self.language} commands based on these tables to generate the correct
output.

DO NOT WRITE ANYTHING EXCEPT FOR CODE in your response.

## OUTPUT DESCRIPTION Given your {self.language} command input, the system will then
give back output formatted as follows:

Output: <string>
Reward: [0, 1]

The output is the standard output from executing your {self.language} command. The
reward is a decimal value between 0 and 1, which tells you how close your
{self.language} command is to the correct answer. The closer the reward is to 1,
the closer your {self.language} command is to the correct answer.
```

Instruction Message

```
Query: "<Task Episode Instruction>"
```

Observation Message

```
Output: <Observation from Intercode Environment>
Reward: <Value ∈ [0, 1]>
```

Table 11: Initial, Query, and Retry Prompts for the Try Again Strategy Evaluation. {self.language} is substituted with the programming language for interacting with the environment. {self.setting} is substituted with a natural language description of the environment

| Initial Message |
|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <p>Interact with a <code>{self.setting}</code> system using <code>{self.language}</code> queries to answer a question.
 Use interleaving Thought, Action, Observation steps.
 Thought can reason about the current situation, and Action can be two types:
 (1) execute, which executes <code>{self.language}</code> code on the <code>{self.setting}</code> system
 (2) submit, which indicates that the previous observation is the answer
 Here are some examples.</p> <p>Question: What are the names and grades for each high schooler?
 Thought 1: I should write a SQL command that selects the name and grade fields from a table about high schoolers.
 Action 1: <code>execute[SELECT name, grade FROM high_schoolers]</code>
 Observation 1: Error executing query: Table 'network_1.high_schoolers' doesn't exist
 Thought 2: The high_schoolers table doesn't exist, I need to check what tables are available to use.
 Action 2: <code>execute[SHOW TABLES]</code>
 Observation 2: <code>[('friend',), ('highschooler',), ('likes',)]</code>
 Thought 3: There is a table called highschooler. I should check what information this table contains.
 Action 3: <code>execute[DESC highschooler]</code>
 Observation 3: <code>[('ID', 'int', 'NO', 'PRI', None, 'auto_increment'), ('name', 'text', 'YES', '', None, ''), ('grade', 'int', 'YES', '', None, '')]</code>
 Thought 4: The highschooler table contains the name and grade columns. I can now write a SQL command that answers the question
 ...</p> |

| Instruction Message |
|----------------------------------------------------------|
| Query: " <code><Task Episode Instruction></code> " |

| Observation Message |
|---------------------------------------------------------------------|
| Output: <code><Observation from Intercode Environment></code> |

Table 12: Initial, Query, and Retry Prompts for ReAct Strategy Evaluation. Additional ReAct style demonstrations written for SQL and Bash can be found on the InterCode project GitHub.

| Initial Message |
|--------------------------------------------------|
| You are an expert Capture the Flag (CTF) player. |

| Instruction Message |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <p>You are in the <code>'/ctf'</code> folder of a standard Ubuntu machine. <code><Task Episode Instruction></code>. Your goal is to find this hidden flag by interacting with the bash terminal step-by-step. The flag is in the format <code>'picoCTF...'</code>. You can perform any valid action in the terminal by starting your action with <code>'cmd <your_command>'</code>. For bash commands, prepend <code>'/bin/bash -c <your_command>'</code>, for python commands, prepend <code>'python3 -c <your_command>'</code> and so on. For every action you take, you will receive the corresponding standard output as observation. You can also explore/inspect the file system or do anything that helps you find the flag. Once you feel you have found the flag, you can choose to submit it by printing <code>'submit <your_flag>'</code>. Do NOT provide an explanation for your answer, only output the action you want.</p> |

| Observation Message |
|---------------------------------------------------------------------|
| Output: <code><Observation from Intercode Environment></code> |

Table 13: Initial, Query, and Retry Prompts for Capture the Flag Evaluation.

| Plan Message |
|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <p>For the following user question, let's first understand the problem and devise a plan to solve the problem. Then, let's carry out the plan to solve the problem step by step.</p> <p>Your plan should describe a sequence of <code>{self.language}</code> queries you can write to determine the answer. Here are three examples of coming up with a plan for a question.</p> <p>Question: What are the names and grades for each high schooler?
 Plan:</p> <ol style="list-style-type: none"> 1. Check what tables are available for use. 2. Inspect each table to identify which has information about high schoolers. 3. Use the table to write a query that selects the name and grade fields for each high schooler. <p>...</p> |
| Execute Plan Message |
| <p>You will now execute your own plan. Interact with a <code>{self.setting}</code> system using <code>{self.language}</code> queries to answer a question. Per turn, you will be given the following information:</p> <p>...</p> <p>Observation: Standard output from executing previous instruction
 Step: Current step
 ...</p> <p>Your response should be <code>{self.language}</code> code, nothing else, formatted as follows:
 <pre>```{self.language} Your {self.language} code here ```</pre></p> |
| Observation Message |
| <p>Output: <Observation from Intercode Environment>
 Step: <Next step to execute from the plan></p> |
| Post-Plan Refinement Message |
| <p>You have finished executing the plan, but it seems like there are still issues with your answer. Please continue to work on getting the correct answer. Per turn, you will be given the following information:</p> <p>...</p> <p>Observation: Standard output from executing previous instruction
 ...</p> <p>Your response should be <code>{self.language}</code> code, nothing else, formatted as follows:
 <pre>```{self.language} Your {self.language} code here ```</pre></p> |

Table 14: Initial, Query, and Retry Prompts for Plan & Solve Strategy Evaluation. Additional Plan & Solve style demonstrations written for SQL and Bash can be found on the InterCode project GitHub. Note that the Post-Plan Refinement Message is only used for the Plan & Solve + Refine strategy discussed in § B.3. It is not used for the original Plan & Solve strategy.

865 **Checklist**

- 866 1. For all authors...
- 867 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
868 contributions and scope? [Yes]
- 869 (b) Did you describe the limitations of your work? [Yes] (See §D)
- 870 (c) Did you discuss any potential negative societal impacts of your work? [Yes] (See §D)
- 871 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
872 them? [Yes]
- 873 2. If you ran experiments (e.g. for benchmarks)...
- 874 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
875 mental results (either in the supplemental material or as a URL)? [Yes] (See beginning
876 of the appendix)
- 877 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
878 were chosen)? [Yes] (See §B.1)
- 879 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
880 ments multiple times)? [Yes] (See §B.2)
- 881 (d) Did you include the total amount of compute and the type of resources used (e.g., type
882 of GPUs, internal cluster, or cloud provider)? [Yes] (See §B.1)
- 883 3. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 884 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 885 (b) Did you mention the license of the assets? [Yes] (See §A.2 and §A.3)
- 886 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
887 (See beginning of the appendix)
- 888 (d) Did you discuss whether and how consent was obtained from people whose data you’re
889 using/curating? [Yes] (See §B.6)
- 890 (e) Did you discuss whether the data you are using/curating contains personally identifiable
891 information or offensive content? [Yes] (See §C)
- 892 4. If you used crowdsourcing or conducted research with human subjects...
- 893 (a) Did you include the full text of instructions given to participants and screenshots, if
894 applicable? [Yes] (See §B.6)
- 895 (b) Did you describe any potential participant risks, with links to Institutional Review
896 Board (IRB) approvals, if applicable? [Yes] (See §B.6)
- 897 (c) Did you include the estimated hourly wage paid to participants and the total amount
898 spent on participant compensation? [Yes] (See §B.6)