



# **Snorkel AI**

## **Terminus - Expert Contributor Onboarding**

---

November 2025

# Project Overview

---

**Terminus** is Snorkel's effort to build a high-quality dataset in the style of [Terminal-Bench](#).

- Terminal-Bench is a benchmark for evaluating how well AI agents can accomplish complex tasks in a terminal environment
- Features multi-step tasks to be completed via a command line interface (CLI)
- Examples include:
  - Compiling and packaging a code repository
  - Downloading a dataset and training a classifier on it
  - Setting up a server

Your role as an EC is to create these tasks, along with an Oracle solution and associated tests that verify its correctness

- These tasks should be **challenging**, targeting a pass rate of <80% from SOTA models (e.g. GPT5 Codex, Claude Code)
- Tasks will be created locally and submitted to a shared GitHub repository where they will undergo CI validation and peer review

# Task Difficulty Ratings

---

The following ratings will be applied to tasks to categorize their difficulty. Ratings are based on the accuracy of GPT-5 and Sonnet 4.5 when attempting to complete the task.

| Easy            | Medium          | Hard           |
|-----------------|-----------------|----------------|
| 60-80% Accuracy | 40-60% Accuracy | < 40% Accuracy |

**NOTE:** Since these ratings are based off model performance, you will not know how a task will be classified until after submitting.

If your task results in a “trivial” rating, you should iterate on the task to increase difficulty to be at least “easy”.

# Task Components - Part 1 (Files)

---

- **Instruction/Task Description (task.yaml)**
  - Clear and self-contained description of the task to be accomplished
  - Includes references to relevant resources necessary for task completion, rules and constraints, and description of success criteria
  - Also contains metadata not available to the agent at runtime
- **Docker Environment**
  - **Dockerfile**
    - Base image that fully sets up environment, including all required tools, resources, and dependencies
    - Should run without privileged mode
  - **Docker-compose.yaml**
    - Config file that defines orchestration for the task
    - Should reference base image from Dockerfile and include any necessary environment variables or mounted resources
- **Oracle Solution (solution.sh)**
  - Step-by-step solution contained within a shell script
  - Reliably and accurately completes the task

# Task Components - Part 2 (Files cont.)

---

- **run-tests.sh**
  - Script to execute and validate end-to-end evaluation of task
  - Coordinates running the agent within the defined environment, invokes test cases or success criteria checks, produces structured output that indicate task completion status
  - Must be deterministic and callable directly from project root
- **Python tests**
  - Series of deterministic Python scripts containing unit tests
  - Check task completion based on the final state of the environment
- **[Optional] Supporting DataFiles**
  - Any additional inputs (e.g. data, config files, etc.) that are required for the task
- **[Optional] Custom Tools**
  - Custom tools that can be used by the agent
  - Can be delivered in multiple ways, from source code files to companion containers

# Task Components - Part 3 (Metadata)

---

- **Pass Rate Difficulty (Model Performance)**
  - Easy/medium/hard rating
  - **Not provided by the EC**, determined by model performance post-submission
- **Difficulty (Time Estimate)**
  - Estimated time that completing the task would take for both a junior and a senior software engineer
- **Task Type**
  - Label for the task according to the standard 9-option taxonomy
  - Captures the primary theme, topic, or activity of the task
- **Task Tags**
  - 3-6 descriptive keyword tags that capture concepts relevant to the task (e.g. text-editing, vim, python, etc.)
  - No pre-existing taxonomy, come up with any relevant tags

# Task Type Taxonomy

---

system-administration

build-and-dependency-management

data-processing

games

software-engineering

machine-learning

debugging

security

scientific-computing

# Example task.yaml

---

```
instruction: |
  Your task is to decompile a Python function in the format of bytecode, which is compiled in CPython 3.8, and save the function in decompiled.py in the same folder as task.yaml.
  The function name in the decompiled.py should be named "func".
  Note that the function is pickled using dill and serialized using base64 in the file called func.serialized.
author_name: anonymous
author_email: anonymous
difficulty: hard
category: software-engineering
tags:
  - software-engineering
  - python
  - bytecode
  - serialization
parser_name: pytest
max_agent_timeout_sec: 360.0
max_test_timeout_sec: 60.0
run_tests_in_same_shell: false
disable_ascinema: false
estimated_duration_sec:
expert_time_estimate_min: 360
junior_time_estimate_min: 720
```

# Task Design Requirements

---

- **Multi-Step**
  - Tasks requiring chaining multiple commands, handling intermediate states
  - Not solvable with a single command or episode of commands
- **Testable**
  - Must be fully specified so that agent can attempt to complete the task without ambiguity
  - Must be able to design unit tests to determine if final state of environment is correct
- **Unique**
  - Tasks should be unique and distinct from all existing tasks in the public Terminal-Bench benchmark
  - Should also be unique from other tasks submitted for this project
- **No Privileged Ops**
  - Tasks must not require root-level privileges or unsafe Docker settings like `--privileged`
- **Standalone**
  - Tasks must run to completion without additional user input after start
  - All parameters must be provided via files, flags, environment variables, or in the task.yaml

# Task Submission Checklist

---

- All behavior checked in the test cases is described in the task instruction.
- All behavior described in the task instruction is checked in the unit tests.
- My test cases have informative docstrings that describe which behavior they check. It is hard for the agent to cheat on my task (e.g. by editing data files, looking inside files for strings that represent solutions, training on the test set, etc.).
- My task.yaml was written by a human.
- My solution.sh was written by a human (with minimal help from a language model).
- If the agent produces structured data (e.g. it is tasked with building an API) the exact schema is described in the task.yaml or a separate file.
- If my task uses external dependencies (e.g. Docker images, pip packages, etc.) their versions are pinned to ensure reproducibility.
- I ran this task using agent with a powerful model (e.g., GPT-5). For failing runs (which are expected for harder tasks!), I've added an analysis below to confirm the task itself is valid. (Hint: tb tasks debug can help)
- I formatted and linted the task (Ruff).

# High-Level Tasking Workflow

---

Submissions should be made through this private GitHub repository: <https://github.com/snorkel-ai/snorkel-tb-tasks>

**Once you are granted access, you should:**

1. Install uv: `curl -LsSf https://astral.sh/uv/install.sh | sh`
2. Clone the repository
3. Create a task in your local environment under tasks directory using the CLI wizard
4. Ensure the oracle solution passes and minimum criteria for difficulty are met
5. Push a branch (e.g. `username/<task-id>`)
6. Create a pull request against the repository and make sure the title starts with "Task:"
7. Iterate on that branch until all CI/Evals pass

We expect submission tasks to take 2-5 hours based on task difficulty.

# Initial Setup and Task Creation

---

## Initial Setup:

1. Install uv: `curl -LsSf https://astral.sh/uv/install.sh | sh`
2. Clone the repo: `git clone https://github.com/snorkel-ai/snorkel-tb-tasks.git`
3. Navigate to the created local directory: `cd snorkel-tb-tasks`

## Task Creation:

1. Run the task creation wizard: `uv run stb tasks create`
2. Follow the steps given in the terminal to instantiate your task
  - a. See the EC Guidelines for more detail

This will create a new folder for your task in your `snorkel-tb-tasks/tasks` directory. The contents of the folder contain the necessary file skeleton to complete the task.

# Completing a Task - Part 1

---

**After you have followed the steps to create your task folder and skeleton, do the following to complete your task:**

1. Edit the created Dockerfile using a text editor to set up your task environment
  - a. Add any dependencies of the task, such as additional required packages
  - b. If you require a multi-container environment or other custom configuration, see [this page](#) for more information on how to customize your docker-compose.yaml
2. Enter your task container in interactive mode: `tb tasks interact -t <task-name>`
3. While interacting with your task container, test your solution idea to make sure that it works as expected
  - a. Once solution is verified, record it and exit the container

# Completing a Task - Part 2

---

## Continued from previous slide:

4. Modify the solution file (solution.sh) with the verified commands from the previous step
  - a. This file will be used by the OracleAgent to ensure the task is solvable
  - b. If you need to run commands that are not possible to run with a bash script (e.g. vim), use a solution.yaml file to configure interactive commands
5. Update the tests/test\_outputs.py file to verify task completion
  - a. Create pytest unit tests to ensure that the task was completed correctly
  - b. If tests require any file dependencies, place them in the tests/ directory
6. Test your task solution passes and meets all the requirements specified in the tests: `tb run --agent oracle --task-id <task-name>`
7. Push a branch (e.g. my-task-1) to the repository
8. Create a PR to submit your task

# Submission Quality Control

---

All submitted tasks are evaluated for quality and accuracy in three ways.

## 1. LLM-as-Judge (LLMaJ)

- a. Programmatic evaluators that run upon task submission
- b. Check for non-deterministic quality indicators (e.g. do the tests cover all intended behavior for the task?)
- c. Should iterate on submission until all of these pass

## 2. Deterministic Continuous Integration (CI)

- a. Programmatic checks that run upon task submission
- b. Check for deterministic adherence to structure and formatting (e.g. do all required files exist?)
- c. Should correct any issues that lead to CI failure

## 3. Manual Peer Review

- a. Once your task has passed both types of programmatic checks, it will be manually checked by an expert peer reviewer
- b. Peer reviewer will leave comments on the task indicating what needs to be added, removed, or corrected

# LLMaJ - Part 1

---

- Behavior in Task Description
  - Checks whether all behaviors asserted by the tests are explicitly described in task.yaml
  - Fails if the tests enforce implementation details that the task does not specify
- Behavior in Tests
  - Checks whether all behaviors described in the task.yaml are actually tested (inverse of above)
  - Fails if one or more specified behaviors from the task is not fully tested
- Informative Test Docstrings
  - Checks for a clear and informative docstring describing each test
  - Fails if tests are missing docstrings or have poorly written ones
- Anti Cheating Measures
  - Checks that tests are executed after agent runs and there is no path to cheat by reading tests
  - Fails if there are ways for an agent to pass the tests by cheating instead of actually implementing functionality
- Structured Data Schema
  - If applicable, checks that any data files implied by the tests are specified in the task.yaml or a separate schema file
  - Fails if any tests imply a data file that is not specified in the task

# LLMaJ - Part 2

---

- Pinned Dependencies
  - Checks if required dependencies are pinned in the Dockerfile for reproducibility
  - Fails if any required dependencies are missing a pinned version
- Typos
  - Checks for typos in filenames, variables, and instructions
  - Fails if any typos are found
- Tests or Solution in Image
  - Checks to make sure that the Dockerfile does not copy tests or a solution into the image
  - Fails if Dockerfile does copy tests or solution into image
- Test Deps in Image
  - Checks to make sure that test dependencies are installed at test time (e.g., in [run-tests.sh](#))
  - Fails if dependencies are installed into the image/during build
- Hardcoded Solution
  - Checks that the solution script writes a full source file implementing required behavior, not a hard-coded answer
  - Fails if the solution script merely echoes a final answer instead of implementing logic
- File Reference Mentioned
  - Checks that the required output file names given by the tests are specified in the task.yaml
  - Fails if the task.yaml does not mention the required output names that the tests check for

# CI Evaluations - Part 1

---

- **Check Test File References**
  - **Operates on:** *test\_outputs.py, solution.yaml, task.yaml*
  - **Checks:** all file references in *test\_outputs.yaml* are mentioned in the *task.yaml* - files that appear in both test and solution files but not in *task.yaml* are flagged
  - **If fails:** either remove the files or add them to the task instructions
- **Validate Task Fields**
  - **Operates on:** *task.yaml*
  - **Checks:** all required fields are present in the *task.yaml*
  - **If fails:** add the missing required field
- **Check Canary String**
  - **Operates on:** *Task.yaml, solution.yaml, Dockerfile, solution.sh, test\_outputs.py*
  - **Checks:** the canary string is present at the top of all required files
  - **If fails:** add the canary string to the top of all required files

# CI Evaluations - Part 2

---

- **Check for Privileged Containers**
  - **Operates on:** *docker-compose.yaml*
  - **Checks:** for privileged containers in the *docker-compose.yaml*
  - **If fails:** remove any privileged container in the *docker-compose.yaml*
- **Check PR Files**
  - **Operates on:** *tasks/*
  - **Checks:** for any files not in the *tasks/* directory
  - **If fails:** delete or move all files not in the *tasks/* directory
- **Python Tests**
  - **Operates on:**
  - **Checks:**
  - **If fails:**

# CI Evals - Part 3

---

- **Check task file sizes**
  - **Operates on:** all files
  - **Checks:** that every file is under 1MB
  - **If fails:** any files over 1MB should be removed or condensed
- **Ruff**
  - **Operates on:** all files
  - **Checks:** for Linting errors
  - **If fails:** fix any Linting errors
- **Check run-tests.sh Sanity**
  - **Operates on:** *run-tests.sh, task.yaml*
  - **Checks:** if the run-tests.sh file uses uv init or has a task.yaml file with global or system-wide keywords
  - **If fails:** update run-tests.sh to contain a uv init or uv venv, or add global or system-wide keywords to the task.yaml

# CI Evaluations - Part 4

---

- **Check Task Absolute Paths**
  - **Operates on:** *task.yaml*
  - **Checks:** that task instructions use absolute paths rather than relative paths
  - **If fails:** update instructions to use absolute paths
- **Check Dockerfile References**
  - **Operates on:** *Dockerfile*
  - **Checks:** if any of solution.yaml, solution.sh, run-tests.sh, or test\_outputs.py are in the Dockerfile
  - **If fails:** remove the forbidden file(s) from the Dockerfile
- **Test Tasks with Multiple Agents**
  - **Operates on:**
  - **Checks:**
  - **If fails:**

## Rate Schedule

---

| Task Type  | Easy  | Medium | Hard  |
|------------|-------|--------|-------|
| Submission | \$140 | \$245  | \$350 |
| Review     | TBD   | TBD    | TBD   |

# Resources

---

- Project Guidelines: <https://docs.google.com/document/d/1FjfHrOqHF-Uv16XEvs29gUlwu9jDxRmli36jsyzYLzk/edit?tab=t.0>
- Training Site: <https://snorkel-ai.github.io/Terminus-EC-Training/>
- GitHub Repo: <https://github.com/snorkel-ai/snorkel-tb-tasks>
- Slack Channel: #ec-terminus-submission