# Snorkel AI
## Terminus - Review/Feedback Session

November 3, 2025

# Common Failures with "Hard" Difficulty

We are seeing some PRs achieving "Hard" difficulty due to the following error modes:

- Testing things that are not described in the task.yaml file
  - Checked by Analysis on Agent Failures
  - E.g. checking if an output file has a very particular, case-sensitive string sequence that wasn't specified

- Not giving long enough timeout to the agent
  - max_agent_timeout_sec in task.yaml determines how long the agent has before it times out
  - If this is set too low, this can create false "hard" difficulty

# Checking for Failures due to Insufficient Task Instructions vs. Normal Agent Limitations

- Sometimes when testing the agent(s), they fail not due to performance limitations or lack of capability, but due to insufficient task instructions given

- When this is the case, need to alter task instructions to create a genuine agent failure

- You can get more detail on this from the steps to the right
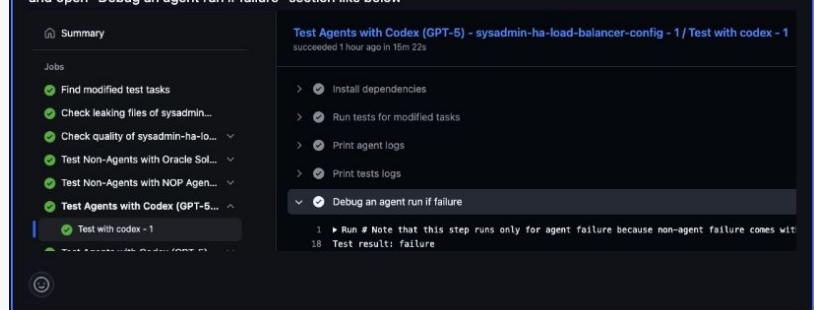
- You should check these logs even when all CI passes!

# Running LLMaJ Locally

- [https://github.com/snorkel-ai/snorkel-tb-tasks/pull/98#issuecomment-3475046582](https://github.com/snorkel-ai/snorkel-tb-tasks/pull/98#issuecomment-3475046582)

- Can run LLMaJ checks locally using: uv run tb tasks check {task_name}

- Requires Anthropic API key which should be delivered shortly

# Over-indexing Checking Output

- [https://github.com/snorkel-ai/snorkel-tb-tasks/pull/112](https://github.com/snorkel-ai/snorkel-tb-tasks/pull/112)

- Checking the output is correct, but you must also try to introduce intermediate milestones that could check the process as well

- In example above, task is just asking for 2 output files that theoretically could be directly generated by LLM
  - Could also ask for a Python script to solve the problem and check if the file was created, runs, has specific functions, etc.

- Need to ask in instructions to specifically create and include any code/scripts - we can't verify that the models actually created code if we don't ask for it to be stored somewhere