

Verification and Validation Report: TPG

Team 3, Tangle
Calvyn Siong
Cyruss Allen Amante
Edward Gao
Richard Li
Mark Angelo Cruz

March 10, 2025

1 Revision History

Date	Version	Notes
3/10/2025	1.0	First write up of VnV Report

2 Symbols, Abbreviations and Acronyms

symbol	description
TPG	Tangled Program Graphs
DNNs	Deep Neural Networks
RL	Reinforcement Learning
SRS	Software Requirement Specification
FR	Functional Requirement
NFR	Non-Functional Requirement
SLN	Solution
VnV	Verification and Validation

Contents

1	Revision History	i
2	Symbols, Abbreviations and Acronyms	ii
3	Functional Requirements Evaluation	1
3.1	MuJoCo Integration	1
3.2	Experiment Visualization	1
3.3	Github Actions CI/CD Pipeline	2
3.4	Software Engineering Practices	2
4	Nonfunctional Requirements Evaluation	3
4.1	Usability	3
4.2	Performance	4
4.3	Operational and Environmental	4
4.4	Maintainability	5
4.5	Security	6
4.6	Compliance	7
5	Unit Testing	8
5.1	Behaviour-Hiding Module	8
5.1.1	RegisterMachine Crossover Tests	8
5.1.2	Team Crossover Tests	9
5.2	MuJoCo Module	10
5.2.1	MuJoCo Environment Test	10
5.2.2	MuJoCo Ant Test	12
5.2.3	MuJoCo Half Cheetah Test	13
5.2.4	MuJoCo Hopper Test	15
5.2.5	MuJoCo Humanoid Standup Test	16
5.2.6	MuJoCo Inverted Double Pendulum Test	16
5.2.7	MuJoCo Inverted Pendulum Test	17
5.2.8	MuJoCo Reacher Test	18
6	Changes Due to Testing	19
6.1	Feedback from Rev 0	19
6.2	Feedback from Usability Testing	19
7	Automated Testing	20

8	Trace to Requirements	21
9	Trace to Modules	22
10	Code Coverage Metrics	22
11	Appendix — Usability Survey	24
11.1	Usability Survey Questions and Answers	24

List of Tables

1	MuJoCo Integration Tests	1
2	Experiment Visualization Tests	1
3	Github Actions CI/CD Pipeline Tests	2
4	Software Engineering Practices Tests	2
5	Usability Tests	3
6	Performance Tests	4
7	Operational and Environmental Tests	4
8	Maintainability Tests	6
9	Security Tests	6
10	Compliance Tests	7

List of Figures

1	Example of a Numerical Computation Test	5
2	Example of a Linter Error	7
3	Example of a Numerical Computation Test	23

This document cohesively summarizes the results of each test as specified in the [VnV Plan](#) documentation.

3 Functional Requirements Evaluation

3.1 MuJoCo Integration

Table 1: MuJoCo Integration Tests

Test Id	Notes	Result
FR-SLN1	When executing the appropriate script, all MuJoCo environments can be run. The best-performing agent within the policy can be visualized using OpenGL or an MP4 file.	Pass
FR-SLN2	MuJoCo environments within the TPG framework can be successfully run within the Digital Research Alliance, enabling research to be conducted by executing experiments.	Pass

3.2 Experiment Visualization

Table 2: Experiment Visualization Tests

Test Id	Notes	Result
FR-SLN3	When an experiment is running or finished training, the best performing policy can be visualized using the TPG CLI tool.	Pass

3.3 Github Actions CI/CD Pipeline

Table 3: Github Actions CI/CD Pipeline Tests

Test Id	Notes	Result
FR-SLN4	Affirmed that the “Build TPG Project” pipeline properly builds the TPG framework with updated code when changes are pushed to any branch.	Pass
FR-SLN5	When the project building pipeline runs properly, the TPG unit test cases are also automatically ran, and the build will only pass if all the unit tests passes .	Pass
FR-SLN6	Tested that linting and Latex compilation pipeline works as expected.	Pass

3.4 Software Engineering Practices

Table 4: Software Engineering Practices Tests

Test Id	Notes	Result
FR-SLN7	Newly added code in the TPG codebase follows Google’s C++ Style Guide and software engineering best practices such as design patterns, and object-oriented design. This includes careful review and consideration of code readability, extendability, maintainability and scalability. A linter has also been implemented to check for such styling as discussed in 3.3 .	Pass

4 Nonfunctional Requirements Evaluation

4.1 Usability

Table 5: Usability Tests

Test Id	Result
NFR-SLN1	Pass
NFR-SLN2	Pass
NFR-SLN3	Pass

For NFR-SLN1, the TPG framework now includes comprehensive documentation across multiple levels. The main README provides installation instructions and a quick start guide, while the CLI tool documentation details commands for evolving, plotting, and replaying policies. Code documentation with inline comments explains classes and methods, and YAML configuration files include detailed parameter descriptions. Usability testing with graduate students confirmed the documentation’s effectiveness, with a rating of 10/10 for installation experience.

For NFR-SLN2, an event-driven logging architecture was implemented to provide real-time, accurate message logging throughout the training process. Metrics are stored in CSV files, capturing timing, selection, replacement, and removal data. The CLI tool provides real-time feedback during experiment execution. Usability testing rated the accuracy and usefulness of system logs at 9/10, with suggestions to clarify CSV file naming conventions and document the meaning of removal and replacement metrics in the wiki.

For NFR-SLN3, the framework now offers extensive customization options for MuJoCo simulation parameters through YAML configuration files. Each environment includes specific parameters (reward weights, health criteria, observation settings) and common simulation parameters (maximum timesteps, model path). The CLI tool integrates these parameters, allowing users to specify different environments and configurations when running experiments. Usability testing confirmed the ease of locating and customizing parameters,

with suggestions to place YAML files within the experiments directory and reduce the number of log files.

4.2 Performance

Table 6: Performance Tests

Test Id	Result
NFR-SLR4	Pass

For NFR-SLN4, test cases within TPG for the experimental environments have been implemented to check for the accuracy of the numerical computations associated during training. Declaration of variables with proper types (e.g. signed long or int, unsigned long or int) has also been taken into consideration to reduce issues in the future for extremely large or small numbers that may overflow. TPG has been comprehensively tested to guarantee that all computations with high numerical precision (e.g. during the runtime of an experiment) are accurate and contain an acceptable tolerance limit of 0.00001. The results were inspected manually by comparing the actual output to the anticipated output, and performing a calculation to check for quantitative error, and if such error meets the requirements for numerical precision.

4.3 Operational and Environmental

Table 7: Operational and Environmental Tests

Test Id	Result
NFR-SLR6	Pass

For NFR-SLN6, TPG now supports contributions from macOS, Windows, and Linux developers. Previously, only Linux was supported because TPG

```

TEST_CASE("Mujoco_Ant_v4 Reset Function", "[reset]") {
    std::unordered_map<std::string, std::any> params = createDefaultParams();
    Mujoco_Ant_v4 ant(params);
    std::mt19937 rng(1234);

    ant.step_ = 50;

    std::vector<double> qpos = {0.5, 0.8, -0.3};
    std::vector<double> qvel = {0.1, -0.05, 0.05};
    ant.set_state(qpos, qvel);

    std::vector<double> obs(ant.obs_size_, 1.0);
    ant.get_obs(obs);

    ant.reset(rng);

    REQUIRE(ant.step_ == 0);

    for (size_t i = 0; i < ant.state_.size(); i++) {
        REQUIRE(ant.state_[i] == Catch::Approx(0.0).margin(1e-2));
    }
}

```

Figure 1: Example of a Numerical Computation Test

used SCons for C++ builds and Linux-specific dependencies from [requirements.txt](#). With VSCode Dev Containers, a Linux development environment is automatically launched for all developers, ensuring a standardized setup. Simply follow the [Wiki](#) instructions to download all necessary Linux dependencies and build the C++ code reliably. Onboarding on a Macbook has been reduced from 2 weeks to just 5 minutes.

4.4 Maintainability

Table 8: Maintainability Tests

Test Id	Result
NFR-SLR7	Pass

To satisfy the testing requirements for NFR-SLR7 - establishing a secure and robust repository management system, the team has implemented checks to ensure the repository prevents unauthorized access and defective code integration. The repository where our team is working on GitHub (whereas the base TPG repository is based in Gitlab), and access is controlled through a combination of two-factor authentication (2FA), and a main branch that is protected to ensure that merge requests can only be performed after the [TPG project GitHub workflow](#) action pipeline has successfully completed. This pipeline validates the building and testing process, ensuring that only code that passes all checks can be merged into the main branch. Any critical build errors or warnings, create blocking pull request conversations that must be resolved before merging. There is also a specific [GitHub workflow](#) that is used to automatically pull changes from GitLab, eliminating the need for manual merging and risk of human error.

4.5 Security

Table 9: Security Tests

Test Id	Result
NFR-SLR8	Pass

For NFR-SLR8, the .csv, .txt, .png and .mp4 files that are generated within Classic Control and MuJoCo experiments are ignored by Git when making commits to the public repositories in GitHub and GitLab to reduce chance of oversharing sensitive data. Currently, none of these files generate sensitive data, but to follow best practice and to keep the repository at a clean state,

these are not recognized when synchronizing code to each respective repository. Additionally, the team has also manually checked all stored .csv, .txt, .png and .mp4 files along with others that may contain textual information to see if data within them are sensitive and must be kept private.

4.6 Compliance

Table 10: Compliance Tests

Test Id	Result
NFR-SLR9	Pass

The modified codebase is successfully analyzed using Clang-Tidy and Clang-Format within the CI/CD pipeline. Code change discussions take place through pull request conversations made to the main branch. All errors and warnings are generated based on the C++ Style Guidelines. Any critical errors found during the linting process create blocking pull request conversations that must be resolved before merging into the main branch.

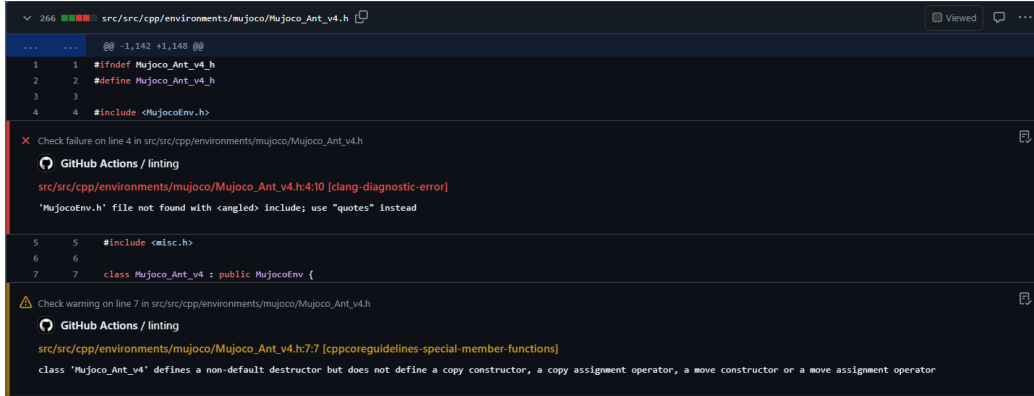


Figure 2: Example of a Linter Error

5 Unit Testing

5.1 Behaviour-Hiding Module

5.1.1 RegisterMachine Crossover Tests

Type: Automatic, Functional

Initial State: The TPG and `RegisterMachine` objects are initialized with default parameters and state.

Test Case Derivation: The expected behavior is derived from the correct crossover functionality, chunk splitting, and recombination of `RegisterMachine` objects, ensuring valid instruction patterns and segment lengths.

Test Procedure: The test will be performed as follows:

- **Basic Crossover Functionality Test:**
 - **Input:** Two parent `RegisterMachine` objects.
 - **Output:** Two child `RegisterMachine` objects with valid instructions and actions.
 - **Test Derivation:** Verifies that crossover produces children with reasonable sizes, valid actions, and instruction patterns derived from both parents.
- **Chunk Splitting and Recombination Test:**
 - **Input:** Two parent `RegisterMachine` objects with predefined instruction sequences.
 - **Output:** Two child `RegisterMachine` objects with instruction counts and patterns derived from both parents.
 - **Test Derivation:** Ensures that crossover produces children with valid instruction counts and different instruction patterns.
- **Crossover Constraints Test:**
 - **Input:** Two parent `RegisterMachine` objects with predefined instruction sequences.

- **Output:** Two child `RegisterMachine` objects adhering to crossover constraints.
- **Test Derivation:** Verifies that crossover points, segment lengths, and resulting program lengths adhere to predefined constraints (`dcmax`, `lsmax`, `dsmax`, `lmin`, `lmax`).

5.1.2 Team Crossover Tests

Type: Automatic, Functional

Initial State: The `TPG` and `team` objects are initialized with default parameters and state.

Test Case Derivation: The expected behavior is derived from the correct crossover functionality of `team` objects, ensuring valid team sizes, atomic program preservation, and adherence to team size limits.

Test Procedure: The test will be performed as follows:

- **Single Program Teams - Linear Crossover Test:**
 - **Input:** Two parent teams with single programs.
 - **Output:** A child team with one program.
 - **Test Derivation:** Verifies that crossover produces a child team with a single program and valid atomic count.
- **Multi-Program Teams - Team Crossover Test:**
 - **Input:** Two parent teams with multiple programs.
 - **Output:** A child team with programs derived from both parents.
 - **Test Derivation:** Ensures that crossover produces a child team with a valid size (within `max_team_size`) and at least one atomic program.
- **Atomic Program Preservation Test:**
 - **Input:** Two parent teams with atomic and non-atomic programs.
 - **Output:** A child team with at least one atomic program.

- **Test Derivation:** Verifies that crossover preserves atomic programs in the child team.
- **Team Size Limits Test:**
 - **Input:** Two parent teams with the maximum number of programs.
 - **Output:** A child team with a size within `max_team_size`.
 - **Test Derivation:** Ensures that crossover produces a child team adhering to the predefined team size limit.

5.2 MuJoCo Module

5.2.1 MuJoCo Environment Test

Type: Automatic, Functional

Initial State: The MuJoCo environment is initialized using the `MockMujocoEnv` class with the appropriate model path.

Test Case Derivation: The expected behavior is derived from the correct initialization, state setting, and simulation execution of the MuJoCo environment.

Test Procedure: The test will be performed as follows:

- **Simulation Initialization Test:**
 - **Input:** Model path determined by the `determine_tpg_env()` function.
 - **Output:** Successful initialization of the MuJoCo environment.
 - **Test Derivation:** Verifies that the `initialize_simulation()` function correctly initializes the MuJoCo environment, ensuring that the model (`m_`) and data (`d_`) pointers are not null.
- **Set State Test:**
 - **Input:** Position vector `qpos` set to `{0.5, 0.5, ...}` and velocity vector `qvel` set to `{0.1, 0.1, ...}`.
 - **Output:** Updated state in the MuJoCo environment.

- **Test Derivation:** Ensures that the `set_state()` function correctly updates the position and velocity states in the MuJoCo environment, verifying that `d_->qpos` and `d_->qvel` match the input values.
- **Do Simulation Test:**
 - **Input:** Control vector `control` set to `{0.2, 0.2, ...}` and a step count of 5.
 - **Output:** Updated control values in the MuJoCo environment.
 - **Test Derivation:** Confirms that the `do_simulation()` function correctly applies the control inputs and updates the simulation state, ensuring that `d_->ctrl` matches the input control values.

Test Cases

Test Case 1: Simulation Initialization

- **Description:** Tests the initialization of the MuJoCo simulation environment.
- **Steps:**
 1. Create a `MockMujocoEnv` object with the model path.
 2. Call `initialize_simulation()`.
 3. Verify that `m_` and `d_` are not null.

Test Case 2: Set State

- **Description:** Tests the ability to set the state of the MuJoCo environment.
- **Steps:**
 1. Create a `MockMujocoEnv` object and initialize the simulation.
 2. Set `qpos` to `{0.5, 0.5, ...}` and `qvel` to `{0.1, 0.1, ...}`.
 3. Call `set_state(qpos, qvel)`.
 4. Verify that `d_->qpos` and `d_->qvel` match the input values.

Test Case 3: Do Simulation

- **Description:** Tests the execution of a simulation step with control inputs.
- **Steps:**
 1. Create a `MockMujocoEnv` object and initialize the simulation.
 2. Set `control` to `{0.2, 0.2, ...}`.
 3. Call `do_simulation(control, 5)`.
 4. Verify that `d_->ctrl` matches the input control values.

5.2.2 MuJoCo Ant Test

Type: Automatic, Functional

Initial State: The `Mujoco_Ant_v4` environment is initialized.

Test Case Derivation: The expected value is based on the logic that the environment should be terminal when the step count reaches 200, as per the environment's design.

Test Procedure: The test will be performed as follows:

- **Healthy Reward Test:**
 - **Input:** None.
 - **Output:** Returns `healthy_reward_`.
 - **Test Derivation:** Verifies that the `healthy_reward()` function correctly returns the predefined `healthy_reward_` value.
- **Control Cost Test:**
 - **Input:** Action vector `{0.1, -0.1, 0.2, 0.3}`.
 - **Output:** Calculated control cost.
 - **Test Derivation:** Ensures the `control_cost()` function computes the cost using `control_cost_weight_` and the squared sum of action values.
- **Contact Cost Test:**
 - **Input:** None.

- **Output:** Non-negative contact cost.
- **Test Derivation:** Confirms that the `contact_cost()` function always returns a non-negative value.
- **Is Healthy Test:**
 - **Input:** Modify `qpos[2]` to test health conditions.
 - **Output:** Boolean indicating health status.
 - **Test Derivation:** Validates that `is_healthy()` returns `true` when `qpos[2]` is within the healthy range and `false` otherwise.
- **Simulation Step Test:**
 - **Input:** Action vector `{0.1, -0.1, 0.2, 0.3}`.
 - **Output:** Finite reward and incremented `step_`.
 - **Test Derivation:** Checks that `sim_step()` processes actions correctly, updating the environment state and returning a valid reward.
- **Get Observation Test:**
 - **Input:** None.
 - **Output:** Non-zero observation vector.
 - **Test Derivation:** Ensures `get_obs()` reflects the current state of the environment in the observation vector.
- **Reset Function Test:**
 - **Input:** Random number generator.
 - **Output:** Reinitialized environment state.
 - **Test Derivation:** Verifies that `reset()` brings the environment back to its initial state, setting `step_` to 0 and state values close to zero.

5.2.3 MuJoCo Half Cheetah Test

Type: Automatic, Functional

Initial State: The `Mujoco_Half_Cheetah_v4` environment is initialized with default parameters.

Test Case Derivation: The expected behavior is derived from the correct initialization, terminal condition, control cost calculation, simulation step execution, and reset functionality of the environment.

Test Procedure: The test will be performed as follows:

- **Initialization Test:**
 - **Input:** Default parameters.
 - **Output:** Correct initialization of environment variables.
 - **Test Derivation:** Verifies that `n_eval_train_`, `n_eval_validation_`, `n_eval_test_`, and `max_step_` are set correctly.
- **Terminal Condition Test:**
 - **Input:** Step count set to 200.
 - **Output:** `terminal()` returns `true`.
 - **Test Derivation:** Ensures the environment terminates when the step count reaches 200.
- **Control Cost Test:**
 - **Input:** Action vector `{0.1, -0.1, 0.2}`.
 - **Output:** Calculated control cost.
 - **Test Derivation:** Confirms that `control_cost()` computes the cost using the squared sum of action values.
- **Simulation Step Test:**
 - **Input:** Action vector `{0.1, -0.1, 0.2}`.
 - **Output:** Finite reward and incremented step count.
 - **Test Derivation:** Verifies that `sim_step()` processes actions correctly and updates the step count.
- **Reset Function Test:**
 - **Input:** Random number generator and modified state.

- **Output:** Reinitialized environment state.
- **Test Derivation:** Ensures `reset()` resets the step count and state values to initial conditions.

5.2.4 MuJoCo Hopper Test

Type: Automatic, Functional

Initial State: The `MujocoHopper_v4` environment is initialized with default parameters.

Test Case Derivation: The expected behavior is derived from the correct initialization, terminal condition, healthy reward, control cost calculation, health check, simulation step execution, observation retrieval, and reset functionality of the environment.

Test Procedure: The test will be performed as follows:

- **Initialization Test:** Similar to Half Cheetah or Ant test, but the input is default parameters.
- **Terminal Condition Test:** Similar to Half Cheetah or Ant test, but the input includes modifying `qpos[1]` to test the healthy z range and step count.
- **Healthy Reward Test:** Similar to Half Cheetah or Ant test, but the input is none, and the output is `healthy_reward_`.
- **Control Cost Test:** Similar to Half Cheetah or Ant test, but the input is action vector `{0.1, -0.1, 0.2}`.
- **Is Healthy Test:** Similar to Half Cheetah or Ant test, but the input includes modifying `qpos[1]` and `qpos[2]` to test the healthy z range and angle range.
- **Simulation Step Test:** Similar to Half Cheetah or Ant test, but the input is action vector `{0.1, -0.1, 0.2}`.
- **Get Observation Test:** Similar to Half Cheetah or Ant test, but the input includes manually setting `qpos` and `qvel` to non-zero values.
- **Reset Function Test:** Similar to Half Cheetah or Ant test, but the input includes modifying `qpos`, `qvel`, and step count before resetting.

5.2.5 MuJoCo Humanoid Standup Test

Type: Automatic, Functional

Initial State: The `Mujoco_Humanoid_Standup_v4` environment is initialized with default parameters.

Test Case Derivation: The expected behavior is derived from the correct initialization, terminal condition, simulation step execution, observation retrieval, and reset functionality of the environment.

Test Procedure: The test will be performed as follows:

- **Initialization Test:** Similar to Hopper or Half Cheetah test, but the input is default parameters.
- **Terminal Condition Test:** Similar to Hopper or Half Cheetah test, but the input is step count set to 200.
- **Simulation Step Test:** Similar to Hopper or Half Cheetah test, but the input is action vector $\{0.1, -0.1, 0.2\}$.
- **Get Observation Test:** Similar to Hopper or Half Cheetah test, but the input includes verifying non-zero observation values.
- **Reset Function Test:** Similar to Hopper or Half Cheetah test, but the input includes modifying `qpos`, `qvel`, and step count before resetting.

5.2.6 MuJoCo Inverted Double Pendulum Test

Type: Automatic, Functional

Initial State: The `Mujoco_Inverted_Double_Pendulum_v4` environment is initialized with default parameters.

Test Case Derivation: The expected behavior is derived from the correct initialization, terminal condition, simulation step execution, observation retrieval, and reset functionality of the environment.

Test Procedure: The test will be performed as follows:

- **Initialization Test:** Similar to Humanoid Standup or Hopper test, but the input is default parameters.

- **Terminal Condition Test:** Similar to Humanoid Standup or Hopper test, but the input includes modifying `site_xpos[2]` to test the terminal threshold and step count.
- **Simulation Step Test:** Similar to Humanoid Standup or Hopper test, but the input is action vector $\{0.1\}$.
- **Get Observation Test:** Similar to Humanoid Standup or Hopper test, but the input includes manually setting `qpos` and `qvel` to non-zero values.
- **Reset Function Test:** Similar to Humanoid Standup or Hopper test, but the input includes modifying `qpos`, `qvel`, and step count before resetting.

5.2.7 MuJoCo Inverted Pendulum Test

Type: Automatic, Functional

Initial State: The `Mujoco_Inverted_Pendulum_v4` environment is initialized with default parameters.

Test Case Derivation: The expected behavior is derived from the correct initialization, terminal condition, simulation step execution, observation retrieval, and reset functionality of the environment.

Test Procedure: The test will be performed as follows:

- **Initialization Test:** Similar to Inverted Double Pendulum or Humanoid Standup test, but the input is default parameters.
- **Terminal Condition Test:** Similar to Inverted Double Pendulum or Humanoid Standup test, but the input includes modifying `qpos[1]` to test the terminal threshold and step count.
- **Simulation Step Test:** Similar to Inverted Double Pendulum or Humanoid Standup test, but the input is action vector $\{0.1\}$ and the expected reward is `1.0`.
- **Get Observation Test:** Similar to Inverted Double Pendulum or Humanoid Standup test, but the input includes manually setting `qpos` and `qvel` to non-zero values and verifying the observation vector.

- **Reset Function Test:** Similar to Inverted Double Pendulum or Humanoid Standup test, but the input includes modifying `qpos`, `qvel`, and step count before resetting.

5.2.8 MuJoCo Reacher Test

Type: Automatic, Functional

Initial State: The `Mujoco_Reacher_v4` environment is initialized with default parameters.

Test Case Derivation: The expected behavior is derived from the correct initialization, terminal condition, control cost calculation, distance retrieval, simulation step execution, observation retrieval, and reset functionality of the environment.

Test Procedure: The test will be performed as follows:

- **Initialization Test:** Similar to Inverted Pendulum or Inverted Double Pendulum test, but the input is default parameters.
- **Terminal Condition Test:** Similar to Inverted Pendulum or Inverted Double Pendulum test, but the input is step count set to 200.
- **Control Cost Test:** Similar to Hopper or Half Cheetah test, but the input is action vector $\{0.1, -0.1\}$.
- **Get Distance Test:** Unique to Reacher, the input is none, and the output is a distance vector of size 2.
- **Simulation Step Test:** Similar to Inverted Pendulum or Inverted Double Pendulum test, but the input is action vector $\{0.1, -0.1\}$.
- **Get Observation Test:** Similar to Inverted Pendulum or Inverted Double Pendulum test, but the input includes verifying non-zero observation values.
- **Reset Function Test:** Similar to Inverted Pendulum or Inverted Double Pendulum test, but the input includes modifying `qpos`, `qvel`, and step count before resetting.

6 Changes Due to Testing

6.1 Feedback from Rev 0

The feedback given by the instructor and teaching assistant during Revision 0 was essential in guiding the next steps as the team looks toward the final demonstration. Emphasis was placed on ensuring that usability testing was executed systematically rather than in the more ad-hoc manner initially planned by the team. Some additional changes to be made include ensuring that unit testing and benchmarking of the implemented environments are cohesively executed and investigating whether the integration of deployment within the DRA is possible.

6.2 Feedback from Usability Testing

A usability testing session was conducted with a core contributor from Dr. Kelly’s research group, who is also a primary contributor to the TPG framework. Although the TPG framework is Linux-based, the test user operates on MacOS—a platform that traditionally presents challenges in onboarding and setup.

The testing session revealed several positive aspects of the system. The installation process was rated a 10, with the documentation proving clear and straightforward. Executing the simulation environment using MuJoCo was also rated a 10, as the user found it very easy to follow, particularly appreciating the guidance provided for running multiple experiments.

Regarding system feedback, the logging messages received a rating of 9 in terms of accuracy and utility. However, the user suggested improvements in the clarity of naming for CSV files, specifically by outlining in the wiki what “removal” and “replacement” refer to. For the plotting functionality, feedback indicated that replicating the approach of the previous state (via a script such as `tpg-plot-stats.sh` with a “-for all” flag) would help retrieve comprehensive plots. Additionally, providing an option to save visualizations as PDFs, in addition to PNGs, was recommended.

The user also raised concerns about supporting flexible workflows with experiment parameterization. For instance, when experimenting with environments like the “Inverted Pendulum,” he might create two different YAML files to test various hyperparameters. Currently, the repository supports only a

single experiment instance at a time. This limitation means that when reviewing experiment results, he is forced to navigate through the file system to locate the associated hyperparameters file, a process that he finds inconvenient. Addressing this could involve enhancing the repository to support multiple experiment instances simultaneously, thereby streamlining both the experimentation process and the retrieval of results.

A significant challenge highlighted during the session was the difficulty MacOS users face with a Linux-centric codebase. Previously, a current grad student hosted and maintained a Docker image to assist onboarding, but it proved cumbersome. In response, we introduced Visual Studio Dev Container support to automate the creation of a Linux-based development environment within VS Code. This change aims to simplify the onboarding process for MacOS users, ensuring a more seamless setup.

7 Automated Testing

As a result of the team’s conversion from building the project using SCons to CMake, automated testing became significantly easier to execute and debug. To run any automated tests within a developer’s local environment, a developer can simply execute a command to build the project. This not only compiles everything but also runs all automated tests. If a developer wishes to run only the tests, they must navigate to the directory where the tests were already compiled (typically `/build/tests`). From there, the command `ctest` can be entered into the command prompt. Similar to the compilation process, all automated tests are executed once this command is run.

From the repository’s point of view, tests are executed using GitHub Actions or GitLab CI (depending on which repository is being viewed). Both linting and compilation are performed using the same commands that would be executed within a developer’s local environment. These tests run when a new pull request is made to the main branch, ensuring that all tests pass before merging. The compiler workflow is also executed after merging into the main branch to ensure no errors or unintended changes in code behaviour have occurred. If any test or workflow fails, the logs of the workflow can be reviewed, providing a detailed summary of the reason for failure. This not only allows for easier debugging but also resolves the “works on my machine”

issue.

8 Trace to Requirements

Please refer to the [SRS documentation](#) (Tangle, 2024) for detailed information of each requirement.

Req. ID	System Test ID
FR-1	FR-SLN1 , FR-SLN2
FR-2	FR-SLN3
FR-3	FR-SLN4
FR-4	FR-SLN4
FR-5	FR-SLN4
FR-6	FR-SLN5
FR-7	FR-SLN4
FR-8	FR-SLN2
UH-E1	NFR-SLN1
UH-E2	NFR-SLN2
UH-PI1	NFR-SLN3
UH-L1	NFR-SLN1
UH-L2	NFR-SLN1
UH-UP1	NFR-SLN1
UH-UP2	NFR-SLN1
UH-A1	NFR-SLN1
PR-PA	NFR-SLN4
OE-EPE	NFR-SLN6
OE-OSR	NFR-SLN1
MS-M2	NFR-SLN1
MS-S1	NFR-SLN1 , NFR-SLN7
MS-S2	NFR-SLN6
MS-A2	NFR-SLN6

SR-A1	NFR-SLN7
SR-A2	NFR-SLN7
SR-I1	NFR-SLN7
SR-P1	NFR-SLN7 , NFR-SLN8
SR-AU1	NFR-SLN2

9 Trace to Modules

Please refer to the [Module Guide documentation](#) (Tangle, 2025) for detailed information of each module.

Req. ID	System Test ID
M1	N/A
M2	FR-SLN3 , FR-SLN5 , NFR-SLN1 , NFR-SLN2 , NFR-SLN4
M3	FR-SLN1 , FR-SLN2 , FR-SLN3 , FR-SLN5 , NFR-SLN3 , NFR-SLN4
M4	FR-SLN3 , FR-SLN5
M5	NFR-SLN2 , FR-SLN5 , NFR-SLN6 , NFR-SLN7 , NFR-SLN8
M6	FR-SLN4 , FR-SLN5 , NFR-SLN6

10 Code Coverage Metrics

The image below displays code coverage report metrics for TPG. This is specific to the files that our team have unit tested. The code coverage for the existing TPG module are lower as there was a large amount of existing prewritten code, and the tests that have been written are specific to only a couple crucial TPG functions. The MuJoCo module which we have implemented ourselves has a much higher code coverage, as it was solely written by our team.

File	Lines	Exec	Coverage
src/engine/RegisterMachine.cc	412	389	94.4%
src/engine/TPG.cc	1262	345	27.3%
src/engine/team.cc	573	185	32.2%
src/environments/mujoco/MujocoEnv.cc	156	142	91.0%
src/environments/mujoco/Mujoco_Ant_v4.cc	124	112	90.3%
src/environments/mujoco/Mujoco_Half_Cheetah_v4.cc	98	87	88.8%
src/environments/mujoco/Mujoco_Hopper_v4.cc	112	103	92.0%
src/environments/mujoco/Mujoco_Humanoid_Standup_v4.cc	143	118	82.5%
src/environments/mujoco/Mujoco_Inverted_Double_Pendulum_v4.cc	87	82	94.3%
src/environments/mujoco/Mujoco_Inverted_Pendulum_v4.cc	76	72	94.7%
src/environments/mujoco/Mujoco_Reacher_v4.cc	92	85	92.4%
src/experiments/api_client.cc	103	37	35.9%

Figure 3: Example of a Numerical Computation Test

References

- Team 3 Tangle. System requirements specification. <https://github.com/TPGEngine/tpg/blob/main/docs/SRS/SRS.pdf>, 2024.
- Team 3 Tangle. System requirements specification. <https://github.com/TPGEngine/tpg/blob/main/docs/Design/SoftArchitecture/MG.pdf>, 2025.

11 Appendix — Usability Survey

11.1 Usability Survey Questions and Answers

1. What operating system do you use?

- MacOS

2. On a scale from 1-10 (higher means better), how would you rate your installation experience through the documentation?

- 10 — It was pretty simple to follow and understand.

3. On a scale from 1-10 (higher means easier), how easy would you say it was to execute a simulation environment using MuJoCo?

- 10 — Very easy to use. The documentation provided clear instructions, including details on running multiple experiments.

4. On a scale from 1-10 (higher means better), how accurate and useful are the messages logged by the system?

- 9 — Some suggestions included providing more clarity on the naming of CSV files. It was recommended that the wiki should explain what is meant by “removal” and “replacement.” Additionally, for plotting, it was suggested to replicate the approach of the previous setup (e.g., using a script like `tpg-plot-stats.sh` with a `-for all` flag) and offer an option to save figures as PDF instead of just PNG.

5. On a scale from 1-10 (higher means easier), how easy would you say it was to locate and customize parameters for MuJoCo?

- The experience was rated positively. However, the current workflow does not fully support flexible experimentation. For example, when testing an environment like the “Inverted Pendulum” to improve TPG’s performance, the user might create two different YAML files for different hyperparameter configurations. Presently, the repository only supports a single experiment instance at a time, and when reviewing the experiment results, the user has to navigate through the file system

to locate the associated hyperparameter file. A workflow that accommodates multiple experiment instances would improve usability.

6. If applicable, on a scale from 1-10 (higher means easier), how easy was it to implement changes to the code?

- The user mentioned a willingness to implement code changes related to “Recursive Forecasting” and indicated that he would provide further feedback on the process, since the branch he was working on has not been merged yet and there’d be quite a lot of merge conflicts.

Appendix — Reflection

The information in this section will be used to evaluate the team members on the graduate attribute of Reflection.

The purpose of reflection questions is to give you a chance to assess your own learning and that of your group as a whole, and to find ways to improve in the future. Reflection is an important part of the learning process. Reflection is also an essential component of a successful software development process.

Reflections are most interesting and useful when they’re honest, even if the stories they tell are imperfect. You will be marked based on your depth of thought and analysis, and not based on the content of the reflections themselves. Thus, for full marks we encourage you to answer openly and honestly and to avoid simply writing “what you think the evaluator wants to hear.”

Please answer the following questions. Some questions can be answered on the team level, but where appropriate, each team member should write their own response:

1. What went well while writing this deliverable?

One part that went well for this deliverable is that valuating the functional and non functional requirements was relatively simple, as we were able to trace it back to our VNV plan and SRS report, which was relatively straightforward due to the clear traceability between the VnV Plan, the SRS report, and the implemented tests. The team was able to match crucial requirements to relevant test cases, ensuring that all

functional requirements, such as MuJoCo integration, experiment visualization, and CI/CD pipeline functionality, were validated. For example, the functional requirements for MuJoCo integration (FR-SLN1 and FR-SLN2) were tested by running the environments and verifying their compatibility with the TPG framework, while the CI/CD pipeline requirements (FR-SLN4, FR-SLN5, and FR-SLN6) were validated through automated builds, unit tests, and linting checks. This goes for non-functional requirements, such as usability, maintainability, and compliance, which were evaluated through proper testing - including the implementation of VSCode Dev Containers to standardize the development environment and the integration of Clang-Tidy and Clang-Format to enforce coding standards. The team was able to confidently verify that the fundamental requirements were met, ensuring the project's reliability and functionality. Our team believes that we were able to provide clear evidence of the system's compliance with the specified requirements while writing this deliverable.

2. What pain points did you experience during this deliverable, and how did you resolve them?

Writing out the specifications for the unit tests and revising the similar section in VnV plan, while mapping unit tests to their specific modules was relatively challenging. It was challenging because it required a deeper understanding of both the TPG project's structure and the expected behavior of each module.

Another significant challenge was developing code coverage metrics for C++ for the project. Due to mismatching libraries and compiler versions, generating accurate code coverage reports proved to be difficult. The team encountered issues with compatibility between the coverage tools and the build system, which often resulted in errors generating code coverage reports.

3. Which parts of this document stemmed from speaking to your client(s) or a proxy (e.g. your peers)? Which ones were not, and why?

The usability testing section was directly informed by conversations with our client's graduate students, who provided valuable feedback on documentation clarity, logging functionality, and parameter customization. Their ratings (10/10 for installation experience, 9/10 for logging

accuracy) and specific suggestions (clarifying CSV file naming, relocating YAML files) shaped our assessment of the system’s usability.

The MuJoCo integration and performance sections were based on technical evaluations conducted in collaboration with Dr. Kelly’s team, who provided expertise on expected agent behaviors and numerical precision requirements. Their involvement ensured our testing addressed research-relevant concerns rather than just technical functionality.

Sections covering CI/CD pipeline and software engineering practices were developed primarily through our team’s technical assessment, with minimal client input. This was appropriate as these components primarily serve developer needs rather than research objectives, and our team had sufficient expertise to evaluate them independently.

4. In what ways was the Verification and Validation (VnV) Plan different from the activities that were actually conducted for VnV? If there were differences, what changes required the modification in the plan? Why did these changes occur? Would you be able to anticipate these changes in future projects? If there weren’t any differences, how was your team able to clearly predict a feasible amount of effort and the right tasks needed to build the evidence that demonstrates the required quality? (It is expected that most teams will have had to deviate from their original VnV Plan.)

Our VnV activities deviated from the original plan in several key areas. First, we initially planned for comprehensive automated testing across all MuJoCo environments, but discovered that visual inspection by domain experts was more effective for validating agent behaviors. The complexity of reinforcement learning outcomes made it difficult to define automated pass/fail criteria, requiring us to rely more heavily on manual validation.

Second, our usability testing became more structured than originally planned. Following feedback from Revision 0, we implemented formal surveys and specific task-based evaluations rather than the ad-hoc approach initially outlined. This change improved the quality of feedback and provided more actionable insights.

Third, we underestimated the effort required for cross-platform compatibility testing. The diversity of development environments among

researchers necessitated more extensive testing than anticipated, leading us to implement containerization solutions that weren't in the original plan.

These changes occurred primarily due to our initial unfamiliarity with the research domain and underestimation of the complexity of validating reinforcement learning systems. In future projects, we would anticipate similar challenges by consulting domain experts earlier in the planning process and allocating more resources to areas requiring specialized knowledge.