

\*Text highlighted in Green is new (introduced by the framework). Uncolored text was in the participant's original requirements.\*

### 1. TITLE: Real-Time Monitoring for Human Oversight

SYSTEM\_REQUIREMENT: The system shall provide real-time monitoring data to the human overseer to ensure transparency in decision-making and robot performance.

#### IMPLEMENTATION\_DETAILS:

- \* Display location of Center of Gravity
- \* Provide current in ampere as an indicator of torque and resistance
- \* Display robot temperature
- \* Show linear velocity including current and target values
- \* Provide binary foot contacts to determine stability
- \* Show confidence scores of the RL policy for action selection.

RATIONALE: Providing real-time monitoring data ensures that the overseer can interpret system behavior efficiently, detect anomalies early, and respond appropriately.

### 2. TITLE: Intuitive Information Presentation for Human Oversight

SYSTEM\_REQUIREMENT: The system shall present real-time information in a manner that aligns with the expertise and cognitive load of the overseer.

#### IMPLEMENTATION\_DETAILS:

- \* Diagram of the robot model with highlighted visual information for CoG and foot contacts
- \* Real-time dashboard with visual indicators for key state variables such as velocity, current and RL policy confidence score
- \* Simplified color-coded alerts indicating system performance (green for normal, yellow for caution, red for critical)
- \* Timeline visualization of past actions and predicted trajectory
- \* Configurable verbosity levels (high-detail logs for engineers summaries for operators).

RATIONALE: Ensuring information is clearly presented and easily interpretable improves decision-making and reduces cognitive overload for the overseer.

### 3. TITLE: Automated Alert System with Thresholds

SYSTEM\_REQUIREMENT: The system shall trigger alerts based on predefined conditions and notify the overseer using multiple communication methods.

#### IMPLEMENTATION\_DETAILS:

- \* Trigger alert when roll/pitch exceeds  $\pm 15^\circ$  (Instability detected)
- \* Trigger alert when joint failure lasts  $> 50$  ms (Joint failure—compensating)
- \* Trigger alert when foot contact inconsistency is detected (Slip detected—adjusting)
- \* Trigger alert when linear velocity deviation exceeds 0.2 m/s (Unexpected speed variation)
- \* Trigger alert when linear velocity is below 0.5 m/s
- \* Trigger alert when temperature exceeds an experimental threshold
- \* Trigger alert when current exceeds an experimental threshold
- \* Trigger alert when RL policy confidence score drops below 80%
- \* Trigger alert when unexpected external force is detected
- \* Alerts are presented via UI notifications, auditory signals, and log entries.

RATIONALE: Establishing automated alerts ensures quick responses to anomalies, improving safety and system reliability.

#### 4. TITLE: Human Overseer Training and System Familiarization

SYSTEM\_REQUIREMENT: The system shall provide tools and training resources to ensure that the overseer has adequate expertise before deployment in high-risk environments.

IMPLEMENTATION\_DETAILS:

- \* Pre-deployment training modules on expected vs. anomalous behaviors
- \* Interactive simulations demonstrating RL decision-making
- \* Anomaly detection system highlighting risky policy behaviors
- \* Instant halt button for emergency intervention
- \* Ability to switch robot from autonomous mode to remote-controlled mode.

RATIONALE: Ensuring human overseers are well-trained reduces the risk of misinterpretation and improves the ability to respond to critical situations.

#### 5. TITLE: Demonstration of Expected and Failure Behaviors

SYSTEM\_REQUIREMENT: The system shall present examples of expected, degraded, and failure behaviors in an interpretable manner.

IMPLEMENTATION\_DETAILS:

- \* Color-coded evaluation of key variables during trial runs
- \* Comparison of optimally planned motion trajectory vs. real-time trajectory
- \* Categorization of failures into environmental, mechanical, or policy-related categories.

RATIONALE: Providing clear demonstrations of different behavioral outcomes helps overseers better understand and predict robot performance in various conditions.

#### 6. TITLE: Active Human Involvement in Training Process

SYSTEM\_REQUIREMENT: The system shall allow human overseers to actively participate in the RL training process by providing feedback and curating training data.

IMPLEMENTATION\_DETAILS:

- \* Reward-shaping interface allowing human feedback to adjust policy weighting in real-time
- \* Remote-controlled mode enabling intervention in critical scenarios
- \* Post-run evaluation allowing overseers to flag and annotate undesired behaviors
- \* Scenario generation where operators identify challenging situations to be added to the training curriculum
- \* Training data curation where overseers help label and select high-quality data
- \* Rejection sampling allowing operators to discard specific training episodes containing unwanted behaviors.

RATIONALE: Enabling human oversight in the training process increases transparency and ensures that the RL system aligns with operational expectations.

#### 7. TITLE: Comprehensive Logging of System Interactions

SYSTEM\_REQUIREMENT: The system shall record and log critical interactions, decisions, and system state transitions for review and analysis.

IMPLEMENTATION\_DETAILS:

- \* Time-stamped logs for all state transitions, actions, and outcomes
- \* Policy confidence scores recorded for post-analysis
- \* Video and trajectory logs for external review
- \* Performance statistics including maximum and average values for key variables such as temperature, current, velocity, and tilt.

RATIONALE: Comprehensive logging supports transparency, troubleshooting, and accountability in RL-based decision-making.

#### 8. TITLE: Post-Mortem Analysis and Scenario Replay

SYSTEM\_REQUIREMENT: The system shall enable detailed post-mortem analysis, including the ability to replay scenarios to trace decision-making processes.

IMPLEMENTATION\_DETAILS:

- \* Digital-twin configuration of the robot for complete scenario archiving
- \* Sampling of logs of key sensor inputs, decisions, and corresponding actions to balance feasibility and fairness.

RATIONALE: Providing tools for post-mortem analysis ensures that decision-making can be reviewed, understood, and improved upon.

#### 9. TITLE: Bias Detection in RL Policies

SYSTEM\_REQUIREMENT: The system must include a framework for detecting biases in reinforcement learning policies by analyzing action distributions across different terrains.

IMPLEMENTATION\_DETAILS:

- \* RL policies are stress-tested across diverse terrains including surface topology and texture
- \* RL policies are stress-tested across diverse environmental conditions such as rain, high temperatures, obstructing weather like sandstorms, and nighttime
- \* Domain randomization ensures exposure to different environmental variables
- \* Bias detection framework analyzes action distribution across different terrains.

RATIONALE: To ensure fairness, the system must identify and mitigate any biases in learned policies that may favor certain environments over others.

#### 10. TITLE: Ensuring Equitable Performance Across Different Terrains

SYSTEM\_REQUIREMENT: The reinforcement learning system must ensure that robots perform equitably across diverse terrains without mechanical or physical justification for disparity.

IMPLEMENTATION\_DETAILS:

- \* Robot should maintain similar success rates regardless of terrain type
- \* Reinforcement learning policy should adapt equally well to all five test terrains

- \*\* flat ground

- \*\* memory foam

- \*\* mulch

- \*\* lawn

- \*\* and hiking trails

- \* Develop domain-specific fairness metrics for locomotion measuring stability, energy efficiency, and success rates across environments

- \* Ensure equal training exposure across diverse environments

- \* Conduct regular cross-testing to validate that improvements in one environment do not degrade performance in others
  - \* Include feedback from diverse users deploying the robot in varied environmental contexts.
- RATIONALE: To prevent biases that could lead to systematic underperformance in specific environmental contexts, ensuring fair and reliable operation across different terrains.

#### 11. TITLE: Metrics for Bias Identification in Decision-Making

SYSTEM\_REQUIREMENT: The system must use multiple methods to identify biases in the agent's decision-making process across different environmental contexts.

##### IMPLEMENTATION\_DETAILS:

- \* Perform performance disparity analysis using statistical comparisons of success rates, completion times, and energy efficiency
- \* Implement action preference mapping to identify unwarranted movement biases
- \* Categorize failure instances to detect clustering around specific environmental features
- \* Correlate internal confidence measures with actual performance across conditions
- \* Use counterfactual testing to isolate the impact of environmental features on decision-making
- \* Measure behavioral consistency to assess uniform application of learned policies across different contexts.

RATIONALE: Identifying and quantifying biases ensures that RL policies do not develop environment-specific preferences that could impact fairness and reliability.

#### 12. TITLE: Ensuring Diversity and Representativeness in Training Data

SYSTEM\_REQUIREMENT: The training data must be assessed for diversity and representativeness to prevent biases or blind spots in learned behaviors.

##### IMPLEMENTATION\_DETAILS:

- \* Create detailed taxonomies of environmental conditions and track their coverage in training data
- \* Introduce systematic variations in environmental conditions to test model robustness
- \* Search for and document challenging scenarios at the boundaries of the operational envelope
- \* Implement continuous data collection systems to log novel environmental conditions encountered during deployment
- \* Analyze statistical distributions of key environmental features such as incline angles
- \* obstacle densities
- \* and surface compliance
- \* Set minimum training data requirements across all five terrain types with subcategories for variations within each type.

RATIONALE: Ensuring diverse and representative training data reduces the likelihood of biases and ensures robust, fair decision-making across a variety of environments.