CoRL: Environment Creation and Management Focused on System Integration*

Justin D. Merrick[†] and Benjamin K. Heiner[‡] *Air Force Research Laboratory, Wright-Patterson AFB, OH, 45433*

Cameron Long[§], Brian Stieber[¶] and Steve Fierro[∥] *Toyon Research Corporation, Goleta, CA, 93117*

Vardaan Gangal**

Jacobs, Dayton, OH, 45432

Madison Blake^{††} Heron Systems, Alexandria, VA, 22314

Joshua Blackburn^{‡‡} STR, Beavercreek, OH, 45431

Existing reinforcement learning environment libraries use monolithic environment classes, provide shallow methods for altering agent observation and action spaces, and/or are tied to a specific simulation environment. The Core Reinforcement Learning library (CoRL) is a modular, composable, and hyper-configurable environment creation tool. It allows minute control over agent observations, rewards, and done conditions through the use of easy-to-read configuration files, pydantic [1] validators, and a functor design pattern. Using integration pathways allows agents to be quickly implemented in new simulation environments, encourages rapid exploration, and enables transition of knowledge from low-fidelity to high-fidelity simulations. Natively multi-agent design and integration with Ray/RLLib [2] at release allow for easy scalability of agent complexity and computing power. The code is publicly released and available at https://github.com/act3-ace/CoRL.

I. Introduction

Existing Reinforcement Learning (RL) software libraries can be roughly broken into three categories and their combinations: environment API, algorithm implementation, and integration. Integration handles the communication between the previous two. The introduction of the OpenAI Gym environment API [3] led to a proliferation of software packages utilizing this API [4]. Despite this popularity, the Gym API has significant limitations. Most importantly, it lacks effective tooling for environment creation and alteration and does not natively support environments with multiple agents, which has led to the creation of multi-agent specific environment APIs such as PettingZoo [5]. The Core Reinforcement Learning library (CoRL) straddles the environment API and integration categories, providing a modular, composable, and hyper-configurable suite of tools for environment creation §§. CoRL allows users to rapidly explore complex design spaces with minimal re-integration. The contributions of CoRL are as follows:

- 1) A suite of environment and agent creation tools,
- 2) A hyper-configurable environment enabling rapid exploration for task-based RL,

^{*}Approved for Public Release. Case Numbers: APRS-RYZ-2023-01-00006

[†]Deputy Branch Chief, Autonomous Controls Branch, 21300 8th St.

[‡]Behavior Development Lead AACO, Autonomy Capability Team 3, 2241 Avionics Circle.

[§]Analyst, 6800 Cortona Dr

[¶]Deputy Director, Algorithm/AI Solutions, 6800 Cortona Dr

Senior Analyst, 6800 Cortona Dr

^{**}Junior Autonomy Engineer,1415 Research Park Dr

^{††}Senior Staff Engineer, 209 Madison St.

^{‡‡}Lead Engineer, 2611 Commons Blvd, Ste 150

^{§§}Code available at: https://github.com/act3-ace/CoRL

- 3) A flexible framework for developing these environments and the agents that populate them,
- 4) Tooling to validate configuration files and give useful feedback when files are mis-configured, and
- 5) Integration pathways enabling a dramatic reduction in development time to transition agents into a new simulation.

The remainder of this paper compares CoRL to existing RL libraries (Section II), expands upon the design of CoRL and how it achieves its environment creation capabilities (Section III), presents case studies highlighting the composable and modular nature of CoRL (Section IV), and discusses its potential impacts and future improvements (Section V).

II. Related Works

Gym [3] is perhaps the best known environment API library currently in use. Its quantity and variety of environments and its common sense API have made it ubiquitous. This has also resulted in a large number of RL algorithms and RL algorithm libraries being designed for use with Gym. Stable Baselines [6], Tianshou [7], and RLLib [2] are just a couple of the many libraries which have implementations that are utilize this API. Gym organizes each environment into a single class structure, implementing four main methods and six attributes. Of principle importance are the step method, which takes an action as an argument and runs one time step of the environment's dynamics, returning the next observation, the reward for the time step, whether the episode is complete, and an information dictionary that varies per environment; the reset method, which sets the environment to an initial state and returns an initial observation and an information dictionary; the action_space attribute which describes the type, cardinality, and bounds of the actions required by the environment; and the observation_space attribute, which returns the type, cardinality, and bounds of the observations output by the environment. CoRL uses this interface in its own Environment class and includes full integration with existing environments that use the Gym API. However, CoRL eschews the monolithic class structure for a more modular framework that allows users to easily modify their environments via YAML configuration files. The wrappers used to modify existing Gym environments have a similar function to these configuration management tools in CoRL. However, the modular and configurable nature of observations, rewards, and dones in CoRL allow complex permutations of a single task, and implementations of new tasks, in a way that is quick, readable, and repeatable.

Petting Zoo [5] is another environment API which extends the capabilities of Gym to multiple agents and introduces the Agent Environment Cycle Game in which agents act and receive reward in sequential manner. CoRL is also natively multi-agent and capable of handling agent death and exposing complex reward attribution. Unified Distributed Environment (UDE) [8] is also multi-agent capable and seeks to improve on the modularity of simulation environments through its bridge component. This is similar in concept to the CoRL Simulator class, but UDE does not seek to create a configurable and composable environment and instead is merely providing an integration tool. UDE also has environment virtualization capabilities, similar to RLLib, which CoRL does not seek to emulate. Gymnasium [9] is a fork of Gym with a slightly different API that differentiates between "terminated", meaning an agent met its goal or violated a constraint, and "truncated", meaning an episode met its time limit without reaching its goal or violating a constraint. These semantics are mirrored in CoRL through the use of a DoneStatusCode and a dictionary of Done functors.

The Unity ML-Agents Toolkit [10] spans all three categories of library, but in a more closed ecosystem. It uses a similar structure to and implements many of the features of CoRL. The Agents in the ML-Agents Learning Environment are analogous to CoRL Platforms, while the Behaviors are analogous to instances of the Policy class. Sensors exist in both CoRL and Unity ML-Agents in a very similar implementation. This level of modularity is also extended to the done state and the rewards, allowing access to the rewards and done states of individual agents. The differences between Unity ML-Agents and CoRL lie in the ability of CoRL to use any simulation environment, at the cost of some integration work, and in the Glue structure of CoRL being much more flexible and composable than the Sensor structure in Unity ML-Agents. CoRL also does not attempt to implement any RL algorithms, relying on other libraries to implement those in a manner consistent with the Environment class being used.

The DeepMind Control Suite (dm_control) [11] is a popular environment creation tool that pairs the MuJoCo [12] physics library with a Python front-end that allows user to assemble new environments and tasks in ways very similar to CoRL. The XML specifications of the environment used in dm_control provide precise control over the physics and objects in an environment, while CoRL's YAML configuration files provide the same precise control and composability over the observations, rewards, actions, and dones. dm_control provides similar precise control over rewards as well by the implementation of their tolerance() function. CoRL takes a different approach in implementing some simple reward classes and providing a framework for the user to implement their own reward functions. The API utilized by dm_control is different from that of Gym, which CoRL utilizes. The return from the step function returns information about the type of episode termination, similar to that of Gymnasium, which CoRL implements with its

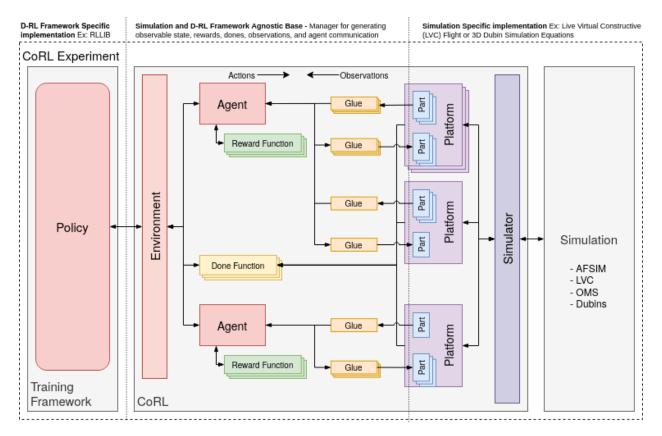


Fig. 1 A high Level view of CoRL showing the default information flows. (from left to right) The Training Framework is external to CoRL and must be compatible with the Environment class. The Environment class houses and manages all other classes. The Agent instances house Reward Functors that utilize the observations provided by Glues to construct the observation passed to the Training Framework. Glues utilize the measurements of the Sensors of Platforms or pass the controls to the Platform's Controller(s). Done Functors evaluate Platform state and report the done status to the Environment. Platforms are managed by the Simulator, which passes data to and from the Simulation, which is external to CoRL.

DoneStatusCode. The focus of dm_control is on the creation of a simulation, while the focus of CoRL is on the management and manipulation of the outputs of that simulation. The slightly different API used by dm_control means that fewer algorithm libraries support its use or do so through the use of wrappers which enforce conformity with the Gym API.

SuperSuit [13] is an integration tool between environments and algorithms that implements a number of more complex wrappers to Gym environments, such as frame stacking or delaying observations. All of these wrappers use the Gym or Gymnasium API and their functions can be or are replicated in CoRL.

III. Design

CoRL is designed to be modular, composable, and hyper-configurable. It accomplishes these goals through the use of pydantic[1] validators which specify the construction of YAML configuration files, a stateful functor design pattern for Glues, Rewards, and Termination/Goal Critiera (i.e. Dones), an Agent and Platform abstraction which allow hierarchies of agents to emerge naturally, and a pluggable simulator interface. An environment class is used to manage each of these constructs and ensure final compatibility with the training algorithm's interface. Fig. 1 shows a high level view of CoRL.

A. Validators and Functors

Each class in CoRL has its own validator. These serve three purposes. The first is to provide typing requirements for the class's attributes. These serve to guide a user in constructing a YAML configuration file for the class. Examples of this pairing between validator and configuration file are shown in Section IV. Second, the validator is used to perform functional checks on the inputs provided. Third, the validators are used to interpret the configuration arguments to generate values and complex datatypes that enable the composability of CoRL Glues, Rewards, and Dones. This last purpose is accomplished through the use of the Episode Parameter Provider (EPP), the Reference Store (Section III.D), and CoRL's unit library. The EPP and Reference store supply values from a distribution to the functors. The unit library automatically converts units based on the configuration arguments. Specifically, the units of a functor must match the type (e.g. velocity, distance) of unit used in any wrapped or referenced classes. In addition to implementing a number of base and derived physical units, CoRL also implements and enforces fractions and percentages, as well as a None unit type.

The stateful functor design pattern used for CoRL's Glues, Rewards, and Dones combines naturally with the validators to allow a single configuration file entry to direct the composition of multiple functors. These create a directed acyclic graph of functors, which is used to reduce the number of instances of identical functors. Additionally, this functor design allows the user to operate at their preferred level of detail. A single functor can be created to do a complex manipulation of platform state, or a chain of functors can be used to accomplish the same calculation, allowing reuse. The statefulness of CoRL functors also allows communication to flow between the three main types of functors (Glues, Rewards, and Dones). Each of these classes and their interactions is explained in greater detail in Sections III.E and III.F.

B. Simulator

The Simulator class allows the CoRL Environment class (Section III.H) to modularly interact with a variety of simulations of varying complexity with minimal reintegration. At its core a Simulator class looks very similar to a Gym Environment, implementing a reset and step method to reset or advance the simulation dynamics; however, the Simulator has attributes for frame_rate, sim_time, and platforms. These attributes allow for the creation, destruction, or unresponsiveness of Platforms (Section III.C) to be communicated to the Environment. The similarity to the Gym API allows for environment dynamics to be implemented directly in the Simulator class for simple problems or where no external simulation is present or necessary.

The minimal reintegration is achieved through a plugin library of Platform Parts. This allows for a rapid and principled transition from virtual simulation, to bench testing, to real-world testing. A Simulator class is defined for a virtual simulation, with a separate Simulator defined for both bench testing and real-world testing. By implementing the same types of Parts for each Simulator, the engineer can be confident that the measurements that the Glues create are consistent across these implementations.

C. Platforms

Platforms are an important piece of the modularity of CoRL. A platform is an abstraction of a piece of the environment that has a measurable or modifiable state. In the simplest implementations, this can be the entire environment. CoRL uses this implementation for default interaction with existing Gym environments. A more physical interpretation might be in an orbital dynamics problem where a platform is defined as a spacecraft. In this example, the Platform class can be used to define the core spacecraft states such as its position and velocity, as well as how to extract them from the Simulator class.

In all cases, Platforms have some combination of Sensors and Controllers, which are collectively called Parts. Parts dictate how the observation or modification of Platform states are achieved, with Sensors observing the state of the platform for downstream use by Glues and Controllers allow modifying the controllable states of the Platform in order to affect its future state, as determined by the dynamics of the Simulator. Fig. 2 shows how this sequence is used in the step method of the environment to construct the observation for a single agent occupying a single platform.

In order to minimize re-integration efforts in the case where the simulation dynamics change or the platform changes, Parts have properties which define the limits, cardinality, and units of the space in which they are valid, similar to the spaces provided in Gym. Additionally, Parts are added to a plug-in library. This library registers a Part with a specific group name and conditions that must be met to invoke the Part. For example, a condition might require a specific Platform type or Simulator class be used. The functionality behind these constructs is that Parts can be added to an Agent's configuration file using the group name. That group name is then referenced with the plug-in library and the

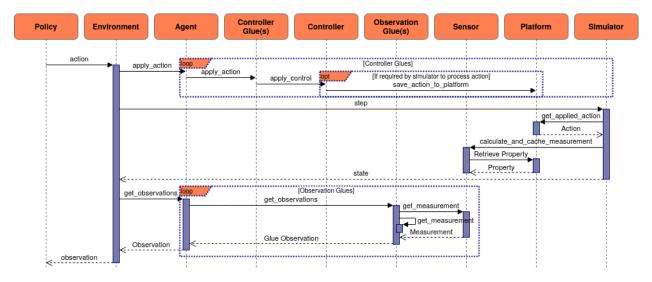


Fig. 2 Sequence diagram for creation of an observation in the step function of CoRL's environment

conditions for using the Parts in that group are checked until a match is found. Then, if the Simulator or Platform is changed, the Agent configuration file need not change, as the plug-in library search will select the appropriate Part to attach to the Platform or Platforms the Agent is attached to. This modularity of Parts also creates a natural pathway to dealing with asynchronous or missing measurement problems in a localized fashion. CoRL Sensors implement a value hold by default in the case of missing or malformed measurements.

D. Episode Parameter Provider and Reference Store

The Episode Parameter Provider (EPP) enables the environment to modify its parameters, such as the initial state, on a per-episode basis throughout the training process. These modifications can happen using two separate mechanisms. First, the parameters returned by the EPP are actually distributions rather than values, which the environment then samples to get the parameter value. This allows the user to interject randomness into the training process by using different values each episode. CoRL natively provides parameters that implement constant, uniform, truncated Gaussian, and discrete choice distributions. Furthermore, there is an abstract interface that allows users to implement other distributions by extending the base parameter class.

Second, the EPP contains hooks that allow it to modify the distribution of parameters as training progresses. When defining the original parameter distributions, the user can specify updaters on the distribution hyperparameters, for example to increment a hyperparameter by a fixed step size. At the end of each training iteration, the environment passes the training result to the episode parameter provider. This allows it to modify its internal data, such as calling the updater on the parameter distributions that it will return to the environment. By creating subclass implementations of the episode parameter provider, the user is able to implement various forms of curriculum learning, such as domain randomization [14, 15] or automatic domain randomization [16].

As an example use of the EPP, consider an extension to Cartpole [17] where the height of the pole could be varied for each episode. In order to solve Cartpole with an arbitrary height, the training might use the episode parameter provider to specify this height as a uniform parameter with some bounds. Furthermore, the upper bound on the distribution could grow larger using an updater as part of a curriculum learning scheme.

The EPP is able to provide parameters to glues, rewards, dones, simulators, and simulator platform initialization. If multiple elements require the same parameter (such as a done termination value and a reward that computes distance from that termination value), the parameter can be put into the Reference Store. Then the relevant objects can all reference the same value after the parameter has been sampled by the environment. The EPP is implemented as a Ray Actor [18, 19], meaning CoRL has a strong dependency on Ray.

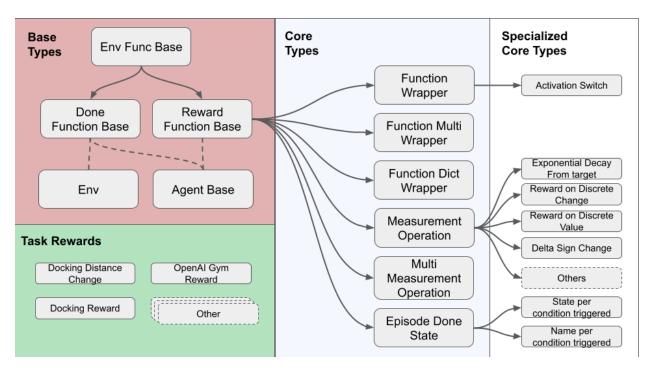


Fig. 3 CoRL Reward types: (1) Base types - general API shared by all rewards, (2) Task Rewards - Specific rewards utilized by the environments in Section IV, (3) Core Types - General task independent rewards, (4) Specialized Core Types - Specific reward implementations that provide shaping or advanced logic for activation.

E. Glues

Glues are a flexible integration class used to communicate and transform information from the Platforms (and their Parts) to the Agents (Section III.G) and vice versa. When communicating from Platform to Agent, the glue instances construct the observation that the Agent will use to calculate the rewards and done, and which the Policy will use to determine what action(s) to take. When communicating from Agent to Platform, the glue instances construct the action that the Platform, and ultimately the Simulator, can use to propagate the dynamics.

Due to their composable nature, Glues can vary greatly in complexity and may use the environment state, platform state and sensors, reference store parameters, and other glues in their calculations. Several base types of Glues are included in CoRL. For the most basic functionality, the ObserveSensor and ControllerGlue Glues connect Platform Sensors and Controllers to the Agent and the Policy. To aid in Glue composition, wrapper Glues that can wrap one or more other Glues and, in the case of wrapping multiple Glues, store them in list or dictionaries, are also included. Other basic mathematical functions such as unit vector transformation, norm calculation, projection, and differencing are included. Examples Glues and their compositions are shown in Section IV.

To take advantage of the directed acyclic graph created by wrapped Glues (and eventually by Rewards and Dones, Section III.F, as well), Glues may optionally use an Extractor class. The extractor simplifies that configuration YAML and is used to topologically sort the dependencies of the CoRL functors. This reduces memory requirements and improves computation speed by eliminating identical instances of a Glue.

F. Rewards and Dones

Dones are a class of functors that test whether their criteria have been met and, if so, mark their scope done and issue a status code to indicate whether this constitutes a win, partial win, draw, partial loss, or loss. The scope of a particular Done can be a single agent or the environment as a whole. These different scopes allow for a variety of done conditions to hold. For instance, the default behavior in CoRL does not end an episode until each agent has had a Done evaluate True. Through a configuration file, this can be changed to end an episode as soon as any agent is done. Furthermore, with environment Dones, called Shared Dones and which have a slightly different call signature, the episode can be terminated regardless of the done status of particular agents.

Similar to Dones, Rewards are a class of functors that, when evaluated, produce a scalar value for a particular

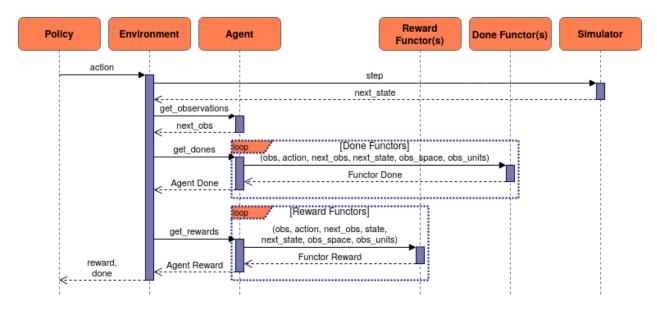


Fig. 4 Sequence diagram for creation of rewards and dones in the step function of CoRL's environment

component of an agent's reward. The rewards for each agent are eventually summed in the Environment to produce the scalar reward that the Gym API expects, but, by calculating the component rewards individually, and reporting them to Tensorboard, analysis of an agent's general behaviors can be conducted analyzed across training, similar to the process proposed in [20]. Rewards are specifically calculated after Dones in order to allow rewards to utilize the done status code created by any Done functor that evaluates to True.

As with all functors in CoRL, both Rewards and Dones are composable and configurable. Fig. 3 shows some examples of Rewards and Dones and how they can compose. Fig. 4 shows a sequence diagram of the processing of the Dones and Rewards as implemented in CoRL's Environment class (Section III.H). Similar to Glues, Rewards and Dones may also take advantage of the graph of functors through the Extractor class.

G. Agents and Policies

The CoRL Agent class implements an agent as described in [21]. Each instance of an agent may partially or fully control any number of platforms and contains a dictionary of Glues, which serve as the agent's sensors and actuators. The instance also contains a dictionary of Rewards, which the agent uses to communicate with a Policy defined in its configuration. The Policy serves to map the Glue observations to an action.

The structure of a Policy class must be dictated by the algorithm implementation library and the environment API that CoRL is integrating with. While integration with RLLib is included in the CoRL release, a Policy implementation is not limited to that library, nor must it be compatible with RL in general. Examples of a random action policy and a scripted policy are included in the CoRL release.

This construction of an agent as a separate entity from a Platform or a Policy allows a natural progression from single-agent to multi-agent environments, both competitive and cooperative. It also allows for shared Policys among Agents, multiple Agents utilizing the same Platform, and a robust means of developing hierarchies of Agents. Furthermore, it minimizes the re-integration necessary for transferring agents from learning environments to inference environments, or from simple dynamics to complex dynamics.

H. Environment

At the opposite end of CoRL from the Simulator is the Environment class, which integrates the Agents, Platforms, Simulator, EPP, Reference Store, Glues, Dones, and Rewards into a cohesive whole that complies with the interfaces expected by the training algorithm. Included in the CoRL release is a multi-agent environment that complies with the the RLLib MultiAgentEnv API, a multi-agent version of the Gym API. While the configuration files specify how the environment is to be created, the Environment class processes these configurations to initialize the component parts of the environment.

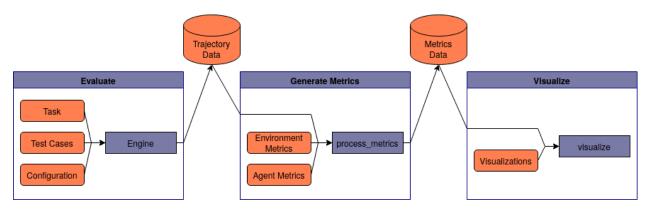


Fig. 5 The Evaluation Framework

The environment is also responsible for maintaining any global Parameters and Dones that are specified in the configuration, such as Parameters that are applicable to multiple agents, or Dones that deal with the overall state of the simulation as opposed to any particular Agent or Platform. A simple example of this is a hard environment horizon. If the episode is terminated after a fixed number of steps, the meta nature of this done condition means it is best placed in the environment. In addition, the Environment also checks the observation(s) obtained from the Agent(s) to ensure that the measurements lie within the space defined by the Glues. This sanity check can be configured to check every step, or it can be used to spot-check observations throughout training.

I. Visualization/Debugging Tools

CoRL, through its Ray/RLLib integration, utilizes Tensorboard [22] to log and track metrics and generate visualizations for them. Within CoRL, Tensorboard has been configured to log and track all rewards, dones, and reference store parameters. Additionally, the DoneStatusCode and simulator reset parameters are also logged. All of these measurements are automatically aggregated to maximum, minimum and mean. All parameters are also logged to a CSV file for ease of visualization outside of Tensorboard.

The configuration of RLLib, the Environment, the Agents, and the Policies are also saved for each trial in order to assist in debugging or determining the differences between runs. These are stored in a JSON file. As some parameters may change over the course of training, especially those modified by an EPP, the information logged to Tensorboard, along with the configuration of the agent and environment for each iteration are logged in a JSON file.

J. Evaluation Framework

After a policy has been trained, the Evaluation Framework can be used to investigate the policy's performance. The evaluation framework follows the same CoRL design practices, utilizing YAML configuration files and a functor design pattern. The framework has been constructed to enable three core processes: Evaluate, Generate Metrics, and Visualize. Each process utilizes separate YAML configuration files to supply the appropriate inputs and each may be run separately or in series as an end-to-end pipeline.

The Evaluate process executes rollouts from a user-defined set of initial conditions. To achieve this, the evaluation framework captures CoRL agent and environment definitions and uses them with an engine to execute rollouts. The engine is specific to the training framework used. At release, CoRL implements an RLLib engine that is compatible with all Policy classes included in CoRL. The trajectory resulting from a rollout is contained in an EpisodeArtifact class which stores agent information such as observations, rewards, dones, and platform state, as well as environment state information. An EpisodeArtifact is generated and saved to disk for each member of the set of initial conditions.

The Generate Metrics process computes metrics from the trajectory data collected and stored to disk by the Evaluate process. Metric classes using the same composable functor design pattern as Glues, Rewards, and Dones define the quantities to be computed. These quantities can be computed directly from a trajectory, can be aggregated across all trajectories, or can be aggregated over quantities generated by other metrics. Some example metrics are the number of times an agent successfully solves the task and the proportion of total reward each reward functor contributed aggregated across a trajectory. The evaluation framework provides a layer of abstraction for the type of quantity returned from a metric functor, categorizing them into two types, terminal and non-terminal. A terminal metric can be categorized

as one which can be represented as a non-container type (e.g. float, string, integer). A non-terminal metric returns a container type (e.g. dictionary, list). This distinction is made to assist in compatibility with non-Python visualizations. A handful of common abstract metrics are implemented at release, but a user can craft more specific metrics to suit their needs.

The Visualize process creates visualizations from the computed metrics. A generic interface is provided that can be implemented to provide a custom visualization. There are two visualizations provided, a visualization that prints the metrics computed in a table to stdout and another that creates HTML plots of the metrics. Multiple visualizations can be included in a single configuration file for the Visualize process.

IV. Case Study

Two case studies are included to illustrate the design principles of CoRL. The first uses the OpenAI Gym CartPole environment as a simulation and uses CoRL's GymSimulator and its associated classes to demonstrate the modularity and composability of CoRL Glues, Rewards, and Dones. The second is an abstraction of the spacecraft docking problem in which the degrees of freedom have been limited to one dimension for simplicity. This case study shows how a simulation environment can be integrated into CoRL without first being set up as a Gym environment.

A. CartPole and OpenAI Gym

Several variations on an agent designed to solve the CartPole-v1 Gym environment are included in CoRL to provide concrete examples of the Gym integration and Glue composability discussed in Section III. The simplest example is shown in Listing 1. Here, a trainable agent (Line 1) is configured with two parts, one sensor and one controller (Lines 3-6). The Sensor_State notation in the configuration file will trigger a search of the plug-in library for a Sensor with that name, which will yield a sensor that reports the Gym environment state as output by the step function. A baseline EPP is included which does not have any parameters (Line 7). Glues are added for each of the parts (Lines 8-24), each specifying which part they are connecting to. The ObserveSensor Glue is also configured to not normalize the observations pulled from the Sensor_State sensor. Finally, Dones and Rewards are added from other files. Both baseline_dones.yml and baseline_rewards.yml contain a simple done and reward that report the done and reward provided by the Gym environment.

From this basic design it is simple to add complexity to the observation. For instance, Listing 2 shows the YAML configuration for adding a Glue which takes the zero-th index of the state and subtracts it from a target value, in this case zero, to return the difference between that state and the target value. Lines 5-8 are functionally identical to Lines 16-20 of Listing 1, but in this case, the TargetValueDifference Glue has wrapped the ObserveSensor Glue in a dictionary with the key "sensor" (Line 3-4). This configuration also includes information on the resulting observation's units ("N/A") and the minimum and maximum measurements expected and allowed, which are used to construct the Glue's observation space (Lines 16-20).

Complexity can also be added to the rewards. Listing 3 shows the configuration for a Reward which uses the Glue from Listing 2 and grants an exponentially decaying reward for how far the observation of the Glue is from a target value. This example makes use of the observation extractors discussed in Section III.F to specify which values, spaces, and units to use (Lines 5-7). In this case, the Glue from Listing 2 is used. While that observation is already a difference to a target value, the modularity of the Rewards and Glues allow the same measurement to be constructed entirely in the reward, or for a different target value than zero to be rewarded. In this case, the default for the target_value field is zero and it is not included. Lines 9-12 specify shaping parameters for the exponential decay as well as how to handle situations where the observation at time t+1 is farther away than the observation at time t.

Other variations on the CartPole agent in Listing 1 are included in the CoRL repository which highlight additional composability in the Glues and Rewards. Additionally, examples of non-operable agents, repeated observations, and agents without the ability to affect the environment are demonstrated. The integration of Gym wrappers and the use of keyword arguments passed to Gym environments is also demonstrated.

B. 1D Docking

The 1D Docking problem is implemented with the same structure as used in [23, 24] but represented as a simple linear ODE,

$$\dot{\mathbf{x}} = A\mathbf{x} + B\mathbf{u} \tag{1}$$

Listing 1 CoRL agent file for CartPole-v1 Gym environment

```
"agent": "corl.agents.base_agent.TrainableBaseAgent"
   "config": {
       "parts": [
           {"part": "Controller_Gym"},
           {"part": "Sensor_State"},
       ],
       "episode_parameter_provider": !include baseline_epp.yml,
       "glues": [
           {
                "functor": "corl.glues.common.controller_glue.ControllerGlue",
                "config": {
11
                    "controller": "Controller_Gym",
                },
13
           },
                "functor": "corl.glues.common.observe_sensor.ObserveSensor",
                "config":{
                    "sensor": "Sensor_State",
                    "normalization": {
                      "enabled": False
                    }
21
                }
           },
23
       ],
       "dones": !include baseline_dones.yml,
25
       "rewards": !include baseline_rewards.yml,
26
   }
27
```

Listing 2 Example of Glue composability using the TargetValueDifference and ObserveSensor Glues

```
{
       "functor": "corl.glues.common.target_value_difference.TargetValueDifference",
2
       "wrapped": {
            "sensor": {
                "functor": "corl.glues.common.observe_sensor.ObserveSensor",
                "config":{
                    "sensor": "Sensor_State".
                    "normalization": {"enabled": False}
           },
       },
11
       "config":{
           "target_value": 0,
13
            "index": 0,
            "unit": N/A,
15
            "limit": {
                "minimum": -5000,
                "maximum": 5000,
                "unit": N/A
            }
       },
21
   }
22
```

Listing 3 Example of Reward complexity using the ExponentialDecayFromTargetValue Reward

```
{
       "name": "OpenAIGymExtractorReward".
2
       "functor":
       "corl.rewards.exponential_decay_from_target_value.ExponentialDecayFromTargetValue",
       "config": {
           "observation": {
               "fields": ["ObserveSensor_Sensor_StateDiff", "direct_observation_diff"]
           "index": 0,
           "eps": 1, # THIS IS RADIANS
           "reward_scale": .000000000001,
10
           "is_closer": true,
           "closer_tolerance": 10,
       }
13
   }
14
```

Listing 4 reset method of 1D Docking Simulator

```
def reset(self, config):
       config = self.get_reset_validator(**config)
2
       self.clock = 0.0
       # construct entities ("Gets the platform object associated with each simulation
        → entity.")
       self.sim_entities = {} # pylint: disable=attribute-defined-outside-init
       for agent_id, agent_config in self.config.agent_configs.items():
           agent_reset_config = config.platforms.get(agent_id, {})
           self.sim_entities[agent_id] = Deputy1D(name=agent_id, **agent_reset_config)
       # construct platforms ("Gets the correct backend simulation entity for each agent.")
11
       sim_platforms = {}
12
       for agent_id, entity in self.sim_entities.items():
13
           agent_config = self.config.agent_configs[agent_id]
14
           sim_platforms[agent_id] = Docking1dPlatform(platform_name=agent_id,

¬ platform=entity, parts_list=agent_config.parts_list)

       self._state = BaseSimulatorState(sim_platforms=sim_platforms, sim_time=self.clock)
17
       self.update_sensor_measurements()
       return self._state
```

where the state $\mathbf{x} = [x, \dot{x}]$, the position and velocity of the docking craft, the control $\mathbf{u} = [T]$, the thrust of the docking craft, and

$$A = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \quad B = \begin{bmatrix} 0 \\ \frac{1}{m} \end{bmatrix},\tag{2}$$

the state and input matrices.

This implementation uses the concept of an Entity class (Deputy1d) which is roughly analogous to a Platform in CoRL and a Dynamics class (Docking1dDynamics) which calculates \dot{x} . The task of implementing this simulation in CoRL then becomes the creation of a Simulator, a Platform, and several Parts. Implementing Docking1dSimulator requires creating concrete methods for several abstract methods defined in the CoRL BaseSimulator, most importantly a reset and a step method. Listing 4 shows the reset method. Some of this method, particularly Lines 2-3 and 17-19, will be typical of most CoRL Simulator classes. Lines 6-10, however, interface directly with the simulation by constructing and storing the Entity class(es) associated with the simulation while Lines 12-15 pair the Entity class(es) with CoRL Platform(s), in particular the Docking1dPlatform created for this environment.

The Docking1DPlatform requires a Deputy1D Entity in its validator, and uses the properties of this Entity as its own properties to pass along the state of the environment. Docking1DPlatform also stores the actions given to it as shown in Fig. 2. This allows the step function of the Simulator class, shown in Listing 5 to retrieve these applied actions (Line 3), apply them to the appropriate Entity classes (Lines 4-5) and then update the sensor measurements on each Platform with the updated state of the Entity (Line 7).

Sensor measurements are calculated from Sensor classes which access the sensor and velocity properties of their parent Platform. These sensors are added to the plug-in library and associated with Docking1dSimulator and Docking1dPlatform. Finally, a Controller is also created which wraps the Docking1dPlatform methods to save and get actions to the abstract Controller class methods. This Part is similarly added to the plug-in library.

Now a working interface between the simulation and the CoRL Environment and Agent classes has been made. It can be configured with YAML files similar to those in Section IV.A. A baseline configuration is included in the CoRL release. More detail on the Simulator, Platform, Part interfaces can be found in the source code or documentation (Footnote §§).

Listing 5 step method of 1D Docking Simulator

```
def step(self):
    for agent_id, platform in self._state.sim_platforms.items():
        action = np.array(platform.get_applied_action(), dtype=np.float32)
        entity = self.sim_entities[agent_id]
        entity.step(action=action, step_size=self.config.step_size)
        platform.sim_time = self.clock
        self.update_sensor_measurements()
        self.clock += self.config.step_size
        self._state.sim_time = self.clock
        return self._state
```

V. Conclusions

CoRL is a modular, composable, and hyper-configurable environment creation tool. While some existing tools rely on monolithic implementations, CoRL allows for minute control of an agent's observations and actions with modification of YAML configuration files. Other tools are tied to specific simulation environments, CoRL allows for new simulations to be integrated with minimal difficulty. Still others are tied only to a reinforcement learning paradigm, CoRL can be used with RL agents trained using any framework, but can also utilize the distributed computing of Ray to evaluate non-RL agents with ease.

At release, CoRL provides useful examples and common implementations of many of its classes, but can be easily extended to more advanced features. Section III.D discussed using the EPP to do curriculum learning, but that framework could be expanded further to include automatic domain randomization. The EPP could also be upgraded to allow the reference store to provide information to any object in the environment, currently it only allows for functors to reference it. The existing unit conversion system is very rigid. Using an existing unit conversion library could allow users to define their own units. RLLib [2] is incompatible with the DeepMind Control Suite [11], but implementing a Simulator, Platform, and basic Glues could allow its popular baseline tasks to be manipulated by CoRL. CoRL is being actively maintained and used for a variety of research tasks at the Air Force Research Laboratory.

References

- [1] Covin, S., "pydantic,", 2021. URL https://github.com/pydantic/pydantic.
- [2] Liang, E., Liaw, R., Nishihara, R., Moritz, P., Fox, R., Goldberg, K., Gonzalez, J., Jordan, M., and Stoica, I., "RLlib: Abstractions for Distributed Reinforcement Learning," *Proceedings of the 35th International Conference on Machine Learning*, Proceedings of Machine Learning Research, Vol. 80, edited by J. Dy and A. Krause, PMLR, 2018, pp. 3053–3062. URL https://proceedings.mlr.press/v80/liang18b.html.
- [3] Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W., "OpenAI Gym,", 2016.
- [4] "OpenAI/Gym: Dependency Graph,", 2022. URL https://github.com/openai/gym/network/dependents.
- [5] Terry, J. K., Black, B., Grammel, N., Jayakumar, M., Hari, A., Sulivan, R., Santos, L., Perez, R., Horsch, C., Dieffendahl, C., Williams, N. L., Lokesh, Y., Sullivan, R., and Ravi, P., "PettingZoo: Gym for Multi-Agent Reinforcement Learning," arXiv preprint arXiv:2009.14471, 2020.
- [6] Raffin, A., Hill, A., Gleave, A., Kanervisto, A., Ernestus, M., and Dormann, N., "Stable-Baselines3: Reliable Reinforcement Learning Implementations," *Journal of Machine Learning Research*, Vol. 22, No. 268, 2021, pp. 1–8. URL http://jmlr.org/papers/v22/20-1364.html.
- [7] Weng, J., Chen, H., Yan, D., You, K., Duburcq, A., Zhang, M., Su, Y., Su, H., and Zhu, J., "Tianshou: A Highly Modularized Deep Reinforcement Learning Library," *Journal of Machine Learning Research*, Vol. 23, No. 267, 2022, pp. 1–6. URL http://jmlr.org/papers/v23/21-1127.html.
- [8] La, W. G., Muralidhara, S., Kong, L., and Nichat, P., "Unified Distributed Environment,", 2022.
- [9] The Farama Foundation, "Gymnasium,", 2022. URL https://github.com/Farama-Foundation/Gymnasium.

- [10] Juliani, A., Berges, V.-P., Teng, E., Cohen, A., Harper, J., Elion, C., Goy, C., Gao, Y., Henry, H., Mattar, M., and Lange, D., "Unity: A General Platform for Intelligent Agents,", 2018. https://doi.org/10.48550/ARXIV.1809.02627, URL https://arxiv.org/abs/1809.02627.
- [11] Tunyasuvunakool, S., Muldal, A., Doron, Y., Liu, S., Bohez, S., Merel, J., Erez, T., Lillicrap, T., Heess, N., and Tassa, Y., "dm_control: Software and tasks for continuous control," *Software Impacts*, Vol. 6, 2020, p. 100022.
- [12] Todorov, E., Erez, T., and Tassa, Y., "MuJoCo: A physics engine for model-based control," 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems, IEEE, 2012, pp. 5026–5033. https://doi.org/10.1109/IROS.2012.6386109.
- [13] Terry, J. K., Black, B., and Hari, A., "SuperSuit: Simple Microwrappers for Reinforcement Learning Environments," *arXiv* preprint arXiv:2008.08932, 2020.
- [14] Tobin, J., Fong, R., Ray, A., Schneider, J., Zaremba, W., and Abbeel, P., "Domain randomization for transferring deep neural networks from simulation to the real world," 2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), 2017, pp. 23–30. https://doi.org/10.1109/IROS.2017.8202133.
- [15] Peng, X. B., Andrychowicz, M., Zaremba, W., and Abbeel, P., "Sim-to-Real Transfer of Robotic Control with Dynamics Randomization," 2018 IEEE International Conference on Robotics and Automation (ICRA), 2018, pp. 3803–3810. https://doi.org/10.1109/ICRA.2018.8460528.
- [16] OpenAI, Akkaya, I., Andrychowicz, M., Chociej, M., Litwin, M., McGrew, B., Petron, A., Paino, A., Plappert, M., Powell, G., Ribas, R., Schneider, J., Tezak, N., Tworek, J., Welinder, P., Weng, L., Yuan, Q., Zaremba, W., and Zhang, L., "Solving Rubik's Cube with a Robot Hand,", 2019. https://doi.org/10.48550/ARXIV.1910.07113, URL https://arxiv.org/abs/1910.07113.
- [17] Barto, A. G., Sutton, R. S., and Anderson, C. W., "Neuronlike adaptive elements that can solve difficult learning control problems," *IEEE Transactions on Systems, Man, and Cybernetics*, Vol. SMC-13, 1983, pp. 834–846.
- [18] Team, R., "Ray v2 Architecture," Tech. rep., Ray Project, October 2022. Available at https://docs.google.com/document/d/1tBw9A4j62ruI5omIJbMxly-la5w4q_TjyJgJL_jN2fI/preview#.
- [19] Moritz, P., Nishihara, R., Wang, S., Tumanov, A., Liaw, R., Liang, E., Elibol, M., Yang, Z., Paul, W., Jordan, M. I., et al., "Ray: A distributed framework for emerging {AI} applications," 13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18), 2018, pp. 561–577.
- [20] MacGlashan, J., Archer, E., Devlic, A., Seno, T., Sherstan, C., Wurman, P. R., and Stone, P., "Value Function Decomposition for Iterative Design of Reinforcement Learning Agents," arXiv preprint arXiv:2206.13901, 2022.
- [21] Russell, S., Russell, S., Norvig, P., and Davis, E., *Artificial Intelligence: A Modern Approach*, Prentice Hall series in artificial intelligence, Prentice Hall, 2010. URL https://books.google.com/books?id=8jZBksh-bUMC.
- [22] Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and Zheng, X., "TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems,", 2015. URL https://www.tensorflow.org/, software available from tensorflow.org.
- [23] Cunningham, J., Ravaioli, U. J., Dunlap, K., and Hobbs, K. L., "safe-autonomy-dynamics,", 2022. URL https://github.com/act3-ace/safe-autonomy-dynamics.
- [24] Ravaioli, U. J., Cunningham, J., McCarroll, J., Gangal, V., Dunlap, K., and Hobbs, K. L., "Safe Reinforcement Learning Benchmark Environments for Aerospace Control Systems," 2022 IEEE Aerospace Conference (AERO), 2022, pp. 1–20. https://doi.org/10.1109/AERO53065.2022.9843750.