

SHARPIE: A Modular Framework for Reinforcement Learning and Human-AI Interaction Experiments

Anonymous authors
Paper under double-blind review

Keywords: Human-in-the-loop, Interactive RL, Hybrid Intelligence, Experimental Platforms

Summary

Reinforcement Learning (RL) offers a general approach for modeling and training AI agents, including human-AI interaction scenarios. In this paper, we propose SHARPIE (Shared Human-AI Reinforcement Learning Platform for Interactive Experiments) to address the need for a generic framework to support experiments with RL agents and humans. Its modular design consists of a versatile wrapper for RL environments and algorithm libraries, a participant-facing web interface, logging utilities, deployment on popular cloud and participant recruitment platforms. The platform is based on a generic interface for human-RL interactions that aims to standardize the field of study on RL in human contexts.

Contribution(s)

1. We present a first-of-its-kind, general-purpose, standardized open-source experimentation platform for research on human-centered reinforcement learning, bridging the gap in existing work by supporting simultaneous multi-agent and multi-human experimentation with configurable communication channels in various modalities.
Context: Prior frameworks provide task-specific end solutions rather than configurable modules, and mainly support either multi-agent or multi-human scenarios, and restrict the type of possible interactions among them.
2. The framework is generic, modular and applicable to a wide range of RL algorithms, environments, and interaction paradigms at a limited additional implementation cost using only minimal wrapper extensions to Gymnasium environments. Moreover, it seamlessly integrates with crowdworker platforms.
Context: The modularity allows using the framework for research on a wide range of research questions. Crowdworker platform integration has proven invaluable in non-RL machine learning domains but has so far seen lacking support in RL platforms involving humans.
3. We evaluate the flexibility and functional capabilities of the platform using four representative case studies that illustrate different ways humans can be integrated into the RL process: through evaluative feedback, action intervention, environment modification, and with natural language instructions.
Context: These representative case studies validates the high compatibility of the versatile wrapper that our platform provides with diverse RL algorithms and environments.
4. We conduct several performance and scalability tests to demonstrate that our framework supports seamless distributed interaction across geographically distant human participants, with end-to-end latency within 0.21 seconds.
Context: Low-latency communication is paramount to several real-time synchronized collaborative human-RL tasks, but challenging due to fundamental networking limitations.

SHARPIE: A Modular Framework for Reinforcement Learning and Human-AI Interaction Experiments

Anonymous authors

Paper under double-blind review

Abstract

1 Reinforcement Learning (RL) offers a general approach for modeling and training AI
2 agents, potentially including human-AI interaction scenarios. In this paper, we pro-
3 pose SHARPIE (Shared **H**uman-AI Reinforcement Learning Platform for **I**nteractive
4 **E**xperiments) to address the need for a generic framework to support experiments with
5 RL agents and humans. Its modular design consists of a versatile wrapper for RL en-
6 vironments and algorithm libraries, a participant-facing web interface, logging utili-
7 ties, deployment on popular cloud and participant recruitment platforms. It empowers
8 researchers to study a wide variety of research questions related to the interaction be-
9 tween humans and RL agents, including those related to interactive reward specification
10 and learning, learning from human feedback, action delegation, preference elicitation,
11 user-modeling, and human-AI teaming. The platform is based on a generic interface
12 for human-RL interactions that aims to standardize the field of study on RL in human
13 contexts. To demonstrate ease-of-use and versatility, we replicate four illustrative ex-
14 periments from the multi-agent RL literature and report on the platform’s performance
15 under stress testing to showcase its robustness.

16 **1 Introduction**

17 Reinforcement learning (RL) agents learn through environmental interaction to maximise cumula-
18 tive long-term reward (Sutton & Barto, 1998). In recent years, this paradigm has expanded beyond
19 pure autonomous discovery to include scenarios involving humans, under the umbrella term of Re-
20 inforcement Learning from Human Feedback (Kaufmann et al., 2025). As available compute in-
21 creases, RL methods are being deployed in increasingly complex domains where humans are central
22 to the loop. RL agents interact with humans in a wide variety of ways: in some problem settings,
23 RL agents observe humans to achieve a common goal or act on their behalf (Natarajan et al., 2010);
24 in other settings, RL agents communicate with humans to provide services or support their decision
25 making (Zhao et al. (2021)); in other settings, human feedback can also be used to evaluate the value
26 of agent behavior (Knox & Stone, 2009; Christiano et al., 2017), to communicate preferences (Bryk
27 et al., 2022), to correct behavior (Saunders et al., 2018), and to guide exploration (Guan et al., 2021;
28 Torne et al., 2023).

29 These different types of interaction highlight the need to incorporate humans in the training and
30 evaluation of RL agents for various aims including efficiency, safety, personalization and alignment
31 purposes. However, implementing these experiments in a standardized manner presents specific in-
32 frastructural challenges. Connecting local RL training loops to human participants requires building
33 interfaces, managing state synchronization, and implementing data logging services. Consequently,
34 current research on RL involving humans frequently relies on custom, ad-hoc point-solutions engi-
35 neered for individual studies. While popular RL libraries such as Gymnasium (Towers et al., 2024),
36 PettingZoo (Terry et al., 2021), JaxMARL (Rutherford et al., 2024), and MO-Gymnasium (Alegre
37 et al., 2022) provide established benchmarks in their related domains, experimentation involving

38 humans in RL has thus far seen a lack of standardization and has relied on study-specific platforms
39 and tools. While some efforts have been made, these have been lacking in their integrative scope,
40 and have suffered from support challenges (Taylor et al., 2023).

41 To address this infrastructural gap, we propose SHARPIE (Shared Human-AI Reinforcement-
42 Learning Platform for Interactive Experiments). This platform provides a reusable toolset designed
43 to standardize the design, deployment, and analysis of human-RL interactions. By moving away
44 from rigid point-solutions and adopting a modular platform, researchers can more systematically
45 develop more experimental setups. What is more, by first-class support for multi-agent and multi-
46 human interaction scenarios, we hope to facilitate the study of interactions beyond the basic and
47 unidirectional patterns we largely see today (Baraka et al., 2026) in order to study effectively ad-
48 dress problems which humans nor agent can solve independently (Akata et al., 2020; Carroll et al.,
49 2019).

50 The proposed platform provides versatile wrappers for popular single-agent RL and multi-agent RL
51 environments and algorithms, supports configurable communication channels between humans and
52 RL agents in various modalities, and offering logging services, deployment utilities, and participant
53 recruitment platforms integration. The platform aims to empower researchers to address research
54 questions on the interaction of RL agents and humans, focusing on interactive reward specifica-
55 tion and learning, learning from human feedback, action delegation, preference elicitation, user-
56 modeling, and human-AI teaming. By providing a technical platform in which these paradigms can
57 be easily studied, compared, and combined, we hope to increase rigor and contribute to setting the
58 standard for human-RL experiments.

59 Our main contributions in this work are outlined below:

- 60 • SHARPIE introduces a potential standard for RL-based human-agent interactions in a multi-
61 human/multi-agent setting in the same way that the Gymnasium API has become this standard
62 for fully simulated RL environments.
- 63 • SHARPIE supports experiments with diverse RL applications via minimal Gymnasium extensions
64 and native integration of crowdworker platforms.
- 65 • four representative case studies are benchmarked to illustrate different modalities human partici-
66 pants can be integrated into the RL learning process.
- 67 • Several performance and scalability tests to demonstrate that our framework supports seamless
68 distributed interaction across different numbers and geographically distant human participants,
69 with end-to-end latency within 0.21 seconds.

70 The rest of the manuscript is organized as follows. Section 2 summarizes the related work for
71 the study. Design principles of SHARPIE framework are explained in Section 3. While Section
72 4 provides the qualitative evaluation of SHARPIE with 4 use cases, quantitative analysis for its
73 performance is presented in Section 5. Finally, the implications of the study along with possible
74 future research directions are discussed in Section 6.

75 2 Related Work

76 There are several software platforms for conducting behavioral experiments (Peirce, 2007;
77 De Leeuw, 2015). None of these offers any infrastructure for any type of artificially intelligent
78 entities, whether learning or not. Various platforms and software packages have been developed for
79 the deployment of RL (Gauci et al., 2019; Albers et al., 2022; Zhu et al., 2024). However, these are
80 not focused on controlled experiments on the interaction between RL agents and humans.

81 Several successful platforms for the development of human-in-the-loop RL algorithms and evalua-
82 tion environments have been made over the past decades, emphasizing a different number of objec-
83 tives, different numbers of agents, different kinds of tasks and different programming languages. We
84 do not intend to review all software packages in RL and here we focus only on the related work that

85 particularly deals with RL platforms that explicitly target the study of RL in a context that involves
86 humans.

87 With recent advances in Reinforcement Learning from Human Feedback (RLHF), platforms such
88 as (Christiano et al., 2017; Yuan et al., 2024) have been developed. However, the interfaces for
89 human interaction provided by these platforms are only capable of handling restricted feedback, and
90 treat humans mostly as a source of information rather than an interaction partner. Along similar
91 lines, Qian et al. (2025) proposed UserRL, an approach for training agentic large language models
92 based on human input with RL in various settings. This work is limited to a *one-to-one* text-based
93 interaction with LLMs whereas we focus on the multi-human multi-agent experiments involving
94 various modalities, including reward annotation, action guidance, as well as textual communication.

95 With a specific focus on user studies, Knierim et al. (2024a) presents a framework in human teach-
96 ing RL agents setting with audio communication. Their codebase also facilitates a multi-participant
97 remote Wizard-of-Oz approach (Green & Wei-Haas, 1985), where one participant can play the role
98 of a teacher and another plays the role of a learning agent, addressing additional requirements for
99 early iterative testing of interactive algorithmic settings, where usability and human behavior are
100 used as a baseline, target, or data collection / calibration step. From an algorithmic perspective,
101 although an alternative modality is proposed in this study (namely prosodic voice cues), the inter-
102 actions are designed to be unidirectional, with human teachers guiding RL-agents with a restricted
103 communication modality.

104 Closer to the scope of SHARPIE, the HIPPO-GYM platform (Taylor et al., 2023) provides a frame-
105 work that involves human interactions in a RL setting. While similar in its main objectives as well
106 as the needs it addresses, this framework is limited to single-agent tasks and provides only specific
107 interactions in those tasks (i.e. humans teaching RL agents through various forms of feedback).

108 On the multi-agent front, previous studies in multi-agent learning such as MAgent (Zheng et al.,
109 2018) and MOMALand (Felten et al., 2024) are capable of handling multi-agent experiments. MA-
110 gent supports a large population of agents, ranging from hundreds to millions, whereas it is possible
111 to train the agents for multiple objectives in MOMALand. However, these packages do not yet sup-
112 port human interaction in any way. SHARPIE specifically targets the gap of easily setting up exper-
113 iments with (multiple) human participants and RL agents by interfacing with existing environment
114 and algorithm libraries.

115 Finally, the closest work of SHARPIE that we have identified in the existing literature is Interactive
116 Gym (McDonald, 2024) which is a framework in early stage of development that enables human-
117 AI interaction in multi-agent RL settings. This framework currently provides two different setups:
118 single human with multiple agents in a browser based execution and any number of human and AI
119 agents in a client-server based execution of experiments. With its current structure, Interactive Gym
120 seems to be specific to the Gymnasium environments where human(s) can take the control of RL
121 agents in a game setup which can require more effort to deal with a diverse set of scenarios including
122 various type of human-RL agent interactions. SHARPIE, on the other hand, aims to provide a more
123 generic and extendable design that is able to efficiently facilitate such scenarios.

124 All in all, there are several existing approaches and subsequently experimental platforms that have
125 looked at specific restricted human-RL scenarios, but do not focus on offering the platform as a ver-
126 satile tool for research in this area. Additionally, we identified a few similar platforms that partially
127 overlap with the aims and scope of SHARPIE, and offer a solid ground to build a unified platform.
128 We envision SHARPIE will empower human-AI researchers to deploy interactive agents and col-
129 lect data on a wide variety of human-in-the-loop RL scenarios in *rich, real-time and potentially*
130 *multi-agent* settings.

131 **3 SHARPIE Framework**

132 SHARPIE is a Python-based web framework that aims to provide a versatile wrapper around pop-
133 ular Reinforcement Learning a) environments, b) algorithms, and c) methodologies. For the first

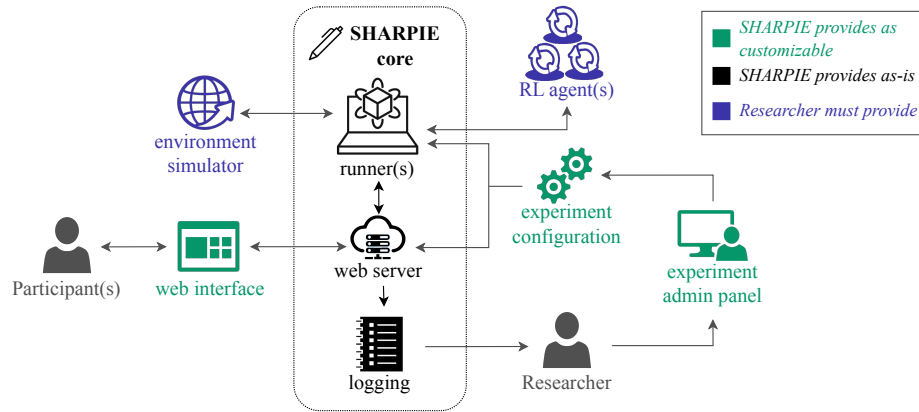


Figure 1: High-level SHARPIE architecture.

134 part, it can encapsulate any existing environment that follows the conventional Gymnasium API
 135 (reset, step, render, etc.) which encompasses most of the existing RL platforms (see Appendix A).
 136 Most importantly, the ambition of SHARPIE is not to be tightly integrated with any one particular
 137 environment or library, but rather to be compatible with many.

138 As illustrated in Figure 1, the library can be decomposed into three main families of modules. In
 139 black are the built-in components, responsible for scheduling episodes and communication between
 140 agent(s), human(s), the learning algorithm(s) and the environment. These provide the backbone of
 141 the library and are largely transparent to any but the most curious of users. In green are components
 142 that may be subjected to modification to suit a particular experiment, for instance, with arrow keys
 143 and dynamic performance feedback or text-based instructions followed by colour-coded answers.
 144 Finally, the parts of SHARPIE that the experimenters will be most manipulating are depicted in
 145 blue. These correspond mainly to the agents and their environment, although in most cases a simple
 146 wrapper around a gym- or pettingzoo-compatible environment is enough.

147 In terms of customizable components, the front-end User Interface (UI) is web-based and multi-
 148 modal. This allows human users to interact with the environment and other (RL or human) agents.
 149 Additionally, RL agents may be able to request action delegation, explicitly or as prompted by their
 150 learning algorithm, to further increase the range of interaction scenarios that SHARPIE can help
 151 streamline.

152 The front-end also provides various complementary utilities to further smoothen out the experimen-
 153 tal processes: from (a)synchronous evaluations of a learning agent to scheduling and management
 154 on long-term data storage. In the first case, this takes many forms, such as a preference elicitation
 155 module to visualize and rank trajectories, actions, and policies. In the second case, this includes all
 156 the logging facilities required in a large-scale RL experiment, especially one involving human par-
 157 ticipants. Indeed, it is of paramount importance to be able to securely and robustly store any relevant
 158 data, as restarting such an experiment from scratch might be, at best, long and costly and, at worse,
 159 practically impossible. To this end, we devised a data model capable of handling the minutiae of
 160 multiple human studies by linking experiments, agents and time steps in a single SQL database (see
 161 Figure 3).

162 In addition, SHARPIE is designed to support multi-modal communication channels to allow com-
 163 munication between agents (again, RL or human). Depending on said agents, these channels may
 164 be used for e.g. coordination, teaching, or overriding. At its core, SHARPIE is providing json-like
 165 communication channels allowing experiments to transmit any type of information between agents,
 166 environments and the learning algorithm. Considering the most common types of data being used
 167 in RL research some types are handled natively by the library. Heavy types such as large byte ar-

Table 1: Use cases included in functional validation. The last columns indicate the numbers of lines of code (LoC) and database (DB) entries to integrate with SHARPIE.

Algorithm	Paradigm	Description	LoC	DB
SayCan (Ichter et al., 2023)	1 messages	Language instructions, learned affordances	17	5
Dagger (Ross et al., 2011)	2 actions	Queries demonstrations for visited states	65	5
SP/BC (Carroll et al., 2019)	3 environment	Plays cooperatively with a human	29	6
TAMER (Knox & Stone, 2009)	4 rewards	Learns from human evaluative feedback	31	5

168 rays are additionally subjected to compression to reduce the overhead brought about by network
 169 communication, in case the participant has limited bandwidth.

170 Finally, the library is designed with utilities to deploy to a cloud server, a private machine, or a local
 171 host. Given the widespread support for Python, this encompasses any desktop operating system
 172 (e.g. Linux, Mac, Windows) with Python 3 installed, including remote web-servers or platforms
 173 such as Amazon Web Services (AWS). This provides an abstraction between the machines that the
 174 researchers and participants are using and where the experiment is actually running. Such a feature
 175 is essential in studies involving numerous participants with varying locations and hardware. Dele-
 176 gation is done through the runners, allowing replicates of an experiment to be performed in parallel
 177 by instantiating multiple such runners on any machine from local servers to service providers.

178 4 Functional Validation via Case Studies

179 We validate the expressiveness and flexibil-
 180 ity of the proposed framework by replicat-
 181 ing multiple interactive scenarios involving RL
 182 agents and human participants. To evaluate
 183 whether the framework is sufficiently expres-
 184 sive to implement a wide range of such scen-
 185 arios, we consider a taxonomy of interaction
 186 paradigms. Figure 2 displays the different inter-
 187 action paradigms based on how humans, agents
 188 and the environment can interact with one an-
 189 other.

190 In the first interaction paradigm, RL agents and
 191 humans can send *messages*. These messages
 192 can be used to give instructions, ask for help,
 193 explain the agent internal state, give action ad-
 194 vice and can have various modalities including
 195 text, images, etc. In the second paradigm, hu-
 196 mans can select and override RL agent actions.
 197 This paradigm includes shared control, learning from demonstrations etc. In the third paradigm,
 198 humans act as agents in the environment. This includes shared tasks, human-RL-teaming etc. In the
 199 fourth paradigm, the human specifies or alters an environment reward to model various evaluative
 200 feedback and reward shaping setups. As a final note on these paradigms, we want to highlight that
 201 they are not mutually exclusive and are often combined.

202 Table 1 lists a representative study that involves human participants for each of the above-mentioned
 203 paradigms. We evaluate SHARPIE’s functional expressiveness by using it to implement a replication
 204 for each of these studies. To assess the framework’s flexibility, we record the number of lines of code
 205 (LoC)¹ and the number of experiment configuration database (DB) entries. In our analysis, we focus
 206 on the overhead imposed by creating a human-participant study and therefore only include code and

¹Excluding docstrings, comments, and empty lines

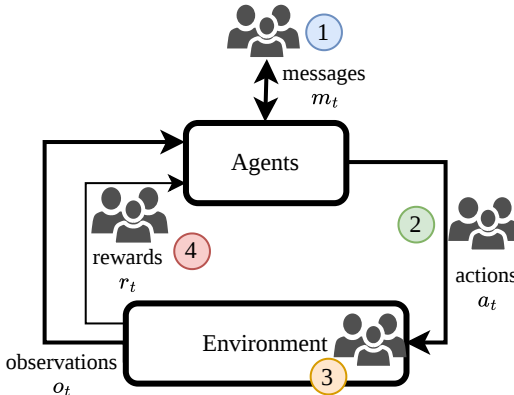


Figure 2: Taxonomy of agent-human-environment interactions.

207 configurations to correctly route policy decisions, participant inputs, agent outputs, and environment
208 renderings to the appropriate location. That is, we explicitly do not incorporate agent training and
209 environment simulation logic as this is out of scope for our framework. The remainder of this section
210 consists of detailed descriptions of these replications per use case.

211 **SayCan** by [Ichter et al. \(2023\)](#) implements the first paradigm, where a human participant and
212 RL agents interact by exchanging messages containing instructions. Specifically, the human user
213 provides high-level, unstructured natural language instructions, which the agent then grounds into
214 executable skills using learned affordances. The primary challenge in this paradigm is message
215 passing: standard RL environments and agent APIs typically expect fixed-dimensional numerical
216 vectors (e.g., standard Gym spaces). Routing variable-length string payloads from a user-facing
217 web interface to the agent’s policy typically requires custom serialization and custom API wrappers.
218 SHARPIE natively abstracts this exchange. We successfully connected a participant-facing text-
219 input interface to a SayCan policy interacting with a table top pick-and-place environment using
220 only 17 LoC and 5 DB entries. This demonstrates SHARPIE’s capability to support experiments
221 involving messages between humans and RL agents.

222 **Dagger** ([Ross et al., 2011](#)) represents the second interaction paradigm, where a human expert
223 provides corrective action labels for states visited by the agent’s current policy. To evaluate the
224 framework’s support for this paradigm, we implemented an early DAgger-based imitation learning
225 experiment using the *Super Mario Bros* environment ([Kauten, 2018](#)). A substantial proportion of the
226 43 LoC implementation deals with migrating the environment from an old version of OpenAI-Gym
227 library and translating participant inputs to the environment’s action space. While the environment
228 uses a bitmask-based input system to represent concurrent button presses (e.g., ‘right’ + ‘jump’),
229 we found that mapping these inputs to the agent’s action space is straightforward using SHARPIE
230 utilities. At the same time, SHARPIE’s policy interface enables the interception of the control loop
231 to implement DAgger’s “query” phase for the expert policy. Here, our implementation routes the
232 current observation to the participant, while the update method ensures that the resulting transition
233 is correctly captured and buffered for the training phase. Overall, complete integration is done with
234 65 LoC and 5 DB entries.

235 **Overcooked** exemplifies paradigm three, where humans and RL agents act simultaneously within
236 a shared environment. In the original study, a human participant plays cooperatively with an agent
237 trained via either Self-Play or Behavioral Cloning (SP/BC) in the Overcooked environment ([Carroll
238 et al., 2019](#)). The primary infrastructural challenge in this paradigm is real-time synchronization:
239 the system must reliably stream environment renderings to the participant-facing web interface,
240 capture participant keystrokes, query the RL policy, merge these separate actions, and transition the
241 environment. Using SHARPIE, we replicated this functionality using 29 LoC and 6 DB entries,
242 which highlights SHARPIE’s capacity to support experiments in which users and agents operate in
243 realtime in the same environment. Assuming access to the original agents, which currently rely on
244 deprecated code², one could fully replicate the user study performed in the original paper or even
245 scale it up through crowd-sourcing.

246 **TAMER** represents the fourth paradigm, where human participants specify an immediate reward
247 via evaluative feedback. In the TAMER framework ([Knox & Stone, 2009](#)), the participants observe
248 the agent act and provide an asynchronous positive or negative label to its decisions. These labels are
249 interpreted as immediate rewards and used to train a policy during interactions with the participant
250 and the environment. Implementing this in a standard RL codebase often requires building com-
251 plex threading or networking logic to inject asynchronous UI inputs into the agent’s synchronous
252 training loop. SHARPIE’s functionalities reduce this overhead significantly. We successfully cap-
253 tured and routed the human’s real-time evaluative feedback for the MountainCar environment to the
254 agent’s reward signal using only 31 LoC and 5 DB entries. This demonstrates that SHARPIE can

²[https://github.com/HumanCompatibleAI/overcooked_ai?tab=readme-ov-file#
deprecated-behavior-cloning-and-reinforcement-learning](https://github.com/HumanCompatibleAI/overcooked_ai?tab=readme-ov-file#deprecated-behavior-cloning-and-reinforcement-learning)

255 elegantly handle the timing and synchronization challenges inherent to interactive reward shaping at
256 low implementation cost.

257 5 Scalability and Performance Profiling

258 We now evaluate SHARPIE’s scalability and performance characteristics to ensure it meets the rig-
259 orous demands of interaction studies involving possibly multiple agents and humans. In our bench-
260 marks we include four aspects that might affect reliability of experiment results. First, we assess the
261 scalability wrt the number of human participants interacting with a single environment simultane-
262 ously to give an idea on the overhead imposed by the SHARPIE framework. Second, we assess the
263 scalability with respect to the number of RL agents to assess the degree to which SHARPIE can
264 support experiments involving multiple RL agents. Third, we assess performance in geographically
265 distributed experimentation scenarios. Although SHARPIE’s architecture supports geographically
266 distributed experiments due to its client-server setup with the environment simulation performed at
267 a single server, this architecture may negatively affect performance in scenarios where participants’
268 connection is of low quality. Finally, we test how the framework is affected by environments with
269 reasonably high computational demands, including environments producing high-resolution render-
270 ings. We detail these scenarios below for a target of 30 frames per second (FPS). All tests have been
271 conducted using the same machine with one Intel® Core™ Ultra 9 (22 cores) and 96 GB of RAM.

272 5.1 Human participant scalability

273 This experiment assesses SHARPIE’s capacity to host large-scale participant studies. We simulate
274 $n \in \{1, 10, 50, 100, 250\}$ simultaneous participants interacting simultaneously with a single real-
275 time environment instance, in a single experiment, hosted at the same machine. To isolate SHARPIE
276 overhead from environment-specific computation, we use a no-op environment which discards the N
277 received actions, and returns a 64×64 RGB observation containing random noise at each time step
278 to simulate a visual payload per timestep. Having all participants interact with the same environment
279 instance is challenging since all inputs need to be processed simultaneously before transitioning the
280 environment to the next state. Moreover, this single-experiment setting is the most relevant since
281 scalability in settings with multiple environment instances can be easily achieved by hosting separate
282 experiments on separate SHARPIE web server instances, i.e. by scaling horizontally with n .

283 We report system performance via several metrics in Table 2. First, we consider actual FPS as
284 an indicator of performance for real-time interactions, averaged across participants throughout an
285 episode. Second, we report the median and 95th percentile (P95) round-trip times (RTT) as the
286 latency between a participant sending an action and receiving its associated observation. Third, we
287 report throughput as a measure of capacity of the system independent of the individual participant
288 experience. We also recorded the frequency of socket-level errors or dropped packets which was 0
289 in all experiments.

290 From Table 2 we see that the system maintains nominal performance up to $n = 100$ concurrent par-
291 ticipants, with average FPS remaining close to the target of 30 and median RTT showing negligible
292 variance across number of participants ($\Delta \approx 0.1\text{ms}$), although P95 latency increases substantially. At
293 $n = 250$, however, the system has reached a performance bottleneck and average FPS degrades by
294 67% relative to the single-participant baseline, and Median RTT increases to 68.66 ms. This non-
295 linear degradation indicates the saturation point of the single-instance server architecture under the
296 tested hardware configuration. Notably, across all tested scales, the system recorded zero commu-
297 nication errors, indicating that while throughput limits affect fluidity and will affect experimental
298 results, the underlying state synchronization remains intact.

299 5.2 Multi-agent scalability

300 Similarly, we evaluate SHARPIE’s scalability in the number of RL agents. We simulate $n \in$
301 $\{1, 10, 50, 100, 250\}$ concurrent random RL agents interacting in a single no-op environment in-

Table 2: Performance metrics by participant count, results in parenthesis are relative to $N = 1$ baseline.

n participants	FPS	median RTT (ms)	p95 RTT (ms)	Throughput (msg/s)
1	32.49 (100%)	20.12 (1.00 \times)	22.54 (1.00 \times)	50.78 (1.0 \times)
10	33.11 (102%)	19.82 (0.99 \times)	22.33 (0.99 \times)	480.67 (9.5 \times)
50	29.69 (91%)	19.90 (0.99 \times)	59.30 (2.63 \times)	2244.50 (44.2 \times)
100	22.98 (71%)	27.88 (1.39 \times)	71.25 (3.16 \times)	3056.07 (60.2 \times)
250	10.64 (33%)	68.66 (3.41 \times)	130.06 (5.77 \times)	3750.71 (73.9 \times)

302 stance. This environment discards all actions and samples a random 64×64 noise image at each
 303 timestep. The metrics reported in 3 show that the system maintains nominal performance in all cases.
 304 We hypothesize that this due to SHARPIEs architecture, in which all RL policies run in the same
 305 process as the environment simulation loop to minimize data transfer and the overhead imposed by
 306 SHARPIE consists mostly of routing and logging of messages.

Table 3: Performance metrics for a variable number of RL agents, results in parenthesis are relative to $N = 1$ baseline.

n RL agents	FPS	Throughput (msg/s)
1	32.33 (100%)	51.48 (1.00 \times)
10	32.40 (100.2%)	54.40 (1.05 \times)
50	32.39 (100.2%)	48.84 (0.94 \times)
100	32.43 (100.3%)	49.53 (0.96 \times)
250	33.08 (102.3%)	50.85 (0.98 \times)

307 5.3 Network Quality and Geographic Distribution

308 To assess the suitability of our server-client architecture in supporting experiments involving geo-
 309 graphically distant participants and RL agents, we simulate the impact of network-induced latency
 310 on system performance. Specifically, we emulate various network latencies to assess FPS and RTT
 311 in the following settings: a SHARPIE webserver instance running on the participants’ localhost
 312 (<0.1 ms), on the participants network (<10 ms), in the participants country (<20 ms), on the par-
 313 ticipants continent (50 ms) and world-wide (200 ms). The results reported in 4 show that, with a
 314 target of 30 FPS, SHARPIE does not suffer loss when running on the same LAN. This is perfect for
 315 hosting interactive real-time experiments. in the same university as the server running SHARPIE.
 316 For an experiment requiring participants from a cloud source platform, the system will be able to run
 317 at least 26 FPS and provide a smooth interaction. However, for regional or world-wide interactions,
 318 it is more suited for experiments with less real-time requirements.

Table 4: Performance stability analysis under variable network conditions, results in parenthesis are relative to localhost baseline.

Trial / Condition	Avg FPS (rel.)	Median RTT (rel.)	P95 RTT (rel.)	Throughput (rel.)
Localhost	32.39 (100%)	20.29 (1.00 \times)	21.91 (1.00 \times)	52.43 (1.0 \times)
LAN	33.08 (102%)	20.13 (0.99 \times)	22.00 (1.00 \times)	53.37 (1.0 \times)
National	26.77 (83%)	26.90 (1.33 \times)	30.34 (1.38 \times)	44.32 (0.85 \times)
Regional	14.62 (45%)	57.88 (2.85 \times)	60.58 (2.76 \times)	26.14 (0.5 \times)
Global	4.59 (14%)	208.72 (10.29 \times)	211.13 (9.64 \times)	8.83 (0.2 \times)

319 5.4 Environment and Rendering Complexity

320 Finally, we test the throughput of the framework when handling high definition rendering or highly
 321 demanding payloads. To illustrate that, we evaluate rendering different image sizes with random
 322 noise as well as testing a Starcraft II environment. The results in Table 5 show that SHARPIE can
 323 handle up to 256x256 payloads at 30FPS which will comfortably serve most use cases. Scalability
 324 beyond this size remains challenging and will require additional measures such as compression
 325 currently not implemented in SHARPIE.

Table 5: Performance metrics for a single participant under varying image size, results in parenthesis are relative to the smallest (64x64).

Image size	Avg FPS (rel.)	Median RTT (rel.)	P95 RTT (rel.)	Throughput (rel.)
64x64	33.12 (100%)	20.28 (100%)	22.22 (100%)	64.34 (100%)
128x128	33.17 (100.2%)	20.13 (99.4%)	22.67 (102.1%)	53.03 (82.5%)
256x256	30.37 (91.7%)	22.30 (110.0%)	27.89 (125.6%)	49.59 (77.1%)
512x512	12.10 (36.6%)	72.80 (359%)	77.17 (347%)	22.12 (34.4%)
1024x1024	3.43 (10.4%)	282.63 (1394%)	326.91 (1470%)	6.64 (10.3%)

326 6 Discussion & Future Work

327 We have presented and motivated the design and implementation of a framework to accelerate hu-
 328 man-AI interaction studies, with a particular focus on studies for multi-agent tasks involving hu-
 329 mans. As stated, the scope of this project is not to implement alternative environment and algorithms
 330 for such studies, but rather to integrate with existing solutions seamlessly and provide utilities to ease
 331 experimentation with participants. This is done by empowering researchers through a modular soft-
 332 ware interface relying, in part, on the de facto standard set by Gymnasium. With SHARPIE we aim
 333 to provide an easily integrable framework that researchers can use to painlessly set up experiments
 334 involving both human and artificial agents. Our hope is that in turn such an architecture lays the
 335 foundation for a standard for the interaction between human and artificial agents.

336 This modular approach also allows for numerous directions of improvement in terms of interoper-
 337 ability and features. In the former case, we aim to provide an increasing number of ready-made,
 338 supported plugins to handle a large part of the existing work on environments, libraries, and deploy-
 339 ment options. In the latter case, we plan to widen the scope of possible human-agent interactions
 340 by incorporating additional modalities such as audio or video (Christofi & Baraka, 2024; Knierim
 341 et al., 2024b). The resulting real-time and multimodal communication between users and agents
 342 would allow the study of rich and fine-grained communication protocols, and support experiments
 343 involving pitch and tone, nonverbal cues, etc. Finally, we envision a hosted version of SHARPIE
 344 that can be used for outreach, education and user literacy purposes.

345 References

- 346 Zeynep Akata, Dan Balliet, Maarten De Rijke, Frank Dignum, Virginia Dignum, Guszt Eiben,
 347 Antske Fokkens, Davide Grossi, Koen Hindriks, Holger Hoos, et al. A research agenda for hybrid
 348 intelligence: augmenting human intellect with collaborative, adaptive, responsible, and explain-
 349 able artificial intelligence. *Computer*, 53(8):18–28, 2020.
- 350 Nele Albers, Mark A Neerinx, and Willem-Paul Brinkman. Addressing people’s current and future
 351 states in a reinforcement learning algorithm for persuading to quit smoking and to be physically
 352 active. *Plos one*, 17(12):e0277295, 2022.
- 353 Lucas N Alegre, Florian Felten, El-Ghazali Talbi, Grégoire Danoy, Ann Nowé, Ana LC Bazzan, and
 354 Bruno C da Silva. Mo-gym: A library of multi-objective reinforcement learning environments. In

- 355 *Proceedings of the 34th Benelux Conference on Artificial Intelligence BNAIC/Benelearn*, volume
356 2022, pp. 2, 2022.
- 357 Kim Baraka, Ifrah Idrees, Taylor Kessler Faulkner, Erdem Biyik, Serena Booth, Mohamed
358 Chetouani, Daniel H. Grollman, Akanksha Saran, Emmanuel Senft, Silvia Tulli, Anna-Lisa
359 Vollmer, Antonio Andriella, Helen Beierling, Tiffany Horter, Jens Kober, Isaac Sheidlower,
360 Matthew E. Taylor, Sanne van Waveren, and Xuesu Xiao. Human-interactive robot learning:
361 Definition, challenges, and recommendations. *J. Hum.-Robot Interact.*, 15(2), February 2026.
362 DOI: 10.1145/3779297. URL [https://doi-org.vu-nl.idm.oclc.org/10.1145/](https://doi-org.vu-nl.idm.oclc.org/10.1145/3779297)
363 [3779297](https://doi-org.vu-nl.idm.oclc.org/10.1145/3779297).
- 364 Erdem Biyık, Aditi Talati, and Dorsa Sadigh. Aprel: A library for active preference-based reward
365 learning algorithms. In *2022 17th ACM/IEEE International Conference on Human-Robot Inter-*
366 *action (HRI)*, pp. 613–617. IEEE, 2022.
- 367 G Brockman. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 368 Micah Carroll, Rohin Shah, Mark K. Ho, Thomas L. Griffiths, Sanjit A. Seshia, Pieter Abbeel,
369 and Anca Dragan. *On the utility of learning about humans for human-AI coordination*. Curran
370 Associates Inc., Red Hook, NY, USA, 2019.
- 371 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep
372 reinforcement learning from human preferences. *Advances in neural information processing sys-*
373 *tems*, 30, 2017.
- 374 Konstantinos Christofi and Kim Baraka. Uncovering patterns in humans that teach robots through
375 demonstrations and feedback. In *Companion of the 2024 ACM/IEEE International Conference*
376 *on Human-Robot Interaction*, pp. 332–336, 2024.
- 377 Joshua R De Leeuw. jspsych: A javascript library for creating behavioral experiments in a web
378 browser. *Behavior research methods*, 47:1–12, 2015.
- 379 DeepMind, Igor Babuschkin, Kate Baumli, Alison Bell, Surya Bhupatiraju, Jake Bruce, Peter
380 Buchlovsky, David Budden, Trevor Cai, Aidan Clark, Ivo Danihelka, Antoine Dedieu, Clau-
381 dio Fantacci, Jonathan Godwin, Chris Jones, Ross Hemsley, Tom Hennigan, Matteo Hessel,
382 Shaobo Hou, Steven Kapturowski, Thomas Keck, Iurii Kemaev, Michael King, Markus Kunesch,
383 Lena Martens, Hamza Merzic, Vladimir Mikulik, Tamara Norman, George Papamakarios, John
384 Quan, Roman Ring, Francisco Ruiz, Alvaro Sanchez, Laurent Sartran, Rosalia Schneider, Eren
385 Sezener, Stephen Spencer, Srivatsan Srinivasan, Miloš Stanojević, Wojciech Stokowiec, Luyu
386 Wang, Guangyao Zhou, and Fabio Viola. The DeepMind JAX Ecosystem, 2020. URL [http:](http://github.com/deepmind)
387 [//github.com/deepmind](http://github.com/deepmind).
- 388 Florian Felten, Lucas N. Alegre, Ann Nowé, Ana L. C. Bazzan, El Ghazali Talbi, Grégoire Danoy,
389 and Bruno Castro da Silva. A toolkit for reliable benchmarking and research in multi-objective
390 reinforcement learning. In *Proceedings of the 37th Conference on Neural Information Processing*
391 *Systems (NeurIPS 2023)*, 2023.
- 392 Florian Felten, Umut Ucak, Hicham Azmani, Gao Peng, Willem Röpke, Hendrik Baier, Patrick
393 Mannion, Diederik M Roijers, Jordan K Terry, El-Ghazali Talbi, et al. Momaland: A set of
394 benchmarks for multi-objective multi-agent reinforcement learning. In *Multi-objective Decision*
395 *Making Workshop at ECAI 2024*, 2024.
- 396 Jason Gauci, Edoardo Conti, Yitao Liang, Kittipat Virochsiri, Yuchen He, Zachary Kaden, Vivek
397 Narayanan, Xiaohui Ye, Zhengxing Chen, and Scott Fujimoto. Horizon: Facebook’s open source
398 applied reinforcement learning platform. In *ICML 2019 Workshop on RL4RealLife*, 2019. URL
399 <https://openreview.net/forum?id=SylQKinLi4>.
- 400 Kevin Godin-Dubois, Karine Miras, and Anna V Kononova. AMaze: A benchmark generator for
401 sighted maze-navigating agents. *Journal of Open Source Software*, pp. in press, 2024.

- 402 Paul Green and Lisa Wei-Haas. The rapid development of user interfaces: Experience with the
403 wizard of oz method. *Proceedings of the Human Factors Society Annual Meeting*, 29(5):470–
404 474, 1985. DOI: 10.1177/154193128502900515. URL [https://doi.org/10.1177/
405 154193128502900515](https://doi.org/10.1177/154193128502900515).
- 406 Lin Guan, Mudit Verma, Suna Sihang Guo, Ruohan Zhang, and Subbarao Kambhampati. Widening
407 the pipeline in human-guided reinforcement learning with explanation and context-aware data
408 augmentation. *Advances in Neural Information Processing Systems*, 34:21885–21897, 2021.
- 409 Shengyi Huang, Rousslan Fernand Julien Dossa, Chang Ye, Jeff Braga, Dipam Chakraborty, Ki-
410 nal Mehta, and João G.M. Araújo. Cleanrl: High-quality single-file implementations of deep
411 reinforcement learning algorithms. *Journal of Machine Learning Research*, 23(274):1–18, 2022.
412 URL <http://jmlr.org/papers/v23/21-1342.html>.
- 413 Brian Ichter, Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog,
414 Daniel Ho, Julian Ibarz, Alex Irpan, Eric Jang, Ryan Julian, Dmitry Kalashnikov, Sergey Levine,
415 Yao Lu, Carolina Parada, Kanishka Rao, Pierre Sermanet, Alexander T Toshev, Vincent Van-
416 houecke, Fei Xia, Ted Xiao, Peng Xu, Mengyuan Yan, Noah Brown, Michael Ahn, Omar Cortes,
417 Nicolas Sievers, Clayton Tan, Sichun Xu, Diego Reyes, Jarek Rettinghouse, Jornell Quiambao,
418 Peter Pastor, Linda Luu, Kuang-Huei Lee, Yuheng Kuang, Sally Jesmonth, Nikhil J. Joshi, Kyle
419 Jeffrey, Rosario Jauregui Ruano, Jasmine Hsu, Keerthana Gopalakrishnan, Byron David, Andy
420 Zeng, and Chuyuan Kelly Fu. Do as i can, not as i say: Grounding language in robotic affordances.
421 In Karen Liu, Dana Kulic, and Jeff Ichnowski (eds.), *Proceedings of The 6th Conference on Robot
422 Learning*, volume 205 of *Proceedings of Machine Learning Research*, pp. 287–318. PMLR, 14–
423 18 Dec 2023. URL <https://proceedings.mlr.press/v205/ichter23a.html>.
- 424 Timo Kaufmann, Paul Weng, Viktor Bengs, and Eyke Hüllermeier. A survey of reinforcement
425 learning from human feedback. *Transactions on Machine Learning Research*, 2025. ISSN 2835-
426 8856. URL <https://openreview.net/forum?id=f7OkTurx4b>. Survey Certification.
- 427 Christian Kauten. Super Mario Bros for OpenAI Gym. GitHub, 2018. URL [https://github.
428 com/Kautenja/gym-super-mario-bros](https://github.com/Kautenja/gym-super-mario-bros).
- 429 Matilda Knierim, Sahil Jain, Murat Han Aydoğın, Kenneth Mitra, Kush Desai, Akanksha Saran, and
430 Kim Baraka. Prosody as a teaching signal for agent learning: Exploratory studies and algorithmic
431 implications. *arXiv preprint arXiv:2410.23554*, 2024a.
- 432 Matilda Knierim, Sahil Jain, Murat Han Aydoğın, Kenneth D Mitra, Kush Desai, Akanksha Saran,
433 and Kim Baraka. Leveraging prosody as an informative teaching signal for agent learning: Ex-
434 ploratory studies and algorithmic implications. In *Proceedings of the 26th International Confer-
435 ence on Multimodal Interaction*, pp. 95–123, 2024b.
- 436 W Bradley Knox and Peter Stone. Interactively shaping agents via human reinforcement: The tamer
437 framework. In *Proceedings of the fifth international conference on Knowledge capture*, pp. 9–16,
438 2009.
- 439 Marc Lanctot, Edward Lockhart, Jean-Baptiste Lespiau, Vinicius Zambaldi, Satyaki Upadhyay,
440 Julien Pérolat, Sriram Srinivasan, Finbarr Timbers, Karl Tuyls, Shayegan Omidshafiei, Daniel
441 Hennes, Dustin Morrill, Paul Muller, Timo Ewalds, Ryan Faulkner, János Kramár, Bart De
442 Vylder, Brennan Saeta, James Bradbury, David Ding, Sebastian Borgeaud, Matthew Lai, Julian
443 Schrittwieser, Thomas Anthony, Edward Hughes, Ivo Danihelka, and Jonah Ryan-Davis. Open-
444 Spiel: A framework for reinforcement learning in games. *CoRR*, abs/1908.09453, 2019. URL
445 <http://arxiv.org/abs/1908.09453>.
- 446 Eric Liang, Richard Liaw, Robert Nishihara, Philipp Moritz, Roy Fox, Ken Goldberg, Joseph Gon-
447 zalez, Michael Jordan, and Ion Stoica. RLlib: Abstractions for distributed reinforcement learning.
448 In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on
449 Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 3053–3062.

- 450 PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/liang18b.html>.
451
- 452 Chase McDonald. Interactive gym. [https://github.com/chasemcd/](https://github.com/chasemcd/interactive-gym)
453 [interactive-gym](https://github.com/chasemcd/interactive-gym), 2024.
- 454 Sriraam Natarajan, Gautam Kunapuli, Kshitij Judah, Prasad Tadepalli, Kristian Kersting, and Jude
455 Shavlik. Multi-agent inverse reinforcement learning. In *2010 ninth international conference on*
456 *machine learning and applications*, pp. 395–400. IEEE, 2010.
- 457 Jonathan W Peirce. Psychopy—psychophysics software in python. *Journal of neuroscience meth-*
458 *ods*, 162(1-2):8–13, 2007.
- 459 Cheng Qian, Zuxin Liu, Akshara Prabhakar, Jielin Qiu, Zhiwei Liu, Haolin Chen, Shirley Kokane,
460 Heng Ji, Weiran Yao, Shelby Heinecke, et al. Userrl: Training interactive user-centric agent via
461 reinforcement learning. *arXiv preprint arXiv:2509.19736*, 2025.
- 462 Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah
463 Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of*
464 *Machine Learning Research*, 22(268):1–8, 2021. URL [http://jmlr.org/papers/v22/](http://jmlr.org/papers/v22/20-1364.html)
465 [20-1364.html](http://jmlr.org/papers/v22/20-1364.html).
- 466 Stephane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and struc-
467 tured prediction to no-regret online learning. In Geoffrey Gordon, David Dunson, and Miroslav
468 Dudík (eds.), *Proceedings of the Fourteenth International Conference on Artificial Intelligence*
469 *and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pp. 627–635, Fort Laud-
470 erdale, FL, USA, 11–13 Apr 2011. PMLR. URL [https://proceedings.mlr.press/](https://proceedings.mlr.press/v15/ross11a.html)
471 [v15/ross11a.html](https://proceedings.mlr.press/v15/ross11a.html).
- 472 Alexander Rutherford, Benjamin Ellis, Matteo Gallici, Jonathan Cook, Andrei Lupu, Garðar Ing-
473 varsson, Timon Willi, Akbir Khan, Christian Schroeder de Witt, Alexandra Souly, Saptarashmi
474 Bandyopadhyay, Mikayel Samvelyan, Minqi Jiang, Robert Lange, Shimon Whiteson, Bruno Lac-
475 erda, Nick Hawes, Tim Rocktäschel, Chris Lu, and Jakob Foerster. Jaxmarl: Multi-agent rl
476 environments and algorithms in jax. In *Proceedings of the 23rd International Conference on Au-*
477 *tonomous Agents and Multiagent Systems*, AAMAS ’24, pp. 2444–2446, Richland, SC, 2024. In-
478 ternational Foundation for Autonomous Agents and Multiagent Systems. ISBN 9798400704864.
- 479 William Saunders, Girish Sastry, Andreas Stuhlmüller, and Owain Evans. Trial without error: To-
480 wards safe reinforcement learning via human intervention. In *Proceedings of the 17th Interna-*
481 *tional Conference on Autonomous Agents and MultiAgent Systems*, AAMAS ’18, pp. 2067–2069,
482 Richland, SC, 2018. International Foundation for Autonomous Agents and Multiagent Systems.
- 483 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press,
484 1998.
- 485 Matthew E Taylor, Nicholas Nissen, Yuan Wang, and Neda Navidi. Improving reinforcement learn-
486 ing with human assistance: an argument for human subject studies with hippo gym. *Neural*
487 *Computing and Applications*, 35(32):23429–23439, 2023.
- 488 Jordan Terry, Benjamin Black, Nathaniel Grammel, Mario Jayakumar, Ananth Hari, Ryan Sullivan,
489 Luis S Santos, Clemens Dieffendahl, Caroline Horsch, Rodrigo Perez-Vicente, et al. Petting-
490 zoo: Gym for multi-agent reinforcement learning. *Advances in Neural Information Processing*
491 *Systems*, 34:15032–15043, 2021.
- 492 Marcel Torne, Max Balsells, Zihan Wang, Samedh Desai, Tao Chen, Pulkit Agrawal, and Abhishek
493 Gupta. Breadcrumbs to the goal: goal-conditioned exploration from human-in-the-loop feedback.
494 In *Proceedings of the 37th International Conference on Neural Information Processing Systems*,
495 pp. 63222–63258, 2023.

- 496 Mark Towers, Ariel Kwiatkowski, Jordan Terry, John U Balis, Gianluca De Cola, Tristan Deleu,
497 Manuel Goulao, Andreas Kallinteris, Markus Krimmel, Arjun KG, et al. Gymnasium: A standard
498 interface for reinforcement learning environments. *arXiv preprint arXiv:2407.17032*, 2024.
- 499 Yifu Yuan, Jianye Hao, Yi Ma, Zibin Dong, Hebin Liang, Jinyi Liu, Zhixin Feng, Kai Zhao, and
500 Yan Zheng. Uni-RLHF: Universal platform and benchmark suite for reinforcement learning with
501 diverse human feedback. In *The Twelfth International Conference on Learning Representations*,
502 *ICLR*, 2024. URL <https://openreview.net/forum?id=WesY0H9ghM>.
- 503 Yangyang Zhao, Zhenyu Wang, Changxi Zhu, and Shihan Wang. Efficient dialogue complementary
504 policy learning via deep q-network policy and episodic memory policy. In *Proceedings of the*
505 *2021 Conference on Empirical Methods in Natural Language Processing*, pp. 4311–4323, 2021.
- 506 Lianmin Zheng, Jiacheng Yang, Han Cai, Ming Zhou, Weinan Zhang, Jun Wang, and Yong Yu.
507 Magent: A many-agent reinforcement learning platform for artificial collective intelligence. In
508 *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- 509 Zheqing Zhu, Rodrigo de Salvo Braz, Jalaj Bhandari, Daniel Jiang, Yi Wan, Yonathan Efroni, Liyuan
510 Wang, Ruiyang Xu, Hongbo Guo, Alex Nikulkov, et al. Pearl: A production-ready reinforcement
511 learning agent. *Journal of Machine Learning Research*, 25(273):1–30, 2024.

512
513
514

Supplementary Materials

The following content was not necessarily subject to peer review.

Name	Description	Language	Target	API
Amaze ¹	Maze navigation environments	Python	Research	Gymnasium+
CleanRL ²	Algorithm library	Python/Torch	Research/Production	Gymnasium
Gymnasium ³	De-facto standard environments	Python	Research	Gymnasium
HIPPO-Gym ⁴	Human Input Parsing	Python	Research	Gymnasium+
Interactive-Gym	Interactive web-based experiment platform	Python	Research	Gymnasium+
JaxMARL ⁵	Multi-agent environments	JAX	Research	Gymnasium+
Momaland ⁶	Multi-agent multi-objective environments	Python	Research	Gymnasium
MO-Gymnasium ⁷	Multi-objective environments	Python	Research	Gymnasium+
MORL Baselines ⁸	Multi-objective algorithm library	Python	Research	Gymnasium+
OpenSpiel ⁹	Game (theory) environments/algorithms	Python	Research	Gymnasium
Pearl ¹⁰	Algorithm library	Python	Production	Gymnasium
PettingZoo ¹¹	Multi-agent environments	Python	Research	Gymnasium+
Ray RLlib ¹²	Algorithm library	Python/Torch	Production	Gymnasium
RLax ¹³	Algorithm library	JAX	Research	Gymnasium
Stable-baselines3 ¹⁴	Algorithm library	Python/Torch	Research/Production	Gymnasium
Uni-RLHF ¹⁵	Annotation tool	Python	Annotation	Gymnasium+

Table 6: Comparison of RL Packages, APIs, and Applications. + indicates minor alterations to the Gymnasium API such as a vectorial reward for multi-objective, or vectorial action for multi-agent RL. References are: 1- [Godin-Dubois et al. \(2024\)](#); 2- [Huang et al. \(2022\)](#); 3- [Brockman \(2016\)](#); 4- [Taylor et al. \(2023\)](#); 5- [Rutherford et al. \(2024\)](#); 6- [Felten et al. \(2024\)](#); 7- [Alegre et al. \(2022\)](#); 8- [Felten et al. \(2023\)](#); 9- [Lanctot et al. \(2019\)](#); 10- [Zhu et al. \(2024\)](#); 11- [Terry et al. \(2021\)](#); 12- [Liang et al. \(2018\)](#); 13- [DeepMind et al. \(2020\)](#); 14- [Raffin et al. \(2021\)](#); 15- [Yuan et al. \(2024\)](#).

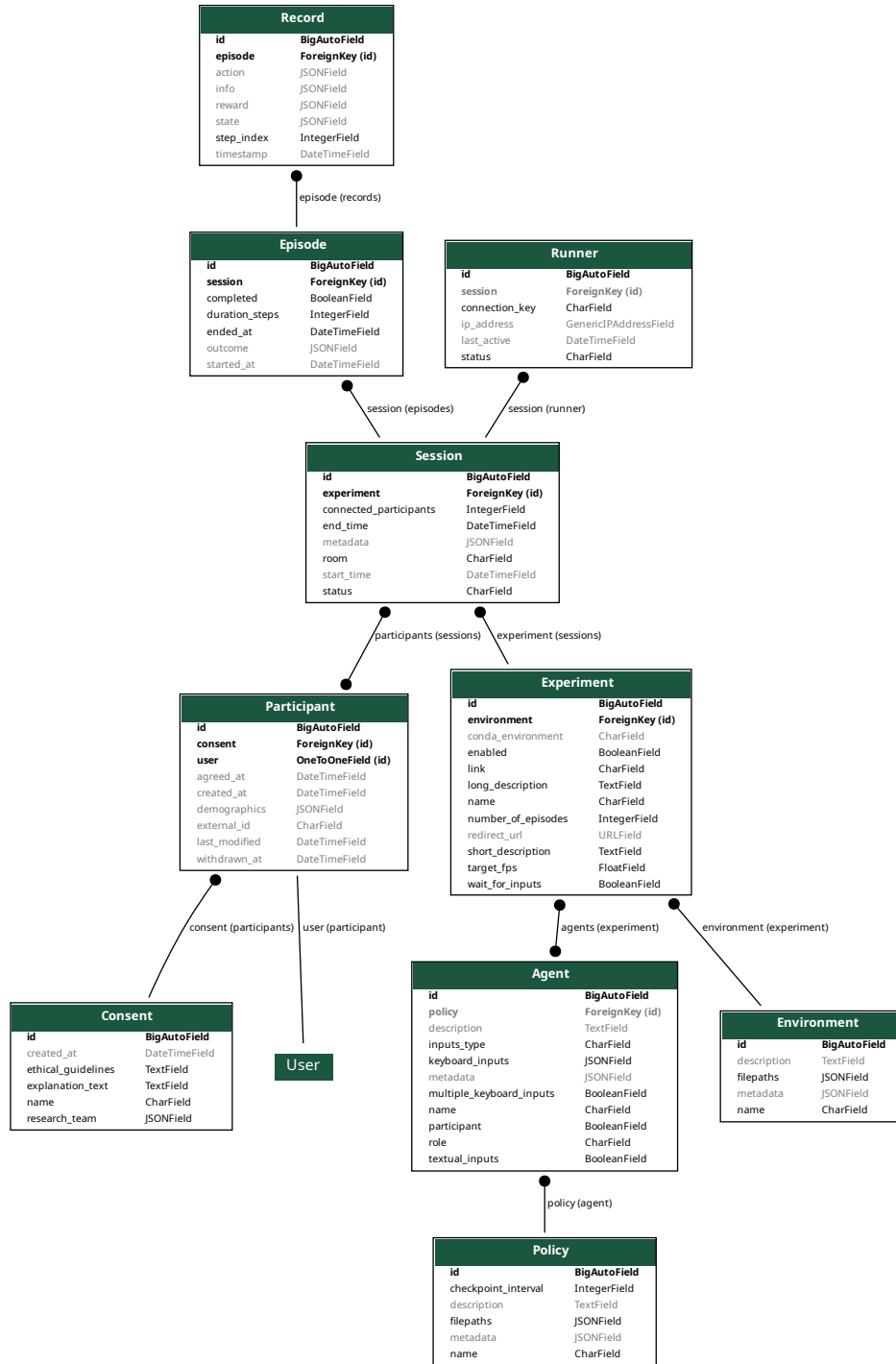


Figure 3: Overview of SHARPIE's data model