
NiceWebRL: a Python library for human subject experiments with reinforcement learning environments

Anonymous Author(s)
Affiliation
Address
email

Abstract

1 We present NiceWebRL, a research tool that enables researchers to use machine
2 reinforcement learning (RL) environments for online human subject experiments.
3 NiceWebRL is a Python library that allows any Jax-based environment to be
4 transformed into an online interface, supporting both single-agent and multi-agent
5 environments. As such, NiceWebRL enables AI researchers to compare their
6 algorithms to human performance, cognitive scientists to test ML algorithms as
7 theories for human cognition, and multi-agent researchers to develop algorithms
8 for human-AI collaboration. We showcase NiceWebRL with 3 case studies that
9 demonstrate its potential to help develop Human-like AI, Human-compatible AI,
10 and Human-assistive AI. In the first case study (Human-like AI), NiceWebRL
11 enables the development of a novel RL model of cognition. Here, NiceWebRL
12 facilitates testing this model against human participants in both a grid world and
13 Craftax, a 2D Minecraft domain. In our second case study (Human-compatible AI),
14 NiceWebRL enables the development of a novel multi-agent RL algorithm that
15 can generalize to human partners in the Overcooked domain. Finally, in our third
16 case study (Human-assistive AI), we show how NiceWebRL can allow researchers
17 to study how an LLM can assist humans on complex tasks in XLand-Minigrid,
18 an environment with millions of hierarchical tasks. The library is available at
19 <https://anonymous.4open.science/r/nicewebml-28BF>.

20 1 Introduction

21 Recent progress in Artificial Intelligence (AI) increasingly motivates researchers to study AI in the
22 context of human behavior. Some ML researchers aim to improve AI systems by comparing them to
23 humans, since humans can still provide an upper bound on the performance our systems can hope
24 to achieve [26, 31]. For example, Minecraft [14, 15] remains a challenging exploration problem
25 for machines but a fun exploration adventure for people [18, 10]. In cognitive science, there is
26 increased interest in asking whether these machines are human-like [23, 36, 6]. Even if they are not,
27 cognitive scientists are interested in using them as the basis for building human-like machines [23, 8].
28 In multi-agent RL, many researchers are interested in whether these agents can act as adaptive
29 partners to humans across the wide range of social settings they might be deployed in [5, 29]. This is
30 increasingly relevant with LLMs, as they possess superhuman knowledge and are improving in their
31 “reasoning” abilities [12]. A natural question is how well they can combine their prior knowledge
32 with environment perception to assist us in completing complex tasks. Collectively, these advances
33 could have a wide impact—ranging from robotics [33] to education [19] to healthcare [30]. Clearly,
34 many are interested in building human-like, human-compatible, and human-assistive AI.

35 Despite this interest in human-centered AI development, pursuing research that integrates human
36 subject experiments with modern ML libraries is currently a cumbersome process. To run experiments
37 with many participants, researchers leverage the internet to get large sample sizes [13]. Thus, most
38 infrastructure is written in the web’s programming language: JavaScript [11, 17, 13, 9]. Machine
39 learning code, on the other hand, relies heavily on Python for model development [1, 28, 4] and

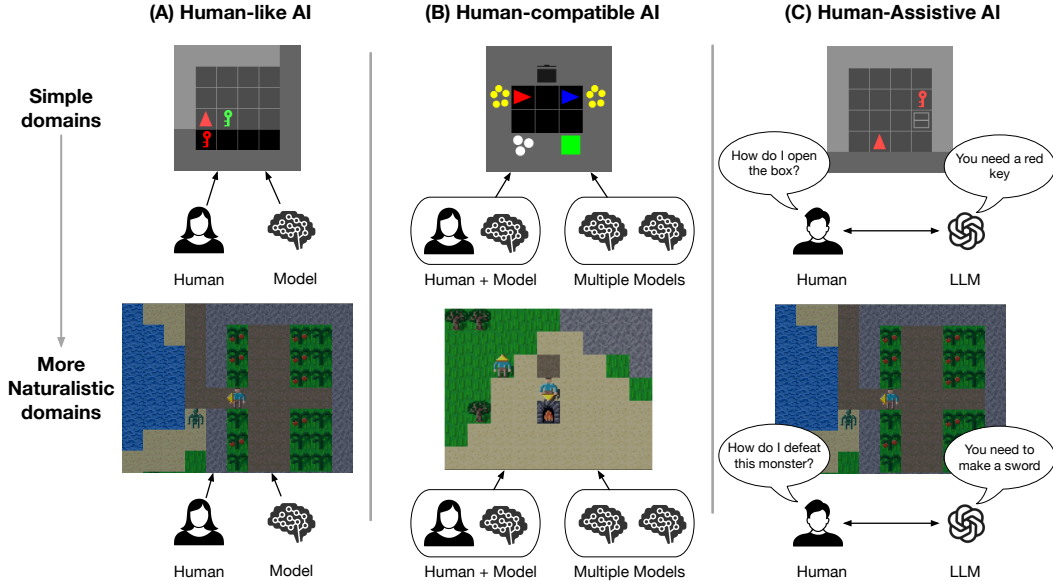


Figure 1: **NiceWebRL is a meta-environment that enables the use of Jax-based environments to develop Human-like, Human-compatible, and Human-assistive AI.** (A) Researchers can compare how humans and AI complete tasks to evaluate if AI behaves in human-like ways. (B) They can study how AI coordinates with humans during task completion to assess if the AI has learned human-compatible social behaviors. (C) They can also integrate Large Language Models into their experiments to evaluate how effectively they combine prior knowledge with environmental perception to assist participants. Importantly, NiceWebRL enables findings to generalize across potentially more complex domains.

40 variants of C for developing fast environments for simulation [22, 34]. Leveraging ML models or
 41 environments for human subject experiments currently requires setting up domain-specific server-
 42 client configurations that integrate Javascript, Python, and sometimes C. Doing this for each domain
 43 makes the process even more cumbersome.

44 To address this challenge, we present NiceWebRL: a research tool that lets researchers leverage ML
 45 environments for human subject experiments (see Figure 1). Integrating Python with JavaScript
 46 requires maintaining a connection between a remote Python-based server and a local Javascript-based
 47 client. This distance can cause latency issues when running online experiments. To circumvent
 48 this challenge, NiceWebRL exploits Jax [4]—a high-performance numerical computing library—to
 49 precompute and cache environment dynamics for arbitrary Jax-based environments. NiceWebRL then
 50 acts as a meta-environment for researchers to use arbitrary Jax-based environments in their human
 51 subject experiments. Critically, NiceWebRL allows researchers to program experiments entirely in
 52 Python by integrating with NiceGUI¹—a library that enables web developers to specify advanced
 53 Graphical User Interface (GUI) components entirely in Python.

54 1.1 Case study 1: Developing Human-like AI with NiceWebRL

55 **Developing a novel Deep RL cognitive science model with NiceWebRL [8].** A central question
 56 in cognitive science is how people represent the environment to enable generalization to new tasks.
 57 Successor features (SFs) are a mechanism for how an agent can cache expectations of what it will
 58 see when pursuing a policy [2]. Recent ML research has shown that SFs enable agents to repurpose
 59 policies for new tasks [3, 7]. Later, cognitive scientists showed that SFs also explained how people
 60 reuse prior policies for new tasks [32]. However, the behavioral work was done in a small grid-world
 61 with 13 states. Carvalho et al. [8] studied whether SFs could explain how humans reuse behaviors
 62 for new tasks in 2 more complex domains: a maze gridworld (Figure 2 A) and Craftax [25], a 2D
 63 minecraft domain (Figure 2 D). Across both domains, they set up training tasks where a test object
 64 was visible from along the optimal training paths (e.g. top-right corner of Figure 2 A). SFs could not
 65 generalize here. People could; however, when people reused a training path to a novel goal, their
 66 response times suggested that they were using a caching-based solution rather than something more

¹<https://github.com/zauberzeug/nicegui/>

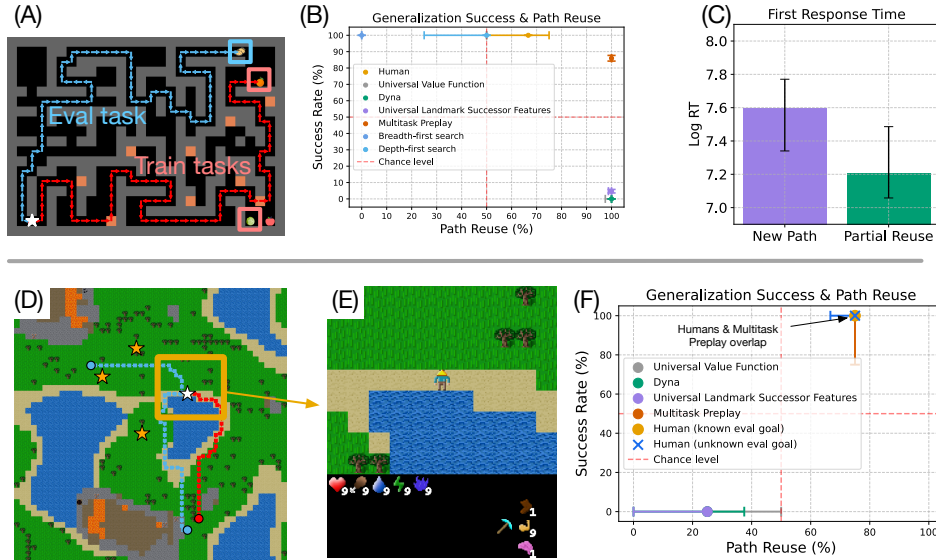


Figure 2: **Case study 1: NiceWebRL enabled the development of a novel Deep RL cognitive model that generalizes to new tasks with the same qualitative behaviors as Humans across multiple domains.** (A, D): a gridworld and 2D minecraft environment that both Human participants and Deep RL models learned in. (B, F): behavioral results studying the same phenomena across the two domains. NiceWebRL enabled developing a Deep RL cognitive model that could both (1) generalize to novel test goals in ways not permitted by previous methods, while (2) doing so with a similar suboptimal path reuse strategy that humans tend to exhibit. Figures reproduced from [8].

67 flexible—but expensive—like planning at decision-time (Figure 2 B-C). They developed an algorithm
 68 termed Multitask Preplay which *preemptively* learns solutions for unpursued tasks nearby training
 69 tasks by augmenting experience replay with small amounts of counterfactual simulation. They found
 70 this algorithm both better accounted for the response times people exhibited and better predicted
 71 how they would reuse prior behaviors. These results generalized to Craftax, where participants and
 72 models had to navigate from partial observations of a large world with many objects (Figure 2 E-F).

73 **Role of NiceWebRL.** First, NiceWebRL enabled comparing human behavior to advanced Deep RL
 74 algorithms including Successor Features. Second, NiceWebRL enabled using the same infrastructure
 75 to study both Deep RL algorithms and human behavior in two domains of increasing complexity:
 76 a gridworld and Craftax. This helped to ensure that findings were generalizable. Finally, NiceWeb-
 77 rL enabled the measurement of response times. This helped to adjudicate between theories that
 78 predict more or less computation at decision-time.

79 1.2 Case study 2: Developing Human-compatible AI with NiceWebRL

80 **Developing a novel Multi-agent reinforcement learning (MARL) algorithm for coordinating**
 81 **with humans using NiceWebRL [20].** One central question in MARL is how we can develop MARL
 82 agents that can generalize to human partners without human training data. One current benchmark for
 83 human-compatible AI is the Overcooked domain [5] where agents must coordinate on basic cooking
 84 tasks. The state-of-the art algorithm is “Efficient End-to-End Training” [E3T; 35], a “Self Play”
 85 algorithm that plays with—and tries to predict the actions of—a noisy variant of itself. Jha et al.
 86 [20] developed a novel algorithm, Cross-Environment Cooperation (CEC), where an agent plays only
 87 against itself but across millions of different procedurally-generated environments (Figure 3 A). They
 88 found that while E3T [35] was able to “succeed” on more episodes when collaborating with humans
 89 (Figure 3 B), humans gave that agent a lower rating than CEC (Figure 3 C). The authors asked
 90 participants questions about their subjective experience using a Likert scale [24] and found that CEC
 91 was rated as more “adaptive” and “human-like”—despite succeeding *less* than E3T. When the authors
 92 analyzed game trajectories, they found CEC would collide less with humans across environments.

93 **Role of NiceWebRL.** First, NiceWebRL enabled comparing multiple MARL algorithms in their
 94 ability to generalize to human partners. Second, NiceWebRL enabled the researchers to collect
 95 feedback from participants after every environment interaction stage using the “Feedback” stage

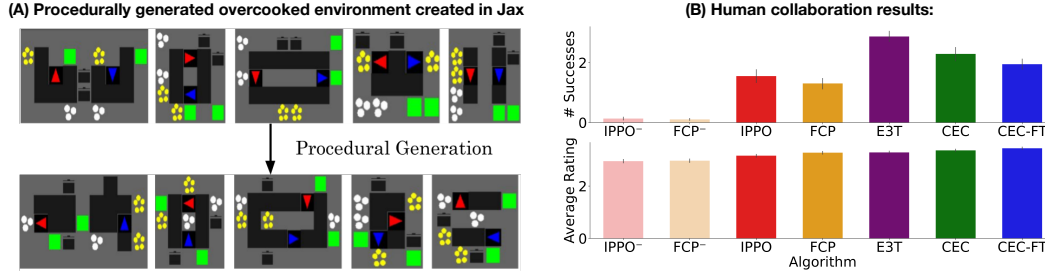


Figure 3: **Case study 2: NiceWebRL enabled the development of a novel MARL algorithm that is more compatible with novel human partners.** (A) A procedurally-generated environment used to design a novel MARL algorithm: Cross-Environment Cooperation (CEC). Prior work had agents learn with diverse sets of agents. CEC has a *single agent* play itself across millions of procedurally-generated environments. (B) While CEC *succeeded less* than other methods when collaborating with humans (top), it succeeded in ways that were *most favorable* to humans (bottom). Analysis suggested showed prior agents succeeded in less collaborative ways. Figures reproduced from [20].

96 object available in NiceWebRL coupled with NiceGUI’s data collection GUI elements. This provided
 97 an easy way to get feedback from participants while agent-interaction data was fresh in their memories.
 98 Third, NiceWebRL stores all environment interactions so participant episodes could be analyzed
 99 post-hoc. This enabled the researchers to analyze trajectories by participants and agents to determine
 100 what qualitative behaviors (such as colliding) were different between the different MARL algorithms.

101 **1.3 Case study 3: Developing Human-assistive AI with NiceWebRL**

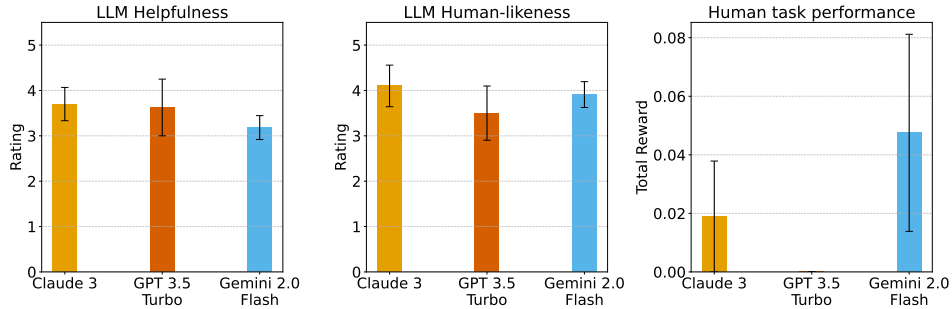


Figure 4: **Case study 3: Proof-of-concept experiment showing that NiceWebRL enables comparing how different LLMs can assist people in completing tasks.** We had Claude 3 Opus, GPT 3.5 Turbo, or Gemini 2.0 Flash act as assistants for people completing tasks in the XLand-Minigrid domain [27]. Each plot is showing the mean and standard error for 10 subjects per model.

102 **Developing an LLM-assistant for sequential-decision making tasks in a virtual environment.**
 103 We created a simple proof-of-concept experiment where people had to interact with tasks from the
 104 Xland-Minigrid domain [27]. To require assistance, they were given no task information but could
 105 ask an (anonymous) LLM assistant for help. We present more experiment details in§D.

106 **Role of NiceWebRL.** NiceWebRL enabled using an existing ML domain to develop an experiment
 107 that studied how an LLM could assist people on long-horizon tasks. Additionally, the Feedback Stage
 108 object in NiceWebRL enabled collecting feedback from participants about the LLMs at the end of the
 109 experiment. This could also be done after every episode or during episodes.

References

- 110
- 111 [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu
112 Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. {TensorFlow}: a system for
113 {Large-Scale} machine learning. In *12th USENIX symposium on operating systems design and*
114 *implementation (OSDI 16)*, pages 265–283, 2016.
- 115 [2] André Barreto, Will Dabney, Rémi Munos, Jonathan J Hunt, Tom Schaul, Hado P van Hasselt,
116 and David Silver. Successor features for transfer in reinforcement learning. *Advances in Neural*
117 *Information Processing Systems*, 30, 2017.
- 118 [3] Andre Barreto, Diana Borsa, John Quan, Tom Schaul, David Silver, Matteo Hessel, Daniel
119 Mankowitz, Augustin Zidek, and Remi Munos. Transfer in deep reinforcement learning using
120 successor features and generalised policy improvement. In *International Conference on Machine*
121 *Learning*, pages 501–510. PMLR, 2018.
- 122 [4] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal
123 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and
124 Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL
125 <http://github.com/google/jax>.
- 126 [5] Micah Carroll, Rohin Shah, Mark K Ho, Tom Griffiths, Sanjit Seshia, Pieter Abbeel, and Anca
127 Dragan. On the utility of learning about humans for human-ai coordination. *Advances in neural*
128 *information processing systems*, 32, 2019.
- 129 [6] Wilka Carvalho and Andrew Lampinen. Naturalistic computational cognitive science: Towards
130 generalizable models and theories that capture the full range of natural behavior. *arXiv preprint*
131 *arXiv:2502.20349*, 2025.
- 132 [7] Wilka Carvalho, Andre Saraiva, Angelos Filos, Andrew Lampinen, Loic Matthey, Richard L
133 Lewis, Honglak Lee, Satinder Singh, Danilo Jimenez Rezende, and Daniel Zoran. Combining
134 behaviors with the successor features keyboard. *Advances in Neural Information Processing*
135 *Systems*, 36, 2024.
- 136 [8] Wilka Carvalho, Sam Hall-McMaster, Honglak Lee, and Sam Gershman. Preemptive solving of
137 future problems: Multitask preplay in humans and machines. *arXiv preprint*, 2025.
- 138 [9] Joshua R De Leeuw. jspsych: A javascript library for creating behavioral experiments in a web
139 browser. *Behavior research methods*, 47:1–12, 2015.
- 140 [10] Yuqing Du, Eliza Kosoy, Alyssa Dayan, Maria Rufova, Pieter Abbeel, and Alison Gopnik. What
141 can ai learn from human exploration? intrinsically-motivated humans and agents in open-world
142 exploration. In *Neurips 2023 workshop: Information-theoretic principles in cognitive systems*,
143 2023.
- 144 [11] Holger Finger, Caspar Goeke, Dorena Diekamp, Kai Standvoß, and Peter König. Labvanced: a
145 unified javascript framework for online studies. In *International conference on computational*
146 *social science (cologne)*, pages 1–3. University of Osnabrück Cologne, 2017.
- 147 [12] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
148 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in
149 llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- 150 [13] Todd M Gureckis, Jay Martin, John McDonnell, Alexander S Rich, Doug Markant, Anna
151 Coenen, David Halpern, Jessica B Hamrick, and Patricia Chan. psiturk: An open-source
152 framework for conducting replicable behavioral experiments online. *Behavior research methods*,
153 48:829–842, 2016.
- 154 [14] William H Guss, Brandon Houghton, Nicholay Topin, Phillip Wang, Cayden Codel, Manuela
155 Veloso, and Ruslan Salakhutdinov. Minerl: A large-scale dataset of minecraft demonstrations.
156 *arXiv preprint arXiv:1907.13440*, 2019.
- 157 [15] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains
158 through world models, 2023. URL <https://arxiv.org/abs/2301.04104>, 2023.

- 159 [16] Assaf Hallak, Dotan Di Castro, and Shie Mannor. Contextual markov decision processes. *arXiv*
160 *preprint arXiv:1502.02259*, 2015.
- 161 [17] Felix Henninger, Yury Shevchenko, Ulf K Mertens, Pascal J Kieslich, and Benjamin E Hilbig.
162 *lab.js: A free, open, online study builder. Behavior Research Methods*, pages 1–18, 2021.
- 163 [18] Larissa Hjorth, Ingrid Richardson, Hugh Davies, and William Balmford. *Exploring Minecraft:*
164 *Ethnographies of play and creativity*. Springer Nature, 2021.
- 165 [19] Wayne Holmes and Ilkka Tuomi. State of the art and practice in ai in education. *European*
166 *journal of education*, 57(4):542–570, 2022.
- 167 [20] Kunal Jha, Wilka Carvalho, Yancheng Liang, Simon S Du, Max Kleiman-Weiner, and Natasha
168 Jaques. Cross-environment cooperation enables zero-shot multi-agent coordination. In *International*
169 *Conference on Machine Learning (ICML)*, 2025.
- 170 [21] Leslie Pack Kaelbling, Michael L Littman, and Anthony R Cassandra. Planning and acting in
171 partially observable stochastic domains. *Artificial intelligence*, 101(1-2):99–134, 1998.
- 172 [22] Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Matt
173 Deitke, Kiana Ehsani, Daniel Gordon, Yuke Zhu, et al. Ai2-thor: An interactive 3d environment
174 for visual ai. *arXiv preprint arXiv:1712.05474*, 2017.
- 175 [23] Brenden M Lake, Tomer D Ullman, Joshua B Tenenbaum, and Samuel J Gershman. Building
176 machines that learn and think like people. *Behavioral and brain sciences*, 40:e253, 2017.
- 177 [24] Rensis Likert. A technique for the measurement of attitudes. *Archives of psychology*, 1932.
- 178 [25] Michael Matthews, Michael Beukman, Benjamin Ellis, Mikayel Samvelyan, Matthew Jackson,
179 Samuel Coward, and Jakob Foerster. Craftax: A lightning-fast benchmark for open-ended
180 reinforcement learning. In *International Conference on Machine Learning (ICML)*, 2024.
- 181 [26] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G
182 Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al.
183 Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- 184 [27] Alexander Nikulin, Vladislav Kurenkov, Ilya Zisman, Artem Agarkov, Viacheslav Sini, and
185 Sergey Kolesnikov. XLand-minigrid: Scalable meta-reinforcement learning environments in
186 JAX. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.
- 187 [28] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
188 Trevor Killeen, Zeming Lin, Nicolas Gribonval, Zachary Devito, et al. Pytorch: An imperative
189 style, high-performance deep learning library. In *Advances in neural information processing*
190 *systems*, pages 8024–8035, 2019.
- 191 [29] Stuart Russell. Human-compatible artificial intelligence., 2022.
- 192 [30] Mohammed Yousef Shaheen. Applications of artificial intelligence (ai) in healthcare: A review.
193 *ScienceOpen Preprints*, 2021.
- 194 [31] Adaptive Agent Team, Jakob Bauer, Kate Baumli, Satinder Baveja, Feryal Behbahani, Avishkar
195 Bhoopchand, Nathalie Bradley-Schmiege, Michael Chang, Natalie Clay, Adrian Collister, et al.
196 Human-timescale adaptation in an open-ended task space. *arXiv preprint arXiv:2301.07608*,
197 2023.
- 198 [32] Momchil S Tomov, Eric Schulz, and Samuel J Gershman. Multi-task reinforcement learning in
199 humans. *Nature Human Behaviour*, 5(6):764–773, 2021.
- 200 [33] Sai H Vemprala, Rogerio Bonatti, Arthur Buckner, and Ashish Kapoor. Chatgpt for robotics:
201 Design principles and model abilities. *Ieee Access*, 2024.
- 202 [34] Tom Ward, Andrew Bolt, Nik Hemmings, Simon Carter, Manuel Sanchez, Ricardo Barreira,
203 Seb Noury, Keith Anderson, Jay Lemmon, Jonathan Coe, et al. Using unity to help solve
204 intelligence. *arXiv preprint arXiv:2011.09294*, 2020.

- 205 [35] Xue Yan, Jiaxian Guo, Xingzhou Lou, Jun Wang, Haifeng Zhang, and Yali Du. An efficient end-
206 to-end training approach for zero-shot human-ai coordination. *Advances in Neural Information*
207 *Processing Systems*, 36:2636–2658, 2023.
- 208 [36] Lance Ying, Katherine M Collins, Lionel Wong, Ilya Sucholutsky, Ryan Liu, Adrian Weller,
209 Tianmin Shu, Thomas L Griffiths, and Joshua B Tenenbaum. On benchmarking human-like
210 intelligence in machines. *arXiv preprint arXiv:2502.20502*, 2025.

211 A Formal domain description

212 Jax functions automatically compile to a fixed behavior when they receive their first input data. As
 213 such, if one wants different domain functionality across different *contexts* (e.g. training vs. testing),
 214 the domain’s functions typically need a “`env_parameter`” argument. Thus, Jax-based domains are
 215 naturally formulated as Partially Observable Contextual Markov Decision Processes (POCMDPs)
 216 $\mathcal{M}_c = \langle S, \mathcal{A}, \mathcal{X}, \mathcal{C}, \rho, P, R, O \rangle$ [16, 21]. Here, S denotes the environment state space, \mathcal{A} denotes
 217 its action space, \mathcal{X} denotes (potentially partial) observations of the environment, and \mathcal{C} denotes a
 218 space of contexts that an MDP can be in. `env_parameter` then corresponds to an MDP’s context
 219 $c \in \mathcal{C}$. It can be used to augment the initial state distribution $\rho_c(s_0)$ (e.g. having an agent start in
 220 different states in different contexts), the transition probabilities, $P_c(s'|s, a)$ (e.g. an agent’s speed
 221 or strength can be changed in different contexts), the reward function $R_c(s)$ (e.g. different objects
 222 can be rewarded in different contexts), or the observation function $O_c(s)$ (e.g. objects can take on
 223 different colors in different contexts).

224 An episode proceeds as follows. An initial state $s_0 \in S$ is sampled from the initial state distribution
 225 $\rho_c(s_0)$. When an agent takes an action $a \in \mathcal{A}$ in state $s \in S$, the next state s' is sampled according
 226 to a next state distribution $s' \sim P_c(\cdot|s, a)$. The agent then receives an observation $x' = O_c(s')$ and
 227 reward $r' = R_c(o)$. Note that c is typically fixed within an episode.

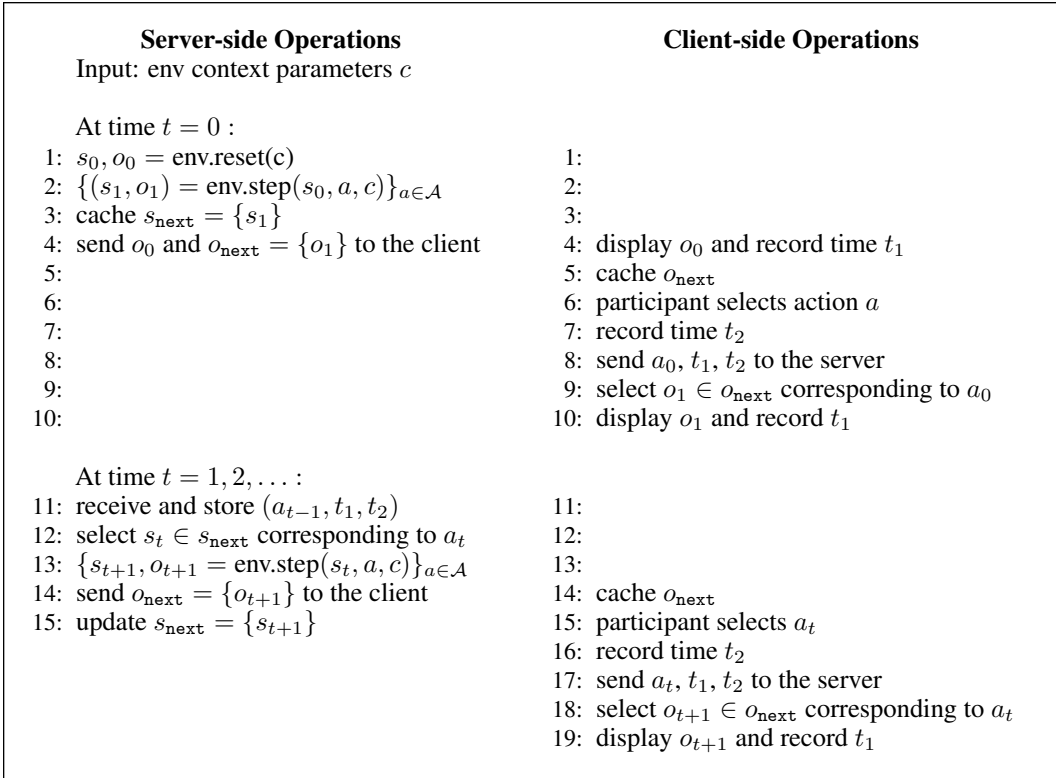


Figure 5: **Server-client Human-Environment Interaction Protocol**. Note that we omit displaying reward r due to space constraints.

228 Let `env` be the programmatic object representing a domain. In our library (and many RL libraries),
 229 $s_0, o_0 = \text{env.step}(c)$ essentially plays the role of sampling from the initial state distribution
 230 and computing the corresponding observation for the agent. The standard practice is to have
 231 $s_{t+1}, o_{t+1}, r_{t+1} = \text{env.step}(s_t, a_t, c)$ implement (a) sampling a new state (b) computing the corre-
 232 sponding reward, (c) computing the observation that an agent will get.

233 B Descriptions of stage types

234 Currently, there are three basic stage classes, though more can easily be added.

- 235 1. **Stage**: used to display instructions or information to a participant.
- 236 2. **FeedbackStage**: used to collect information from participants. Typically involves an
- 237 interactive screen that does *not* interact with the environment.
- 238 3. **EnvStage**: used to interact with an environment. It takes as input an environment and
- 239 environment parameters. We describe how NiceWebRL uses this abstraction to have a
- 240 remote server-side program display images to one’s local web-browser client in Figure 5.
- 241 We present examples of each in Figure 6.

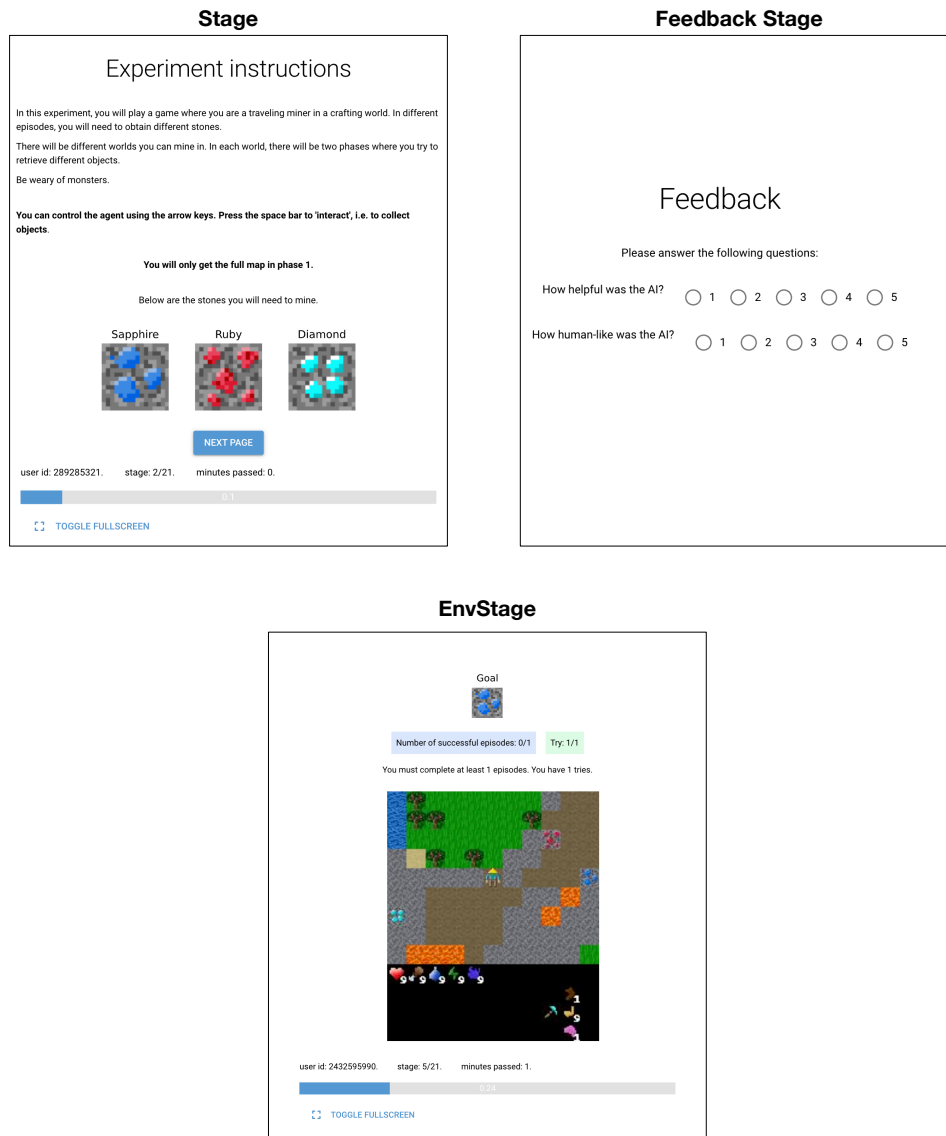


Figure 6: Examples of different kinds of stages.

242 C Computing resources

243 For details on case study 1 or 2, please see [8] or [20], respectively. For case study 3, experiments

244 were conducted using computing infrastructure from the fly.io platform with the “performance-2x”

245 configuration. This is a machine with 4GB of RAM. The machine had no GPU. Even in this setting,

246 Jax’s compilation features provide a significant speed up to environment computation.

247 **D Human subject experiment details**

248 Our study is approved by the University IRB. All subjects were recruited with <https://www.cloudresearch.com/> and provided informed consent. We provide the consent form in the GitHub
249 example. Participants were compensated \$4 for completing the task. The average task completion
250 time was 23.33 minutes. At the beginning of each experiment, the participants provided demographic
251 information (age and gender, coded as male or female).
252

253 At the beginning of the experiment, users were randomly assigned either Claude 3 Opus, GPT 3.5
254 Turbo, or Gemini 2.0 Flash. We set up our server so that it would interact with the LLMs via API
255 calls². Importantly, the LLM assistants were given text descriptions of the ground truth environment
256 state including information on: (1) the goal of the episode (2) the locations and identities of all objects
257 in the environment (3) the rules of the environment (e.g. how objects interact when combined). In
258 principle, this can enable them to help users figure out the goal to maximize task reward. In this
259 proof-of-concept, each participant completed 3 episodes, where a new task was sampled per episode.
260 After all 3 episodes, participants answered two questions on a 5-point scale, “How helpful was the
261 AI?” and “How human-like was the AI?”. We collect data from 30 participants via CloudResearch.
262 For this proof-of-concept, we used older models with cheap API calls. We don’t expect that these
263 results are representative of what is possible with frontier models.

²we provide an example of how to set up a web experiment with a local LLM in our examples folder