

## A FORMAL DOMAIN DESCRIPTION

Jax functions automatically compile to a fixed behavior when they receive their first input data. As such, if one wants different domain functionality across different *contexts* (e.g. training vs. testing), the domain’s functions typically need a “`env_parameter`” argument. Thus, Jax-based domains are naturally formulated as Partially Observable Contextual Markov Decision Processes (POCMDPs)  $\mathcal{M}_c = \langle S, \mathcal{A}, \mathcal{X}, \mathcal{C}, \rho, P, R, O \rangle$  (Hallak et al., 2015; Kaelbling et al., 1998). Here,  $S$  denotes the environment state space,  $\mathcal{A}$  denotes its action space,  $\mathcal{X}$  denotes (potentially partial) observations of the environment, and  $\mathcal{C}$  denotes a space of contexts that an MDP can be in. `env_parameter` then corresponds to an MDP’s context  $c \in \mathcal{C}$ . It can be used to augment the initial state distribution  $\rho_c(s_0)$  (e.g. having an agent start in different states in different contexts), the transition probabilities,  $P_c(s'|s, a)$  (e.g. an agent’s speed or strength can be changed in different contexts), the reward function  $R_c(s)$  (e.g. different objects can be rewarded in different contexts), or the observation function  $O_c(s)$  (e.g. objects can take on different colors in different contexts).

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

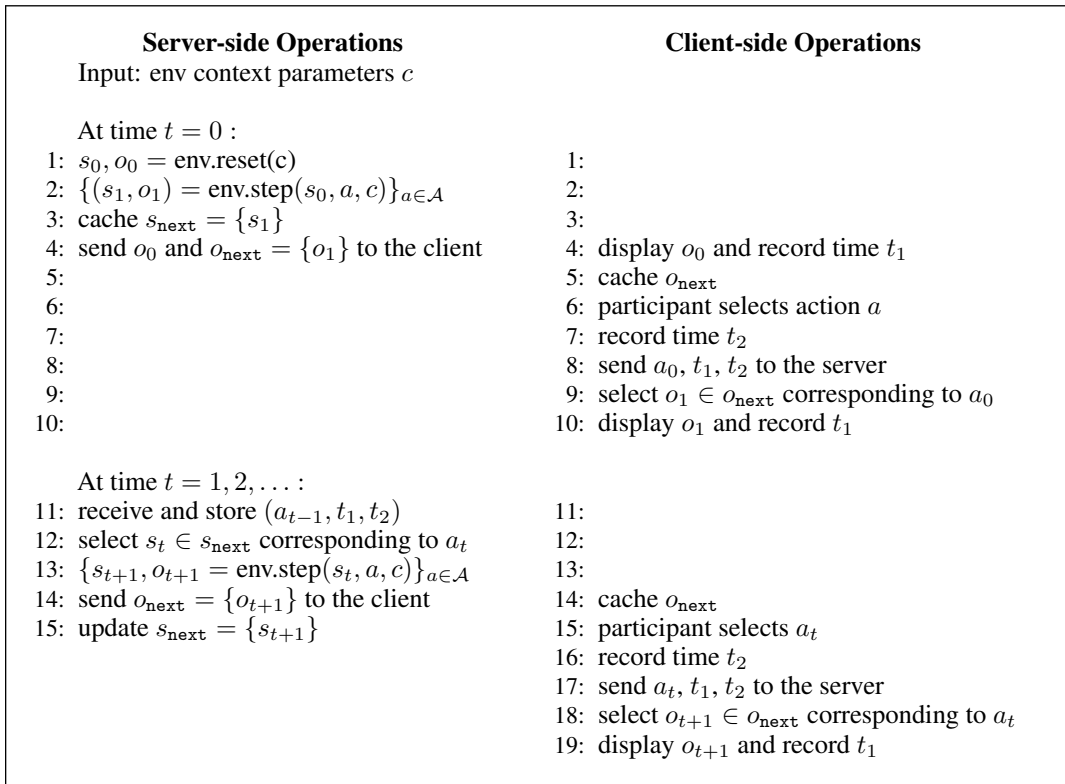


Figure 9: **Server-client Human-Environment Interaction Protocol**. Note that we omit displaying reward  $r$  due to space constraints.

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

## B DESCRIPTIONS OF STAGE TYPES

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

1. Stage: used to display instructions or information to a participant.
2. FeedbackStage: used to collect information from participants. Typically involves an interactive screen that does *not* interact with the environment.
3. EnvStage: used to interact with an environment. It takes as input an environment and environment parameters. We describe how NiceWebRL uses this abstraction to have a remote server-side program display images to one’s local web-browser client in figure 9.

We present examples of each in figure 10.

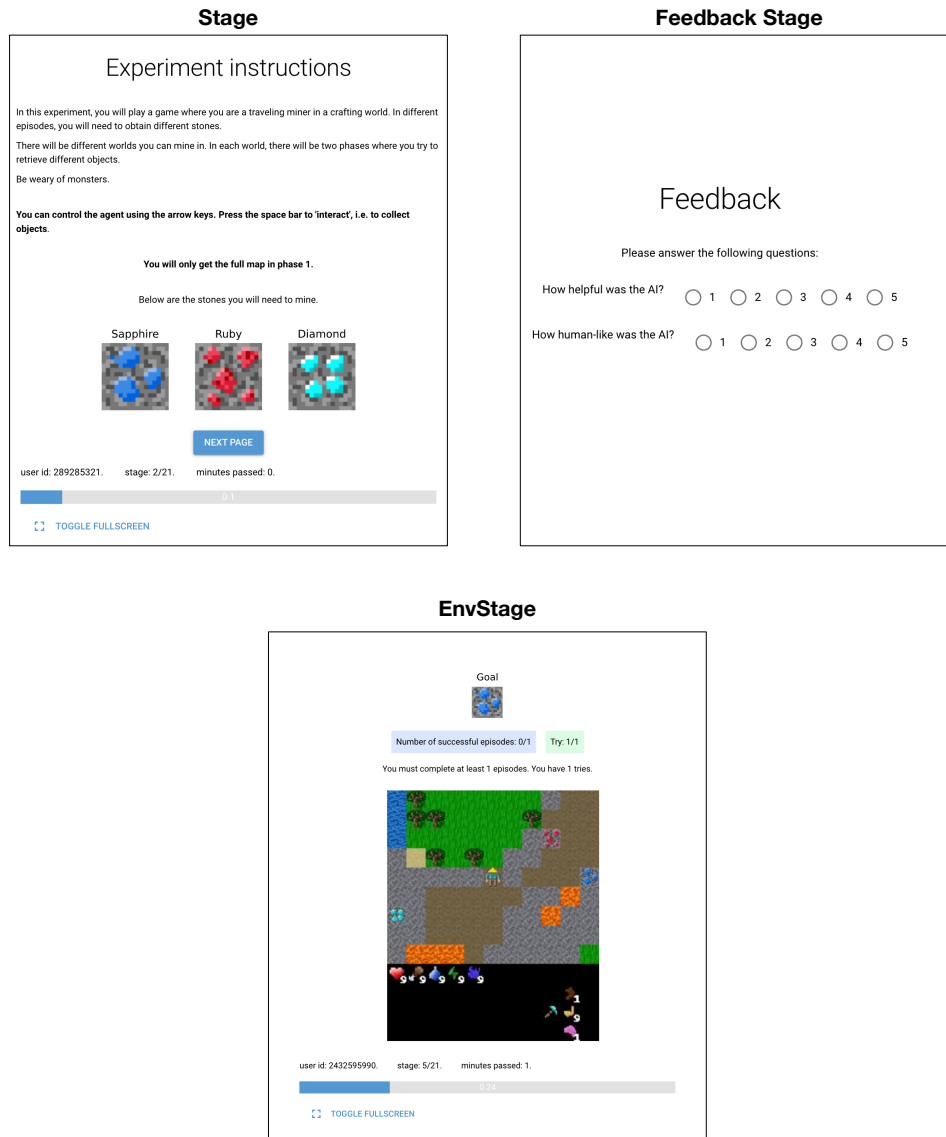


Figure 10: Examples of different kinds of stages.

## C ADDITIONAL RESULTS

## D COMPUTING RESOURCES

For details on case study 1 or 2, please see (Carvalho et al., 2025) or (Jha et al., 2025), respectively. For case study 3, experiments were conducted using computing infrastructure from the [fly.io](https://fly.io) platform with the “performance-2x” configuration. This is a machine with 4GB of RAM. The

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

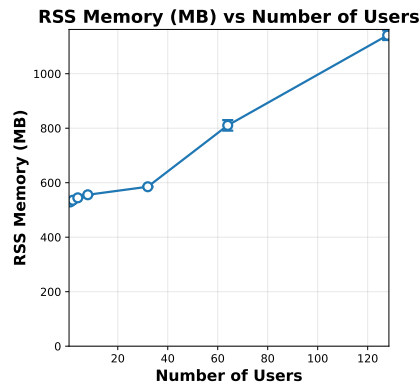


Figure 11: Memory usage vs. number of users

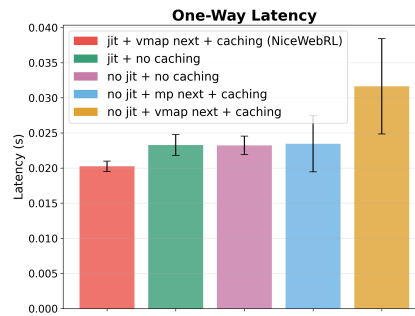


Figure 12: Round Trip Latency

machine had no GPU. Even in this setting, Jax’s compilation features provide a significant speed up to environment computation.

## E HUMAN SUBJECT EXPERIMENT DETAILS

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 example. Participants were compensated \$4 for completing the task. The average task completion time was 23.33 minutes. At the beginning of each experiment, the participants provided demographic information (age and gender, coded as male or female).