1 We thank the reviewers for their interest in the paper and their constructive feedback. In the following, we address the
2 outstanding comments and outline our plan for improving the paper.

3 **Comparison to Facebook's Horizon platform (Reviewer 3).** Face-
4 book's Horizon platform plays an orthogonal and complementary role
5 to Park. In Figure 1, we show where Horizon and Park reside in a
6 typical RL workflow diagram. As shown, Horizon implements several
7 RL algorithms and exposes APIs to developers to integrate RL training
8 and maintenance into their systems. By contrast, Park is a research
9 platform that provides a common interface to connect to a wide range
10 of computer system environments. In fact, Park's primary contribution
11 is this environment suite that can serve as benchmarks for RL research
12 in computer systems. In principle, a developer can use Horizon to
13 try out existing RL algorithms or develop new RL algorithms for the
14 system environments in Park. We will add this comparison to Horizon
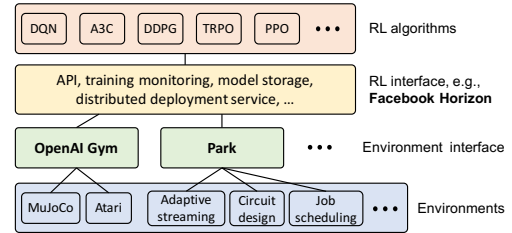15 to the introduction (where we compare with OpenAI Gym) to better position the Park platform.



**Figure 1:** RL workflow. Park uses a common interface to connect to twelve computer system environments.

16 **"Correctness" of system MDP (Reviewer 3).** Thank you for raising this point. Formulating the MDP is an important,
17 problem-specific step in applying RL to any system. In fact, some of Park's systems have multiple possible MDP
18 formulations, which can affect the efficiency of RL training (e.g., for the Tensorflow device placement problem, an
19 "incremental improvement" formulation (Addanki et al., 2019) is more efficient than the "one-shot" MDP (Mirhoseini
20 et al., 2018)). While the main contribution of Park is not defining the MDP for each system environment, we thought
21 carefully to pick an appropriate MDP formulation for each problem. Our guiding principle was to provide the RL agent
22 with all the information and actions available to existing baselines schemes in that environment, such that the agent could
23 at least express existing human-engineered policies. In most cases, the MDP formulations are straightforward and their
24 correctness is self-explanatory. However, some are more subtle (e.g., the Spark scheduling and TF device placement),
25 and in these cases we adopt the formulations from prior work. In the final version of the paper, we will discuss the MDP
26 formulation for each system in more detail in the appendix, pointing out any subtleties with reference to prior work.

27 **How extensible is Park's interface? (Reviewer 3).** We will use lines of code added for each system to quantify how
28 much work it takes to add new environments with Park's interface. For example, the adaptive video streaming environment
29 uses 324 additional lines of code for Park interaction. To add a new environment, as we commented in §4, the developer
30 only needs to specify the state and action space, and make them accessible through RPCs. We will add the additional lines
31 of code for each environment as a reference in the appendix.

32 **Details of baseline schemes (Reviewer 3).** In our experiment, we compare against existing heuristics specifically
33 designed for each environment. In the appendix, we describe the details for the compared heuristics, the overview of their
34 implementations and the intuitions for how they work. Due to space limit, we only provide the names of these baseline
35 schemes in Figure 4. We will add a forward pointer to the appendix in the experiment section for the details.

36 **More detailed discussions on the experimental results (Reviewer 2 and 3).** We agree with Reviewer 2 about trimming
37 §4 (mostly the "Real System Interaction Loop" part) and adding more details from the appendix to §5. As suggested,
38 we will highlight and comment on the results where the RL policies are not stable. We will also describe our model
39 tuning strategy, a form of grid search, which is performed separately for each environment. We believe that the instability
40 observed in some of the environments are due to the fundamental challenges discussed in the paper, not poorly tuned
41 training procedure. For example, the policy in Figure 4(h) is unable to converge smoothly partially due to the variance
42 caused by the cache arrival input sequence (i.e., the challenge in §3.2). Also, to clarify, the horizontal lines in Figure 4
43 correspond to fixed, non-learning baseline schemes in each system. These policies are not trained. We include their
44 performance level as a constant line for reference in Figure 4.

45 **Standardized interface for testing environments (Reviewer 1).** In our experiments, we test the RL agent on unseen
46 scenarios by loading a different set of traces. For example, the adaptive video streaming environment loads unseen
47 network traces and the Tensorflow device placement environment tests with unseen models (i.e., unseen computation
48 graphs). The users have the flexibility to create different scenarios to test the robustness of a learning model. However, we
49 agree that standardizing the training/testing scenarios (e.g., fixing which traces to use for training/testing) would provide
50 a better common ground for comparison. We will provide flags to launch default training and testing scenarios in our
51 implementation.

52 **Elaborate on future directions (Reviewer 1 and 3).** In the conclusion section, we will highlight the future research
53 direction by summarizing the challenges discussed in §3. In addition, the rightmost column of Table 2 outlined the
54 challenges for each environment. We will make them stand out more when we discuss the table in §4.

55 **Adding reference (Reviewer 2).** We will add the citation for adaptive video streaming. Thanks for pointing this out!