

1 We appreciate the reviewers’ thoughtful comments. Due to the space limit, we only address the most salient con-
2 cerns/questions in the reviews. We will do our best to address the reviewers’ other points in the paper’s next version.

3 **Global responses to some common concerns among reviewers**

4 *Relationship to prior work (R2, R3):* While mirror descent has been explored in RL settings before, our paper stands out
5 from all existing work in making three key innovations: (i) using a hybrid representation, a combination of a program
6 selected from a generic programming language and a neural network, for policies; (ii) implementing the projection
7 operator using combinatorial program synthesis; and (iii) a thorough analysis of the consequences of doing gradient
8 descent only in the neural component of hybrid policies, which introduces bias in the gradient operator.

9 To maintain focus on our main contributions, we used relatively basic versions of the two key components of the IPPG
10 meta-algorithm: mirror descent and program synthesis. However, each of these components can be improved further in
11 future work. In particular, like some of the papers pointed out by R2, one can use extensions of mirror descent such
12 as composite objective mirror descent [Duchi et al., JMLR11]. One can also implement the projection operator using
13 learning-accelerated program synthesis techniques, such as those pointed out by R3. We will include further discussion
14 of these possible directions in the final version.

15 *Incorrect link to code and reproducibility (R1, R3):* We most sincerely apologize for the error! We accidentally created
16 two Dropbox folders during submission and uploaded the code to the incorrect folder (due to over-excitement with
17 our paper :-)). We have verified that the code is now available at the included link. We have also provided a Docker
18 container with all dependencies and the TORCS simulator installed, to ease reproducibility.

19 **Responses to other specific questions/concerns**

20 *R1 on why projection is easier than direct search in program space:* Intuitively, this is because we frame projection as an
21 imitation learning (IL) problem to imitate a neural net, in contrast to the full RL problem of searching over programmatic
22 space. It is widely acknowledged that IL tends to be easier than general RL, which needs to solve planning under
23 possibly long horizon (cf. Deeply AggreVaTeD paper by [Sun et al., ICML18]). This intuition is corroborated by
24 related work on programmatic RL [Verma et al., ICML18], which shows that direct search over programs often fails to
25 meet basic performance objectives (for example having a TORCS car finish a lap in any amount of time).

26 *R1 on our benchmarks and baselines being simple:* We respectfully disagree regarding the difficulty level of the
27 TORCS task, our primary benchmark. TORCS has several race tracks, many of which are challenging domains for RL
28 due to the continuous action space and long horizon. TORCS is also among the few standard benchmarks that allow
29 studying generalization properties, i.e., where the training and testing environment can differ. Of the popular policy
30 gradient algorithms, only DDPG has been shown to complete laps in TORCS in the tracks that we considered. Two
31 other algorithms, TRPO and PPO, cannot find policies capable of completing laps in most of these tracks [Cheng et al.,
32 ICML19]. In Table 1 in the paper, we show that IPPG learns to drive on some tracks (Ruudskogen and Alpine-2) where
33 even DDPG fails to learn a sensible policy. Thus, we do compare with (and show an improvement over) state-of-the-art
34 deep RL. Finally, we compare IPPG with VIPER and NDPS, the most relevant techniques for programmatic RL.

35 *R3 on assumption for Theorem 4.2:* This simplifying assumption is contained to Section 4.2 of the paper, for the purpose
36 of providing a clean finite-sample analysis under vanilla policy gradient. The subsequent Section 4.3 addresses the
37 setting when this assumption does not hold. Our analysis and experiments do not require this assumption in general.

38 *R3 on initial policy:* IPPG can be initialized with a neural policy, learned for example via DDPG, and thus can be made
39 to learn “from scratch”. This amounts to running the program synthesis projection as a pre-training step (Lines 4-6 in
40 Alg. 1). On most tracks, this gives results comparable to those we have presented. However, note that IPPG with a (very
41 simple) hand-crafted prior can sometimes finish tracks on which DDPG fails. Naturally, this is not possible when we
42 invoke the program synthesis projection during pre-training. Also, using a prior results in “safer” training (Figure 5).

43 *R3 on experimental setup:* In principle, IPPG and its theoretical analysis do not depend on the particular policy gradient
44 approach. However, there is a huge empirical difference between different policy gradient algorithms on TORCS. PPO
45 and TRPO actually cannot find policies capable of completing laps in most TORCS tracks we considered (also noted in
46 some previous papers). For this reason, we did not try using them as a component of IPPG in this paper.

47 *R3 on experimental results & reproducibility:* We have added a video for the TORCS benchmark, comparing IPPG and
48 NDPS at 4 snapshots during training, to the Dropbox folder for the code. The video shows that IPPG is safer during
49 training, shows better iterative improvements, and recovers more gracefully from crashes than the NDPS agent. The
50 tracks are deterministic, but the initial state (starting position) of the car in each race is chosen randomly from a fixed
51 set. While we did not calculate confidence bounds in the submitted version, a key benefit of hybrid policies (over deep
52 RL) is that they lead to lower-variance learning [Cheng et al., ICML19], and as a result, IPPG is likely to have tighter
53 confidence bounds than DDPG. We will substantiate this point using concrete numbers in the final version of the paper.