

1 We warmly thank the reviewers for their careful reading of the paper and for their feedback which will help us to
2 significantly improve the final version. We agree with the reviewers that we can present our work more clearly, that
3 some additional experiments will provide valuable information and that we can give a more accurate perspective about
4 our future work and the current limitations of AlphaNPI.

5 **1 Clarity**

6 We thank the reviewers for their insightful suggestions to improve the clarity of our paper.

7 In Section 3, we will focus on clarifying the key principles of the algorithm, and the way the hierarchies of programs
8 are exploited. We will add a short description of Algorithms 1 and 2, the MCTS simulation and tree search, to improve
9 readability without necessarily referring to the appendix. Besides, as requested by Reviewer #2, we have already
10 worked on a new figure that will illustrate better the main mechanisms of AlphaNPI.

11 We will also add a paragraph to recap with more details the ways the π^{mcts} output is computed and used. In a nutshell,
12 for each observation e , the π^{mcts} vector, computed with AlphaZero via a guided MCTS simulation, corresponds to a
13 distribution over actions. This is assumed to be better than the network output for this same observation, where better
14 means closer to the one generated by an optimal policy. Therefore, for all observations in the generated traces, we
15 minimize a distance between the distribution represented by π^{mcts} and the current network output distribution to adjust
16 the model weights. The exact computation of the π^{mcts} vector can be found in Appendix A.5, which we will improve.

17 The comments of the reviewers also helped us to realize that the role of M_{prog} could be better explained. We will
18 clarify the training and use of program embeddings, and the slight differences between our work and the original NPI.

19 Finally, Reviewer #1 comments that, compared to the work of Cai et al., a previous extension of NPI, we have traded
20 strong supervision of one kind for strong supervision of another. However, we would like to clarify that with our method,
21 the sizes of instances are limited but random during training. There is no ordering of instances, as the curriculum uses
22 only the hierarchy of programs. This hierarchy is not explicit in Cai et al. and in Reed & de Freitas, but it could be
23 easily inferred from the observed execution traces. In our work, traces are not given in advance, they are generated
24 via interactions between the system and the environment. Since execution traces somehow require the problem to
25 be already solved, replacing them by a hierarchy of programs is, in our opinion, a significantly weaker supervision
26 requirement. In the final version of the paper, we will make this supervision difference clearer and more explicit.

27 **2 Additional experiments**

28 We agree with Reviewer #3 that training curves will be a useful addition to our paper. We generated a set of curves to
29 include in the final version, which show how the progress on non-elementary programs starts once a good performance
30 has been reached with all the programs of lower level in the hierarchy. It can be observed that, for the BUBBLESORT
31 program, the performance converges towards almost 100% of success after approximately 250 iterations, which
32 corresponds to the training on $250 \times 20 = 5000$ traces. In the original NPI work, the network was trained on 1216
33 BUBBLESORT execution traces, which represents much more data as each trace contains also sub-traces corresponding
34 to all the executed sub-programs, while in our setting traces are limited to the actions of a single program. Yet, our
35 results exhibit better generalization properties. In future work, we will estimate more precisely the sample efficiency
36 and generalization of AlphaNPI.

37 Regarding additional baselines, we will include the results of Cai et al. and Xiao et al.

38 **3 Future work**

39 The necessary verifiability of post-conditions is a current limitation of our work, as in some contexts post-conditions may
40 not be easy to check. However, without post-conditions, new sub-programs would have to be defined during training,
41 and the discovery of relevant sub-programs is a very difficult challenge. Removing the requirement of post-conditions
42 while keeping the global program easily interpretable would probably be a very difficult extension to our work. Instead,
43 in immediate future work we intend to remove the need for the hierarchy of program levels, resulting in the need for
44 more exploration and a less immediate curriculum learning strategy. We intend to use curriculum learning based on
45 learning progress (LP) to overcome the need for an explicit hierarchy given in advance. With this strategy, the agent
46 should be able to focus its training on low-level programs which are easier to learn and thus provide high LP before
47 attempting training on higher-level ones.