

1 **All reviewers - Comparison with Ellis 2018.** A crucial component of Ellis 2018 is the ability to incrementally render  
2 specifications inferred from the input image and measure similarity with that image. This similarity is used as a  
3 Bayesian likelihood for specifications sampling. However, this is not applicable to our setting as there is no direct  
4 equivalent to the incremental rendering process. Our approach offers more generality for neuro-symbolic systems.  
5 **Action/perception primitives.** In Section 2.2, we opted for an abstract definition of action and perception primitives as  
6 these are domain dependent. For instance, Vizdoom actions include "move forward", "turn left" and "attack". Vizdoom  
7 perceptions include "is there a monster?" or "is a revenant in target?". We will add a detailed list in the paper.

8 **Reviewer 1 - Discarding single tokens instead of entire traces.** This is a promising direction for future work to  
9 improve sample efficiency. Yet there is one fundamental difference with Barman 2010. Angelic programming allows  
10 for uncertainty in the program, while the approach proposed by the reviewer introduces uncertainty in the specifications.  
11 We agree that there is a conceptual connection, but it is not obvious to us whether both problems are equivalent.  
12 **Comparison to Chen 2019: Execution-guided neural program synthesis.** We do not compare because the synthesis  
13 tasks are different. To cite Chen 2019 (about demo2program) "Sun 2018 propose to synthesize the program from  
14 demonstration videos, which can be viewed as sequences of states. In such a problem, all intermediate states can be  
15 extracted from the videos. On the contrary, in the input-output program synthesis problem studied in our work, the input  
16 to the program synthesizer provides only the initial state and the final state." Chen 2019 infers likely intermediate states  
17 in order to simplify synthesis, while in our case they are part of the specification. This makes the tasks incomparable.  
18 **Related work.** We will refer to Solar-Lezama 2008 for the concept of program sketching, and improve the dichotomy  
19 between "rule-based" and "statistical" methods. We will correct the reference to DeepCoder.  
20 **Progressively decrease cost bound until problem is unsat.** We have implemented an alternative decremental scheme,  
21 preliminary results do not show a noticeable performance gap. We thank the reviewer for the literature references.

	Accuracy (%)	Karel			ViZDoom		
22 PLANS (best configuration)	91.6	53.9	34.2	88.0	65.5	58.8	
Decrement cost bound (best configuration)	92.1	53.8	34.5	87.5	64.8	57.8	

23 **Reviewer 2 - Rigidity of the approach.** The filtering technique developed to deal with missclassified perceptions  
24 primitives does not make domain specific assumptions and is potentially applicable to similar neuro-symbolic systems.  
25 **Cost of iterative search over multiple attributes.** A possible solution to this problem is to directly encode the cost  
26 function over programs in the solver. This is not possible in Rosette, that we chose for its flexibility in encoding DSLs.  
27 Future work could attempt to build on different solvers integrating cost functions and compare performance.  
28 **Impact to the NeurIPS community.** We believe that developing efficient neuro-symbolic systems is crucial to the  
29 machine learning community. Dealing with specification uncertainty is one of the major challenges involved.  
30 **Softmax scores vs. ensemble methods.** Our method is orthogonal to the choice of uncertainty measure. Ensemble tech-  
31 niques provide a good measure of uncertainty, but there is a high computational cost involved with training sufficiently  
32 many members. This is why we followed prior work using softmax response, which yields good performance.

33 **Reviewer 3 - No thresholds for Karel.** There is no need for filtering as the network reliably infers the specifications.  
34 **Perception primitives.** No semantic understanding of the perception primitives is required for our method. Besides,  
35 our filtering technique introduces some resistance to noisy training labels as a side effect. Concerning scalability, our  
36 empirical observation is that program length and control-flow depth matter more than number of primitives.  
37 **Encoding of actions/perceptions into I/O specifications.** High-level semantics of the encoding are indeed intuitive.  
38 The technical challenge lies in implementing these semantics accurately in the Rosette solver based on Racket.  
39 **Technical contributions.** The ability to deal with uncertainty in specifications has been repeatedly described as a  
40 crucial aspect of neuro-symbolic systems, since traditional symbolic solvers are extremely vulnerable to noise. We  
41 addressed the key technical challenge of dealing with misspecifications in a generic way.  
42 **Generality of the action/perception framework.** It is a very natural way of describing agent-environment interaction.  
43 An important direction for future work is to generalize to arbitrary combination of perceptions in conditionals.  
44 **Advantages of removing program supervision.** (i) action/perception labels are much easier to acquire. (ii) it removes  
45 the need to train on several different demonstrations of the same program. (iii) it drastically reduces training cost.

46 **Reviewer 4 - Experimental comparison with Raychev 2016.** This is problematic, because the algorithm in Raychev  
47 2016 deeply modifies the traditional program synthesis process. It was only evaluated on bit-stream programs, and has  
48 not been integrated in a general purpose solver yet. In contrast, our filtering algorithms build on top of state-of-the-art  
49 solvers, and are simple to implement. We have contacted the authors of Raychev 2016 and agree on additional  
50 advantages of our approach in a neuro-symbolic setting (i) it focuses on satisfying the most confident samples, which  
51 makes finding the correct program more likely. (ii) sorting the samples by confidence reduces overall time complexity.  
52 **Program size - scalability.** Programs are composed of up to 43 tokens, and we observe balanced time measurements  
53 between neural network inference and symbolic solver calls. We believe that parallelism inside solver calls could help  
54 scaling to larger programs with more control-flow statements. This would be fair as the network runs on a GPU.