1 **(R1) Extension to POMDP domains.** We appreciate the insightful question. We are indeed planning to generalize
2 RPN to POMDP domains. The POMDP setting is challenging for both classical planning and modern learning-based
3 planning mainly because it requires *active sensing* steps [19]. Consider the example goal `InHand(Apple)` when the
4 apple is in a closed fridge. To learn to plan for such a goal within the RPN framework, we may assume that the
5 video demonstration includes sensing intermediate goals, e.g., a human demonstrator searches for the apple at a few
6 locations including the fridge. Given such demonstration data, we can train an RPN to predict these sensing goals (e.g.,
7 `IsOpen(Fridge)`) as the *preconditions* of achieving the final goal. However, more principled approaches [19] may
8 require extending RPN to explicitly reason about uncertainty of the state, which is out of the scope of this paper. We
9 will revise the current running example to avoid possible confusion as suggested by the reviewer.

10 **(R1) Task with goal `IsClose(Fridge) AND InHand(Apple).`** We agree with the reviewer that this example presents
11 a challenging case, where solving the precondition `IsOpen(Fridge)` of the final goal `InHand(Apple)` undoes the
12 subgoal `IsClose(Fridge)`, which may cause the planner to oscillate between the two subgoals. Common solutions in
13 classical planning require considering the global planning solution to resolve the correct subgoal order [1]. However,
14 RPN does not assume a global planning setting. Solutions within the RPN framework can be either 1) annotate the
15 dependency IsClose(Fridge) -> InHand(Apple), which leads the planner to solve `InHand(Apple)` first regardless the
16 state of `IsClose(Fridge)`, or 2) consider the final task goal as a subgoal block instead of individual subgoals. We
17 will revise the paper to include a detailed discussion of this challenge.

18 **(R1) Source of supervision.** We note that the form and the granularity of the supervisions in this work are very similar
19 to that of instructional videos [36], which are extensively used in the video understanding community. Our eventual
20 goal is to scale RPNs to real-world tasks by learning from these instructional video datasets.

21 **(R2) Role of learning and difference from classical symbolic planners.** While RPN is inspired by classical planners,
22 the additional assumptions and the extra information required by classical planners prohibit a fair comparison to RPN.
23 As noted in the main text (line 5, 24, 91, 122-127), classical symbolic planners search for plans based on *planning*
24 *domains*, which define high-level actions as the symbolic rules of the task domains. The main challenge for symbolic
25 planners is thus not only to estimate entity states but even more to hand-define such rules for real-world domains at scale.
26 The key contribution of our paper is a novel neural network architecture that learns to plan from task demonstrations
27 without the need of a planning domain or explicit state estimation. Each component of the network directly takes
28 high-dimensional observations as input and encapsulates a part of the overall regression planning process (as described
29 in Sec. 4 of the paper). Furthermore, RPN works with an incomplete predicate set. For example, to plan in the kitchen
30 domain, a classical planner would require predicates such as `Graspable` and `Activatable` to define necessary rules,
31 whereas RPN learns to plan without these predicates (Appendix 3.2 includes a list of predicates used in the experiments).
32 We will revise the paper to further clarify these differences and the advantages of RPN.

33 **(R2) Application in more "realistic" settings.** First, we designed our evaluation environment to emulate a realistic
34 setting: our Kitchen 3D environment features visual scenes rendering meshes scanned from physical objects and full 6-
35 DoF object pose variations. Second, while the chosen task domains facilitate symbolic definitions, this type of domains
36 is widely adopted by prior works on real-world robot planning, especially in task and motion planning [12,13,19].
37 However, while these works assume access to symbolic planning domains, our paper presents a data-driven method that
38 learns to perform task planning without a planning domain. Furthermore, we note that because RPN represents the
39 planning steps with generic function approximators, we may use other representations that exhibit similar compositional
40 structures. As noted in the paper (line 345), we plan to extend RPN to work with hybrid symbolic-geometric task goals,
41 e.g., `On(A, B)` with a specified 6D pose. In such setting, the generated plans could be fed to an off-the-shelf motion
42 planner to enable task and motion planning. We will update the paper to highlight the relevance.

43 **(R4) Comparison to Rocktäschel and Riedel.** We thank the reviewer for the highly relevant pointer. Rocktäschel and
44 Riedel [2] proposes a neural theorem prover (NTP) that works with incomplete knowledge bases: it uses distributed
45 symbolic representations to resolve non-exact symbol matches and derive new logical rules. As noted by the reviewer,
46 NTP is related to our approach in that, similar to our regression planning mechanism, it employs a learnable backward-
47 chaining algorithm that recursively decomposes a proving goal into subgoals, and derives new proving states until the
48 terms in the state are supported by facts in the knowledge base. Besides the difference in the observation spaces, the
49 main conceptual difference between NTP and RPN is that the goal of NTP is to predict if a statement is true, whereas
50 RPN focuses on the correctness of the proof (plan) – whether the next intermediate goal in the plan is reachable and can
51 ultimately lead to the final goal. Hence, in principle, NTP may derive the correct answer from a wrong proof, whereas a
52 wrong plan (proof) would lead to failure in robotic tasks. On the other hand, RPN and NTP are complementary in many
53 aspects. For example, RPN may employ a distributed symbolic representation to handle task goals specified with an
54 open vocabulary. We will include a discussion on the relevance and will consider the notations in [2].

55 **(R4) Details.** $e_t$ denotes the entity features at time step $t$. We will incorporate the comments in the next revision.

56 [1] D. Chapman, "Planning for conjunctive goals," *Artificial intelligence*, vol. 32, no. 3, pp. 333–377, 1987.
57 [2] T. Rocktäschel and S. Riedel, "End-to-end differentiable proving," in *Advances in Neural Information Processing Systems*, pp. 3788–3800, 2017.