1 We thank the reviewers for their thoughtful comments! We are encouraged that all reviewers voted to accept, finding
2 the paper to be creative and well-motivated [**R1**], extremely clear / well-written [**R1 R2 R3**], and that the technical
3 contributions of the paper – a novel model for language-conditioned instruction following – was judged to have
4 medium-high [**R1**] to high [**R2**] significance. We address specific concerns below (and will incorporate all feedback).

5 [**R1**] **Comparison to prior work.** Thanks for pointing out the relevant ablated comparisons in [4] and [7], which
6 achieve val-unseen success rates (SR) of 42% and 43.6% respectively. With additional hyperparameter tuning, we found
7 the results of our model increased to 33% SR / 30% SPL on val-unseen (up from 31% / 27%). Compared to these works,
8 we use less powerful ResNet-34 image features (instead of ResNet-152) and our model is much less dependent on the
9 structure of the navigation graph (see L36-41 below). We will update Table 2 with these comparisons instead of [1].

10 [**R1**] **Impact of agent heading.** Excellent point for which we will add discussion. Heading indicates important cues
11 about a trajectory which the model can leverage in the belief state. For instance, is is unlikely for an agent following the
12 true path to turn 180 degrees midway through (unless this is commanded by the instruction). Similarly, without knowing
13 heading, the model can represent instructions such as 'go past the table' but not 'go past with the table on your left'.

14 [**R1**] **Are the mapper-filter experiments in Sec 5.2 produced by training without the policy? Are the filter and**
15 **semantic map parameters updated while training with the policy in Sec 5.3?** Yes, to both. In Sec 5.2 we train
16 without the policy, using only the filter loss (KL-divergence between the predicted belief $b_{1:T}$ and the true state $s^*_{1:T}$). In
17 Sec 5.3 we train concurrently with both the filter loss and the policy loss (cross-entropy loss to maximize the likelihood
18 of the ground-truth target action). In both cases we train all parameters end-to-end (except for the pretrained CNN). We
19 have verified that the stand-alone performance of the filter is not unduly impacted by the addition of the policy.

20 [**R1**] **Is the true state sequence $s^*_{1:T}$ (L245) always the human trajectory, or does it include the exploration that**
21 **is done by the agent during training?** Yes it is always the ground-truth (human) trajectory. Recall that the filter
22 subcomponent is simply trying to estimate the path of an ideal agent following the instructions, given a partially-
23 observed semantic map $\mathcal{M}$. However, the agent's exploration does determine the content of the map (i.e., the map
24 contains the observations from the agent's explorations in the current episode). We will clarify.

25 [**R1**] **No 'direct' access to the instruction and map** We will add this qualification at L238.

26 [**R2**] **"Middling" performance compared to SoTA.** As in the paper, we freely admit that the proposed method does
27 not approach SoTA performance, but we agree with the reviewer guidelines in that *"Solid, technical papers that explore*
28 *new territory or point out new directions for research are preferable to papers that advance the state of the art, but*
29 *only incrementally."* Moreover, our novel formulation for instruction-conditioned goal-identification in VLN as state
30 tracking provides increased inspectability of agent beliefs and reasons over 2D space rather than just navigable nodes.
31 As such, it may be of broader interest outside the VLN task.

32 [**R2**] **Lack indications that this architecture will lead to SoTA.** The performance of the model improved with
33 additional hyperparameter tuning (see L5-9 above). Given the rebuttal period, it was impractical to incorporate the
34 orthogonal additions like RL training, data augmentation from a trained speaker model, or trajectory re-ranking which
35 have lead to much of the advancement at VLN.

36 [**R2**] **The action space of the agent... requires significant domain and task knowledge in the form of the structure**
37 **of the navigation graph.** We would like to clarify that the navigation graph is defined by the R2R dataset [1]. Previous
38 works [2-8] use LSTMs that operate directly over the nodes in this graph. Our approach is far less reliant on this
39 structure than prior work (at test time, our mapper and filter operate on arbitrarily discretized space, with no knowledge
40 of the navigation graph). We anticipate this could be an advantage for sim-to-real transfer (i.e., in real robot scenarios
41 where a navigation graph is not provided, and could be non-trivial to generate).

42 [**R3**] **How does the state rejection mechanism compare to the one of Weib et el. arxiv.org/abs/1511.06458. Can**
43 **one think of the language component as a refined prior?** We are not sure what the reviewer is referring by a "state
44 rejection mechanism" in our model, but are happy to hear more details in an updated review. With regard to Weib et al. -
45 we do not need to employ such rejection or importance sampling techniques because we use a histogram representation
46 for the posterior, not a set of weighted samples as in particle-based Bayesian state tracking. At a high-level, we agree
47 that language does define the set of expected observation and actions an agent should take throughout the trajectory.

48 [**R3**] **At L134 what are XY?** $X$ and $Y$ are the initial spatial dimensions of the semantic map in grid cells. In
49 experiments this dimension is 96 (refer supplementary). In practice the map could be dynamically resized if necessary.

50 [**R3**] **What does 'strong baseline' in abstract refer to?** We are referring to the LingUNet model and the Hand-coded
51 baseline that exploits dataset biases (Sec 5.2). Note that existing models [2-8] do not explicitly model the goal location,
52 and are thus not capable of predicting the goal location from a provided sequence of observations.