

1 We would like to thank all the reviewers for their insightful comments. All reviewers confirm the contribution of our
2 paper in providing a novel framework for improving the generalisation of VLN tasks in addition to a sound formulation
3 with wider range of applications. Please see below for specific responses.

4 **Response to Reviewer #1:**

5 **Relation to Mixup:** Mixup is not directly applicable to VLN since (1) it is sequential in nature, (2) an interpolation of
6 state-action from one trajectory to another may lead to catastrophic difference in the objective. Our approach intervenes
7 in the visual features to simulate the agent’s behaviour in a counterfactual environment, where the agent still has to
8 follow the same instruction and sequence of actions.

9 **Difference with [34]:** It is the closest approach to ours and general differences are highlighted in L89-93. Specifically
10 for counterfactual distribution learning, similar to theirs, our approach creates counterfactual samples that are hard for
11 the agent to follow with minimum interventions. However our approach is different to [34] in (1) ours is formulated
12 to minimise the expected difference of the model on observations and the intervened samples while [34] resorted to
13 importance sampling that could have a large variance (in fact it was not as successful in VLN); (2) ours is sequential;
14 (3) our instructions have actions as opposed to vision and language only (4) we had a model of a speaker to incorporate.

15 **Response to Reviewer #2:**

16 **Prior distribution and random noise:** The hyper-parameters of the prior distribution have been selected through a
17 simple grid search before applying the counterfactual distribution learning. Additionally, Environment Dropout [11]
18 can be considered as a random noise injection which has been reported in Table 1 and 2.

19 **Improvements:** Generally, the reported improvements are indeed significant (around 4% in success rate and SPL),
20 considering the fact that R2R task has been explored extensively in recent years so that even large-scale self-supervised
21 pre-training [44] gains less than 4% improvement. By adding an insightful though simple counterfactual learning
22 process, we gain further 1 – 2% improvement on top of the improvement gained using the prior (2 – 3%), which is a
23 significant improvement on top of a model that enjoys a better generalisation than other baselines.

24 **Counterfactual Learning:** For the counterfactual distribution learning, we sample two pairs of real trajectories and
25 only use language instruction of the first one. We will clarify this in the camera-ready version. In L167, we are
26 approximating the marginalisation of \mathbf{u} by adjusting the variable generated from the prior, instead of relying on costly
27 methods like MCMC or a variational lower bound. It is also worth mentioning that as stated in L180 and according
28 to Eq. 9, minimum edit (intervention) happens when \mathbf{u} is close to one. In addition, for L194, please note that as
29 stated in L178, $p(\mathbf{u} | \tau, \mathbf{c}) \propto p(\mathbf{u})\tilde{\pi}_\theta(\tau | \mathbf{c}, \mathbf{u})$. Therefore, $p(\mathbf{u} | \tau, \mathbf{c})$ is maximised when \mathbf{u} is close to one. Rather than
30 integrating \mathbf{u} out, we choose the most likely sample corresponding to a counterfactual. It should be noted that Figure
31 2.a in supplements is revealing the general SCM in VLN. Starting from a prior distribution, this variable is learnt
32 conditioned on the inventions on the real observation and then is used for the counterfactual generation (Figure. 2.b).

33 **State variable:** In our experiments, s_t is the hidden state of the RNN model and the parameters are omitted for brevity.

34 **Training setting:** We freeze the speaker model during the counterfactual distribution learning. Additionally, as stated
35 in L231, following the same setting as [11] and [8] for a fair comparison, we use teacher-forcing during IL.

36 **Response to Reviewer #3:**

37 **Computational cost:** Thanks for mentioning an interesting point. The computational cost of our proposed method
38 highly depends on the number of iterations for learning the counterfactual distribution. Higher values result in better
39 approximations of the exogenous variable while increasing the training time. We undertook several experiments for
40 finding a point of best trade-off and found that even few iterations (5 as reported in Section 3.1 in the supplements)
41 could contribute to a good understanding of the variable and therefore better results.

42 **Sampling for RL:** As indicated in Algorithm. 1, after finding the optimum value of \mathbf{u} using $\{\tau, \mathbf{c}, \tau'\}$ and Eq. 11, at
43 each step of the policy rollout, we intervene the observation using the learned \mathbf{u} and τ' .

44 **Response to Reviewer #4:**

45 **Computational cost:** Please refer to the response to reviewer #3.

46 **The effects of counterfactual data:** We have demonstrated some trajectories revealing the difference between the
47 baseline and our model over the same starting point and the same instruction in supplements. In addition, using our
48 approach we observe a significant improvement in the test set that highlights the positive effect of our approach.

49 **State variable:** Please refer the response to reviewer #2. We will clarify in the camera-ready.

50 **Equation 11:** Thanks for pointing out the typo. We will replace a with a_t in the camera-ready version.

51 **Algorithm 1:** In fact, we use the optimised \mathbf{u} to generate a new τ based on Eq. 9 for IL gain calculation. We haven’t
52 repeated the counterfactual generation for brevity of the algorithm.

53 **Equation 4:** As is mentioned in L182-185, for simplicity and efficiency, \mathbf{u} is marginalised out from the Eq. 4 and is
54 approximated using Eq. 11.