

1 **[REVIEWER #1] Long-term dependencies and generative likelihood:** We computed the estimated negative log
2 probability on a test set of Bouncing ball dataset using importance sampling with the learned encoders (500 sam-
3 plings). The results are RSSM: 3.518 and HRSSM: 3.601. Although it is true that we could expect the proposed
4 model to capture better long-term dependency, we would like to emphasize that our main focus in this work is ob-
5 taining *structure* which is interpretable, stochastic, and temporally hierarchical. As such, our main metric is also
6 the quality of the structure, not the accuracy of the generation like language modeling. We are actually satisfying
7 with the fact that we obtain such structure without sacrificing the likelihood performance (of course, it would have
8 been even better with better likelihood performance.) This seems somewhat similar to the fact that discrete latent
9 variable generative models usually do not perform better than its continuous counterpart, but more interpretable.
10 It is in fact not clear whether this learned structure should improve the likelihood
11 performance as well because, unlike other architecture learning problems like
12 NAS, our problem imposes particular temporal hierarchy structures, and thus
13 severely constraints the model space. **Investigation over latent spaces:** We
14 agree. We will add more analysis such as Fig 1 where the temporal abstraction
15 latent z has the subsequence-level context such as *color*, *direction* and *length*. Fig
16 2 shows how the subsequence are generated from the same z (but different s, m)
17 and the *velocity* of ball and the *length* of subsequence can be varied.

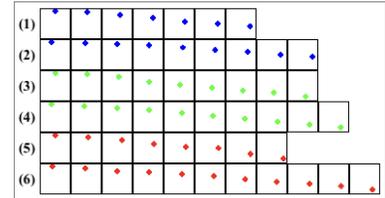


Figure 1: Subseq. with different z

18 **[REVIEWER #2] Interpretable results without color change:** Yes. Because a random color is selected at bouncing
19 from a color set of size 3 including the current color, actually it does not change its color with probability $1/3$. Even in
20 this case, we observe the model cut the segment at bouncing the walls. **Quantitative comparison to a baseline:** This
21 is described above (L1-13) with some explanations. **Maze without action-conditioning:** Yes, we initially also trained
22 the model without action-conditioning and it showed a very similar segmentation result. Due to space limitation, we
23 couldn't include this result, but we will add this result in the Appendix of the camera-ready. We agree that actions
24 and observations can both affect the structure of the segmentation, particularly in a complex way if they are not
25 consistent each other. **RSSM baseline:** RSSM is a single layer version of our model. It is implemented by using
26 the same architecture as the observation-level of our HRSSM and by removing conditioning on z and m (L118-119
27 in the paper). **Training variation:** The training curve changes rather highly during the early stages of the training
28 as the model searches for a stable temporal structure from $q(M|X)$. It would have been more stable if the temporal
29 structure was given or fixed like VHRED instead of learning it as we do. **Prior over segments:** Our prior is designed
30 to regularize $q(M|X)$ to avoid *over-* and *under-segmentation*. This is done not by explicitly changing M or $q(M|X)$
31 but by regularizing it through KL and generation terms. That is, if $q(M|X)$ assigns segments that exceeds the limit
32 defined by the prior, it will lead to lower ELBO via the KL term. We, however, agree that more explicitly controlling the
33 posterior class with the segment limits is worth to try. **Independence of M binary indicators:** We agree that giving
34 more structure and conditioning to $q(m_t|X)$ is also worth to try. Nonetheless, we would like to say that m_t is not fully
35 independent but is independent conditionally after observing X . Thus, we believe that, although it is somewhat indirect,
36 each m_t can still see the global temporal dependencies by observing the full sequence X .

37 **[REVIEWER #3] Forcing the model to produce a new subsequence at each time step:** Although this is used
38 during training, at test time, we only use the UPDATE operation at the z level without generating subsequences. In
39 experiments, we show that this sequence of z 's are good abstract representation of the future and we believe that this is
40 a principled way of performing jumpy imagination in the sense that the formulation is based on the (recurrent) abstract
41 state-space transition model. Also, like other works on imagination-based planning or RL agents, we assume that the
42 test time environment is similar to the training environment so that the learned environment model is useful at test
43 time. So, we do not expect the model to generalize when the test time environment is much different from training
44 distribution. That is, in a new environment, we should not rely on empirically-learned (inductive) imagination until
45 it collects and completes learning from the new environment. **Importance of stochasticity** In our experiment on the
46 navigation task, we observed that we can actually rollout multiple future imaginations from the same state by sampling
47 multiple rollouts. This would not be possible in HMRNN as its rollout is deterministic. This result is obtained when we
48 train the model without action-conditioning. With action-conditioning, we found the model uncertainty reduces and not
49 generate various futures, which is what can be expected. We will add this result in the camera-ready. **Outside image
50 domain:** We agree that applying the proposed model to other domains like text or speech signal would be interesting.
51 For this work, as our main focus was to apply it for agent learning, we put our priority on the image domain and the
52 planning task. **RL experiments which are directly comparable with prior work.** Due to limited time and as our
53 focus was on the structure learning which can also be evaluated by the planning agent without RL, in this work we could
54 not highly prioritize the RL experiment. We, however, agree that it will be helpful in making the paper more complete.
55 **Gumbel softmax:** We used temperature annealing (from $\tau = 1.0$ to $\tau = 0.5$).
56 Without annealing, it was unable to train the model as the entropy of $q(M|X)$
57 didn't properly decrease and showed meaningless random segmentation.

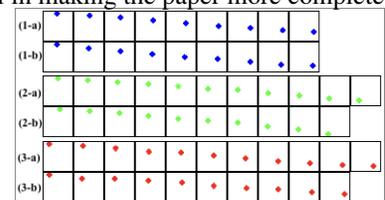


Figure 2: Subseq with same z