

1 We thank the reviewers for their thoughtful comments and suggestions. We performed several new experiments and  
 2 analyses to address the comments and will make the suggested changes to the main text. We also thank all reviewers for  
 3 taking the time to point out minor errors. Below, we address the reviewers’ comments individually.

4 **R1, R6: Additional analyses/ablations for  $\mathcal{L}_{\text{sparse}}$  and  $\mathcal{L}_{\text{sep}}$ .** We agree with Reviewer 1 that much of the novelty of  
 5 our work lies in the losses and training approach. We performed new analyses to show that  $\mathcal{L}_{\text{sparse}}$  and  $\mathcal{L}_{\text{sep}}$  are crucial  
 6 to the performance and stability of the model, both in terms of video metrics (especially FVD, Fig. A) and coordinate  
 7 tracking accuracy (Fig. B), on which downstream tasks depend. We will add these analyses to the main text.

8 **R1: Temporal consistency and “jumping” keypoints.** We initially experimented with using predictions from the  
 9 dynamics model as “prior” for the keypoint detector, but achieved better performance without enforcing temporal  
 10 consistency explicitly. Keypoints can indeed “jump” between frames, but we show in a new analysis (Fig. D) that the  
 11 VRNN partially smooths over such jumps: We displaced the location of one keypoint by  $0.5 \times$  image width in the  
 12 direction of the image center for one frame (Basketball dataset). The keypoint location inferred by the VRNN jumps by  
 13 less than  $0.5 \times$  image width in the perturbed frame and quickly recovers. Jumping thus seems to be a minor issue.

14 **R1: Did you observe training issues when combining a large  $K$  with  $\mathcal{L}_{\text{sep}}$ ?** Note that the optimal  $\sigma_{\text{sep}}$  (spatial  
 15 Gaussian radius of  $\mathcal{L}_{\text{sep}}$ ) is very small ( $\sigma_{\text{sep}} = 2 \times 10^{-3} \times$  image width for Human3.6M). At this  $\sigma_{\text{sep}}$ , the loss does not  
 16 interfere with initial training even for large  $K$ , but still prevents keypoints from collapsing onto the same image feature.

17 **R1: What is the size of the feature vector in CNN-VRNN?** We made sure to match the size of the feature vectors  
 18 of the models, such that the CNN-VRNN had  $K \times 3$  dimension at the narrowest point. Therefore, in principle, the  
 19 CNN-VRNN had the capacity to exactly recapitulate the Struct-VRNN structure.

20 **R1: Usefulness of KP structure for RL.** Our claim has since been confirmed by Kulkarni et al. (arXiv 1906.11883v1).

21 **R5: How is spatial structure imposed and why is it not sensitive to initialization?** See Jakob et al. [12] for how  
 22 the keypoint detector imposes spatial structure. A naive application of [12] to video indeed suffers from sensitivity to  
 23 initialization (see Figs. A and B, “no  $\mathcal{L}_{\text{sparse}}/\mathcal{L}_{\text{sep}}$  loss”). By adding  $\mathcal{L}_{\text{sparse}}$  and  $\mathcal{L}_{\text{sep}}$ , we achieve high robustness.

24 **R5, R6: Comparison to adversarial methods.** We note that we do compare to an adversarial method (“EPVA-GAN”,  
 25 Fig. 3, bottom right). A GAN loss could also be added to our model as a complementary objective; this is orthogonal to  
 26 our contributions. We agree that comparison to SAVP would be interesting, but we could not obtain results in time for  
 27 the rebuttal. We will include them in the final paper.

28 **R5: Why train keypoint detector and dynamics model separately?** We initially tried to train the model jointly  
 29 ( $\varphi^{\text{det}} \rightarrow \text{VRNN} \rightarrow \varphi^{\text{rec}}$ ), but found that the model learned an unstructured latent code, rather than spatially meaningful  
 30 keypoints. Presumably it was easier for  $\varphi^{\text{rec}}$  to reconstruct the image from an unstructured code, than for the VRNN to  
 31 learn the keypoint structure. Isolating the keypoint detector from the dynamics model solves this problem.

32 **R6: Why not apply B.o.M. sampling and  $\mathcal{L}_{\text{sparse}}$  to CNN-VRNN?** We did apply both to CNN-VRNN, but this yields  
 33 no gains because sample evaluation and sparsity are less meaningful in an unstructured space than in keypoint space.

34 **R6, R7: Is sample diversity an advantage? Are all samples good?** We agree with Reviewers 6 and 7 that we need  
 35 to expand the discussion of sample diversity. Fig. E below shows that even the samples with the lowest VGG cosine  
 36 similarity to ground truth are of high visual quality. For videos, see Sections 2 and 3 on the anonymous website (link  
 37 in original submission). We will add more examples and videos to the final paper. We emphasize that frame-wise  
 38 similarity to GT (e.g. VGG sim, PSNR and SSIM) does not meaningfully measure video prediction quality. For real  
 39 data, at test time, there is no single “ground truth”. Instead, there is an astronomical number of plausible futures that  
 40 are all consistent with the conditioning frames. We believe that most previous models dramatically underestimate this  
 41 diversity; our model comes closer to it. This is backed up by FVD, which is designed to measure sample diversity.

42 **R7: More fine-grained evaluation of object tracking.** We performed a new analysis of per-object tracking perform-  
 43 ance on Basketball (Fig. C). We identified two different failure modes: The basketball (yellow traces) has relatively  
 44 large tracking errors across all 10 model initializations, presumably because the dynamics of the ball are hard to learn.  
 45 On the other hand, Player 3 (pink) is tracked well in some and poorly in other model initializations, presumably because  
 46 the keypoint detector fails to recognize the light-colored object. We will describe these failure modes in the main text.

