

1 We thank all reviewers for their constructive feedback and for their time in creating well thought out reviews.

2 Below we address all raised concerns, namely we perform ablation studies of adding (i) 2nd-order ODEs and (ii) BNNs;
3 (iii) address more complex experiments and comparisons; and (iv) discuss the role of the KL and regularisation.

4 **A new 1st-order baseline:** We tested a new ODE¹VAE variant where the latent space is governed by 1st-order ODE
5 system. ODE¹VAE is similar to the NeuralODE [Chen et al 2018], except for having BNNs, and for NeuralODE placing
6 a variational distribution on initial value $q(\mathbf{x}_0)$, while ODE¹VAE models the posterior over full trajectory $q(\mathbf{x}_{0:T})$.

7 **[R1,R3] ODE¹VAE vs ODE²VAE:** We
8 performed a new comparison study of
9 ODE¹VAE against ODE²VAE on bounc-
10 ing balls dataset. The experimental setup
11 is kept the same, except that the number
12 of convolutional filters is reduced so that
13 the impact of differential function choice
14 becomes more apparent. Table 1 shows the
15 resulting MSE over 10 frame ahead pre-
16 dictions. Note that ODE²VAE models the
17 acceleration $\dot{\mathbf{v}}_t = \mathbf{f}(\mathbf{s}_t, \mathbf{v}_t) : \mathbb{R}^{2d} \rightarrow \mathbb{R}^d$

18 whereas 1st-order systems learn $\dot{\mathbf{z}}_t = \mathbf{f}(\mathbf{z}_t) : \mathbb{R}^d \rightarrow \mathbb{R}^d$. Results show that the 2nd-order dynamics results in far better
19 accuracy, even if the first order dynamics has more flops ($d = 50$). We will include ablation studies in the paper.

20 **[R1,R2] NN vs BNN:** Table 1 shows comparable performance of BNNs and NNs on bouncing balls. In order to
21 demonstrate the benefit of using a BNN, we repeat the CMU walking experiment with a NN differential function. The
22 MSE achieved by ODE²VAE-NN over three test sequences is 9.96, whereas ODE²VAE-BNN error improves to 9.43.

23 **[R2] Learning of BNNs:** Learning BNN is performed via mean-field variational approximation (simultaneously with
24 variational inference of the whole ODE²VAE model), where each weight and bias component has its own mean and
25 shares a global variance parameter. The ODE solver used in our experiments is fixed step Runge-Kutta for both NN and
26 BNN systems; hence NFEs are also the same.

27 **[R1] Comprehensive experiments:** Our model is suitable for sequential datasets, of which we demonstrated good
28 performance on motion capture data, bouncing balls experiments and on rotating MNIST. Conventional image datasets
29 such as CIFAR-10 or Celeb are not directly applicable for our model as they do not have an immediate dynamic
30 dimension. In this work we proposed the theoretical foundations of latent differential equations, and in future we intend
31 to explore video prediction application as separate work due to its daunting scope and complexity.

32 **[R2] Comparison to moving MNIST:** Moving MNIST is a dataset of digits bouncing off the walls of a box. Physical
33 interaction rules in bouncing balls dataset is more complicated because balls collide with each other, as well. In that
34 sense, inferring the dynamics in bouncing balls dataset is more challenging. On the other hand, MNIST dataset possibly
35 requires more powerful decoders, which we will consider as part of future work on video prediction.

36 **[R3] Missing NeuralODE baseline in rotating MNIST and bouncing balls:** While the public NeuralODE imple-
37 mentation worked as expected in the CMU walking experiments, we were unable to get NeuralODE model to work in
38 BOUNCING BALLS and ROTATING MNIST datasets. We included ConvNet architectures and tried these experiments
39 numerous times with different encoder/decoder hyperparameters and initialisations; however we always got fully black
40 frames as reconstructions. We believe the ODE¹VAE results instead to be informative enough to demonstrate inherent
41 limitations of 1st-order models, such as NeuralODE.

42 **[R1] Regularisation parameters:** The β and γ parameters weigh the regularising KL terms to be comparable to the
43 weight of the likelihood term (see e.g. "Fixing the Broken ELBO" paper). We choose to fix $\beta = |q|/|\mathbb{W}|$ to the ratio
44 between the latent space dimensionality q and number of weight parameters of the differential function $|\mathbb{W}|$, in order to
45 counter-balance the penalties. We chose $\gamma = 0.001$ by cross validation from [0,0.1,0.01,...0.00001].

46 **[R2] ODE²VAE-KL variant:** As correctly pointed out by the reviewer, all consecutive triplets in a sequence are
47 encoded. We then compute the KL divergence between encoder distributions and the state distributions induced by
48 ODE integration. This way, the entire sequence (rather than only the initial values) is utilized for encoder training.

49 **[R3] Long-term forecasting:** Long-term forecasting of non-linear dynamical systems requires an almost perfect
50 underlying dynamics model for the trajectories not to deviate. We regard "long-term" forecasting to be up around
51 20 frames ahead in bouncing balls, multiple cycles of walking, or a full rotation of MNIST numbers. We found out
52 empirically that NeuralODE can not forecast sufficiently, while the GPPVAE model interpolates states over time with
53 an RBF kernel with little extrapolation capability.

Table 1: Comparison of neural network (NN) and Bayesian neural network (BNN) ODE's with different latent dimensionalities on BOUNCING BALL experiment. Adding 2nd order momentum achieves superior performance, while BNN's have a smaller impact.

Model	Latent dimensions d		Test MSE	
	1st-order state	2nd-order momentum	NN	BNN
ODE ¹ VAE	25	-	45	43
	50	-	36	35
ODE ² VAE	25	25	26	27