

1 We thank reviewers for the relevant comments. We first address general questions and then give brief individual answers.

2 **On the necessity of GRU-Bayes and link to filtering methods.** GRU-ODE and GRU-Bayes have complementary
3 roles and should be used together. GRU-ODE integrates the dynamics of the hidden process $h(t)$ in time. It thus
4 provides the future estimated distribution of the observations. It computes our **belief** about the unobserved future time
5 series observations. Those projected distributions vary smoothly as they are driven by an ODE. But, as soon as any new
6 observation is available, our belief about the state of the process should be updated **instantaneously**. GRU-Bayes is
7 responsible for this and compares the prediction from GRU-ODE and the actual observation. This update necessarily
8 has to be discrete to accommodate for the new information arriving in packets (or sporadically). See figures 1 and 2
9 for illustration of this dynamic. This also motivates why sporadic measurements cannot be fed directly to the ODE
10 : observing a sample should make us update our belief immediately. A useful analogy is the Kalman filter, which
11 consists in a **prediction** and an **update** phase. In prediction phase, it propagates in time the predictions about the
12 distribution of the state (analog to GRU-ODE). In update phase, it computes a new estimate for the state probability
13 distribution conditioned on the new information (analog to GRU-Bayes). Yet, in contrast to the (extended) Kalman
14 filter, our approach is able to learn complex non-linear dynamics for the hidden process and is computationally cheaper.

15 **More related works context.** Other recent works have previously investigated the relationship between deep neural
16 nets and differential equations. They mainly focused on deriving better deep architectures motivated by the stability of
17 the corresponding partial differential equations (Ruthotto and Haber, *arXiv* 2018; Chang et al., *ICLR* 2019). Despite
18 their ODE motivation, those approaches aim at designing new **discrete** architectures and don't explore neural networks
19 parametrized ODEs as such.

20 **Point processes** (Mei and Eisner, *NeurIPS* 2017; Gunawardana et al., *NeurIPS* 2011) are intrinsically continuous as
21 they focus on **time-to-event** modelling. Continuous-time Bayesian networks (Nodelman et al., *UAI* 2002) address a
22 related problem where they frame events as state transitions. In contrast, our work aims at modelling continuous-time
23 **real-valued** measurements, not only events. Yet, our continuous modelling of the latent process $h(t)$ allows to easily
24 **jointly** model a continuous intensity function (e.g. by modulating the intensity function of a Poisson process with $h(t)$).
25 This joint modeling of continuous measurements and events was left for future work.

26 **Reviewer 2.** Some assumptions have to be made about the conditional distribution of the observations. However, this is
27 not very restrictive as a broad range of distributions have a tractable KL (e.g. Poisson, Exponential, ...).

28 The main difference with GRU-Discretized-Bayes resides in the continuity of the latent process (will be detailed in the
29 Appendix). We then considered this experiment as a continuity ablation study. Note that the continuity prior can be
30 tuned by rescaling time accordingly. Elongating time by a factor 2 would lead to a Lipschitz-4 prior.

31 Unlike sporadic measurements, the continuous measurements can indeed be directly fed to the GRU-ODE as suggested
32 at line 89 and would considerably improve the prediction of the model as vital signs are very strong predictors in critical
33 care. In this work we didn't include them (1) to compare equally with the other baselines which are not capable of
34 handling continuous inputs, (2) we re-use the same variable subsets as used in (Che et al. *Scientific Reports* 2018). Still,
35 we warmly welcome this suggestion from the reviewer as this would further highlight the capabilities of our model.

36 For fair comparison, all compared models were *made* probabilistic if not already (i.e. they all output log-variance).

37 **Reviewer 4.** We added the following (at 1.264) : “We assessed the performance of the models by creating 5 different
38 folds, each consisting of training (70%), validation (20%) and a left out test set (10%). Those folds are **reused** across
39 compared architectures. For each fold, the models were trained with various hyper-parameters on the train set then
40 evaluated on the validation set to select the best ones. We then retrain the model on the train data and assess performance
41 on the test set. We do this 5 times and report the mean and stdev of our test set performances.” Test samples were **never**
42 used in training nor in tuning hyper-parameters, but for performance reporting **only**. The same train, val and test indices
43 were used to compare all models and we used the **same number of hyper-parameters combinations** for all methods.

44 Minibatch is performed by jointly integrating all time series in the batch between unique observation times over the
45 whole batch. At each unique time, we only update (with GRU-Bayes) the time series actually having an observation at
46 that time and leave the other untouched. We then proceed with joint time integration until the next unique time.

47 **Reviewer 5.** A basic MLP does not enjoy the properties stated in sections 2.2 and in first paragraph of section 2.3. Still
48 we investigated an other ODE parametrization : *GRU-minimal*, as described in Appendix G.

49 A important strength of our approach is to be able to predict future observations at **any** point in time. Yet it can also
50 be easily used for time series classification or regression (e.g. using the last latent $h(T)$ to feed a classifier). If those
51 labels are widely used in the literature, we are convinced vitals forecasting is just as important in medical practice as the
52 health status of a patient is complex and cannot usually be fully captured by discrete labels.