

# SELF-SUPERVISED INFERENCE IN STATE-SPACE MODELS

**David Ruhe**

AI4Science, AMLab, Anton Pannekoek Institute  
University of Amsterdam, The Netherlands  
d.ruhe@uva.nl

**Patrick Forré**

AI4Science, AMLab  
University of Amsterdam, The Netherlands  
p.d.forre@uva.nl

## ABSTRACT

We perform approximate inference in state-space models with nonlinear state transitions. Without parameterizing a generative model, we apply Bayesian update formulas using a local linearity approximation parameterized by neural networks. This comes accompanied by a maximum likelihood objective that requires no supervision via uncorrupt observations or ground truth latent states. The optimization backpropagates through a recursion similar to the classical Kalman filter and smoother. Additionally, using an approximate conditional independence, we can perform smoothing without having to parameterize a separate model. In scientific applications, domain knowledge can give a linear approximation of the latent transition maps, which we can easily incorporate into our model. Usage of such domain knowledge is reflected in excellent results (despite our model’s simplicity) on the chaotic Lorenz system compared to fully supervised and variational inference methods. Finally, we show competitive results on an audio denoising experiment.

## 1 INTRODUCTION

Many sequential processes in industry and research involve noisy measurements that describe latent dynamics. A state-space model is a type of graphical model that effectively represents such noise-afflicted data (Bishop, 2006). The joint distribution is assumed to factorize according to a directed graph that encodes the dependency between variables using conditional probabilities. One is usually interested in performing inference, meaning to obtain reasonable estimates of the posterior distribution of the latent states or uncorrupt measurements. Approaches involving sampling (Neal et al., 2011), variational inference (Kingma & Welling, 2013), or belief propagation (Koller & Friedman, 2009) have been proposed before. Assuming a hidden Markov process (Koller & Friedman, 2009), the celebrated Kalman filter and smoother (Kalman, 1960; Rauch et al., 1965) are classical approaches to solving the posterior inference problem. However, the Markov assumption, together with linear Gaussian transition and emission probabilities, limit their flexibility. We present filtering and smoothing methods that are related to the classical Kalman filter updates but are augmented with flexible function estimators without using a constrained graphical model. By noting that the filtering and smoothing recursions can be back-propagated through, these estimators can be trained with a principled maximum-likelihood objective reminiscent of the `noise2noise` objective (Lehtinen et al., 2018; Laine et al., 2019). By using a locally linear transition distribution, the posterior distribution remains tractable despite the use of non-linear function estimators. Further, we show how a *linearized smoothing* procedure can be applied directly to the filtering distributions, discarding the need to train a separate model for smoothing.

To verify what is claimed, we perform three experiments. (1) A linear dynamics filtering experiment, where we show how our models approximate the optimal solution with sufficient data. We also report that including expert knowledge can yield better estimates of latent states. (2) A more challenging chaotic Lorenz smoothing experiment that shows how our models perform on par with recently proposed supervised models. (3) An audio denoising experiment that uses real-world noise showing practical applicability of the methods.

Our contributions can be summarized as follows.

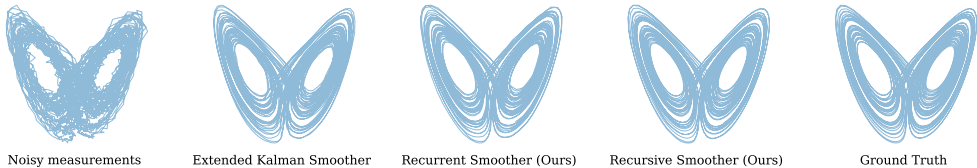


Figure 1: Best viewed on screen. Qualitative results of our work. To the noisy measurements (1st from left), we apply an extended Kalman smoother (2nd). From the noisy measurements, we learn a recurrent model that does slightly better (3rd). Our recursive model combines expert knowledge with inference (4th), yielding the best result. Ground truth provided for comparison (5th).

1. We show that the posterior inference distribution of a state-space model is tractable while parameter estimation is performed by neural networks. This means that we can apply the classical recursive Bayesian updates, akin to the Kalman filter and smoother, with mild assumptions on the generative process.
2. Our proposed method is optimized using maximum likelihood in a self-supervised manner. That is, ground truth values of states and measurements are not assumed to be available for training. Still, despite our model’s simplicity, our experiments show that it performs better or on par with several baselines.
3. We show that the model can be combined with prior knowledge about the transition and emission probabilities, allowing for better applicability in low data regimes and incentivizing the model to provide more interpretable estimates of the latent states.
4. A linearized smoothing approach is presented that does not require explicit additional parameterization and learning of the smoothing distribution.

## 2 RELATED WORK

Becker et al. (2019) provide a detailed discussion of recent related work, which we build on here and in table 1. An early method that extends the earlier introduced Kalman filter by allowing nonlinear transitions and emissions is the Extended Kalman filter (Ljung, 1979). It is limited due to the naive approach to locally linearize the transition and emission distributions. Furthermore, the transition and emission mechanisms are usually assumed to be known, or estimated with Expectation Maximization (Moon, 1996). More flexible methods that combine deep learning with variational inference include Black Box Variational Inference (Archer et al., 2015), Structured Inference Networks (Krishnan et al., 2017), Kalman Variational Autoencoder (Fraccaro et al., 2017), Deep Variational Bayes Filters (Karl et al., 2017), Variational Sequential Monte Carlo (Naesseth et al., 2018) and Disentangled Sequential Autoencoder (Yingzhen & Mandt, 2018). However, the lower-bound objective makes the approach less scalable and accurate (see also Becker et al. (2019)). Furthermore, all of the above methods explicitly assume a graphical model, imposing a strong but potentially harmful inductive bias. The BackpropKF (Haarnoja et al., 2016) and Recurrent Kalman Network (Becker et al., 2019) move away from variational inference and borrow Bayesian filtering techniques from the Kalman filter. We follow this direction but do not require supervision through ground truth latent states or uncorrupt emissions. Satorras et al. (2019) combine Kalman filters through message passing with graph neural networks to perform hybrid inference. We perform some of their experiments by also incorporating expert knowledge. However, contrary to their approach, we do not need supervision. Finally, concurrently to this work, Revach et al. (2021) develop KalmanNet. It proposes similar techniques but evaluates them in a supervised manner. The authors, however, do suggest that an unsupervised approach can also be feasible. Additionally, we more explicitly state what generative assumptions are required, then target the posterior distribution of interest, and develop the model and objective function from there. Moreover, the current paper includes linearized smoothing (section 6), parameterized smoothing (appendix A), and the recurrent model (appendix C). We also denote theoretical guarantees under the `noise2noise` objective.

	scalable	state est.	uncertainty	noise	dir. opt.	self-sup.
Ljung (1979)	✓/×	✓	✓	✓	×	×
Hochreiter et al. (1997)	✓	✓	✓/×	✓	✓	×
Cho et al. (2014)	✓	✓	✓/×	✓	✓	×
Wahlström et al. (2015)	✓	✓	✓/×	×	✓	×
Watter et al. (2015)	✓	×	✓	✓	×	✓
Archer et al. (2015)	✓/×	×	✓	✓	×	✓
Krishnan et al. (2017)	✓/×	×	✓	✓	×	✓
Fracaro et al. (2017)	✓/×	×	✓	✓	×	✓
Karl et al. (2017)	✓	×	✓	✓	×	✓
Naeseth et al. (2018)	✓	×	✓	✓	×	✓
Yingzhen et al. (2018)	✓	×	✓	×	×	✓
Rangapuram et al. (2018)	✓/×	✓ (1D)	✓	×	✓	✓
Doerr et al. (2018)	×	✓	✓	✓	✓	✓
Satorras et al. (2019)	✓	✓	×	✓	✓	×
Haarnoja et al. (2016)	✓	✓	✓	✓	✓	×
Becker et al. (2019)	✓	✓	✓	✓	✓	×
<b>Ours</b>	✓/×	✓	✓	✓	✓	✓

Table 1: We compare whether algorithms are scalable, state estimation can be performed, models provide uncertainty estimates, noisy or missing data can be handled, optimization is performed directly and if supervision is required. “✓/×” means that it depends on the parameterization.

### 3 GENERATIVE MODEL ASSUMPTIONS

In this section, we explicitly state the model’s generative process assumptions. First, we assume that we can measure (at least) one run of (noise-afflicted) sequential data  $\mathbf{y}_{0:K} := (\mathbf{y}_0, \dots, \mathbf{y}_K)$ , where each  $\mathbf{y}_k \in \mathbb{R}^M$ ,  $k = 0, \dots, K$ . We abbreviate:  $\mathbf{y}_{l:k} := (\mathbf{y}_l, \dots, \mathbf{y}_k)$  and  $\mathbf{y}_{<k} := \mathbf{y}_{0:k-1}$  and  $\mathbf{y}_{\leq k} := \mathbf{y}_{0:k}$  and  $\mathbf{y}_{-k} := (\mathbf{y}_{0:k-1}, \mathbf{y}_{k+1:K})$ . We then assume that  $\mathbf{y}_{0:K}$  is the result of some possibly non-linear probabilistic latent dynamics, i.e., of a distribution  $p(\mathbf{x}_{0:K})$ , whose variables are given by  $\mathbf{x}_{0:K} := (\mathbf{x}_0, \dots, \mathbf{x}_K)$  with  $\mathbf{x}_k \in \mathbb{R}^N$ . Each  $\mathbf{y}_k$  is assumed to be drawn from some shared noisy emission probability  $p(\mathbf{y}_k | \mathbf{x}_k)$ . The joint probability is then assumed to factorize as:

$$p(\mathbf{y}_{0:K}, \mathbf{x}_{0:K}) = p(\mathbf{x}_{0:K}) \prod_{k=0}^K p(\mathbf{y}_k | \mathbf{x}_k). \quad (1)$$

Further implicit assumptions about the generative model are imposed via inference model choices (see section 7). Note that this factorization encodes several conditional independences like

$$\mathbf{y}_k \perp\!\!\!\perp (\mathbf{y}_{-k}, \mathbf{x}_{-k}) | \mathbf{x}_k. \quad (2)$$

Typical models that follow these assumptions are linear dynamical systems, hidden Markov models, but also nonlinear state-space models with higher-order Markov chains in latent space like presented in fig. 2.

In contrast to other approaches (e.g., Krishnan et al. (2015); Johnson et al. (2016); Krishnan et al. (2017)) where one tries to model the latent dynamics with transition probabilities  $p(\mathbf{x}_k | \mathbf{x}_{k-1})$  and possibly non-linear emission probabilities  $p(\mathbf{y}_k | \mathbf{x}_k)$ , we go the other way around. We assume that all the non-linear dynamics are captured inside the latent distribution  $p(\mathbf{x}_{0:K})$ , where at this point we make no further assumption about its factorization, and the emission probabilities are (well-approximated with) a linear Gaussian noise model:

$$p(\mathbf{y}_k | \mathbf{x}_k) = \mathcal{N}(\mathbf{y}_k | \mathbf{H}\mathbf{x}_k, \mathbf{R}), \quad (3)$$

where the matrix  $\mathbf{H}$  represents the measurement device and  $\mathbf{R}$  is the covariance matrix of the independent additive noise. We make a brief argument why this assumption is not too restrictive. First, if one is interested in denoising corrupted measurements, any nonlinear emission can be captured directly inside the latent states  $\mathbf{x}_k$ . To see this, let  $\mathbf{z}_k \in \mathbb{R}^{N-M}$  denote non-emitted state variables. We then put  $\mathbf{x}_k := [\mathbf{y}_k \quad \mathbf{z}_k]^\top$ , where  $\mathbf{y}_k$  is computed by applying the nonlinear emission to  $\mathbf{z}_k$ . We thus include the measurements in the modeled latent state  $\mathbf{x}_k$ . Then we can put  $\mathbf{H} := [\mathbf{I}_M \quad \mathbf{0}_{M \times (N-M)}]$ . Second, techniques proposed by Laine et al. (2019) allow for non-Gaussian noise models, relaxing the need for assumption eq. (3). Third, we can locally linearize the

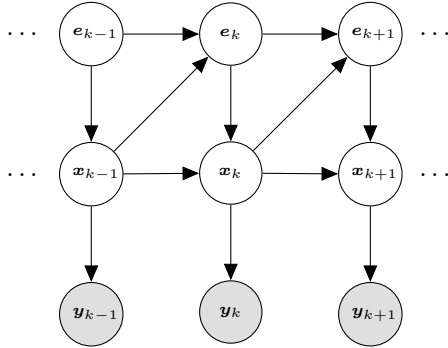


Figure 2: State-space model with deeper latent structure.

emission (Ljung, 1979). Finally, industrial or academic applications include cases where emissions are (sparse) Gaussian measurements and the challenging nonlinear dynamics occur in latent space. Examples can be found in MRI imaging (Lustig et al., 2007) and radio astronomy (Thompson et al., 2017).

#### 4 PARAMETERIZATION

In this section, we show how we parameterize the inference model. A lot of the paper’s work relies on established Bayesian filtering machinery. However, for completeness, we like to prove how all the update steps remain valid while using neural networks for function estimation.

Given our noisy measurements  $\mathbf{y}_{0:K} = (\mathbf{y}_0, \dots, \mathbf{y}_K)$  we want to find good estimates for the latent states  $\mathbf{x}_{0:K} = (\mathbf{x}_0, \dots, \mathbf{x}_K)$ , which generated  $\mathbf{y}_{0:K}$ . For this, we want to infer the marginal conditional distributions  $p(\mathbf{x}_k | \mathbf{y}_{\leq k})$  or  $p(\mathbf{x}_k | \mathbf{y}_{< k})$  (for forecasting), for an online inference approach (*filtering*); and  $p(\mathbf{x}_k | \mathbf{y}_{0:K})$  or  $p(\mathbf{x}_k | \mathbf{y}_{-k})$ , for a full inference approach (*smoothing*). In the main body of the paper, we only consider filtering. Smoothing can be performed similarly, which is detailed in the supplementary material (appendix A).

We start with the following advantageous parameterization:

$$p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{y}_{< k}) = \mathcal{N}(\mathbf{x}_k | \hat{\mathbf{F}}_{k|<k} \mathbf{x}_{k-1} + \hat{\mathbf{e}}_{k|<k}, \hat{\mathbf{Q}}_{k|<k}), \quad (4)$$

where  $\hat{\mathbf{F}}_{k|<k} := \hat{\mathbf{F}}_{k|<k}(\mathbf{y}_{<k})$ ,  $\hat{\mathbf{e}}_{k|<k} := \hat{\mathbf{e}}_{k|<k}(\mathbf{y}_{<k})$  and  $\hat{\mathbf{Q}}_{k|<k} := \hat{\mathbf{Q}}_{k|<k}(\mathbf{y}_{<k})$  are parameterized with neural networks. Next, we have available

$$p(\mathbf{x}_{k-1} | \mathbf{y}_{\leq(k-1)}) = \mathcal{N}(\mathbf{x}_{k-1} | \hat{\mathbf{x}}_{(k-1)|\leq(k-1)}, \hat{\mathbf{P}}_{(k-1)|\leq(k-1)}), \quad (5)$$

i.e., the previous time-step’s conditional marginal distribution of interest. For  $k = 1$ , this is some initialization. Otherwise, it is the result of the procedure we are currently describing. We use this distribution to evaluate the marginalization

$$\begin{aligned} p(\mathbf{x}_k | \mathbf{y}_{<k}) &= \int p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{y}_{<k}) p(\mathbf{x}_{k-1} | \mathbf{y}_{<k}) d\mathbf{x}_{k-1} \\ &= \mathcal{N}(\mathbf{x}_k | \hat{\mathbf{x}}_{k|<k}, \hat{\mathbf{P}}_{k|<k}), \end{aligned} \quad (6)$$

with

$$\hat{\mathbf{x}}_{k|<k}(\mathbf{y}_{<k}) = \hat{\mathbf{F}}_{k|<k} \hat{\mathbf{x}}_{k-1|\leq k-1} + \hat{\mathbf{e}}_{k|<k}, \quad \hat{\mathbf{P}}_{k|<k}(\mathbf{y}_{<k}) = \hat{\mathbf{F}}_{k|<k} \hat{\mathbf{P}}_{k-1|\leq k-1} \hat{\mathbf{F}}_{k|<k}^\top + \hat{\mathbf{Q}}_{k|<k}. \quad (7)$$

Note that the distributions under the integral eq. (6) are jointly Gaussian only because of the parameterization eq. (4). Hence, we can evaluate the integral analytically.

Finally, to obtain the conditional  $p(\mathbf{x}_k | \mathbf{y}_{\leq k}) = p(\mathbf{x}_k | \mathbf{y}_k, \mathbf{y}_{<k})$  we use Bayes’ rule:

$$p(\mathbf{x}_k | \mathbf{y}_k, \mathbf{y}_{<k}) = \frac{p(\mathbf{y}_k | \mathbf{x}_k, \mathbf{y}_{<k}) \cdot p(\mathbf{x}_k | \mathbf{y}_{<k})}{p(\mathbf{y}_k | \mathbf{y}_{<k})} \stackrel{\text{eq. (2)}}{=} \frac{p(\mathbf{y}_k | \mathbf{x}_k)}{p(\mathbf{y}_k | \mathbf{y}_{<k})} \cdot p(\mathbf{x}_k | \mathbf{y}_{<k}). \quad (8)$$

Equation (3) and the result eq. (6) allow us to also get an analytic expression for

$$p(\mathbf{x}_k | \mathbf{y}_{\leq k}) = \mathcal{N}(\mathbf{x}_k | \hat{\mathbf{x}}_{k|\leq k}, \hat{\mathbf{P}}_{k|\leq k}) \quad (9)$$

with the following abbreviations:

$$\hat{\mathbf{x}}_{k|\leq k}(\mathbf{y}_{\leq k}) := \hat{\mathbf{x}}_{k|<k} + \hat{\mathbf{K}}_k (\mathbf{y}_k - \mathbf{H} \hat{\mathbf{x}}_{k|<k}), \quad \hat{\mathbf{P}}_{k|\leq k}(\mathbf{y}_{\leq k}) := \hat{\mathbf{P}}_{k|<k} - \hat{\mathbf{K}}_k \mathbf{H} \hat{\mathbf{P}}_{k|<k}. \quad (10)$$

We introduce the *Kalman gain matrix* similar to the classical formulas:

$$\hat{\mathbf{K}}_k := \hat{\mathbf{P}}_{k|<k} \mathbf{H}^\top (\mathbf{H} \hat{\mathbf{P}}_{k|<k} \mathbf{H}^\top + \mathbf{R})^{-1}. \quad (11)$$

Note that taking the matrix inverse at this place in eq. (11) is more efficient than in the standard Gaussian formulas (for reference presented in appendix B) if  $M \leq N$ , which holds for our experiments.

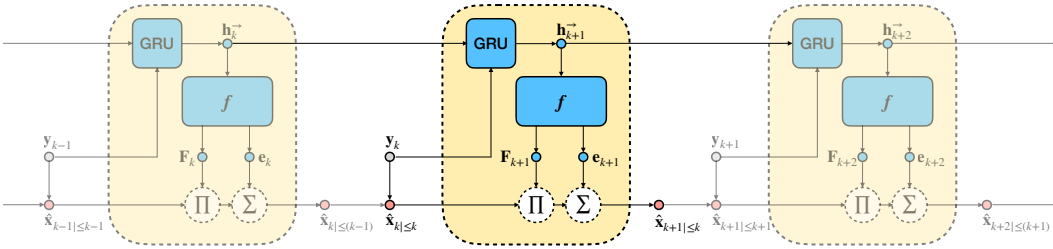


Figure 3: Our recursive model visualized. Data-point  $\mathbf{y}_k$  is fed, together with hidden state  $\mathbf{h}_k^{\rightarrow}$ , into a GRU unit. The new hidden state  $\mathbf{h}_{k+1}^{\rightarrow}$  is decoded into multiplicative component  $\mathbf{F}_{k+1}$  and additive component  $\mathbf{e}_{k+1}$ . Using these, the previous posterior mean  $\hat{\mathbf{x}}_{k|\leq k}$  is transformed into the prior estimate for time-step  $k+1$   $\hat{\mathbf{x}}_{k+1|\leq k}$ .  $\mathbf{y}_{k+1}$  is used to obtain posterior mean  $\hat{\mathbf{x}}_{k+1|\leq(k+1)}$ .

This completes the recursion: we can use eq. (9) for a new time-step  $k+1$  by plugging it back into eq. (6). We have shown how estimating a local linear transition using neural networks in eq. (4) ensures that all the recursive update steps from the Kalman filter analytically hold without specifying and estimating a generative model.

We note that we could also have parameterized

$$p(\mathbf{x}_k | \mathbf{y}_{<k}) = \mathcal{N}(\mathbf{x}_k | \hat{\mathbf{x}}_{k|<k}(\mathbf{y}_{<k}), \hat{\mathbf{P}}_{k|<k}(\mathbf{y}_{<k})) \quad (12)$$

with  $\hat{\mathbf{x}}_{k|<k}(\mathbf{y}_{<k})$  and  $\hat{\mathbf{P}}_{k|<k}(\mathbf{y}_{<k})$  directly estimated by a neural network. This has the advantage that we do not rely on a local linear transition model. However, it also means that we are estimating  $\mathbf{x}_k$  without any form of temporal regularization. Additionally, it is harder to incorporate prior knowledge about the transitions maps into such a model. Nonetheless, we detail this parameterization further in the supplementary material (appendix C) and include its performance in our experiments in section 8.

To conclude the section, we like to discuss some of the limitations of the approach. 1. The Gaussianity assumption of eq. (4) ensures but also restricts eq. (8) and eq. (9) to these forms. That is, we make a direct assumption about the form of the posterior  $p(\mathbf{x}_k | \mathbf{y}_{<k})$ . Defending our case, we like to point out that methods such as variational inference (Krishnan et al., 2017) or posterior regularization (Ganchev et al., 2010) also make assumptions (e.g., mean-field Gaussian) about the posterior. 2. Since we did not explicitly specify a factorization of  $p(\mathbf{x}_{0:K})$ , we cannot ensure that the distributions we obtain from the above procedure form a posterior to the ground truth generative process. This does not mean, however, that we cannot perform accurate inference. Arguably, not making explicit assumptions about the generative process is preferred to making wrong assumptions and using those for modeling, which can be the case for variational auto-encoders. 3. The local linearity assumption is justifiable if the length between time-steps is sufficiently small. However, note that the model is more flexible than directly parameterizing eq. (12) (see appendix C) since it *reduces* to that case by putting  $\mathbf{F}_k := \mathbf{0}$  for all  $k$ .

## 5 FITTING

We have shown in the previous section how parameterization of a local linear transition model leads to recursive estimation of  $p(\mathbf{x}_k | \mathbf{y}_{<k})$  and  $p(\mathbf{x}_k | \mathbf{y}_{\leq k})$  for all  $k$  using classical Bayesian filtering formulas. The inference is only effective if the estimates  $\hat{\mathbf{F}}_{k|<k}$ ,  $\hat{\mathbf{e}}_{k|<k}$  and  $\hat{\mathbf{Q}}_{k|<k}$  from eq. (4) are accurate. We can use the parameterization  $p(\mathbf{x}_k | \mathbf{y}_{<k})$  of eq. (6), the emission model  $p(\mathbf{y}_k | \mathbf{x}_k) = \mathcal{N}(\mathbf{y}_k | \mathbf{H} \mathbf{x}_k, \mathbf{R})$  from eq. (3), and the factorization from eq. (1) to see that an analytic form of the log-likelihood of the data emerges:

$$\begin{aligned} p(\mathbf{y}_k | \mathbf{y}_{<k}) &= \int p(\mathbf{y}_k | \mathbf{x}_k) p(\mathbf{x}_k | \mathbf{y}_{<k}) d\mathbf{x}_k \\ &= \mathcal{N}(\mathbf{y}_k | \mathbf{H} \hat{\mathbf{x}}_{k|<k}(\mathbf{y}_{<k}), \mathbf{H} \hat{\mathbf{P}}_{k|<k}(\mathbf{y}_{<k}) \mathbf{H}^\top + \mathbf{R}), \end{aligned} \quad (13)$$

$$\log p(\mathbf{y}_{0:K}) = \sum_{k=0}^K \log \mathcal{N} \left( \mathbf{y}_k \mid \mathbf{H} \hat{\mathbf{x}}_{k|<k}(\mathbf{y}_{<k}), \mathbf{H} \hat{\mathbf{P}}_{k|<k}(\mathbf{y}_{<k}) \mathbf{H}^\top + \mathbf{R} \right). \quad (14)$$

If we put  $\hat{\mathbf{y}}_{k|<k} := \mathbf{H} \hat{\mathbf{x}}_{k|<k}(\mathbf{y}_{<k})$  and  $\hat{\mathbf{M}}_{k|<k} := \mathbf{H} \hat{\mathbf{P}}_{k|<k}(\mathbf{y}_{<k}) \mathbf{H}^\top + \mathbf{R}$ , then the maximum-likelihood objective leads to the following loss function, which we can minimize using gradient descent methods w.r.t. all model parameters:

$$\mathcal{L} := \sum_{k=0}^K \left[ (\hat{\mathbf{y}}_{k|<k} - \mathbf{y}_k)^\top \hat{\mathbf{M}}_{k|<k}^{-1} (\hat{\mathbf{y}}_{k|<k} - \mathbf{y}_k) + \log \det \hat{\mathbf{M}}_{k|<k} \right]. \quad (15)$$

Note that each term in the sum above represents a one-step-ahead *self-supervised* error term. We thus minimize the prediction residuals  $\hat{\mathbf{y}}_{k|<k}(\mathbf{y}_{<k}) - \mathbf{y}_k$  in a norm that is inversely scaled with the above covariance matrix, plus a regularizing determinant term, which prevents the covariance matrix from diverging. The arisen loss function is similar to the `noise2noise` (Lehtinen et al., 2018; Krull et al., 2019; Batson & Royer, 2019; Laine et al., 2019) objective from computer vision literature, combined with a locally linear transition model like Becker et al. (2019). We show in appendix D that this objective will yield correct results (meaning estimating the ground-truth  $\mathbf{x}_k$ ) if the noise is independent with  $\mathbb{E}[\mathbf{y}_k \mid \mathbf{H} \mathbf{x}_k] = \mathbf{H} \mathbf{x}_k$ . A similar procedure in the causality literature is given by Schölkopf et al. (2016). An algorithmic presentation of performing inference and fitting is presented in appendix E.

Note that after fitting the parameters to the data, eq. (6) can directly be used to do one-step ahead forecasting. Forecasting an arbitrary number of time steps is possible by plugging the new value  $\hat{\mathbf{x}}_{K+1|K}$  via  $\mathbf{y}_{K+1} := \mathbf{H} \hat{\mathbf{x}}_{K+1|K}$  back into the recurrent model, and so on. This is not a generative model but merely a convenience that we deemed worth mentioning.

## 6 LINEARIZED SMOOTHING

So far, we have only discussed how to perform filtering. Recall that for smoothing, we are instead interested in the quantity  $p(\mathbf{x}_k \mid \mathbf{y}_{0:K})$ . A smoothing strategy highly similar to the methods described earlier can be obtained by explicitly parameterizing such a model, which we detail in the supplementary material. Here, we introduce a *linearized smoothing* procedure. The essential advantage is that no additional model has to be trained, which can be costly. Several algorithms stemming from the Kalman filter literature can be applied, such as the RTS smoother (Rauch et al., 1965) and the two-filter smoother (Kitagawa, 1994). To enable this, we need to assume that the conditional mutual information  $I(\mathbf{x}_{k-1}; \mathbf{y}_{k:K} \mid \mathbf{x}_k, \mathbf{y}_{0:k-1})$  is small for all  $k = 1, \dots, K$ . In other words, we assume that we approximately have the following conditional independences:

$$\mathbf{x}_{k-1} \perp\!\!\!\perp \mathbf{y}_{k:K} \mid (\mathbf{x}_k, \mathbf{y}_{0:k-1}). \quad (16)$$

To explain the motivation for this requirement, consider the model in fig. 2. If the states of  $\mathbf{y}_{0:k-1}$  and  $\mathbf{x}_k$  are known, then the additional information that  $\mathbf{y}_{k:K}$  has about the latent variable  $\mathbf{x}_{k-1}$  would need to be passed along the unblocked deeper paths like  $\mathbf{y}_{k+1} \leftarrow \mathbf{x}_{k+1} \leftarrow \mathbf{e}_{k+1} \leftarrow \mathbf{e}_k \leftarrow \mathbf{x}_{k-1}$ . Then the assumption of small  $I(\mathbf{x}_{k-1}; \mathbf{y}_{k:K} \mid \mathbf{x}_k, \mathbf{y}_{0:k-1})$  can be interpreted as that the deeper paths transport less information than the lower direct paths. If we consider all edges to the  $\mathbf{x}_k$ 's as linear and the edges to the  $\mathbf{e}_k$ 's as non-linear maps, the above could be interpreted as an information-theoretic version of expressing that the non-linear correction terms are small compared to the linear parts in the functional relations between the variables.

We will now show that under the earlier assumptions and eq. (16) we get a Gaussian approximation:  $p(\mathbf{x}_k \mid \mathbf{y}_{0:K}) \approx \mathcal{N}(\mathbf{x}_k \mid \hat{\mathbf{z}}_k, \hat{\mathbf{G}}_k)$ . We will do backward induction with  $\hat{\mathbf{z}}_K := \hat{\mathbf{x}}_{K|K}$  and  $\hat{\mathbf{G}}_K := \hat{\mathbf{P}}_{K|K}$ . To propagate this to previous time steps  $k - 1$  we use the chain rule:

$$p(\mathbf{x}_{k-1} \mid \mathbf{y}_{0:K}) = \int p(\mathbf{x}_{k-1} \mid \mathbf{x}_k, \mathbf{y}_{0:K}) p(\mathbf{x}_k \mid \mathbf{y}_{0:K}) d\mathbf{x}_k, \quad (17)$$

where the second term is known by backward induction and for the first term we make use of the approximate conditional independence from eq. (16) to get

$$p(\mathbf{x}_{k-1} \mid \mathbf{x}_k, \mathbf{y}_{0:K}) = p(\mathbf{x}_{k-1} \mid \mathbf{y}_{k:K}, \mathbf{x}_k, \mathbf{y}_{0:k-1}) \stackrel{\text{eq. (16)}}{\approx} p(\mathbf{x}_{k-1} \mid \mathbf{x}_k, \mathbf{y}_{0:k-1}). \quad (18)$$

The latter was shown to be Gaussian in section 4:

$$p(\mathbf{x}_{k-1}, \mathbf{x}_k \mid \mathbf{y}_{0:k-1}) = \mathcal{N} \left( \begin{bmatrix} \mathbf{x}_{k-1} \\ \mathbf{x}_k \end{bmatrix} \mid \begin{bmatrix} \hat{\mathbf{x}}_{k-1|\leq k-1} \\ \hat{\mathbf{x}}_{k|\leq k} \end{bmatrix}, \begin{bmatrix} \hat{\mathbf{P}}_{k-1|\leq k-1} & \hat{\mathbf{P}}_{k-1|\leq k-1} \hat{\mathbf{F}}_{k|\leq k}^\top \\ \hat{\mathbf{F}}_{k|\leq k} \hat{\mathbf{P}}_{k-1|\leq k-1} & \hat{\mathbf{P}}_{k|\leq k} \end{bmatrix} \right). \quad (19)$$

By use of the usual formulas for Gaussians and the *reverse Kalman gain matrix*  $\hat{\mathbf{J}}_{k-1|k}$  we arrive at the following update formulas,  $k = K, \dots, 1$ , with  $\hat{\mathbf{z}}_K := \hat{\mathbf{x}}_{K|\leq K}$  and  $\hat{\mathbf{G}}_K := \hat{\mathbf{P}}_{K|\leq K}$ :

$$\hat{\mathbf{J}}_{k-1|k} := \hat{\mathbf{P}}_{k-1|\leq k-1} \hat{\mathbf{F}}_{k|\leq k}^\top \hat{\mathbf{P}}_{k|\leq k}^{-1}, \quad (20)$$

$$\hat{\mathbf{G}}_{k-1} := \hat{\mathbf{P}}_{k-1|\leq k-1} + \hat{\mathbf{J}}_{k-1|k} \left( \hat{\mathbf{P}}_{k|\leq k} - \hat{\mathbf{P}}_{k|\leq k} \right) \hat{\mathbf{J}}_{k-1|k}^\top, \quad (21)$$

$$\hat{\mathbf{z}}_{k-1} := \hat{\mathbf{x}}_{k-1|\leq k-1} + \hat{\mathbf{J}}_{k-1|k} (\hat{\mathbf{z}}_k - \hat{\mathbf{x}}_{k|\leq k}). \quad (22)$$

As such, we can perform inference for all  $k = 0, \dots, K$  with  $p(\mathbf{x}_k \mid \mathbf{y}_{0:K}) \approx \mathcal{N}(\mathbf{x}_k \mid \hat{\mathbf{z}}_k, \hat{\mathbf{G}}_k)$ . Algorithmically, the above is presented in appendix E.

## 7 RECURRENT NEURAL NETWORK

Before going into the experiments section, we briefly explain how we specifically estimate the functions  $\hat{\mathbf{F}}_{k|\leq k}(\mathbf{y}_{<k})$ ,  $\hat{\mathbf{e}}_{k|\leq k}(\mathbf{y}_{<k})$  and  $\hat{\mathbf{L}}_{k|\leq k}(\mathbf{y}_{<k})$  that parameterize the transition probability  $p(\mathbf{x}_k \mid \mathbf{y}_{<k}, \mathbf{x}_{k-1})$  (eq. (4)). The choice of the model here implicitly makes further assumptions about the generative model. If we consider neural networks, the temporal nature of the data suggests recurrent neural networks (Graves et al., 2013), convolutional neural networks (Kalchbrenner et al., 2014), or transformer architectures (Vaswani et al., 2017). Additionally, if the data is image-based, one might further make use of convolutions. For our experiments, we use a Gated Recurrent Unit (GRU) network (Cho et al., 2014), that recursively encodes hidden representations. Therefore, we put

$$\hat{\mathbf{Q}}_{k|\leq k} := \mathbf{L}_k \mathbf{L}_k^\top, \quad \begin{bmatrix} \hat{\mathbf{F}}_k \\ \hat{\mathbf{e}}_k \\ \hat{\mathbf{L}}_k \end{bmatrix} := \mathbf{f}(\mathbf{h}_k^\rightarrow), \quad \mathbf{h}_k^\rightarrow := \text{GRU}(\mathbf{y}_{k-1}^\rightarrow, \mathbf{h}_{k-1}^\rightarrow), \quad (23)$$

where  $\hat{\mathbf{L}}_k$  is a Cholesky factor and  $\mathbf{f}$  is a multi-layer perceptron decoder.

## 8 EXPERIMENTS

We perform three experiments, as motivated in section 1. Technical details on the setup of the experiments can be found in the supplementary material (appendix F). We refer to the model detailed in section 4 as the *recursive filter*, as it uses the Bayesian update recursion. For smoothing experiments, we use *recursive smoother*. The model obtained by parameterizing  $p(\mathbf{x}_k \mid \mathbf{y}_{<k})$  directly (eq. (12)) is referred to as the *recurrent filter* or *recurrent smoother*, as it only employs recurrent neural networks (and no Bayesian recursion) to estimate said density directly.

### 8.1 LINEAR DYNAMICS

In the linear Gaussian case, it is known that the Kalman filter will give the optimal solution. Thus, we can get a lower bound on the test loss. In this toy experiment, we simulate particle tracking under linear dynamics and noisy measurements of the location. We use Newtonian physics equations as prior knowledge. We generate trajectories  $\mathcal{T} = \{\mathbf{x}_{0:K}, \mathbf{y}_{0:K}\}$  with  $\mathbf{x}_k \in \mathbb{R}^6$  and  $\mathbf{y}_k \in \mathbb{R}^2$  according to the differential equations:

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & -c & 1 \\ 0 & -\tau c & 0 \end{bmatrix} \begin{bmatrix} p \\ v \\ a \end{bmatrix}. \quad (24)$$

We obtain sparse, noisy measurements  $\mathbf{y}_k = \mathbf{H}\mathbf{x}_k + \mathbf{r}$  with  $\mathbf{r} \sim \mathcal{N}(\mathbf{0}, \mathbf{R})$ .  $\mathbf{H}$  is a selection matrix that returns a two-dimensional position vector. We run this experiment in a filtering setting, i.e.,

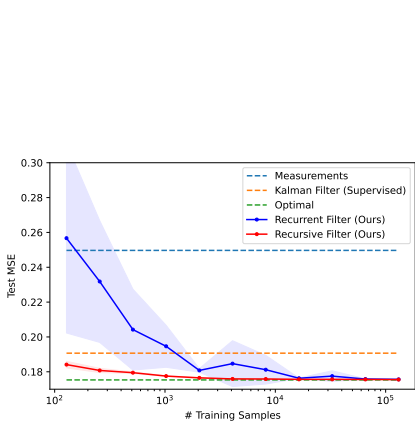


Figure 4: Considers the linear dynamics experiment (section 8.1). The mean squared error on the test set (lower is better) as a function of the number of examples.

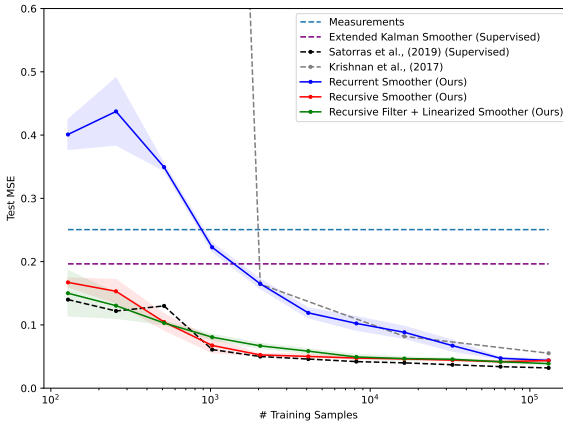


Figure 5: Considers the Lorenz experiment (section 8.2). The mean squared error on the test set (lower is better) as a function of the number of examples used for training.

we only use past observations. We compare against (1) the raw, noisy measurements which inherently deviate from the clean measurements, (2) the Kalman filter solution where we optimized the transition covariance matrix using clean data (hence supervised), (3) the optimal solution, which is a Kalman filter with ground truth parameters performing exact inference. To estimate the transition maps for the Kalman filter, we use the standard Taylor series of  $e^{A \cdot \Delta t}$  up to the first order. Additionally, we use this expert knowledge as an *inductive bias* for the recursive filter’s transition maps.

In fig. 4 we depict the test mean squared error (MSE, lower is better) as a function of the number of training samples. Given enough data, the self-supervised models approximate the optimal solution arbitrarily well. Our recursive model significantly outperforms both the Kalman filter and the inference model in the low-data regime by using incorporated expert knowledge. Additionally, we report that the recursive model’s distance to the ground truth latent states is much closer to the optimal solution than both the inference model and Kalman filter. Specifically, we report average mean squared errors of 0.685 for the inference model, 0.241 for the Kalman filter, **0.161** for the recursive model compared to 0.135 for the optimal Kalman filter. Finally, it is worth noting that the recursive model has much less variance as a function of its initialization.

## 8.2 LORENZ EQUATIONS

We simulate a Lorenz system according to

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} = \begin{bmatrix} -\sigma & \sigma & 0 \\ \rho - x_1 & -1 & 0 \\ x_2 & 0 & -\beta \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}. \tag{25}$$

We have  $\mathbf{H} = \mathbf{I}$ ,  $\mathbf{x} \in \mathbb{R}^3$  and  $\mathbf{y} \in \mathbb{R}^3$ . The Lorenz equations model atmospheric convection and form a classic example of chaos. Therefore, performing inference is much more complex than in the linear case. This time, we perform smoothing (see appendix A) and compare against (1) the raw measurements, (2) a supervised Extended Kalman smoother (Ljung, 1979), (3) the variational inference approach of Krishnan et al. (2017), (4) the supervised recursive model of Satorras et al. (2019). Our models include a recurrent smoother, a recursive smoother, and the recursive filter with linearized smoothing (section 6). Transition maps for the Extended Kalman smoother and the recursive models are obtained by taking a second-order Taylor series of  $e^{A \cdot \Delta t}$ . For the supervised extended Kalman filter, we again optimize its covariance estimate using ground truth data.

In fig. 5 we plot the test mean squared error (MSE, lower is better) as a function of the number of examples available for training. It is clear that our methods approach the ground truth states with



	Whitenoise	Doing the dishes	Dude miaowing	Exercise Bike	Pink noise	Running tap	Combined
Kalman Filter (Supervised)	0.225	0.230	0.232	0.237	0.235	0.230	0.227
Noise2Noise (Lehtinen et al., 2018)	0.327	0.399	0.448	0.430	0.440	0.383	0.526
SIN (Krishnan et al., 2017)	0.297	0.373	0.352	0.348	0.377	0.342	0.343
Recurrent Filter (Ours)	<b>0.102</b>	0.207	0.213	0.200	0.234	0.175	0.181
Recursive Filter (Ours)	0.107	0.206	<b>0.213</b>	0.198	0.232	0.166	<b>0.175</b>
Recursive Filter + RTS Smoother (Ours)	<b>0.100</b>	<b>0.204</b>	0.215	<b>0.197</b>	<b>0.231</b>	<b>0.166</b>	0.176
RKN (Becker et al., 2019, Supervised)	0.121	<b>0.127</b>	<b>0.109</b>	<b>0.105</b>	<b>0.085</b>	<b>0.121</b>	<b>0.125</b>

Table 2: Considers the audio denoising experiment (section 8.3). Test mean squared error (MSE, lower is better) per model and noise subset. Blue numbers indicate second-best performing models.

more data. This is in contrast to the Extended Kalman smoother, which barely outperforms the noisy measurements. We also see that the recursive models significantly outperform the recurrent model in the low-data regime. The recursive filter with linearized smoothing performs comparably to the other models and even better in low-data regimes. We hypothesize that this is because the required assumption for the linearized smoother holds (section 6) and regularizes the model. The variational method of Krishnan et al. (2017) performs poorly in low-data regimes. Most notably, the supervised method of Satorras et al. (2019) outperforms our models only slightly.

### 8.3 AUDIO DENOISING

Next, we test the model on non-fabricated data with less ideal noise characteristics. Specifically, we use the `SpeechCommands` spoken audio dataset (Warden, 2018). Performing inference on spoken audio is challenging, as it arguably requires understanding natural language. To this end, recent progress on synthesizing raw audio has been made (Lakhotia et al., 2021). However, this requires scaling to much larger and more sophisticated neural networks than presented here, which we deem out the current work’s scope. Therefore, we take a subset of the entire dataset, using audio from the classes “tree”, “six”, “eight”, “yes”, and “cat”. We overlay these clean audio sequences  $\{\mathbf{x}_{0:K}^{(1)}, \dots, \mathbf{x}_{0:K}^{(N)}\}$  with various noise classes that the dataset provides. That is, for every noise class  $C$  we obtain a set of noisy sequences  $\{\mathbf{y}_{0:K}^{(1)}, \dots, \mathbf{y}_{0:K}^{(N)}\}_C$ . We also consider a “combined” class in which we sample from the union of the noise sets. The task is to denoise the audio without having access to clean data. We evaluate the models on non-silent parts of the audio, as performance on those sections is the most interesting. Notably, none of these noise classes is Gaussian distributed.

We show the mean squared error on the test set of all models per noise class in table 2. Our models outperform the Kalman filter, Noise2Noise (Lehtinen et al., 2018), and SIN (Krishnan et al., 2017) unsupervised baselines. We suspect that the relatively poor performance of SIN is due to its generative Markov assumption, regularizing the model too strongly. The poor performance of Noise2Noise is due to the fact that it does not use the current measurement  $\mathbf{y}_k$  to infer  $\mathbf{x}_k$ . Like before, note that the Kalman filter is “supervised” as we optimized its covariance matrix using clean data. The supervised RKN (Becker et al., 2019) outperforms our models on most noise classes, but notably not on white noise. Most of these noise classes have temporal structure, making them predictable from past data. This is confirmed by observing these mean squared error values over the course of training. Initially, the values were better than reported in table 2, but the model increasingly fits the noise over time. Thus, although two of the main assumptions about the model (independent Gaussian noise) are violated, we still can denoise effectively. Since the RKN’s targets are denoised (hence “supervised”), it does not have this problem. In practice, obtaining clean data can be challenging.

## 9 CONCLUSION

We presented an advantageous parameterization of an inference procedure for nonlinear state-space models with potentially higher-order latent Markov chains. The inference distribution is split into linear and nonlinear parts, allowing for a recursion akin to the Kalman filter and smoother algorithms. Optimization is performed directly using a maximum-likelihood objective that backpropagates through these recursions. Smoothing can be performed similarly, but we additionally proposed linearized smoothing that can directly be applied to the filtering distributions. Our model is simple and builds on established methods from signal processing. Despite this, results showed that it can perform better or on par with fully supervised or variational inference methods.

## 10 ETHICS STATEMENT

The paper presents a simple method to perform inference using noisy sequential data. Applications can be found throughout society, e.g., tracking particles, denoising images or audio, or estimating system states. While many such examples are for good, there are applications with ethically debatable motivations. Among these could be tracking humans or denoising purposefully corrupted data (e.g., to hide one’s identity).

## 11 REPRODUCIBILITY STATEMENT

We are in the process of releasing code for the current work. For clarity and reproducibility, the presented methods are available as algorithms in the supplementary material. Furthermore, we made explicit wherever we needed to make an approximation or an assumption.

## REFERENCES

- Evan Archer, Il Memming Park, Lars Buesing, John Cunningham, and Liam Paninski. Black box variational inference for state space models. *arXiv preprint arXiv:1511.07367*, 2015.
- Philipp Bader, Sergio Blanes, and Fernando Casas. Computing the matrix exponential with an optimized taylor polynomial approximation. *Mathematics*, 7(12):1174, 2019.
- Joshua Batson and Loïc Royer. Noise2self: Blind denoising by self-supervision. In *ICML*, 2019.
- Philipp Becker, Harit Pandya, Gregor H. W. Gebhardt, Cheng Zhao, C. James Taylor, and Gerhard Neumann. Recurrent kalman networks: Factorized inference in high-dimensional deep feature spaces. In *ICML*, 2019.
- Christopher M Bishop. Pattern recognition. *Machine learning*, 128(9), 2006.
- Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- Andreas Doerr, Christian Daniel, Martin Schiegg, Nguyen-Tuong Duy, Stefan Schaal, Marc Toussaint, and Trimpe Sebastian. Probabilistic recurrent state-space models. In *International Conference on Machine Learning*, pp. 1280–1289. PMLR, 2018.
- Marco Fraccaro, Simon Kamronn, Ulrich Paquet, and Ole Winther. A disentangled recognition and nonlinear dynamics model for unsupervised learning. *arXiv preprint arXiv:1710.05741*, 2017.
- Kuzman Ganchev, Joao Graça, Jennifer Gillenwater, and Ben Taskar. Posterior regularization for structured latent variable models. *The Journal of Machine Learning Research*, 11:2001–2049, 2010.
- Christian Gourieroux, Alain Monfort, and Alain Trognon. Pseudo maximum likelihood methods: Theory. *Econometrica: journal of the Econometric Society*, pp. 681–700, 1984.
- Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. Speech recognition with deep recurrent neural networks. In *2013 IEEE international conference on acoustics, speech and signal processing*, pp. 6645–6649. Ieee, 2013.
- Tuomas Haarnoja, Anurag Ajay, Sergey Levine, and Pieter Abbeel. Backprop KF: learning discriminative deterministic state estimators. In *NeurIPS*, 2016.
- Matthew J. Johnson, David Duvenaud, Alexander B. Wiltschko, Ryan P. Adams, and Sandeep R. Datta. Composing graphical models with neural networks for structured representations and fast inference. In *NeurIPS*, 2016.
- Nal Kalchbrenner, Edward Grefenstette, and Phil Blunsom. A convolutional neural network for modelling sentences. *arXiv preprint arXiv:1404.2188*, 2014.

- Rudolph Emil Kalman. A new approach to linear filtering and prediction problems. 1960.
- Maximilian Karl, Maximilian Soelch, Justin Bayer, and Patrick van der Smagt. Deep variational bayes filters: Unsupervised learning of state space models from raw data. In *ICLR*, 2017.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Genshiro Kitagawa. The two-filter formula for smoothing and an implementation of the gaussian-sum smoother. *Annals of the Institute of Statistical Mathematics*, 46(4):605–623, 1994.
- Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- Rahul G Krishnan, Uri Shalit, and David Sontag. Deep kalman filters. *arXiv preprint arXiv:1511.05121*, 2015.
- Rahul G. Krishnan, Uri Shalit, and David A. Sontag. Structured inference networks for nonlinear state space models. In *AAAI*, pp. 2101–2109. AAAI Press, 2017.
- Alexander Krull, Tim-Oliver Buchholz, and Florian Jug. Noise2void - learning denoising from single noisy images. In *CVPR*, 2019.
- Samuli Laine, Tero Karras, Jaakko Lehtinen, and Timo Aila. High-quality self-supervised deep image denoising. In *NeurIPS*, 2019.
- Kushal Lakhotia, Evgeny Kharitonov, Wei-Ning Hsu, Yossi Adi, Adam Polyak, Benjamin Bolte, Tu-Anh Nguyen, Jade Copet, Alexei Baevski, Adelrahman Mohamed, et al. Generative spoken language modeling from raw audio. *arXiv preprint arXiv:2102.01192*, 2021.
- Jaakko Lehtinen, Jacob Munkberg, Jon Hasselgren, Samuli Laine, Tero Karras, Miika Aittala, and Timo Aila. Noise2noise: Learning image restoration without clean data. In *ICML*, 2018.
- Lennart Ljung. Asymptotic behavior of the extended kalman filter as a parameter estimator for linear systems. *IEEE Transactions on Automatic Control*, 24(1):36–50, 1979.
- Michael Lustig, David Donoho, and John M Pauly. Sparse mri: The application of compressed sensing for rapid mr imaging. *Magnetic Resonance in Medicine: An Official Journal of the International Society for Magnetic Resonance in Medicine*, 58(6):1182–1195, 2007.
- Todd K Moon. The expectation-maximization algorithm. *IEEE Signal processing magazine*, 13(6): 47–60, 1996.
- Christian Naesseth, Scott Linderman, Rajesh Ranganath, and David Blei. Variational sequential monte carlo. In *International conference on artificial intelligence and statistics*, pp. 968–977. PMLR, 2018.
- Radford M Neal et al. Mcmc using hamiltonian dynamics. *Handbook of markov chain monte carlo*, 2(11):2, 2011.
- Syama Sundar Rangapuram, Matthias W Seeger, Jan Gasthaus, Lorenzo Stella, Yuyang Wang, and Tim Januschowski. Deep state space models for time series forecasting. *Advances in neural information processing systems*, 31:7785–7794, 2018.
- Herbert E Rauch, F Tung, and Charlotte T Striebel. Maximum likelihood estimates of linear dynamic systems. *AIAA journal*, 3(8):1445–1450, 1965.
- Guy Revach, Nir Shlezinger, Xiaoyong Ni, Adria Lopez Escoriza, Ruud JG van Sloun, and Yonina C Eldar. Kalmannet: Neural network aided kalman filtering for partially known dynamics. *arXiv preprint arXiv:2107.10043*, 2021.
- Victor Garcia Satorras, Max Welling, and Zeynep Akata. Combining generative and discriminative models for hybrid inference. In *NeurIPS*, 2019.

Bernhard Schölkopf, David W Hogg, Dun Wang, Daniel Foreman-Mackey, Dominik Janzing, Carl-Johann Simon-Gabriel, and Jonas Peters. Modeling confounding by half-sibling regression. *Proceedings of the National Academy of Sciences*, 113(27):7391–7398, 2016.

Richard A Thompson, James M Moran, and George W Swenson Jr. *Interferometry and synthesis in radio astronomy*. Springer Nature, 2017.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pp. 5998–6008, 2017.

Niklas Wahlström, Thomas B Schön, and Marc Peter Deisenroth. From pixels to torques: Policy learning with deep dynamical models. *arXiv preprint arXiv:1502.02251*, 2015.

Pete Warden. Speech commands: A dataset for limited-vocabulary speech recognition. *arXiv preprint arXiv:1804.03209*, 2018.

Manuel Watter, Jost Tobias Springenberg, Joschka Boedecker, and Martin Riedmiller. Embed to control: A locally linear latent dynamics model for control from raw images. *arXiv preprint arXiv:1506.07365*, 2015.

Li Yingzhen and Stephan Mandt. Disentangled sequential autoencoder. In *International Conference on Machine Learning*, pp. 5670–5679. PMLR, 2018.

## A PARAMETERIZED SMOOTHING

In the main body of the paper, we showed how to parameterize the model for recursive estimation of  $p(\mathbf{x}_k | \mathbf{y}_{\leq k})$ ,  $k = 0, \dots, K$ . Additionally, we provided a *linearized smoothing* procedure that yields  $p(\mathbf{x}_k | \mathbf{y}_{1:K})$ . The disadvantage clearly is the linearization. Here, we show how we can recursively estimate  $p(\mathbf{x}_k | \mathbf{y}_{1:K})$  in a similar sense to the filtering setting.

First, we put

$$p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{y}_{-k}) = \mathcal{N}(\mathbf{x}_k | \hat{\mathbf{F}}_{k|-k} \mathbf{x}_{k-1} + \hat{\mathbf{e}}_{k|-k}, \hat{\mathbf{Q}}_{k|-k}), \quad (26)$$

where  $\hat{\mathbf{F}}_{k|-k}(\mathbf{y}_{-k})$  and  $\hat{\mathbf{e}}_{k|-k}(\mathbf{y}_{-k})$  and  $\hat{\mathbf{Q}}_{k|-k}(\mathbf{y}_{-k})$  are two-sided recurrent neural network outputs similar to section 7. We compute the distribution of interest as follows:

$$p(\mathbf{x}_k | \mathbf{y}_{0:K}) = \frac{p(\mathbf{y}_k | \mathbf{x}_k)p(\mathbf{x}_k | \mathbf{y}_{-k})}{p(\mathbf{y}_k | \mathbf{y}_{-k})} \quad (27)$$

$$\propto p(\mathbf{y}_k | \mathbf{x}_k) \int p(\mathbf{x}_k, \mathbf{x}_{k-1} | \mathbf{y}_{-k}) d\mathbf{x}_{k-1} \quad (28)$$

$$= p(\mathbf{y}_k | \mathbf{x}_k) \int p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{y}_{-k}) p(\mathbf{x}_{k-1} | \mathbf{y}_{-k}) d\mathbf{x}_{k-1} \quad (29)$$

$$\stackrel{\mathbf{x}_{k-1} \perp \mathbf{y}_k | \mathbf{y}_{-k}}{\approx} p(\mathbf{y}_k | \mathbf{x}_k) \int p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{y}_{-k}) p(\mathbf{x}_{k-1} | \mathbf{y}_{0:K}) d\mathbf{x}_{k-1} \quad (30)$$

where we make the approximation to compute eq. (27) efficiently and recursively. It is justified if  $I(\mathbf{x}_{k-1}; \mathbf{y}_k | \mathbf{y}_{-k}) < \epsilon$  for small  $\epsilon$ . That is, the additional information that  $\mathbf{y}_k$  conveys about  $\mathbf{x}_{k-1}$  is marginal if we have all other data. Let

$$p(\mathbf{x}_{k-1} | \mathbf{y}_{0:K}) := \mathcal{N}(\mathbf{x}_{k-1} | \hat{\mathbf{x}}_{k-1|0:K}, \hat{\mathbf{P}}_{k-1|0:K}) \quad (31)$$

be previous time-step’s posterior. Then put

$$p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{y}_{-k}) := \mathcal{N}(\mathbf{x}_k | \hat{\mathbf{F}}_{k|-k} \mathbf{x}_{k-1} + \hat{\mathbf{e}}_{k|-k}, \hat{\mathbf{Q}}_{k|-k}). \quad (32)$$

where we left out the arguments for the following quantities estimated by an RNN.

$$\begin{bmatrix} \hat{\mathbf{F}}_{k|-k}(\mathbf{y}_{-k}) \\ \hat{\mathbf{e}}_{k|-k}(\mathbf{y}_{-k}) \\ \hat{\mathbf{Q}}_{k|-k}(\mathbf{y}_{-k}) \end{bmatrix} := \mathbf{f}(\mathbf{h}_k^{\rightarrow}, \mathbf{h}_k^{\leftarrow}) \quad \begin{bmatrix} \mathbf{h}_k^{\rightarrow} \\ \mathbf{h}_k^{\leftarrow} \end{bmatrix} := \begin{bmatrix} \text{GRU}(\mathbf{h}_{k-1}^{\rightarrow}, \mathbf{y}_{k-1}) \\ \text{GRU}(\mathbf{h}_{k+1}^{\leftarrow}, \mathbf{y}_{k+1}) \end{bmatrix} \quad (33)$$

Applying the integral in eq. (30) we get

$$\int p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{y}_{-k}) p(\mathbf{x}_{k-1} | \mathbf{y}_{0:K}) d\mathbf{x}_{k-1} = \mathcal{N}(\mathbf{x}_k | \hat{\mathbf{x}}_{k|-k}, \hat{\mathbf{P}}_{k|-k}) \quad (34)$$

where we put

$$\hat{\mathbf{x}}_{k|-k} := \hat{\mathbf{F}}_{k|-k} \hat{\mathbf{x}}_{k-1|0:K} + \hat{\mathbf{e}}_{k|-k} \quad \hat{\mathbf{P}}_{k|-k} := \hat{\mathbf{F}}_{k|-k} \hat{\mathbf{P}}_{k-1|0:K} \hat{\mathbf{F}}_{k|-k}^\top + \hat{\mathbf{Q}}_{k|-k} \quad (35)$$

A data likelihood can be computed as follows

$$p(\mathbf{y}_k | \mathbf{y}_{-k}) = \int p(\mathbf{y}_k | \mathbf{x}_k) p(\mathbf{x}_k | \mathbf{y}_{-k}) d\mathbf{x}_k \quad (36)$$

$$= \int p(\mathbf{y}_k | \mathbf{x}_k) \int p(\mathbf{x}_{k-1}, \mathbf{x}_k | \mathbf{y}_{-k}) d\mathbf{x}_k d\mathbf{x}_{k-1} \quad (37)$$

$$= \int p(\mathbf{y}_k | \mathbf{x}_k) \int p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{y}_{-k}) p(\mathbf{x}_{k-1} | \mathbf{y}_{-k}) d\mathbf{x}_k d\mathbf{x}_{k-1} \quad (38)$$

$$\stackrel{\mathbf{x}_{k-1} \perp \mathbf{y}_k | \mathbf{y}_{-k}}{\approx} \int p(\mathbf{y}_k | \mathbf{x}_k) \int p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{y}_{-k}) p(\mathbf{x}_{k-1} | \mathbf{y}_{0:K}) d\mathbf{x}_k d\mathbf{x}_{k-1} \quad (39)$$

where we made the same approximation eq. (30) as before. It evaluates to

$$\int p(\mathbf{y}_k | \mathbf{x}_k) p(\mathbf{x}_k | \mathbf{y}_{-k}) d\mathbf{x}_k = \mathcal{N}(\mathbf{y}_k | \hat{\mathbf{y}}_{k|-k}, \hat{\mathbf{M}}_{k|-k}) \quad (40)$$

with

$$\hat{\mathbf{y}}_{k|-k} := \mathbf{H} \hat{\mathbf{x}}_{k|-k} \quad \hat{\mathbf{M}}_{k|-k} := \mathbf{H} \hat{\mathbf{P}}_{k|-k} \mathbf{H}^\top + \mathbf{R} \quad (41)$$

For fitting to the data we now would use a maximum-pseudo-likelihood (Gourieroux et al., 1984) approach by maximizing:

$$\sum_{k=0}^K \log p(\mathbf{y}_k | \mathbf{y}_{-k}) = \sum_{k=0}^K \log \mathcal{N}(\mathbf{y}_k | \mathbf{H} \hat{\mathbf{x}}_{k|-k}(\mathbf{y}_{-k}), \mathbf{H} \hat{\mathbf{P}}_{k|-k}(\mathbf{y}_{-k}) \mathbf{H}^\top + \mathbf{R}), \quad (42)$$

leading to minimizing the following self-supervised loss function:

$$\mathcal{L} := \sum_{k=0}^K \left[ (\hat{\mathbf{y}}_{k|-k} - \mathbf{y}_k)^\top \hat{\mathbf{M}}_{k|-k}^{-1} (\hat{\mathbf{y}}_{k|-k} - \mathbf{y}_k) + \log \det \hat{\mathbf{M}}_{k|-k} \right], \quad (43)$$

where  $\hat{\mathbf{y}}_{k|-k} := \mathbf{H} \hat{\mathbf{x}}_{k|-k}(\mathbf{y}_{-k})$  and  $\hat{\mathbf{M}}_{k|-k} := \mathbf{H} \hat{\mathbf{P}}_{k|-k}(\mathbf{y}_{-k}) \mathbf{H}^\top + \mathbf{R}$ .

## B GAUSSIAN CONDITIONING FORMULAS

Since many of the calculations used in this work are based on the Gaussian conditioning formulas, we provide them here. If

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \mathbf{P}), \quad (44)$$

$$p(\mathbf{y} | \mathbf{x}) = \mathcal{N}(\mathbf{y} | \mathbf{H}\mathbf{x} + \mathbf{b}, \mathbf{R}), \quad (45)$$

then

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{H}\boldsymbol{\mu} + \mathbf{b}, \mathbf{R} + \mathbf{H}\mathbf{P}\mathbf{H}^\top), \quad (46)$$

$$p(\mathbf{x} | \mathbf{y}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\Sigma} [\mathbf{H}^\top \mathbf{R}^{-1} (\mathbf{y} - \mathbf{b}) + \mathbf{P}^{-1} \boldsymbol{\mu}], \boldsymbol{\Sigma}), \quad (47)$$

with

$$\boldsymbol{\Sigma} = (\mathbf{P}^{-1} + \mathbf{H}^\top \mathbf{R}^{-1} \mathbf{H})^{-1}. \quad (48)$$

## C DIRECT PARAMETERIZATION OF $p(\mathbf{x} \mid \mathbf{y}_{<k})$

We show here how to directly parameterize  $p(\mathbf{x} \mid \mathbf{y}_{<k})$  and  $p(\mathbf{x} \mid \mathbf{y}_{-k})$ . This parameterization is referred to as the *recurrent model* in our experiments. The procedure is rather straightforward. For filtering, we put  $p(\mathbf{x}_k \mid \mathbf{y}_{<k}) = \mathcal{N}(\mathbf{x}_k \mid \hat{\mathbf{x}}_{k|<k}(\mathbf{y}_{<k}), \hat{\mathbf{P}}_{k|<k}(\mathbf{y}_{<k}))$ . We model:

$$\hat{\mathbf{x}}_{k|<k} := \mathbf{e}_k, \quad \hat{\mathbf{P}}_{k|<k} := \mathbf{L}_k \mathbf{L}_k^\top, \quad \begin{bmatrix} \mathbf{e}_k \\ \mathbf{L}_k \end{bmatrix} := \mathbf{f}(\mathbf{h}_k^\rightarrow), \quad (49)$$

where  $\mathbf{L}_k$  is a cholesky factor and  $\mathbf{f}$  is a multi-layer perceptron. The argument  $\mathbf{h}_k^\rightarrow$  is recursively given by

$$\mathbf{h}_k^\rightarrow := \text{GRU}(\mathbf{y}_{k-1}, \mathbf{h}_{k-1}^\rightarrow), \quad (50)$$

where we employ a Gated Recurrent Unit (GRU) network (Cho et al., 2014). Note that this parameterization is equivalent to the model described in the main paper with  $\mathbf{F}_k := \mathbf{0}$ .

For smoothing,  $p(\mathbf{x}_k \mid \mathbf{y}_{-k}) = \mathcal{N}(\mathbf{x}_k \mid \hat{\mathbf{x}}_{k|-k}(\mathbf{y}_{-k}), \hat{\mathbf{P}}_{k|-k}(\mathbf{y}_{-k}))$ .

$$\hat{\mathbf{x}}_{k|-k} := \mathbf{e}_k, \quad \hat{\mathbf{P}}_{k|-k} := \mathbf{L}_k \mathbf{L}_k^\top, \quad \begin{bmatrix} \mathbf{e}_k \\ \mathbf{L}_k \end{bmatrix} := \mathbf{f}(\mathbf{h}_k^\rightarrow, \mathbf{h}_k^\leftarrow), \quad (51)$$

where  $\mathbf{L}_k$  is a cholesky factor and  $\mathbf{f}$  is a multi-layer perceptron.

$$\mathbf{h}_k^\rightarrow := \text{GRU}(\mathbf{y}_{k-1}, \mathbf{h}_{k-1}^\rightarrow), \quad \mathbf{h}_k^\leftarrow := \text{GRU}(\mathbf{y}_{k+1}, \mathbf{h}_{k+1}^\leftarrow). \quad (52)$$

Once  $p(\mathbf{x} \mid \mathbf{y}_{<k})$  (filtering) or  $p(\mathbf{x} \mid \mathbf{y}_{-k})$  (smoothing) is obtained, all the procedures for inference and optimization described in the main paper and appendix A remain the same.

## D BIAS-VARIANCE-NOISE DECOMPOSITION OF THE SELF-SUPERVISED GENERALIZATION ERROR

Any estimate  $\hat{\mathbf{x}}_k = \hat{\mathbf{x}}_k(\mathbf{y}_{-k})$  for  $\mathbf{x}_k$  that is not dependent on  $\mathbf{y}_k$  will give us a bias-variance-noise decomposition of the generalization error. Note that this setting covers both the filtering and smoothing case. Define ‘‘optimal model’’  $\hat{\mathbf{y}}_k^* := \mathbb{E}[\mathbf{y}_k \mid \mathbf{y}_{-k}]$ , then under the specified generative model (section 3) we have

$$\mathbb{E} \left[ \|\hat{\mathbf{y}}_k - \mathbf{y}_k\|_2^2 \mid \mathbf{y}_{-k} \right] = \mathbb{E} \left[ \|\hat{\mathbf{y}}_k - \hat{\mathbf{y}}_k^* + \hat{\mathbf{y}}_k^* - \mathbf{y}_k\|_2^2 \mid \mathbf{y}_{-k} \right] \quad (53)$$

$$= \mathbb{E} \left[ \|\hat{\mathbf{y}}_k - \hat{\mathbf{y}}_k^*\|_2^2 \mid \mathbf{y}_{-k} \right] + \text{Var}[\mathbf{y}_k \mid \mathbf{y}_{-k}] + 2\mathbb{E} \left[ (\hat{\mathbf{y}}_k - \hat{\mathbf{y}}_k^*)^\top (\hat{\mathbf{y}}_k^* - \mathbf{y}_k) \mid \mathbf{y}_{-k} \right] \quad (54)$$

$$= \|\hat{\mathbf{y}}_k - \hat{\mathbf{y}}_k^*\|_2^2 + \text{Var}[\mathbf{y}_k \mid \mathbf{y}_{-k}] + 2(\hat{\mathbf{y}}_k - \hat{\mathbf{y}}_k^*)^\top \underbrace{\mathbb{E}[(\hat{\mathbf{y}}_k^* - \mathbf{y}_k) \mid \mathbf{y}_{-k}]}_0 \quad (55)$$

$$\text{Var}[\mathbf{y}_k \mid \mathbf{y}_{-k}] = \mathbb{E} \left[ \|\hat{\mathbf{y}}_k^* - \mathbf{H}\mathbf{x}_k - \mathbf{n}_k\|_2^2 \mid \mathbf{y}_{-k} \right] \quad (56)$$

$$= \mathbb{E} \left[ \|\hat{\mathbf{y}}_k^* - \mathbf{H}\mathbf{x}_k\|_2^2 \mid \mathbf{y}_{-k} \right] + \mathbb{E} \left[ \|\mathbf{n}_k\|_2^2 \mid \mathbf{y}_{-k} \right] + 2\mathbb{E} \left[ \underbrace{(\mathbf{H}\mathbf{x}_k - \hat{\mathbf{y}}_k^*)^\top \mathbf{n}_k}_{0} \mid \mathbf{y}_{-k} \right] \quad (57)$$

$$\mathbf{x}_k, \mathbf{y}_{-k} \perp \mathbf{n}_k \quad \mathbb{E} \left[ \|\hat{\mathbf{y}}_k^* - \mathbf{H}\mathbf{x}_k\|_2^2 \mid \mathbf{y}_{-k} \right] + \text{tr}(\mathbf{R}) \quad (58)$$

Thus,

$$\mathbb{E} \left[ \|\hat{\mathbf{y}}_k - \mathbf{y}_k\|_2^2 \mid \mathbf{y}_{-k} \right] = \|\hat{\mathbf{y}}_k - \hat{\mathbf{y}}_k^*\|_2^2 + \mathbb{E} \left[ \|\hat{\mathbf{y}}_k^* - \mathbf{H}\mathbf{x}_k\|_2^2 \mid \mathbf{y}_{-k} \right] + \text{tr}(\mathbf{R}) \quad (59)$$

Note that  $\hat{\mathbf{y}}_k^* = \mathbb{E}[\mathbf{H}\mathbf{x}_k \mid \mathbf{y}_{-k}] = \mathbf{H}\mathbb{E}[\mathbf{x}_k \mid \mathbf{y}_{-k}]$ . Then, define our model  $\hat{\mathbf{y}}_k := \mathbf{H}\hat{\mathbf{x}}_k$ . The *reducible* part of the error becomes

$$\|\hat{\mathbf{y}}_k^* - \hat{\mathbf{y}}_k\|_2^2 = \|\mathbf{H}(\mathbb{E}[\mathbf{x}_k \mid \mathbf{y}_{-k}] - \hat{\mathbf{x}}_k)\|_2^2 \quad (60)$$

For this reason, the model output  $\hat{\mathbf{x}}_k$  approaches the optimal model under minimization of the self-supervised error (perturbed by  $\mathbf{H}$ ).

Additionally, the model  $\hat{\mathbf{y}}_k$  approaches  $\mathbf{H}\mathbf{x}_k$  (the uncorrupted measurement) under this criterion.

$$\mathbb{E} \left[ \|\hat{\mathbf{y}}_k - \mathbf{y}_k\|^2 \right] = \mathbb{E} \left[ \|\hat{\mathbf{y}}_k - \mathbf{H}\mathbf{x}_k + \mathbf{H}\mathbf{x}_k - \mathbf{y}_k\|^2 \right] \quad (61)$$

$$= \mathbb{E} \left[ \|\hat{\mathbf{y}}_k - \mathbf{H}\mathbf{x}_k\|^2 \right] + \mathbb{E} \left[ \|\mathbf{H}\mathbf{x}_k - \mathbf{y}_k\|^2 \right] + 2 \mathbb{E} \left[ (\hat{\mathbf{y}}_k - \mathbf{H}\mathbf{x}_k)^\top (\mathbf{H}\mathbf{x}_k - \mathbf{y}_k) \right] \quad (62)$$

$$= \mathbb{E} \left[ \|\hat{\mathbf{y}}_k - \mathbf{H}\mathbf{x}_k\|^2 \right] + \mathbb{E} \left[ \|\mathbf{H}\mathbf{x}_k - \mathbf{y}_k\|^2 \right] + 2 \mathbb{E} \left[ (\hat{\mathbf{y}}_k - \mathbf{H}\mathbf{x}_k)^\top \underbrace{\mathbb{E} [\mathbf{H}\mathbf{x}_k - \mathbf{y}_k]}_0 \right] \quad (63)$$

$$= \mathbb{E} \left[ \|\hat{\mathbf{y}}_k - \mathbf{H}\mathbf{x}_k\|^2 \right] + \text{tr}(\mathbf{R}) + \underbrace{\mathbb{E} [\|\mathbf{H}\mathbf{x}_k - \mathbf{y}_k\|^2]}_0 \quad (64)$$

If we then take  $\hat{\mathbf{y}}_k := \mathbf{H}\hat{\mathbf{x}}_k$  then the reducible part of the error approximates the true  $\mathbf{x}_k$  (perturbed by  $\mathbf{H}$ ).

$$\mathbb{E} \left[ \|\hat{\mathbf{y}}_k - \mathbf{H}\mathbf{x}_k\|^2 \right] = \mathbb{E} \left[ \|\mathbf{H}(\hat{\mathbf{x}}_k - \mathbf{x}_k)\|^2 \right] \quad (65)$$

## E ALGORITHMS

**Algorithm 1:** Recursive Filter (Inference)

**input :** Data (time-series)  $\mathbf{y}_{0:K} = (\mathbf{y}_0, \dots, \mathbf{y}_K)$ , emission matrices  $\mathbf{H}$  and  $\mathbf{R}$ , parameters  $\phi$ .  
**output:**

1. For training: Loss value  $\mathcal{L}_K$  and its gradient  $\nabla_{\phi} \mathcal{L}_K$ .
2. For inference:  $\hat{\mathbf{x}}_{k|\leq k}$  and  $\hat{\mathbf{P}}_{k|\leq k}$  for  $k$  in  $0, \dots, K$ . Inference is done via:  

$$p(\mathbf{x}_k | \mathbf{y}_{\leq k}) = \mathcal{N}(\mathbf{x}_k | \hat{\mathbf{x}}_{k|\leq k}, \hat{\mathbf{P}}_{k|\leq k}).$$
3. For linearized smoothing (section 6):  $\hat{\mathbf{F}}_{k|<k}$ ,  $\hat{\mathbf{P}}_{k|<k}$ ,  $\hat{\mathbf{x}}_{k|\leq k}$  and  $\hat{\mathbf{P}}_{k|\leq k}$ , for  $k$  in  $0, \dots, K$ .
4. For forecasting:  $\hat{\mathbf{x}}_{K+1|<(K+1)}$  and  $\hat{\mathbf{P}}_{K+1|<(K+1)}$ . Forecasting is done via:  

$$p(\mathbf{x}_{K+1} | \mathbf{y}_{0:K}) = \mathcal{N}(\mathbf{x}_{K+1} | \hat{\mathbf{x}}_{K+1|<(K+1)}, \hat{\mathbf{P}}_{K+1|<(K+1)}).$$

$\mathbf{h}_0 := \mathbf{0}$

$\mathcal{L}_{-1} := 0$

$\hat{\mathbf{P}}_{0|<0} := \hat{\mathbf{Q}}_0$

$\hat{\mathbf{x}}_{0|<0} := \hat{\mathbf{e}}_0$

**for**  $k$  in  $0, \dots, K$  **do**

$$\hat{\mathbf{B}}_{k|<k} := (\mathbf{H} \hat{\mathbf{P}}_{k|<k} \mathbf{H}^{\top} + \mathbf{R})^{-1}$$

$$\hat{\mathbf{y}}_{k|<k} := \mathbf{H} \hat{\mathbf{x}}_{k|<k}$$

$$\mathcal{L}_k := \mathcal{L}_{k-1} + (\mathbf{y}_k - \hat{\mathbf{y}}_{k|<k})^{\top} \hat{\mathbf{B}}_{k|<k} (\mathbf{y}_k - \hat{\mathbf{y}}_{k|<k}) - \log \det \hat{\mathbf{B}}_{k|<k}$$

$$\hat{\mathbf{K}}_k := \hat{\mathbf{P}}_{k|<k} \mathbf{H}^{\top} \hat{\mathbf{B}}_{k|<k}$$

$$\hat{\mathbf{P}}_{k|\leq k} := \hat{\mathbf{P}}_{k|<k} - \hat{\mathbf{K}}_k \mathbf{H} \hat{\mathbf{P}}_{k|<k}$$

$$\hat{\mathbf{x}}_{k|\leq k} := \hat{\mathbf{x}}_{k|<k} + \hat{\mathbf{K}}_k (\mathbf{y}_k - \hat{\mathbf{y}}_{k|<k})$$

$$\mathbf{h}_{k+1} := \text{GRU}_{\phi}(\mathbf{h}_k, \mathbf{y}_k)$$

$$\begin{bmatrix} \hat{\mathbf{e}}_{k+1|<(k+1)} \\ \hat{\mathbf{F}}_{k+1|<(k+1)} \\ \hat{\mathbf{L}}_{k+1|<(k+1)} \end{bmatrix} := \mathbf{f}_{\phi}(\mathbf{h}_{k+1})$$

$$\hat{\mathbf{Q}}_{k+1|<(k+1)} := \hat{\mathbf{L}}_{k+1|<(k+1)} \hat{\mathbf{L}}_{k+1|<(k+1)}^{\top}$$

$$\hat{\mathbf{P}}_{k+1|<(k+1)} := \hat{\mathbf{F}}_{k+1|<(k+1)} \hat{\mathbf{P}}_{k|\leq k} \hat{\mathbf{F}}_{k+1|<(k+1)}^{\top} + \hat{\mathbf{Q}}_{k+1|<(k+1)}$$

$$\hat{\mathbf{x}}_{k+1|<(k+1)} := \hat{\mathbf{F}}_{k+1|<(k+1)} \hat{\mathbf{x}}_{k|\leq k} + \hat{\mathbf{e}}_{k+1|<(k+1)}$$

**end**

For the training case we use backpropagation through the above loop to compute  $\nabla_{\phi} \mathcal{L}_K$ .

**Algorithm 2:** Recursive Filter (Training)

**input :** Data (time-series)  $\mathbf{y}_{0:K} = (\mathbf{y}_0, \dots, \mathbf{y}_K)$ , emission matrices  $\mathbf{H}$  and  $\mathbf{R}$ , initialized parameters  $\phi_0$ , number of training rounds  $I$ .

**output:** Model parameters  $\phi^*$  for inference at test-time.

**for**  $i$  in  $0, \dots, I$  **do**

    Obtain  $\mathcal{L}_K^{(i)}$  and  $\nabla_{\phi} \mathcal{L}_K^{(i)}$  from algorithm 1.

    Run preferred optimizer step to update parameters  $\phi$  with  $\nabla_{\phi} \mathcal{L}_K^{(i)}$  (and  $\mathcal{L}_K^{(i)}$ ).

**end**



**Algorithm 3:** Parameterized Recursive Smoother (Inference)

**input** : Data (time-series)  $\mathbf{y}_{0:K} = (\mathbf{y}_0, \dots, \mathbf{y}_K)$ , emission matrices  $\mathbf{H}$  and  $\mathbf{R}$ , initialized parameters  $\phi_0$

**output:**

1. Loss value  $\mathcal{L}_K$  and its gradient w.r.t. all model parameters  $\nabla_{\phi}(\mathcal{L}_K)$  for training.
2. For all  $k$  in  $0, \dots, K$ :  $\hat{\mathbf{x}}_{k|0:K}, \hat{\mathbf{P}}_{k|0:K}$ . These can be used for inference through  $p(\mathbf{x}_k | \mathbf{y}_{0:K}) = \mathcal{N}(\hat{\mathbf{x}}_{k|0:K}, \hat{\mathbf{P}}_{k|0:K})$ .

$$\mathbf{h}_0^{\rightarrow} := \mathbf{0}$$

$$\mathbf{h}_{K+1}^{\leftarrow} := \mathbf{0}$$

$$\mathcal{L}_{-1}^{(i)} := 0$$

$$\hat{\mathbf{P}}_{0|-0} := \hat{\mathbf{Q}}_0$$

$$\hat{\mathbf{x}}_{0|-0} := \hat{\mathbf{e}}_0$$

**for**  $k$  in  $0, \dots, K$  **do**

$$\hat{\mathbf{B}}_{k|-k} := (\mathbf{H} \hat{\mathbf{P}}_{k|-k} \mathbf{H}^{\top} + \mathbf{R})^{-1}$$

$$\hat{\mathbf{y}}_{k|-k} := \mathbf{H} \hat{\mathbf{x}}_{k|-k}$$

$$\mathcal{L}_k^{(i)} := \mathcal{L}_{k-1}^{(i)} + (\mathbf{y}_k - \hat{\mathbf{y}}_{k|-k})^{\top} \hat{\mathbf{B}}_{k|-k} (\mathbf{y}_k - \hat{\mathbf{y}}_{k|-k}) - \log \det \hat{\mathbf{B}}_{k|-k}$$

$$\hat{\mathbf{K}}_k := \hat{\mathbf{P}}_{k|-k} \mathbf{H}^{\top} \hat{\mathbf{B}}_{k|-k}$$

$$\hat{\mathbf{P}}_{k|-k} := \hat{\mathbf{P}}_{k|-k} - \hat{\mathbf{K}}_k \mathbf{H} \hat{\mathbf{P}}_{k|-k}$$

$$\hat{\mathbf{x}}_{k|-k} := \hat{\mathbf{x}}_{k|-k} + \hat{\mathbf{K}}_k (\mathbf{y}_k - \hat{\mathbf{y}}_{k|-k})$$

$$\mathbf{h}_{k+1}^{\rightarrow} := \text{GRU}_{\phi}(\mathbf{h}_k^{\rightarrow}, \mathbf{y}_k)$$

$$\mathbf{h}_{k+1}^{\leftarrow} := \text{GRU}_{\phi}(\mathbf{h}_k^{\leftarrow}, \mathbf{y}_{k+1})$$

$$\begin{bmatrix} \hat{\mathbf{e}}_{k+1|-(k+1)} \\ \hat{\mathbf{F}}_{k+1|-(k+1)} \\ \hat{\mathbf{L}}_{k+1|-(k+1)} \end{bmatrix} := \mathbf{f}_{\phi}(\mathbf{h}_{k+1}^{\rightarrow}, \mathbf{h}_{k+1}^{\leftarrow})$$

$$\hat{\mathbf{Q}}_{k+1|-(k+1)} := \hat{\mathbf{L}}_{k+1|-(k+1)} \hat{\mathbf{L}}_{k+1|-(k+1)}^{\top}$$

$$\hat{\mathbf{P}}_{k+1|-(k+1)} := \hat{\mathbf{F}}_{k+1|-(k+1)} \hat{\mathbf{P}}_{k|0:K} \hat{\mathbf{F}}_{k+1|-(k+1)}^{\top} + \hat{\mathbf{Q}}_{k+1|-(k+1)}$$

$$\hat{\mathbf{x}}_{k+1|-(k+1)} := \hat{\mathbf{F}}_{k+1|-(k+1)} \hat{\mathbf{x}}_{k|0:K} + \hat{\mathbf{e}}_{k+1|-(k+1)}$$

**end**

For training case we also use backpropagation through the above loop to compute  $\nabla_{\phi} \mathcal{L}_K$

**Algorithm 4:** Parameterized Recursive Smoother (Training)

**input** : Data (time-series)  $\mathbf{y}_{0:K} = (\mathbf{y}_0, \dots, \mathbf{y}_K)$ , emission matrices  $\mathbf{H}$  and  $\mathbf{R}$ , initialized parameters  $\phi_0$ , number of training rounds  $I$

**output:** Model parameters  $\phi^*$  for inference at test-time.

**for**  $i$  in  $1, \dots, I$  **do**

    Obtain  $\nabla_{\phi} \mathcal{L}_K^{(i)}$  from algorithm 3.

    Run preferred optimizer step.

**end**

**Algorithm 5:** Linearized Smoother

**input :** Values of:  $\hat{\mathbf{F}}_{k|<k}$ ,  $\hat{\mathbf{P}}_{k|<k}$ ,  $\hat{\mathbf{P}}_{k|\leq k}$ ,  $\hat{\mathbf{x}}_{k|\leq k}$  for all  $k = 0, \dots, K$ , obtained from the recursive filter algorithm.

**output:** Linearly smoothed distributions  $p(\mathbf{x}_k | \mathbf{y}_{0:K}) = \mathcal{N}(\mathbf{x}_k | \hat{\mathbf{z}}_k, \hat{\mathbf{G}}_k)$  for all  $k = 0, \dots, K$ .

$$\hat{\mathbf{z}}_K := \hat{\mathbf{x}}_{K|\leq K}$$

$$\hat{\mathbf{G}}_K := \hat{\mathbf{P}}_{K|\leq K}$$

**for**  $k = K, \dots, 1$  **do**

$$\hat{\mathbf{J}}_{k-1|k} := \hat{\mathbf{P}}_{k-1|\leq k-1} \hat{\mathbf{F}}_{k|<k}^\top \hat{\mathbf{P}}_{k|<k}^{-1}$$

$$\hat{\mathbf{G}}_{k-1} := \hat{\mathbf{P}}_{k-1|\leq k-1} + \hat{\mathbf{J}}_{k-1|k} \left( \hat{\mathbf{P}}_{k|\leq k} - \hat{\mathbf{P}}_{k|<k} \right) \hat{\mathbf{J}}_{k-1|k}^\top$$

$$\hat{\mathbf{z}}_{k-1} := \hat{\mathbf{x}}_{k-1|\leq k-1} + \hat{\mathbf{J}}_{k-1|k} (\hat{\mathbf{z}}_k - \hat{\mathbf{x}}_{k|\leq k})$$

**end**

**Algorithm 6:** Recurrent Smoother (Training)

**input :** Training data (time-series)  $\mathbf{y}_{0:K} = (\mathbf{y}_0, \dots, \mathbf{y}_K)$ , emission function  $\mathcal{N}(\mathbf{H}\mathbf{x}_k, \mathbf{R})$ , initialized parameters  $\phi_0$ , number of training iterations  $n$ .

**output:** Optimized parameters  $\phi^*$

**for**  $i$  in  $1$  to  $n$  **do**

$$\mathbf{h}_0 := \mathbf{0}$$

$$\mathbf{h}_{K+1} := \mathbf{0}$$

$$\mathcal{L}^{(i)} := 0$$

**for**  $k$  in  $0$  to  $K$  **do**

$$\mathbf{h}_k^{\rightarrow} := \text{GRU}_\phi(\mathbf{h}_{k-1}^{\rightarrow}, \mathbf{y}_{k-1})$$

$$\mathbf{h}_k^{\leftarrow} := \text{GRU}_\phi(\mathbf{h}_{k-1}^{\leftarrow}, \mathbf{y}_{k+1})$$

$$\begin{bmatrix} \hat{\mathbf{x}}_{k|>k} \\ \hat{\mathbf{L}}_{k|>k} \end{bmatrix} := \mathbf{f}_\phi(\mathbf{h}_k^{\leftarrow}, \mathbf{h}_k^{\rightarrow})$$

$$\hat{\mathbf{P}}_{k|>k} := \hat{\mathbf{L}}_{k|>k} \hat{\mathbf{L}}_{k|>k}^\top$$

$$\hat{\mathbf{y}}_{k|>k} := \mathbf{H} \hat{\mathbf{x}}_{k|>k}$$

$$\hat{\mathbf{B}}_{k|>k} := \left( \mathbf{H} \hat{\mathbf{P}}_{k|>k} \mathbf{H}^\top + \mathbf{R} \right)^{-1}$$

$$\mathcal{L}_k^{(i)} := \mathcal{L}_{k-1}^{(i)} + (\mathbf{y}_k - \hat{\mathbf{y}}_{k|>k})^\top \hat{\mathbf{B}}_{k|>k} (\mathbf{y}_k - \hat{\mathbf{y}}_{k|>k}) - \log \det \hat{\mathbf{B}}_{k|>k}$$

**end**

Compute  $\nabla_\phi \mathcal{L}_K^{(i)}$  and apply SGD step with respect to all model parameters, which amounts to backpropagation through the above calculations.

**end**

*For the filter variant, the steps that involve the backward direction ( $\leftarrow$ ) are left out.*

**Algorithm 7:** Recurrent Smoother (Inference)

**input :** Test data  $\mathbf{y}_{0:K} = (\mathbf{y}_0, \dots, \mathbf{y}_K)$ , trained parameters  $\phi^*$ , emission function  $\mathcal{N}(\mathbf{H}\mathbf{x}, \mathbf{R})$ .

**output:** Inferred posteriors  $p(\mathbf{x}_k | \mathbf{y}_{0:K}) = \mathcal{N}(\mathbf{x}_k | \hat{\mathbf{x}}_{k|0:K}, \hat{\mathbf{P}}_{k|0:K})$  for all  $k$

$\mathbf{h}_0 := \mathbf{0}$

$\mathbf{h}_{K+1} := \mathbf{0}$

**for**  $k$  in 0 to  $K$  **do**

$$\mathbf{h}_k^{\rightarrow} := \text{GRU}_{\phi}(\mathbf{h}_{k-1}^{\rightarrow}, \mathbf{y}_{k-1})$$

$$\mathbf{h}_k^{\leftarrow} := \text{GRU}_{\phi}(\mathbf{h}_{k-1}^{\leftarrow}, \mathbf{y}_{k+1})$$

$$\begin{bmatrix} \hat{\mathbf{x}}_{k|k-k} \\ \hat{\mathbf{L}}_{k|k-k} \end{bmatrix} := \mathbf{f}_{\phi}(\mathbf{h}_k^{\leftarrow}, \mathbf{h}_k^{\rightarrow})$$

$$\hat{\mathbf{P}}_{k|k-k} := \hat{\mathbf{L}}_{k|k-k} \hat{\mathbf{L}}_{k|k-k}^{\top}$$

$$\hat{\mathbf{K}}_k := \hat{\mathbf{P}}_{k|k-k} \mathbf{H}^{\top} \left( \mathbf{H} \hat{\mathbf{P}}_{k|k-k} \mathbf{H}^{\top} + \mathbf{R} \right)^{-1}$$

$$\hat{\mathbf{x}}_{k|0:K} := \hat{\mathbf{x}}_{k|k-k} + \hat{\mathbf{K}}_k (\mathbf{y}_k - \mathbf{H} \hat{\mathbf{x}}_{k|k-k})$$

$$\hat{\mathbf{P}}_{k|0:K} := \hat{\mathbf{P}}_{k|k-k} - \hat{\mathbf{K}}_k \mathbf{H} \hat{\mathbf{P}}_{k|k-k}$$

**end**

*For the filter variant, the steps that involve the backward direction ( $\leftarrow$ ) are left out.*

## F EXPERIMENTS: DETAILS

### F.1 LINEAR DYNAMICS

As specified in the main paper, the dynamics are according to

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & -c & 1 \\ 0 & -\tau c & 0 \end{bmatrix} \begin{bmatrix} p \\ v \\ a \end{bmatrix}. \quad (66)$$

Since these are linear transitions, we can calculate any transition directly using  $\mathbf{x}(t + \Delta t) = e^{\mathbf{A}\Delta t}\mathbf{x}(t)$ .

$$\mathbf{F} := \begin{bmatrix} e^{\mathbf{A}} & 0 \\ 0 & e^{\mathbf{A}} \end{bmatrix} \quad \mathbf{Q} := \begin{bmatrix} \bar{\mathbf{Q}} & 0 \\ 0 & \bar{\mathbf{Q}} \end{bmatrix} \quad (67)$$

We used  $c = 0.06$ ,  $\tau = 0.17$ ,  $\Delta t := 1$  and covariance

$$\bar{\mathbf{Q}} := 0.1^2 \cdot \begin{bmatrix} \frac{1}{3} & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 3 \end{bmatrix} \quad (68)$$

The matrix exponential is computed using Bader et al. (2019). The parameters for the emission distribution:

$$\mathbf{H} := \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix} \quad \mathbf{R} := 0.5^2 \cdot \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad (69)$$

We simulate a  $K := 131,072$  trajectory for training,  $K := 16,384$  trajectory for validation and  $K := 32,768$  for testing. The  $\tilde{\mathbf{F}}$  that is used in the (non-optimal) Kalman filter and recursive model is computed as follows:

$$e^{\mathbf{A}\Delta t} \approx \tilde{\mathbf{F}} := \sum_{n=0}^1 (\Delta t \mathbf{A}^n) / n! \quad (70)$$

### F.2 LORENZ EQUATIONS

We simulate a Lorenz system according to

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} = \begin{bmatrix} -\sigma & \sigma & 0 \\ \rho - x_1 & -1 & 0 \\ x_2 & 0 & -\beta \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}. \quad (71)$$

We integrate the system using  $dt = 0.00001$  and sample it uniformly at  $\Delta t = 0.05$ . We use  $\rho = 28$ ,  $\sigma = 10$ ,  $\beta = 8/3$ . The transition in  $\Delta t$  arbitrary time-steps is linearly approximated by a Taylor expansion and used in the Kalman smoother and recursive models.

$$e^{\mathbf{A}|_{\mathbf{x}_k} \Delta t} \approx \tilde{\mathbf{F}}_k := \sum_{n=0}^2 (\Delta t \mathbf{A}|_{\mathbf{x}_k})^n / n! \quad (72)$$

We simulate  $K := 131,072$  steps for training,  $K := 32,768$  for testing and  $K := 16,384$  for validation. We have  $\mathbf{H} := \mathbf{I}$  and thus  $\mathbf{x} \in \mathbb{R}^3$  and  $\mathbf{y} \in \mathbb{R}^3$ . We use  $\mathbf{R} := 0.5^2 \mathbf{I}$ .

## G PARAMETERIZED SMOOTHING: ALTERNATIVE POSTERIOR EVALUATION

We would like to point out that the distribution  $p(\mathbf{x}_k | \mathbf{y}_{0:K})$  can be obtained without making assumption eq. (30), which we stretch out here. Initial experiments showed that using these calculations the model did not converge as smoothly as when using the ones stated before. However, it could be of interest to further investigate. Returning to the posterior of interest

$$p(\mathbf{x}_k | \mathbf{y}_{0:K}) = \int p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{y}_{0:K}) p(\mathbf{x}_{k-1} | \mathbf{y}_{0:K}) d\mathbf{x}_{k-1} \quad (73)$$

$$= \int \frac{p(\mathbf{y}_k | \mathbf{x}_k)}{p(\mathbf{y}_k | \mathbf{x}_{k-1}, \mathbf{y}_{-k})} p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{y}_{-k}) p(\mathbf{x}_{k-1} | \mathbf{y}_{0:K}) d\mathbf{x}_{k-1}. \quad (74)$$

We have

$$p(\mathbf{x}_{k-1} \mid \mathbf{y}_{0:K}) = \mathcal{N}\left(\mathbf{x}_{k-1} \mid \hat{\mathbf{x}}_{k-1|0:K}(\mathbf{y}_{0:K}), \hat{\mathbf{P}}_{k-1|0:K}(\mathbf{y}_{0:K})\right) \quad (75)$$

as the previous time-step's posterior.

$$p(\mathbf{x}_k \mid \mathbf{x}_{k-1}, \mathbf{y}_{-k}) = \mathcal{N}\left(\mathbf{x}_k \mid \hat{\mathbf{F}}_{k|-k} \mathbf{x}_{k-1} + \hat{\mathbf{e}}_{k|-k}, \hat{\mathbf{Q}}_{k|-k}\right) \quad (76)$$

where  $\hat{\mathbf{F}}_{k|-k}(\mathbf{y}_{-k})$ ,  $\hat{\mathbf{e}}_{k|-k}(\mathbf{y}_{-k})$  and  $\hat{\mathbf{Q}}_{k|-k}(\mathbf{y}_{-k})$  are estimated by a neural network. Combining this with noise model eq. (3) we get:

$$\begin{aligned} p(\mathbf{x}_k \mid \mathbf{x}_{k-1}, \mathbf{y}_{0:K}) \\ = \mathcal{N}\left(\mathbf{x}_k \mid \hat{\mathbf{F}}_{k|-k} \mathbf{x}_{k-1} + \hat{\mathbf{e}}_{k|-k} + \mathbf{K}_k (\mathbf{y}_k - \mathbf{H} (\hat{\mathbf{F}}_{k|-k} \mathbf{x}_{k-1} + \hat{\mathbf{e}}_{k|-k})), \hat{\mathbf{Q}}_{k|-k} - \mathbf{K}_k \mathbf{H} \hat{\mathbf{Q}}_{k|-k}\right) \end{aligned} \quad (77)$$

$$= \mathcal{N}\left(\mathbf{x}_k \mid (\mathbf{I} - \mathbf{K}_n \mathbf{H}) (\hat{\mathbf{F}}_{k|-k} \mathbf{x}_{k-1} + \hat{\mathbf{e}}_{k|-k}) + \mathbf{K}_k \mathbf{y}_k, (\mathbf{I} - \mathbf{K}_n \mathbf{H}) \hat{\mathbf{Q}}_{k|-k}\right), \quad (78)$$

where we directly applied the Woodbury matrix identity to obtain Kalman gain matrix

$$\mathbf{K}_k := \hat{\mathbf{Q}}_{k|-k} \mathbf{H}^\top \left( \mathbf{H} \hat{\mathbf{Q}}_{k|-k} \mathbf{H}^\top + \mathbf{R} \right)^{-1} \quad (79)$$

Then, applying the integral we get:

$$p(\mathbf{x}_k \mid \mathbf{y}_{0:K}) = \mathcal{N}\left(\mathbf{x}_k \mid \hat{\mathbf{x}}_{k|0:K}, \hat{\mathbf{P}}_{k|0:K}\right), \quad (80)$$

with:

$$\hat{\mathbf{P}}_{k|0:K} = (\mathbf{I} - \mathbf{K}_n \mathbf{H}) \hat{\mathbf{F}}_{k|-k} \hat{\mathbf{P}}_{k-1|0:K} \hat{\mathbf{F}}_{k|-k}^\top (\mathbf{I} - \mathbf{K}_n \mathbf{H})^\top + (\mathbf{I} - \mathbf{K}_n \mathbf{H}) \hat{\mathbf{Q}}_{k|-k}, \quad (81)$$

and

$$\hat{\mathbf{x}}_{k|0:K} = \hat{\mathbf{F}}_{k|-k} \hat{\mathbf{x}}_{k-1|0:K} + \hat{\mathbf{e}}_{k|-k} + \mathbf{K}_k (\mathbf{y}_k - \mathbf{H} (\hat{\mathbf{F}}_{k|-k} \hat{\mathbf{x}}_{k-1|0:K} + \hat{\mathbf{e}}_{k|-k})) \quad (82)$$

$$= (\mathbf{I} - \mathbf{K}_n \mathbf{H}) \left( \hat{\mathbf{F}}_{k|-k} \hat{\mathbf{x}}_{k-1|0:K} + \hat{\mathbf{e}}_{k|-k} \right) + \mathbf{K}_k \mathbf{y}_k. \quad (83)$$