

Convergence Rates of Gaussian ODE Filters

Hans Kersting*, T. J. Sullivan†, and Philipp Hennig‡*

Abstract. A recently-introduced class of probabilistic (uncertainty-aware) solvers for ordinary differential equations (ODEs) applies Gaussian (Kalman) filtering to initial value problems. These methods model the true solution x and its first q derivatives *a priori* as a Gauss–Markov process \mathbf{X} , which is then iteratively conditioned on information about \dot{x} . We prove worst-case local convergence rates of order h^{q+1} for a wide range of versions of this Gaussian ODE filter, as well as global convergence rates of order h^q in the case of $q = 1$ and an integrated Brownian motion prior, and analyse how inaccurate information on \dot{x} coming from approximate evaluations of f affects these rates. Moreover, we present explicit formulas for the steady states and show that the posterior confidence intervals are well calibrated in all considered cases that exhibit global convergence—in the sense that they globally contract at the same rate as the truncation error.

Key words. probabilistic numerics, ordinary differential equations, initial value problems, numerical analysis, Gaussian processes, Markov processes

AMS subject classifications. 60G15, 60J70, 62G20, 62M05, 65C20, 65L05

1. Introduction. A solver of an initial value problem (IVP) outputs an approximate solution $\hat{x}: [0, T] \rightarrow \mathbb{R}^d$ of an ordinary differential equation (ODE) with initial condition:

$$(1.1) \quad x^{(1)}(t) := \frac{dx}{dt}(t) = f(x(t)), \quad \forall t \in [0, T], \quad x(0) = x_0 \in \mathbb{R}^d.$$

(Without loss of generality, we simplify the presentation by restricting attention to the autonomous case.) The numerical solution \hat{x} is computed by iteratively collecting information on $x^{(1)}(t)$ by evaluating $f: \mathbb{R}^d \rightarrow \mathbb{R}^d$ at a numerical estimate $\hat{x}(t)$ of $x(t)$ and using these approximate evaluations of the time derivative to extrapolate along the time axis. In other words, the numerical solution (or *estimator*) \hat{x} of the exact solution (or *estimand*) x is calculated based on evaluations of the vector field f (or *data*). Accordingly, we treat \hat{x} itself as an estimator, i.e. a statistic that translates evaluations of f into a probability distribution over $C^1([0, T]; \mathbb{R}^d)$, the space of continuously differentiable functions from $[0, T]$ to \mathbb{R}^d . This probabilistic interpretation of numerical computations of tractable from intractable quantities as statistical inference of latent from observable quantities applies to all numerical problems and has been repeatedly recommended in the past [22, 7, 29, 20, 24]. It employs the language of probability theory to account for the epistemic uncertainty (i.e. limited knowledge) about the accuracy of intermediate and final numerical computations, thereby yielding algorithms which can be more aware of—as well as more robust against—uncertainty over intermediate computational results. Such algorithms can output probability measures, instead of point estimates, over the final quantity of interest. This approach, now called *probabilistic*

*Max-Planck-Institute for Intelligent Systems, Tübingen, Germany (hkersting@tue.mpg.de, phennig@tue.mpg.de).

†Freie Universität Berlin, Germany (t.j.sullivan@fu-berlin.de), and Zuse Institute Berlin, Germany (sullivan@zib.de).

‡Universität of Tübingen, Germany (philipp.hennig@uni-tuebingen.de).

numerics (PN) [9], has in recent years been spelled out for a wide range of numerical tasks, including linear algebra, optimization, integration and differential equations, thereby working towards the long-term goal of a coherent framework to propagate uncertainty through chained computations, as desirable, e.g., in statistical machine learning.

In this paper, we determine the convergence rates of a recent family of PN methods [27, 13, 17, 28] which recast an IVP as a *stochastic filtering problem* [21, Chapter 6], an approach that has been studied in other settings [11], but has not been applied to IVPs before. These methods assume *a priori* that the solution x and its first q derivatives follow a Gauss–Markov process \mathbf{X} that solves a stochastic differential equation (SDE). The evaluations of f at numerical estimates of the true solution can then be regarded as imperfect evaluations of \dot{x} , which can then be used for a Bayesian update of \mathbf{X} . Such recursive updates along the time axis yield an algorithm whose structure resembles that of Gaussian (Kalman) filtering [26, Chapter 4]. These methods add only slight computational overhead compared to classical methods [28] and have been shown to inherit local convergence rates from equivalent classical methods in specific cases [27, 28]. These equivalences (i.e. the equality of the filtering posterior mean and the classical method) are only known to hold in the case of the integrated Brownian motion (IBM) prior and noiseless evaluations of f (in terms of our later notation, the case $R \equiv 0$), as well as under the following restrictions:

Firstly, for $q \in \{1, 2, 3\}$, and if the first step is divided into sub-steps resembling those of Runge–Kutta methods, an equivalence of the posterior mean of the first step of the filter and the explicit Runge–Kutta method of order q was established in [27] (but for $q \in \{2, 3\}$ only in the limit of the initial time of the IBM to $-\infty$). Secondly, it was shown in [28] that, for $q = 1$, the posterior means after each step coincides with the trapezoidal rule, if it takes an additional evaluation of f at the end of each step, known as P(EC)¹. In the same paper, it was shown that for $q = 2$ the filter coincides with a third-order Nordsieck method [18], if the filter is in the steady state, i.e. after the sequence of error covariance matrices has converged. These results neither cover filters with the integrated Ornstein–Uhlenbeck process (IOUP) prior [17] nor non-zero noise models on evaluations of f .

In this paper, we directly prove convergence rates—without first fitting the filter to existing methods—and thereby lift many of the above restrictions on the convergence rates.

1.1. Contribution. Our main results—Theorems 5.2 and 6.7—provide local and global convergence rates of the ODE filter. Theorem 5.2 shows local convergence rates of h^{q+1} without the above-mentioned previous restrictions—i.e. for a generic Gaussian ODE filter for all $q \in \mathbb{N}$, both IBM and IOUP prior, flexible Gaussian initialization (see Assumptions 2 and 3), and arbitrary evaluation noise $R \geq 0$. As a first global convergence result, Theorem 6.7 establishes global convergence rates of h^q in the case of $q = 1$, the IBM prior and all fixed measurement uncertainty models $R \equiv Kh^p$ of order $p \geq 1$ (see Assumption 4). Interestingly, this global rate of the worst-case error is matched by the contraction rate of the posterior confidence intervals precisely for all models with $p \geq 1$, as we show in Proposition 7.1. In the course of this analysis, we also give closed-form expressions for the steady states in the global case, and show that all q modeled derivatives have optimal convergence order given the error on the non-differentiated solution, in the sense that for every additional derivative the rate is reduced by 1—yielding local convergence of h^{q+1-i} for the i^{th} derivative (see Remark 5.3),

and global bound (independent of h) on the error on the derivative for the global convergence (see Remark 6.8).

1.2. Related work on probabilistic ODE solvers. The Gaussian ODE filter can be thought of as a self-consistent Bayesian decision agent who iteratively updates its prior belief \mathbf{X} over $x: [0, T] \rightarrow \mathbb{R}^d$ (and its first q derivatives) with information on \dot{x} from evaluating f .¹ For Gauss–Markov priors, it performs exact Bayesian inference and optimally (with respect to the L^2 -loss) extrapolates along the time axis. Accordingly, all of its computations are deterministic and—due to its restriction to Gaussian distributions—only slightly more expensive than classical solvers. Experiments demonstrating competitive performance are provided in [28, Section 5].

Another line of work (comprising the methods from [3, 6, 31, 15, 1]) introduces probability measures to ODE solvers in a fundamentally different way—by representing the distribution of all numerically possible trajectories with a set of sample paths. To compute these sample paths, [3] draws them from a (Bayesian) Gaussian process (GP) regression; [6, 31, 15] perturb classical estimates after an integration step with a suitably scaled Gaussian noise; and [1] perturbs the classical estimate instead by choosing a stochastic step-size. While [6, 31, 15, 1] can be thought of as (non-Bayesian) ‘stochastic wrappers’ around classical solvers, which produce samples with the same convergence rate, [3] employs—like the filter—GP regression to represent the belief on x . However, [3] also aims for a sample representation of numerical errors and thereby allows a flexible class of (possibly non-Markov) GP priors, from which it iteratively draws samples. A conceptual and experimental comparison between [3] and the filter can be found in [28].

All of the above sampling-based methods can hence represent more expressive, non-Gaussian posteriors (as e.g. desirable for bifurcations), but multiply the computational cost of the underlying method by the number of samples.² The ODE filter is, in contrast, not a perturbation of known methods, but a novel method designed for computational speed and for a robust treatment of intermediate uncertain values (such as the evaluations of f at estimated points). Accordingly, our convergence results concern the convergence rate of the posterior mean to the true solution—while the theoretical results in [3, 6, 31, 15, 1] provide convergence rates of the variance of the non-Gaussian empirical measure of samples (and not for an individual sample).

1.3. Outline. The paper begins with a brief introduction to Gaussian ODE filtering in Section 2. Next, Sections 3 and 4 provide auxiliary bounds on the flow map of the ODE and on intermediate quantities of the filter respectively. With the help of these bounds, sections 5 and 6 establish local and global convergence rates of the filtering mean respectively. Finally, in light of these rates, Section 7 analyses for which measurement noise models the posterior

¹Here, the word ‘Bayesian’ describes the algorithm in the sense that it employs a prior over the quantity of interest and updates it by Bayes rule according to a prespecified measurement model (as also used in [29, 3, 13]). The ODE filter is not Bayesian in the stronger sense of [5], and it remains an open problem to construct a Bayesian solver in this strong sense without restrictive assumptions, as discussed in [32].

²According to [1, Theorem 5.1.] the convergence speed of these Monte Carlo estimators in the mean squared sense is independent of the number of samples, which generated hope that, in many settings, a small number of samples could suffice.

confidence intervals are well-calibrated—followed by a concluding discussion in [Section 8](#).

1.3.1. Proof of main results. The main results are twofold: local and global convergence of the filtering mean in [Sections 5](#) and [6](#) respectively. The basis of the proof consists (besides some tailor-made regularity results on the flow map in [Section 3](#)) of the orders of the prediction errors and measurement residuals (with explicit dependence on the state misalignments; see [\(4.3\)](#)) in [Lemmas 4.1](#) and [4.2](#). Using these orders, local convergence is obtained in [Theorem 5.2](#), by locally bounding the Kalman gains and the state misalignments in [Lemmas 4.3](#) and [5.1](#) respectively. Next, in [Section 6](#), the strategy to prove global convergence rates is to apply a special discrete Grönwall inequality [\[4\]](#). To this end, we provide global bounds over the Kalman gains (using their attractive steady states) in [Proposition 6.2](#) and over the state misalignments in [Lemma 6.4](#), which—via a summarizing bound in [Lemma 6.5](#)—yields the prerequisite for the intended application of this Grönwall inequality in [Lemma 6.6](#). Grönwall’s inequality then implies our global convergence result in [Theorem 6.7](#). In the global convergence part ([Section 6](#)), we track, by [Assumption 4](#), the order p of the measurement noise $R = Kh^p$ through the proof—which reveals that global convergence holds if and only if $p \geq 1$.

1.4. Notation. We will use the notation $[n] = \{0, \dots, n-1\}$. For vectors matrix and vectors, we will use zero-based numbering, e.g. $x = (x_0, \dots, x_{d-1}) \in \mathbb{R}^d$. For a matrix $P \in \mathbb{R}^{n \times m}$ and $(i, j) \in [n] \times [m]$, we will write $P_{i,:} \in \mathbb{R}^{1 \times m}$ for the i^{th} row and $P_{:,j}$ for the j^{th} column of P . A fixed but arbitrary norm on \mathbb{R}^d will be denoted by $\|\cdot\|$.

2. Gaussian ODE filtering. Let $q, d \in \mathbb{N}$. In this section, we define how a Gaussian filter can solve the IVP [\(1.1\)](#). In the following respective subsections, we first explain the choice of prior on x , then describe how the algorithm computes a posterior output from this prior (by defining a numerical integrator Ψ), and add explanations on the measurement noise of the derivative observations.

2.1. Prior on x . In PN, it is common [\[9, Section 3\(a\)\]](#) to put a prior measure on the unknown solution; often, for fast Bayesian inference by linear algebra [\[23, Chapter 2\]](#), this prior is Gaussian. To this end, for $j \in [d]$ and $i \in [q+1]$, we model the unknown i^{th} derivative of the j^{th} dimension of the solution x of the IVP [\(1.1\)](#), denoted by $x_j^{(i)}$, as a draw from a real-valued, one-dimensional GP $X_j^{(i)}$. To enable GP inference in linear time by Kalman filtering [\[8\]](#), we restrict the prior to Markov processes. As shown in [\[8\]](#), a wide class of such Gauss–Markov processes can be equivalently captured by a law of the (strong) solution [\[21, Chapter 5.3.\]](#) of a linear SDE with Gaussian initial condition. Here, for all $j \in [d]$, we define the vector of time derivatives by $\mathbf{X}_j = (X_j^{(0)}, \dots, X_j^{(q)})^\top$, which we define as a q -times integrated stochastic process driven by an underlying Brownian motion with scaling $\sigma_j > 0$,

$$(2.1) \quad d\mathbf{X}_{j,t} = \begin{pmatrix} dX_{j,t}^{(0)} \\ \vdots \\ dX_{j,t}^{(q-1)} \\ dX_{j,t}^{(q)} \end{pmatrix} = \begin{pmatrix} 0 & 1 & 0 \dots & 0 \\ \vdots & \ddots & \ddots & 0 \\ \vdots & & 0 & 1 \\ a_0 & \dots & \dots & a_q \end{pmatrix} \begin{pmatrix} X_{j,t}^{(0)} \\ \vdots \\ X_{j,t}^{(q-1)} \\ X_{j,t}^{(q)} \end{pmatrix} dt + \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \sigma_j \end{pmatrix} dB_{j,t},$$

with initial condition $X_{j,0} \sim \mathcal{N}(m_j(0), P_j(0))$ and independent underlying one-dimensional Brownian motions $\{B_j; j \in [d]\}$ (independent of \mathbf{X}_0). We assume that $\{X_{j,0}; j \in [d]\}$ are independent. Hence, by the independence of the components of the d -dimensional Brownian motion, the components $\{\{\mathbf{X}_{j,t}; 0 \leq t \leq T\}; j \in [d]\}$ of $\{\mathbf{X}_t; 0 \leq t \leq T\}$ are independent.³ The solution of (2.1) is a Gauss–Markov process with mean $m_j: [0, T] \rightarrow \mathbb{R}^{q+1}$ and covariance matrix $P_j: [0, T] \rightarrow \mathbb{R}^{(q+1) \times (q+1)}$ given by

$$(2.2) \quad m_j(h) = A(h)m_j(0), \quad P_j(h) = A(h)P_j(0)A(h)^\top + Q(h, \sigma_j),$$

where the matrices $A(h), Q(h, \sigma_j) \in \mathbb{R}^{(q+1) \times (q+1)}$ yielded by the SDE (2.1) are known in closed form [25, Theorem 2.9.]; we give the explicit expressions in two key cases in (2.3)–(2.5) below. The precise choice of the prior stochastic process \mathbf{X} depends on the choice of $a = (a_0, \dots, a_q) \in \mathbb{R}^{q+1}$ in (2.1). For $a = (0, \dots, 0)$, \mathbf{X}_j is a q -times integrated Brownian motion (IBM) and, for $a = (0, \dots, 0, -\theta)$, $\theta > 0$, a q -times integrated Ornstein–Uhlenbeck process (IOUP). While the algorithm works for every choice of $a \in \mathbb{R}^{q+1}$, the BM and the OUP are two of the most widely used Gauss–Markov processes for modeling stochastic dynamics, which (as examined in [17]) incorporate different prior assumptions with respect to the stationarity of \dot{x} . Hence, we restrict our attention to these two cases. By (2.2), the prior \mathbf{X} enters the algorithm via A and Q , which can be calculated in closed form. For IBM, we obtain

$$(2.3) \quad A^{\text{IBM}}(h)_{i,j} = \mathbb{1}_{i \leq j} \frac{h^{j-i}}{(j-i)!}, \quad Q^{\text{IBM}}(\sigma, h)_{ij} = \frac{\sigma^2 h^{2q+1-i-j}}{(2q+1-i-j)(q-i)!(q-j)!},$$

and for IOUP

$$(2.4) \quad A^{\text{IOUP}}(h)_{ij} = \begin{cases} \mathbb{1}_{j \geq i} \frac{h^{j-i}}{(j-i)!}, & \text{if } j \neq q, \\ \frac{1}{(-\theta)^{q-i}} \sum_{k=q-i}^{\infty} \frac{(-\theta h)^k}{k!}, & \text{if } j = q, \end{cases}$$

$$(2.5) \quad Q^{\text{IOUP}}(h, \sigma)_{ij} = Q^{\text{IBM}}(h, \sigma)_{ij} + \Theta(h^{2q+2-i-j}).$$

(Derivations of (2.4) and (2.5), as well as the precise form of Q^{IOUP} , are presented in the appendix.) Hence, for all $i \in [q+1]$, the difference between the IBM and IOUP prediction of size h of the i^{th} dimension from any state $u \in \mathbb{R}^{q+1}$ is of order h^{q+1-i} , in the sense that

$$(2.6) \quad |(A^{\text{IOUP}}(h)u)_i - (A^{\text{IBM}}(h)u)_i| \leq K|u_q|h^{q+1-i}.$$

Here, and in the sequel, $K > 0$ denotes a constant independent of $h > 0$ which may change from line to line. If a statement holds for both IBM and IOUP, we will omit the dependence of A and Q on them in the following.

2.2. The algorithm. To avoid the introduction of additional indices, we will define the algorithm Ψ for $d = 1$; for statements on the general case of $d \in \mathbb{N}$ we will use the same symbols from (2.8)–(2.12) as vectors over the whole dimension—see e.g. (4.6) for a statement about a general $r \in \mathbb{R}^d$. By the assumed independence of the dimensions of \mathbf{X} , due to (2.1),

³More involved correlation models of $\{\{\mathbf{X}_{j,t}; 0 \leq t \leq T\}; j \in [d]\}$ are straight-forward to incorporate into the SDE (2.1), but seem complicated to analyse. Therefore, we restrict our attention to independent dimensions.

extension to $d \in \mathbb{N}$ amounts to applying Ψ to every dimension independently. Accordingly, we may in many of the below proofs w.l.o.g. assume $d = 1$. Now, as previously spelled out in [13, 28], Bayesian filtering of \mathbf{X} —i.e. iteratively conditioning \mathbf{X} on the information on $X^{(1)}$ from evaluations of f at the mean of the current conditioned $X^{(0)}$ —yields the following numerical method Ψ . Let $\mathbf{m}(t) = (m^{(0)}(t), \dots, m^{(q)}(t))^\top \in \mathbb{R}^{(q+1)}$ be an arbitrary state at some point in time $t \in [0, T]$ (i.e. $m^{(i)}(t)$ is an estimate for $x^{(i)}(t)$), and let $P(t) \in \mathbb{R}^{(q+1) \times (q+1)}$ be its covariance matrix. For $t \in [0, T]$, let the current estimate of $\mathbf{x}(t)$ be a normal distribution $\mathcal{N}(\mathbf{m}(t), P(t))$, i.e. the mean $\mathbf{m}(t) \in \mathbb{R}^{(q+1)}$ represents the best numerical estimate (given data $\{y(h), \dots, y(t)\}$, see (2.10)) and the covariance matrix $P(t) \in \mathbb{R}^{(q+1) \times (q+1)}$ its uncertainty. Then, the ODE filter for the time step $t \rightarrow t + h$ of size $h > 0$ (we assume w.l.o.g. that $T/h \in \mathbb{N}$) is given by

$$(2.7) \quad \Psi_{P(t),h}(\mathbf{m}(t)) := \left(\Psi_{P(t),h}^{(0)}, \dots, \Psi_{P(t),h}^{(q)} \right)^\top (\mathbf{m}(t)) := \mathbf{m}^-(t+h) + \beta(t+h)r(t+h),$$

with predicted mean $\mathbf{m}^-(t+h) \in \mathbb{R}^{q+1}$ given by $\mathbf{m}^-(t+h) := A(h)\mathbf{m}(t)$ (recall (2.2)) and $r(t+h)$ and $\beta(t+h) = (\beta^{(0)}(t+h), \dots, \beta^{(q+1)}(t+h))^\top$ are calculated by

$$(2.8) \quad P^-(t+h) := A(h)P(t)A(h)^\top + Q(h, \sigma) \in \mathbb{R}^{(q+1) \times (q+1)}, \quad (\text{pred. covariance}),$$

$$(2.9) \quad \beta^{(i)}(t+h) := \frac{P^-(t+h)_{i1}}{(P^-(t+h))_{11} + R(t+h)} \in \mathbb{R}, \quad (\text{Kalman gain on } i^{\text{th}} \text{ state}),$$

$$(2.10) \quad y(t+h) := f\left(m^{-(0)}(t+h)\right) \in \mathbb{R}, \quad (\text{measurement/data on } \dot{x}),$$

$$(2.11) \quad r(t+h) := y(t+h) - m^{-(1)}(t+h) \in \mathbb{R}, \quad (\text{innovation/residual}).$$

Here, R denotes the variance of y (the ‘measurement noise’) and captures the squared difference between the data $y(t+h) = f(m^-(t+h))$ that the algorithm actually receives and the idealised data $\dot{x}(t+h) = f(x(t+h))$ that it ‘should’ receive (see Subsection 2.3). Finally, the mean and the covariance matrix are conditioned on the data: $\mathbf{m}(t+h) = \Psi_{P(t),h}(\mathbf{m}(t))$ (updated mean) by (2.7) and

$$(2.12) \quad P(t+h) = P^-(t+h) - \frac{P^-(t+h)_{:,1}P^-(t+h)_{1,:}^\top}{P^-(t+h)_{11} + R(t+h)}, \quad (\text{updated covariance}).$$

The algorithm is iterated by computing $\mathbf{m}(t+2h) := \Psi_{P(t+h),h}(\mathbf{m}(t+h))$ as well as repeating (2.8) and (2.12), with $P(t+h)$ instead of $P(t)$, to obtain $P(t+2h)$. In the following, to avoid notational clutter, the dependence of the above quantities on t , h and σ will be omitted if their values are unambiguous. Parameter adaptation reminiscent of classical methods (e.g. for σ s.t. the added variance per step coincide with standard error estimates) have been explored in [28, Section 4].

Since this filter is essentially an iterative application of Bayes rule (see e.g. [26, Chapter 4]) based on the prior \mathbf{X} on \mathbf{x} specified by (2.1) (entering the algorithm via A and Q) and the measurement model $y \sim \mathcal{N}(\dot{x}, R)$, it remains to detail the latter (recall subsection 2.1 for the choice of prior). Concerning the data generation mechanism for y (2.10), we only consider the maximum-a-posteriori point estimate of $\dot{x}(t)$ given $\mathcal{N}(m^{-(0)}(t), P_{00}^-(t))$; a discussion of

more inclined statistical models for y can be found in [28, Subsection 2.2.]. Next, for lack of such a discussion for R , we will examine different choices of R —which have proved central to the uncertainty quantification of the filter [13] and will turn out to affect global convergence properties in Section 6.

2.3. Measurement noise R . Two sources of uncertainty add to $R(t)$: noise from imprecise knowledge of $x(t)$ and f . Given f , previous integration steps of the filter (as well as an imprecise initial value) inject uncertainty about how close $m^-(t)$ is to $x(t)$ and how close $y = f(m^-(t))$ is to $\dot{x}(t) = f(x(t))$. This uncertainty stems from the discretization error $\|m^{-(0)}(t) - x(t)\|$ and, hence, strictly increases with h . Additionally, there can be uncertainty from a misspecified f , e.g. when f has estimated parameters, or from numerically imprecise evaluations of f , which can be added to R —a functionality which classical solvers do not possess. In this paper, since R naturally depends on h , we analyse the influence of noise of order $R(t) \equiv Kh^p$ (see Assumption 4), $p \geq 0$, on the quality of the solution to illuminate for which orders of noise we can trust the solution to which extent and when we should, instead of decreasing h , rather spend computational budget on specifying or evaluating f more precisely. The explicit dependence of the noise on its order p in h resembles, despite the fundamentally different role of R compared to additive noise in [6, 1], the variable p in [6, Assumption 1] and [1, Assumption 2.2.] in the sense that the analysis highlights how uncertainty of this order can still be modeled without breaking the convergence rates. (Adaptive noise models are computationally feasible [13], but not within the scope of our analysis.)

3. Regularity of flow. Let $q, d \in \mathbb{N}$. Before we proceed to the analysis of Ψ , we provide all below-needed regularity results in this section.

Assumption 1. *The vector field $f \in C^q(\mathbb{R}^d; \mathbb{R}^d)$ and all its derivatives of order up to q are uniformly bounded and globally Lipschitz, i.e. there exists some $L > 0$ such that for all multi-indices $\alpha \in \mathbb{N}_0^d$ with $\sum_i \alpha_i \leq q$,*

$$(3.1) \quad \|D^\alpha f\|_\infty \leq L, \quad \|D^\alpha f(a) - D^\alpha f(b)\| \leq L\|a - b\|.$$

Under Assumption 1, the solution x is—by the Picard-Lindelöf theorem—well-defined and—by differentiating (1.1) q times using the chain rule—in $C^{q+1}([0, T]; \mathbb{R}^d)$. For $i \in [q + 1]$, we denote $\frac{d^i x}{dt^i}$ by $x^{(i)}$. By a bold symbol, we denote the vector of these derivatives: $\mathbf{x} \equiv (x^{(0)}, \dots, x^{(q)})^\top$. In particular, the solution x of (1.1) is denoted by $x^{(0)}$ from now on. Note that, whenever a vector spans multiple derivatives, we will write a bold symbol—such as \mathbf{x} . In analogous notation, for every $i \in [q + 1]$, we denote the flow of the ODE (1.1) by $\Phi^{(0)}$, i.e. $\Phi_t^{(0)}(x_0) \equiv x^{(0)}(t)$, and, for $i \in [q + 1]$, its i^{th} partial derivative with respect to t by $\Phi^{(i)}$, so that $\Phi_t^{(i)}(x_0) \equiv x^{(i)}(t)$. Again, these derivatives are summarized with a bold symbol: $\Phi := (\Phi_t^{(0)}, \dots, \Phi_t^{(q)})^\top$.

Lemma 3.1. *Under Assumption 1, for all $a \in \mathbb{R}^d$ and all $h > 0$,*

$$(3.2) \quad \left\| \Phi_h^{(i)}(a) - \sum_{k=i}^q \frac{h^{k-i}}{(k-i)!} \Phi_0^{(k)}(a) \right\| \leq Kh^{q+1-i}.$$

265 *Proof.* By [Assumption 1](#), $\Phi^{(q+1)}$ exists and is bounded by $\|\Phi^{(q+1)}\| \leq L$, which can be
 266 seen by applying the chain rule q times to both sides of [\(1.1\)](#). Now, applying $\|\Phi^{(q+1)}\| \leq L$
 267 to the term $\Phi_\tau^{(q+1)}(a)$ (for some $\tau \in (0, h)$) in the Lagrange remainder of the $(q - i)^{\text{th}}$ -order
 268 Taylor expansion of $\Phi_h^{(i)}(a)$ yields [\(3.2\)](#). ■

269 **Lemma 3.2.** Under [Assumption 1](#) and for sufficiently small $h > 0$,

$$270 \quad (3.3) \quad \sup_{a \neq b \in \mathbb{R}^d} \frac{\|\Phi_h^{(0)}(a) - \Phi_h^{(0)}(b)\|}{\|a - b\|} \leq 1 + 2Lh.$$

272 *Proof.* Immediate corollary of [\[30, Theorem 2.8.\]](#). ■

273 For global convergence in [Section 6](#), we will need the following modification of the previous
 274 lemma.

275 **Lemma 3.3.** Let $q = 1$. Then, under [Assumption 1](#) and for sufficiently small $h > 0$,

$$276 \quad (3.4) \quad \sup_{a \neq b \in \mathbb{R}^d} \frac{\|\Phi_h(a) - \Phi_h(b)\|_h}{\|a - b\|} \leq 1 + Kh,$$

278 where, given the arbitrary norm $\|\cdot\|$ on \mathbb{R}^d and $h > 0$, the new norm $\|\cdot\|_h$ on $\mathbb{R}^{(q+1) \times d}$ is
 279 defined by

$$280 \quad (3.5) \quad \|a\|_h := \sum_{i=0}^q h^i \|a_{i,:}\|.$$

282 **Remark 3.4.** The necessity of $\|\cdot\|_h$ stems from the fact that—unlike other ODE solvers—
 283 the ODE filter Ψ additionally estimates and uses the first q derivatives in its state $\mathbf{m} \in$
 284 $\mathbb{R}^{(q+1) \times d}$, whose development cannot be bounded in $\|\cdot\|$, but in $\|\cdot\|_h$. The norm $\|\cdot\|_h$ is used
 285 to make rigorous the intuition that the estimates of the solution's time derivative are ‘one
 286 order of h worse per derivative’.

287 *Proof.* We bound the second summand of

$$288 \quad (3.6) \quad \|\Phi_h(a) - \Phi_h(b)\|_h \stackrel{(3.5)}{=} \underbrace{\|\Phi_h^{(0)}(a) - \Phi_h^{(0)}(b)\|}_{\leq (1+2Lh)\|a-b\|, \text{ by (3.3)}} + h \left\| \underbrace{\Phi_h^{(1)}(a)}_{=f(\Phi_h^{(0)}(a))} - \underbrace{\Phi_h^{(1)}(b)}_{=f(\Phi_h^{(0)}(b))} \right\|$$

290 by

$$291 \quad (3.7) \quad \left\| f(\Phi_h^{(0)}(a)) - f(\Phi_h^{(0)}(b)) \right\| \stackrel{(3.1)}{\leq} L \|\Phi_h^{(0)}(a) - \Phi_h^{(0)}(b)\| \stackrel{(3.3)}{\leq} L(1 + 2Lh)\|a - b\|.$$

293 Now, insertion of [\(3.7\)](#) into [\(3.6\)](#) concludes the proof. ■

4. Auxiliary bounds on intermediate quantities. The ODE filter Ψ iteratively computes the filtering mean $\mathbf{m}(nh) = (m^{(0)}(nh), \dots, m^{(q)}(nh))^T \in \mathbb{R}^{(q+1)}$ as well as error covariance matrices $P(nh) \in \mathbb{R}$ on the mesh $\{nh\}_{n=0}^{T/h}$. Ideally, the truncation error over all derivatives, $\varepsilon(nh) := (\varepsilon^{(0)}(nh), \dots, \varepsilon^{(q)}(nh))^T := \mathbf{m}(nh) - \mathbf{x}(nh)$, falls fast as $h \rightarrow 0$ and is captured by the standard deviation $\sqrt{P_{00}(nh)}$. Next, we present a classical worst-case convergence analysis over all f satisfying [Assumption 1](#) in this paper (see [Section 8](#) for a discussion of the desirability and feasibility of an average-case analysis). To this end, we bound the added error of every step by its $\Delta^{(i)}$ intermediate values, defined in [\(2.9\)](#) and [\(2.11\)](#),

$$(4.1) \quad \Delta^{(i)}((n+1)h) := \left\| \Psi_{P(nh),h}^{(i)}(\mathbf{m}(nh)) - \Phi_h^{(i)}(m^{(0)}(nh)) \right\|$$

$$(4.2) \quad \stackrel{(2.7)}{\leq} \underbrace{\left\| (A(h)\mathbf{m}(nh))_i - \Phi_h^{(i)}(m^{(0)}(nh)) \right\|}_{=: \Delta^{-(i)}((n+1)h)} + \left\| \beta^{(i)}(nh) \right\| \|r(nh)\|,$$

and bound these quantities in the order $\Delta^{-(i)}$, r , $\beta^{(i)}$. These bounds will be needed for the local and global convergence analysis in [Sections 5](#) and [6](#). Inconveniently for the analysis⁴, the entries of \mathbf{m} will in general not be aligned in the sense that, $m^{(i)}$ should coincide with the i^{th} derivative in time of the flow Φ starting at $m^{(0)}$; i.e., in general, there is a *state misalignment*

$$(4.3) \quad \delta^{(i)}(nh) := \left\| m^{(i)}(nh) - f^{(i-1)}(m^{(0)}(nh)) \right\| = \left\| m^{(i)}(nh) - \Phi_0^{(i)}(m^{(0)}(nh)) \right\| > 0,$$

for $i \in [q+1]$, where $f^{(i)}: \mathbb{R}^d \rightarrow \mathbb{R}^d$ is recursively defined by $f^{(0)}(a) = a$, $f^{(1)}(a) = f(a)$ and $f^{(i)}(a) = (\nabla_x f^{(i-1)} \cdot f)(a)$.⁵ While $\delta^{(0)} \equiv 0$, some of the following bounds will therefore depend on $\{\delta^{(i)}; i = 1, \dots, q\}$.

Lemma 4.1. *Under [Assumption 1](#), for all $i \in [q+1]$ and all $h > 0$,*

$$(4.4) \quad \Delta^{-(i)}((n+1)h) \leq K \left[1 + \int \left\| m^{(q)}(nh) \right\| \right] h^{q+1-i} + \sum_{k=i}^q \frac{h^{k-i}}{(k-i)!} \delta^{(k)}(nh),$$

where the notation $\int x$ stands for $\int x := \begin{cases} x, & \text{if } \mathbf{X} = \text{IOUP}, \\ 0, & \text{if } \mathbf{X} = \text{IBM}. \end{cases}$

In the following, all statements will be formulated using the above notation; i.e. statements without this notation will hold both for IOUP and IBM. From [Section 6](#) on, the (global) analysis (and all statements) will be restricted to IBM.

Proof. We may assume, as explained in [Subsection 2.2](#), w.l.o.g. that $d = 1$. We apply the triangle inequality to the definition of $\Delta^{-(i)}((n+1)h)$, as defined in [\(4.2\)](#), which, by [\(2.6\)](#),

⁴This theoretical inconvenience is reasonable in practice, since the possibility of $\delta > 0$ facilitates averaging out uncertainties on the estimates $\{m^{(i)}; i = 1, \dots, q\}$ of the derivatives over multiple steps instead of forcing them to obey the imprecise state estimate $m^{(0)}$.

⁵Note that $f^{(i)}(x^{(0)}(t)) = x^{(i+1)}(t)$ [[10](#), Section 2]. Hence, if $\varepsilon(nh) = 0$, then $\delta^{(i)}(nh) = 0$, $\forall i \in [q+1]$.

323 yields

$$\begin{aligned}
 324 \quad \Delta^{-(i)}((n+1)h) &\leq \sum_{k=i}^q \frac{h^{k-i}}{(k-i)!} \delta^{(k)}(nh) + \int K |m^{(q)}(nh)| h^{q+1-i} \Big\{ \\
 325 \quad (4.5) \quad &+ \underbrace{\left[\sum_{l=i}^q \frac{h^{l-i}}{(l-i)!} \Phi_0^{(l)}(m^{(0)}(nh)) - \Phi_h^{(i)}(m^{(0)}(nh)) \right]}_{\leq Kh^{q+1-i}, \text{ by (3.2)}} \Big\}. \quad \blacksquare
 \end{aligned}$$

326
327 **Lemma 4.2.** Under [Assumption 1](#) and for sufficiently small $h > 0$,

$$\begin{aligned}
 328 \quad (4.6) \quad \|r((n+1)h)\| &\leq K \left[1 + \int \|m^{(q)}(nh)\| \right] h^q + \sum_{k=1}^q \frac{h^{k-1}}{(k-1)!} \delta^{(k)}(nh) + K \sum_{k=1}^q \frac{h^k}{k!} \delta^{(k)}(nh). \\
 329
 \end{aligned}$$

330 *Proof.* Again, w.l.o.g. $d = 1$. Recall that, by (2.11), r is implied by the values of $m^{-,(0)}$
331 and $m^{-,(1)}$. By insertion of $m^{-,(i)}((n+1)h) = \sum_{k=i}^q \frac{h^{k-i}}{(k-i)!} m^{(k)}(nh) + \int K |m^{(q)}(nh)| h^{q+1-i}$
332 (due to (2.6) and (2.7)) into the definition (2.11) of $r((n+1)h)$, we obtain the following
333 equality which we then bound by repeated application of the triangle inequality:

$$\begin{aligned}
 334 \quad |r((n+1)h)| &= \left| f \left(\sum_{k=0}^q \frac{h^k}{k!} m^{(k)}(nh) + \int K |m^{(q)}(nh)| h^{q+1} \right) \right. \\
 335 \quad &\quad \left. - \left(\sum_{k=1}^q \frac{h^{k-1}}{(k-1)!} m^{(k)}(nh) + \int K |m^{(q)}(nh)| h^q \right) \right| \\
 336 \quad &\leq \left| f \left(\sum_{k=0}^q \frac{h^k}{k!} m^{(k)}(nh) + \int K |m^{(q)}(nh)| h^{q+1} \right) - \left(\sum_{k=1}^q \frac{h^{k-1}}{(k-1)!} m^{(k)}(nh) \right) \right| \\
 337 \quad &\quad + \int K |m^{(q)}(nh)| h^q \Big\} \\
 338 \quad (4.7) \quad &\stackrel{(4.3)}{\leq} I_1(h) + I_2(h) + I_3(h) + \sum_{k=1}^q \frac{h^{k-1}}{(k-1)!} \delta^{(k)}(nh) + \int K |m^{(q)}(nh)| h^q, \\
 339
 \end{aligned}$$

340 where I_1 , I_2 , and I_3 are defined and bounded as follows, using [Assumption 1](#) and [Lemma 3.1](#):

$$\begin{aligned}
 341 \quad I_1(h) &:= \left| f \left(\sum_{k=0}^q \frac{h^k}{k!} m^{(k)}(nh) + \int K |m^{(q)}(nh)| h^{q+1} \right) - f \left(\sum_{k=0}^q \frac{h^k}{k!} \Phi_0^{(k)}(m^{(0)}(nh)) \right) \right| \\
 342 \quad (4.8) \quad &\leq L \sum_{k=0}^q \frac{h^k}{k!} \delta^{(k)}(nh) + \int LK |m^{(q)}(nh)| h^{q+1}, \\
 343
 \end{aligned}$$

$$\begin{aligned}
 344 \quad I_2(h) &:= \left| f \left(\sum_{k=0}^q \frac{h^k}{k!} \Phi_0^{(k)}(m^{(0)}(nh)) \right) - f \left(\Phi_h^{(0)}(m^{(0)}(nh)) \right) \right| \\
 345 \quad (4.9) \quad &\leq L \left| \sum_{k=0}^q \frac{h^k}{k!} \Phi_0^{(k)}(m^{(0)}(nh)) - \Phi_h^{(0)}(m^{(0)}(nh)) \right| \stackrel{(3.2)}{\leq} Kh^{q+1}, \\
 346
 \end{aligned}$$

and

$$(4.10) \quad I_3(h) := \left| \Phi_h^{(1)}(m^{(0)}(nh)) - \sum_{k=1}^q \frac{h^{k-1}}{(k-1)!} \Phi_0^{(k)}(m^{(0)}(nh)) \right| \stackrel{(3.2)}{\leq} Kh^q.$$

Inserting (4.8), (4.9), and (4.10) into (4.7) (and recalling $\delta^{(0)} = 0$) yields the claim (4.6). ■

To bound the Kalman gains $\beta(nh)$, we first need to assume that the orders of the initial covariance matrices are sufficiently high (matching the latter required orders of the initialization error; see Assumption 3).

Assumption 2. *The entries of the initial covariance matrix $P(0)$ have order of $\|P(0)_{k,l}\| \leq K_0 h^{2q+1-k-l}$, for all $k, l \in [q+1]$. Here, $K_0 > 0$ is a constant independent of h .*

We make this assumption, as well as Assumption 3, explicit (instead of just making the stronger assumption of exact initializations with zero variance), because it highlights how (statistical or numerical) uncertainty on the initial value effects the accuracy of the output of the filter—a novel functionality of PN with the potential to facilitate a management of the computational budget across a computational chain with respect to the respective perturbations from different sources of uncertainty [9, Section 3(d)].

Lemma 4.3. *Under Assumption 2, for all $i \in [q+1]$ and for all $h > 0$, $\|\beta^{(i)}(h)\| \leq Kh^{1-i}$.*

Proof. Again, w.l.o.g. $d = 1$. Application of the orders of A and Q from (2.3)–(2.5), the triangle inequality and Assumption 2 to the definition of P^- in (2.8) yields

$$(4.11) \quad \begin{aligned} |P^-(h)_{k,l}| &\stackrel{(2.8)}{\leq} |[A(h)P(0)A(h)^\top]_{k,l}| + |Q(h, \sigma)_{k,l}| \\ &\stackrel{(2.3)-(2.5)}{\leq} K \left[\sum_{a=k}^q \sum_{b=l}^q |P(0)_{a,b}| h^{a+b-k-l} + h^{2q+1-k-l} \right] \stackrel{\text{Ass. 2}}{\leq} Kh^{2q+1-k-l}. \end{aligned}$$

Now, by recalling that P and Q are (positive semi-definite) covariance matrices, we deduce that $P^-(h)_{1,1} \geq Kh^{2q-1}$. Hence, insertion of these orders into the definition of $\beta^{(i)}$ (2.9) and recalling that $R \geq 0$ concludes the proof. ■

5. Local convergence rates. With the above bounds on intermediate algorithmic quantities (involving state misalignments $\delta^{(i)}$) in place, we only need an additional assumption to proceed—via a bound on $\delta^{(i)}(0)$ —to our first main result on local convergence orders of Ψ .

Assumption 3. *The initial errors on the initial estimate of the i^{th} derivative $m^{(i)}(0)$ satisfy $\|\varepsilon^{(i)}(0)\| = \|m^{(i)}(nh) - x^{(i)}(nh)\| \leq K_0 h^{q+1-i}$. (This assumption is, like Assumption 2, weaker than the standard assumption of exact initializations.)*

Lemma 5.1. *Under Assumptions 1 and 3, for all $i \in [q+1]$ and for all $h > 0$, $\delta^{(i)}(0) \leq Kh^{q+1-i}$.*

Proof. The claim follows, using Assumptions 1 and 3, from

$$(5.1) \quad \delta^{(i)}(0) \leq \underbrace{\|m^{(i)}(0) - x^{(i)}(0)\|}_{=\|\varepsilon^{(i)}(0)\| \leq K_0 h^{q+1-i}} + \underbrace{\left\| f^{(i-1)}(x^{(0)}(0)) - f^{(i-1)}(m^{(0)}(0)) \right\|}_{\leq L \|\varepsilon^{(0)}(0)\| \leq LK_0 h^{q+1}}.$$

Theorem 5.2 (Local Truncation Error). Under *Assumptions 1 to 3* and for sufficiently small $h > 0$,

$$(5.2) \quad \|\varepsilon^{(0)}(h)\| \leq \|\varepsilon(h)\|_h \leq K \left[1 + \int \|m^{(q)}(0)\| \right] h^{q+1}.$$

Proof. By the triangle inequality for $\|\cdot\|_h$ and subsequent application of [Lemma 3.3](#) and [Assumption 3](#) to the second summand of the resulting inequality, we obtain

$$(5.3) \quad \|\varepsilon(h)\|_h \leq \underbrace{\left\| \Psi_{P(0),h}(\mathbf{m}(0)) - \Phi_h(x^{(0)}(0)) \right\|_h}_{=\sum_{i=0}^q h^i \Delta^{(i)}(h), \text{ by (4.1)}} + \underbrace{\left\| \Phi_h(x^{(0)}(0)) - \Phi_h(m^{(0)}(0)) \right\|_h}_{\leq (1+Kh)\|\varepsilon^{(0)}(0)\| \leq Kh^{q+1}}.$$

The remaining bound on $\Delta^{(i)}(h)$, for all $i \in [q+1]$ and sufficiently small $h > 0$, is obtained by insertion of the bounds from [Lemmas 4.1 to 4.3](#) (in the case of $n = 0$), into [\(4.2\)](#):

$$(5.4) \quad \begin{aligned} \Delta^{(i)}(h) &\leq K \left[1 + \int \|m^{(q)}(0)\| \right] h^{q+1-i} + K \sum_{k=0}^q \frac{h^{k-i}}{(k-1)!} \delta^{(k)}(0) \\ &\stackrel{\text{Lemma 5.1}}{\leq} K \left[1 + \int \|m^{(q)}(0)\| \right] h^{q+1-i}. \end{aligned}$$

Insertion of [\(5.4\)](#) into [\(5.3\)](#) and $\|\varepsilon^{(0)}(h)\| \leq \|\varepsilon(h)\|_h$ (by [\(3.5\)](#)) concludes the proof. \blacksquare

Remark 5.3. [Theorem 5.2](#) establishes a bound on the local truncation error after one step—including the additional effect of imprecise initialization under [Assumption 3](#). Moreover, this theorem is stronger than most error bounds on classical solvers in the sense that it—by the definition [\(3.5\)](#) of $\|\cdot\|_h$ —implies bounds of order h^{q+1-i} on the error $\varepsilon^{(i)}(h)$ on the i^{th} derivative for all $i \in [q+1]$, i.e. in particular a classical bound of order h^{q+1} on the truncation error $\varepsilon^{(0)}(h)$ on the solution of [\(1.1\)](#). While for IBM these bounds do not depend on the initialization beyond [Assumption 3](#), the constants for IOUP additionally depend on $\|m^{(q)}(0)\|$ due to the effect of the ‘mean-reverting’ drift on the update step (see [\(2.6\)](#) for the deviation of the prediction of IOUP from IBM). A global analysis for IOUP would therefore require additional assumptions ensuring a global bound on $\|m^{(q)}\|$. Hence, the following first global analysis is restricted to IBM.

6. Global analysis. As argued above, we only consider the case of the IBM prior for the global analysis in this section. Moreover, as for $q \geq 2$ the proof of an analogue of [Proposition 6.2](#) would be very technically involved, we restrict our analysis to $q = 1$ in order to maintain readability of this first global analysis. (See [Section 8](#) for a discussion of these restrictions.) While, for local convergence, all noise models R yielded the same convergence rates in [Theorem 5.2](#), it is unclear how the order of R in h (as described in [Subsection 2.3](#)) affects global convergence rates: E.g., for the limiting case $R \equiv Kh^0$, the steady-state Kalman gains β^∞ would converge to zero (see [\(6.7\)](#) and [\(6.8\)](#) below) for $h \rightarrow 0$, and hence the evaluation of f would not be taken into account—yielding a filter Ψ which assumes that the evaluations of f are equally off, regardless of $h > 0$, and eventually just extrapolates along the prior without being global convergence of the posterior mean \mathbf{m} . For the opposite limiting case $R \equiv \lim_{p \rightarrow \infty} Kh^p \equiv 0$, it has already been shown in [\[28, Proposition 1 and Theorem 1\]](#)

that—in the steady state and for $q = 1, 2$ —the filter Ψ inherits global convergence rates from known multistep methods in Nordsieck form [18]. Therefore we assume a fixed noise model with arbitrary order p .

Assumption 4. *The fixed noise model is chosen to be $R \equiv Kh^p$, for some fixed $p \geq 0$.*

In the following, we analyse how small p can be in order for Ψ to exhibit fast global convergence (cf. the similar role of the order p of additive perturbations in [6, Assumption 1] and [1, Assumption 2.2.]). In light of Theorem 5.2, the highest possible global convergence rate is $\mathcal{O}(h)$ —which will indeed be obtained for all $p \geq 1$ in Theorem 6.7. Since every extrapolation step of Ψ from t to $t+h$ depends not only on the current state, but also on the covariance matrix $P(t)$ —which itself depends on all previous steps— Ψ is neither a single-step nor a multistep method. Therefore, the global analysis begins with an analysis of these covariance matrices, culminating in global bounds on the Kalman gains in Proposition 6.2. It proceeds—via a global bound on the state misalignments in Lemma 6.4—to Lemma 6.6 which, by Grönwall’s inequality, implies our main result Theorem 6.7. Contrary to [28], we do not restrict our theoretical analysis to the steady-state case, but provide our results under the weaker Assumptions 2 and 3 that were already sufficient for local convergence in Theorem 5.2—which is made possible by the bounds (6.12) and (6.13) in Proposition 6.2.

6.1. Global bounds on Kalman gains. Since we will analyse the sequence of covariances using contractions, we first introduce the following lemma.

Lemma 6.1. *Let (\mathcal{X}, d) be a non-empty complete metric space, $T_n: \mathcal{X} \rightarrow \mathcal{X}$, $n \in \mathbb{N}$, a sequence of L_n -Lipschitz continuous contractions with $\sup_n L_n \leq \bar{L} < 1$. Let u_n be the fixed point of T_n , as well-defined by Banach fixed-point theorem (BFT), and let $\lim_{n \rightarrow \infty} u_n = u^* \in \mathcal{X}$. Then, for all $x_0 \in \mathcal{X}$, the recursive sequence $x_n := T_n(x_{n-1})$ converges to u^* .*

Proof. Let $\tilde{u}_0 = u^*$ and $\tilde{u}_n = T_n(\tilde{u}_{n-1})$, for $n \in \mathbb{N}$. Then,

$$d(u^*, x_n) \leq \underbrace{d(u^*, u_n)}_{\rightarrow 0} + \underbrace{d(u_n, \tilde{u}_n)}_{=: a_n} + \underbrace{d(\tilde{u}_n, x_n)}_{= d((T_n \circ \dots \circ T_1)(u^*), (T_n \circ \dots \circ T_1)(x_0)) \leq \bar{L}^n d(u^*, x_0) \rightarrow 0}.$$

It remains to show that $\lim_{n \rightarrow \infty} a_n = 0$. The \bar{L} -Lipschitz continuity of T_n and the triangle inequality yield

$$(6.1) \quad a_n = d(T_n(u_n), T_n(\tilde{u}_{n-1})) \leq \bar{L}[d(u_n, u_{n-1}) + d(u_{n-1}, \tilde{u}_{n-1})] = \bar{L}a_{n-1} + b_{n-1},$$

where $b_n := \bar{L}d(u_{n+1}, u_n) \rightarrow 0$. Now, for all $m \in \mathbb{N}$, let $a_0^{(m)} := a_0$ and $a_n^{(m)} := \bar{L}a_{n-1}^{(m)} + b_m$. By BFT, $\lim_{n \rightarrow \infty} a_n^{(m)} = b_m/(1 - \bar{L})$. Since, for all $m \in \mathbb{N}$, $a_n \leq a_n^{(m)}$ for sufficiently large n , it follows that

$$(6.2) \quad 0 \leq \limsup_{n \rightarrow \infty} a_n \leq \lim_{n \rightarrow \infty} a_n^{(m)} = \frac{b_m}{1 - \bar{L}}, \quad \forall m \in \mathbb{N}.$$

Since the convergent sequence u_n is in particular a Cauchy sequence, $\lim_{m \rightarrow \infty} b_m = 0$ and, hence, $0 \leq \lim_{n \rightarrow \infty} a_n = \limsup_{n \rightarrow \infty} a_n \leq 0$. Hence, $\lim_{n \rightarrow \infty} a_n = 0$. ■

Proposition 6.2. Under *Assumption 2*, the unique (attractive) steady states for the following quantities are

$$(6.3) \quad P_{11}^{-,\infty} := \lim_{n \rightarrow \infty} P_{11}^-(nh) = \frac{1}{2} \left(\sigma^2 h + \sqrt{4\sigma^2 Rh + \sigma^4 h^2} \right),$$

$$(6.4) \quad P_{11}^{\infty} := \lim_{n \rightarrow \infty} P_{01}^-(nh) = \frac{\left(\sigma^2 h + \sqrt{4\sigma^2 Rh + \sigma^4 h^2} \right) R}{\sigma^2 h + \sqrt{4\sigma^2 Rh + \sigma^4 h^2} + 2R},$$

$$(6.5) \quad P_{01}^{-,\infty} := \lim_{n \rightarrow \infty} P_{01}^-(nh) = \frac{\sigma^4 h^2 + (2R + \sigma^2 h) \sqrt{4\sigma^2 Rh + \sigma^4 h^2} + 4R\sigma^2 h}{2(\sigma^2 h + \sqrt{4\sigma^2 Rh + \sigma^4 h^2})} h,$$

$$(6.6) \quad P_{01}^{\infty} := \lim_{n \rightarrow \infty} P_{01}(nh) = \frac{R \sqrt{4R\sigma^2 h + \sigma^4 h^2}}{\sigma^2 h + \sqrt{4\sigma^2 Rh + \sigma^4 h^2}} h,$$

$$(6.7) \quad \beta^{\infty,(0)} := \lim_{n \rightarrow \infty} \beta^{(0)}(nh) = \frac{\sqrt{4R\sigma^2 h + \sigma^4 h^2}}{\sigma^2 h + \sqrt{4\sigma^2 Rh + \sigma^4 h^2}} h, \quad \text{and}$$

$$(6.8) \quad \beta^{\infty,(1)} := \lim_{n \rightarrow \infty} \beta^{(1)}(nh) = \frac{\sigma^2 h + \sqrt{4\sigma^2 Rh + \sigma^4 h^2}}{\sigma^2 h + \sqrt{4\sigma^2 Rh + \sigma^4 h^2} + 2R}.$$

Furthermore, under *Assumptions 2* and *4* and for sufficiently small $h > 0$,

$$(6.9) \quad \max_{n \in [T/h+1]} P_{11}^-(nh) \leq Kh^{1 \wedge \frac{p+1}{2}},$$

$$(6.10) \quad \max_{n \in [T/h+1]} P_{11}(nh) \leq Kh^{p \vee \frac{p+1}{2}},$$

$$(6.11) \quad \max_{n \in [T/h+1]} \|P_{01}(nh)\| \leq Kh^{p+1},$$

$$(6.12) \quad \max_{n \in [T/h+1]} \|\beta^{(0)}(nh)\| \leq Kh, \quad \text{and}$$

$$(6.13) \quad \max_{n \in [T/h+1]} \|1 - \beta^{(1)}(nh)\| \leq Kh^{(p-1) \vee 0}.$$

All of these bounds are sharp—in the sense that they would not hold for any higher order in the exponent of h .

Remark 6.3. The recursions for $P(nh)$ and $P^-(nh)$ given by (2.8) and (2.12) follow a discrete algebraic Riccati equation (DARE)—a topic studied in many related settings [14]. While the asymptotic behavior (6.3) of the completely detectable state $X^{(1)}$ can also be obtained using classical filtering theory [2, Chapter 4.4.], the remaining statements of Proposition 6.2 also concern the undetectable state $X^{(0)}$ and are, to the best of our knowledge, not directly obtainable from existing theory on DAREs or filtering (which makes the following proof necessary). Note that, in the special case of no measurement noise ($R \equiv 0$), (6.7) and (6.8) yield the equivalence of the filter in the steady state with the $P(EC)^1$ implementation of the trapezoidal rule which was previously shown in [28, Proposition 1]. For future research, it would be interesting to examine whether insertion of positive choices of R into (6.7) and (6.8) can reproduce known methods as well.

Proof. Again, w.l.o.g. $d = 1$. We prove the claims in the following order: (6.3), (6.9), (6.4), (6.10), (6.5), (6.7), (6.8), (6.6), (6.13), (6.12), (6.11). The sharpness of these bounds

is shown, directly after they are proved. As a start, for (6.3), we show that $P_{11}^{-,\infty}$ is indeed the unique fixed point of the recursion for $\{P_{11}^{-}(nh)\}_n$ by checking that, if $P_{11}^{-}(nh) = \frac{1}{2}(\sigma^2 h + \sqrt{4\sigma^2 Rh + \sigma^4 h^2})$, then also $P_{11}^{-}((n+1)h) = \frac{1}{2}(\sigma^2 h + \sqrt{4\sigma^2 Rh + \sigma^4 h^2})$:

$$(6.14) \quad P_{11}^{-}(nh) \stackrel{(2.12)}{=} P_{11}^{-}(nh) \left(1 - \frac{P_{11}^{-}(nh)}{P_{11}^{-}(nh) + R}\right) = \frac{(\sigma^2 h + \sqrt{4\sigma^2 Rh + \sigma^4 h^2})R}{\sigma^2 h + \sqrt{4\sigma^2 Rh + \sigma^4 h^2} + 2R}, \quad \text{and}$$

$$(6.15) \quad P_{11}^{-}((n+1)h) = P_{11}^{-}(nh) + \sigma^2 h \stackrel{(6.14)}{=} \frac{1}{2}(\sigma^2 h + \sqrt{4\sigma^2 Rh + \sigma^4 h^2}) = P_{11}^{-}(nh).$$

After combining (6.14) and (6.15), the recursion for P_{11}^{-} is given by

$$(6.16) \quad P_{11}^{-}((n+1)h) = \underbrace{\left(\frac{R}{P_{11}^{-}(nh) + R}\right)}_{=: \alpha(nh)} P_{11}^{-}(nh) + \sigma^2 h =: \tilde{T}(P_{11}^{-}(nh)).$$

Since R and $P_{11}^{-}(nh)$ are positive variances, we know that $\inf_{n \in [T/h+1]} P_{11}^{-}(nh) \geq \sigma^2 h$, and hence $\max_{n \in [T/h+1]} \alpha(nh) \leq R/(\sigma^2 h + R) < 1$. Hence, \tilde{T} is a contraction. By BFT, $P_{11}^{-,\infty}$ is the unique (attractive) fixed point of \tilde{T} , and the sequence $\{|P_{11}^{-}(nh) - P_{11}^{-,\infty}|\}_n$ is strictly decreasing. Since, by (2.12), (2.3) and Assumption 2,

$$(6.17) \quad P_{11}^{-}(h) = P_{11}(0) + \sigma^2 h \leq Kh,$$

we can, using the reverse triangle inequality and the (by BFT) strictly decreasing sequence $\{|P_{11}^{-}(nh) - P_{11}^{-,\infty}|\}_n$, derive (6.9):

$$(6.18) \quad |P_{11}^{-}(nh)| \leq \underbrace{|P_{11}^{-}(nh) - P_{11}^{-,\infty}|}_{\leq |P_{11}^{-}(h) - P_{11}^{-,\infty}|} + |P_{11}^{-,\infty}| \leq \underbrace{P_{11}^{-}(h)}_{\leq Kh} + \underbrace{2P_{11}^{-,\infty}}_{\leq Kh^{1 \wedge \frac{p+1}{2}}, \text{ by (6.3)}} \leq Kh^{1 \wedge \frac{p+1}{2}},$$

which is sharp because it is estimated against the maximum of the initial P_{11}^{-} and the steady state that can both be attained. Recall that, by (6.14), $P_{11}(nh)$ depends continuously on $P_{11}^{-}(nh)$, and, hence, inserting (6.3) into (6.14) yields (6.4)—the necessary computation was already performed in (6.14). Since $P_{11}(nh)$ monotonically increases in $P_{11}^{-}(nh)$ (because the derivative of $P_{11}(nh)$ with respect to $P_{11}^{-}(nh)$ is non-negative for all $P_{11}^{-}(nh)$ due to $R \geq 0$; see (6.14)), we obtain (6.10):

$$(6.18) \quad P_{11}(nh) \stackrel{(6.14)}{\leq} \frac{(\max_n P_{11}^{-}(nh))R}{\max_n P_{11}^{-}(nh) + R} \stackrel{\text{Ass. 4}}{\leq} \frac{Kh^{1 \wedge \frac{p+1}{2}} Kh^p}{Kh^{1 \wedge \frac{p+1}{2}} + Kh^p} \leq \frac{Kh^{(p+1) \wedge \frac{3p+1}{2}}}{Kh^{1 \wedge p}} \\ \leq \begin{cases} Kh^{\frac{p+1}{2}}, & \text{if } p \leq 1, \\ Kh^p, & \text{if } p \geq 1, \end{cases} \leq Kh^{p \vee \frac{p+1}{2}},$$

which is sharp because the steady state (6.9) has these rates. For (6.5), we again first construct the following recursion (from (2.8), (2.12) and (2.3))

$$(6.19) \quad P_{01}^-((n+1)h) = \underbrace{\frac{R}{P_{11}^-(nh) + R}}_{=\alpha(nh)} P_{01}^-(nh) + \underbrace{\left(P_{11}(nh) + \frac{\sigma^2 h}{2}\right)}_{=:g(nh)} h = T_n(P_{01}^-(nh)),$$

where the $\alpha(nh)$ -Lipschitz continuous contractions T_n satisfy the prerequisites of Lemma 6.1, since $\sup_n \alpha(nh) \leq R/(\sigma^2 h + R) < 1$ (due to $\inf_n P_{11}^-(nh) \geq \sigma^2 h$) and the sequence of fixed points $(1 - \alpha(nh))^{-1} g(nh)$ of T_n (well-defined by BFT) converges. Both $\alpha(nh)$ and $g(nh)$ depend continuously on $P_{11}^-(nh)$. Hence, insertion of the limits (6.3) and (6.4) yield

$$(6.20) \quad \lim_{n \rightarrow \infty} (1 - \alpha(nh))^{-1} = \frac{\sigma^2 h + \sqrt{4\sigma^2 R h + \sigma^4 h^2} + 2R}{\sigma^2 h + \sqrt{4\sigma^2 R h + \sigma^4 h^2}}, \quad \text{and}$$

$$(6.21) \quad \lim_{n \rightarrow \infty} g(nh) = \frac{(\sigma^4 h^2 + (2R + \sigma^2 h)\sqrt{4\sigma^2 R h + \sigma^4 h^2} + 4R\sigma^2 h)}{2(\sigma^2 h + \sqrt{4\sigma^2 R h + \sigma^4 h^2} + 2R)} h.$$

Now, application of Lemma 6.1 implies convergence of the recursion (6.19) to the product of these two limits (6.20) and (6.21), i.e. (6.5):

$$(6.22) \quad \begin{aligned} \lim_{n \rightarrow \infty} P_{01}^-(nh) &= \lim_{n \rightarrow \infty} (1 - \alpha(nh))^{-1} \cdot \lim_{n \rightarrow \infty} g(nh) \\ &= \frac{\sigma^4 h^2 + (2R + \sigma^2 h)\sqrt{4\sigma^2 R h + \sigma^4 h^2} + 4R\sigma^2 h}{2(\sigma^2 h + \sqrt{4\sigma^2 R h + \sigma^4 h^2})} h. \end{aligned}$$

For Equations (6.7) and (6.8), we can simply insert Equations (6.3) and (6.5) for $P_{01}^-(nh)$ and $P_{11}^-(nh)$ respectively into their definition (2.9):

$$(6.23) \quad \beta^{\infty, (0)} \stackrel{(2.9)}{=} \frac{P_{01}^{\infty, (0)}}{P_{11}^{\infty, (0)} + R} \stackrel{(6.3), (6.5)}{=} \frac{\sqrt{4R\sigma^2 h + \sigma^4 h^2}}{\sigma^2 h + \sqrt{4R\sigma^2 h + \sigma^4 h^2}} h, \quad \text{and}$$

$$(6.24) \quad \beta^{\infty, (1)} \stackrel{(2.9), (6.3)}{=} \frac{\sigma^2 h + \sqrt{4\sigma^2 R h + \sigma^4 h^2}}{\sigma^2 h + \sqrt{4\sigma^2 R h + \sigma^4 h^2} + 2R}.$$

These steady states (6.7) and (6.8) are again unique and attractive because $\beta^{(0)}(nh)$ and $\beta^{(1)}(nh)$ depend continuously on $P_{11}^-(nh)$ and $P_{01}^-(nh)$. Next, recall that

$$(6.25) \quad P_{01}(nh) \stackrel{(2.12)}{=} \left(1 - \frac{P_{11}^-(nh)}{P_{11}^-(nh) + R}\right) P_{01}^-(nh) = R \frac{P_{01}^-(nh)}{P_{11}^-(nh) + R} \stackrel{(2.9)}{=} R\beta^{(0)}(nh),$$

which, since $P_{01}(nh)$ depends continuously on $\beta^{(0)}(nh)$, implies the unique (attractive) fixed point $P_{01}^\infty(nh) = R\beta^{\infty, (0)}$, which yields (6.6). Now, exploiting (2.9) and $\inf_n P_{11}^-(nh) \geq \sigma^2 h$ yields (6.13):

$$(6.26) \quad \left|1 - \beta^{(1)}(nh)\right| = \frac{R}{P_{11}^-(nh) + R} \leq \frac{R}{\sigma^2 h + R} \stackrel{\text{Ass. 4}}{=} \frac{Kh^p}{Kh + Kh^p} \leq Kh^{(p-1) \vee 0},$$

which is sharp because $\inf_n P_{11}^-(nh) \geq Kh$ is sharp (due to (2.3) and (2.8)). And since, for $\beta^{(0)}$, maximizing over both $P_{01}^-(nh)$ and $P_{11}^-(nh)$ at the same time does not yield a sharp bound (while above in (6.18) and (6.26) the maximization over just one quantity does), we prove (6.12) by inductively showing that

$$(6.27) \quad \left| \beta^{(0)}(nh) \right| \leq \hat{\beta}h, \quad \forall n \in \mathbb{N}, \quad \text{with} \quad \hat{\beta} := \left(\frac{2K_0}{\sigma^2} + \frac{1}{2} \right) \vee 1 > 0,$$

where $K_0 > 0$ is the constant from Assumption 2. The constant $\hat{\beta}$ is independent of n and a possible choice for K in (6.12). The basis ($n = 1$) follows from

$$(6.28) \quad \left| \beta^{(0)}(h) \right| = \frac{|P_{01}^-(h)|}{P_{11}^-(h) + R} \stackrel{(2.8)}{\leq} \frac{|P_{01}(0)| + hP_{11}(0) + \frac{\sigma^2}{2}h^2}{\sigma^2h} \stackrel{\text{Ass. 2}}{\leq} \left(\frac{2K_0}{\sigma^2} + \frac{1}{2} \right)h \leq \hat{\beta}h.$$

In the following inductive step ($n - 1 \rightarrow n$) we, to avoid notational clutter, simply denote $P^-((n - 1)h)_{ij}$ by P_{ij}^- which leaves us—by (2.9), (2.8) and (2.12)—with the following term to bound:

$$(6.29) \quad \left| \beta^{(0)}(nh) \right| = \frac{|P_{01}^-(nh)|}{P_{11}^-(nh) + R} \leq \frac{|P_{01}^-|\alpha(nh) + hP_{11}^-\alpha(nh) + \frac{\sigma^2}{2}h^2}{P_{11}^-\alpha(nh) + \sigma^2h + R},$$

with $\alpha(nh) = \left(1 - \frac{P_{11}^-}{P_{11}^- + R} \right) = \frac{R}{P_{11}^- + R}$. Application of the inductive hypothesis (i.e. $P_{01}^- \leq \hat{\beta}(P_{11}^- + R)$) yields, after some rearrangements, that

$$(6.30) \quad \begin{aligned} \left| \beta^{(0)}(nh) \right| &\leq \frac{\hat{\beta}(P_{11}^- + R)h\alpha(nh) + hP_{11}^-\alpha(nh) + \frac{\sigma^2}{2}h^2}{P_{11}^-\alpha(nh) + \sigma^2h + R} \\ &= \frac{2\hat{\beta}P_{11}^-R + \sigma^2h(P_{11}^- + R) + 2P_{11}^-R + 2\hat{\beta}R^2}{2(P_{11}^-R + \sigma^2h(P_{11}^- + R) + P_{11}^-R + R^2)}h \\ &= \frac{2(\hat{\beta} + 1)\Lambda_1 + \Lambda_2 + 2\hat{\beta}\Lambda_3}{4\Lambda_1 + 2\Lambda_2 + 2\Lambda_3}h, \end{aligned}$$

with $\Lambda_1 := 2P_{11}^-R$, $\Lambda_2 := \sigma^2h(P_{11}^- + R)$, and $\Lambda_3 := R^2$. Now, application of $\hat{\beta} \geq 1$ yields $|\beta^{(0)}(nh)| \leq \hat{\beta}h$, which completes the inductive proof of (6.27). This implies (6.12), which is sharp because it is the order of $\beta^{(0)}$ in the steady state (6.7), for all $p \geq 0$. Now, insertion of (6.12) into (6.25) immediately yields (6.11), which—by (6.25)—inherits the sharpness of (6.12). ■

6.2. Global bounds on state misalignments. Before bounding the added deviation of Ψ from the flow Ψ , a global bound on the state misalignments defined in (4.3) is necessary.

Lemma 6.4. Under Assumptions 1 to 4 and for sufficiently small $h > 0$,

$$(6.31) \quad \max_{n \in [T/h+1]} \delta^{(1)}(nh) \leq \begin{cases} Kh, & \text{if } p \geq 1, \\ K(T), & \text{if } p < 1. \end{cases}$$

Here and in the sequel, $K(T) > 0$ is a constant that depends on T , but not on h .

Proof. For all $n \in [T/h + 1]$, we can estimate

$$\begin{aligned} \delta^{(1)}(nh) &= \left\| m^{(1)}(nh) - f\left(m^{(0)}(nh)\right) \right\| = \left\| \Psi^{(1)}(\mathbf{m}((n-1)h) - f\left(m^{(0)}(nh)\right) \right\| \\ (6.32) \quad &\leq \underbrace{\left\| \Psi^{(1)}(\mathbf{m}((n-1)h) - f\left(m^{-(0)}(nh)\right) \right\|}_{=: J_1(h)} + \underbrace{\left\| f\left(m^{-(0)}(nh)\right) - f\left(m^{(0)}(nh)\right) \right\|}_{=: J_2(h)}, \end{aligned}$$

bound J_1 , by the definition (2.7) of $\Psi^{(1)}(\mathbf{m}((n-1)h)$ as well as the definition (2.11) of $r(nh)$ by,

$$\begin{aligned} J_1(h) &= \left\| m^{(-1)}(nh) + \beta^{(1)}(nh) \left[f\left(m^{-(0)}(nh)\right) - m^{(-1)}(nh) \right] - f\left(m^{-(0)}(nh)\right) \right\| \\ (6.33) \quad &\leq \left\| 1 - \beta^{(1)}(nh) \right\| r(nh) \stackrel{(6.13)}{\leq} Kh^{(p-1) \vee 0} r(nh) \end{aligned}$$

and bound J_2 , by exploiting L -Lipschitz continuity of f , inserting the definition (2.7) of $\Psi^{(0)}(\mathbf{m}((n-1)h)$ and applying (6.12) to $\|\beta^{(0)}(nh)\|$,

$$(6.34) \quad J_2(h) \leq L \left\| m^{(0)}(nh) - m^{-(0)}(nh) \right\| \leq L \left\| \beta^{(0)}(nh) \right\| r(nh) \stackrel{(6.12)}{\leq} Khr(nh).$$

Altogether, after inserting these bounds into (6.32),

$$\begin{aligned} \delta^{(1)}(nh) &\leq \left(Kh^{(p-1) \vee 0} + Kh \right) r(nh) \leq Kh^{((p-1) \vee 0) \wedge 1} r(nh) \\ (6.35) \quad &\stackrel{(4.6)}{\leq} Kh^{(p \vee 1) \wedge 2} + \left(Kh^{((p-1) \vee 0) \wedge 1} + Kh^{(p \vee 1) \wedge 2} \right) \delta^{(1)}((n-1)h) =: \bar{T} \left(\delta^{(1)}((n-1)h) \right). \end{aligned}$$

If $p \geq 1$, by BFT (applicable for sufficiently small $h > 0$ such that $(Kh^{((p-1) \vee 0) \wedge 1} + Kh^{(p \vee 1) \wedge 2}) < 1$ and thereby \bar{T} is a contraction), there exists a unique fixed point δ^∞ of order

$$(6.36) \quad \delta^\infty \leq \frac{Kh^{(p \vee 1) \wedge 2}}{1 - (Kh^{((p-1) \vee 0) \wedge 1} + Kh^{(p \vee 1) \wedge 2})} \leq Kh^{(p \vee 1) \wedge 2}.$$

We proceed with showing by induction that, for all $n \in [T/h]$,

$$(6.37) \quad \delta^{(1)}(nh) \leq \delta^{(1)}(0) \vee 2\delta^\infty.$$

The base case $n = 0$ is trivial. For the inductive step, we distinguish two cases. If $\delta^{(1)}((n-1)h) \leq \delta^\infty$, then $\bar{T}(\delta^{(1)}((n-1)h)) < 2\delta^\infty$, since

$$\bar{T}(\delta^{(1)}((n-1)h)) - \delta^\infty \leq \left| \delta^\infty - \bar{T}(\delta^{(1)}((n-1)h)) \right| < \delta^\infty - \underbrace{\delta^{(1)}((n-1)h)}_{\geq 0} \leq \delta^\infty.$$

In this case,

$$(6.38) \quad \delta^{(1)}(nh) \stackrel{(6.35)}{\leq} \bar{T} \left(\delta^{(1)}((n-1)h) \right) < 2\delta^\infty \leq \delta^{(1)}(0) \vee 2\delta^\infty,$$

where the last inequality follows from the inductive hypothesis. In the other case, namely $\delta^{(1)}((n-1)h) > \delta^\infty$, it follows that

$$\begin{aligned} \delta^{(1)}(nh) - \delta^\infty &\stackrel{(6.35)}{\leq} \bar{T}(\delta^{(1)}((n-1)h)) - \delta^\infty \leq \left| \bar{T}(\delta^{(1)}((n-1)h)) - \delta^\infty \right| \\ &\leq \left| \delta^{(1)}((n-1)h) - \delta^\infty \right| = \delta^{(1)}((n-1)h) - \delta^\infty, \end{aligned} \quad (6.39)$$

which, after adding δ^∞ and applying the inductive hypothesis, completes the inductive step. Hence, (6.37) holds. Since this bound is uniform in n , inserting the orders of $\delta^{(1)}(0)$ from Lemma 5.1 and of δ^∞ from (6.36), concludes the proof for $p \geq 1$.

If $p < 1$, inserting (6.33) and (6.34) into (6.32) yields the (sharp) bound

$$(6.40) \quad \delta^{(1)}(nh) \leq Kh + \underbrace{\left(\left\| 1 - \beta^{(1)}(nh) \right\| + Kh \right)}_{\leq 1, \text{ by (2.9)}} \delta^{(1)}((n-1)h).$$

Now, application of a special discrete Grönwall inequality [4] to (6.40) only yields the (sharp) bound of $K(T) > 0$, claimed in (6.31), that depends on T but not on h . (See Proof 14 for an analogous application of this discrete Grönwall inequality in more detail.) ■

6.3. Prerequisite for discrete Grönwall inequality.

Lemma 6.5. Under Assumptions 1 to 4 and for sufficiently small $h > 0$,

$$(6.41) \quad \max_{n \in [T/h+1]} \left\| \Psi_{P(nh),h}(\mathbf{m}(nh)) - \Phi_h(m^{(0)}(nh)) \right\|_h \leq \begin{cases} Kh^2, & \text{if } p \geq 1, \\ K(T)h, & \text{if } p < 1. \end{cases}$$

Proof. By (3.5), we have

$$(6.42) \quad \left\| \Psi_{P(nh),h}(\mathbf{m}(nh)) - \Phi_h(m^{(0)}(nh)) \right\|_h = S_1(h) + hS_2(h),$$

with $S_1(h)$ and $S_2(h)$ defined and bounded by

$$\begin{aligned} S_1(h) &:= \left\| \Psi_h^{(0)}(\mathbf{m}(nh)) - \Phi_h^{(0)}(m^{(0)}(nh)) \right\| \\ &\stackrel{(4.2)}{\leq} \underbrace{\Delta^{-(0)}((n+1)h)}_{\stackrel{(4.4)}{\leq} Kh^2 + \delta^{(0)}(nh) + h\delta^{(1)}(nh)} + \underbrace{\left\| \beta^{(0)}((n+1)h) \right\|}_{\stackrel{(6.12)}{\leq} Kh} \underbrace{\left\| r((n+1)h) \right\|}_{\stackrel{(4.6)}{\leq} Kh + (1+Kh)\delta^{(1)}(nh)}, \end{aligned} \quad (6.43)$$

and, analogously,

$$\begin{aligned} S_2(h) &:= \left\| \Psi_h^{(1)}(\mathbf{m}(nh)) - \Phi_h^{(1)}(m^{(0)}(nh)) \right\| \\ &\stackrel{(4.2)}{\leq} \underbrace{\Delta^{-(1)}((n+1)h)}_{\stackrel{(4.4)}{\leq} Kh + \delta^{(1)}(nh)} + \underbrace{\left\| \beta^{(1)}((n+1)h) \right\|}_{\stackrel{(2.9)}{\leq} 1} \underbrace{\left\| r((n+1)h) \right\|}_{\stackrel{(4.6)}{\leq} Kh + (1+Kh)\delta^{(1)}(nh)}. \end{aligned} \quad (6.44)$$

642 Insertion of (6.43) and (6.44) into (6.42) yields

$$643 \quad \left\| \Psi_{P(nh),h}(\mathbf{m}(nh)) - \Phi_h(m^{(0)}(nh)) \right\|_h \leq Kh^2 + \delta^{(0)}(nh) + Kh\delta^{(1)}(nh),$$

645 which—after recalling $\delta^{(0)}(nh) = 0$ and applying Lemma 6.4 to $\delta^{(1)}(nh)$ —implies the remain-
646 ing claim for $p < 1$. ■

647 The previous lemma now implies a suitable prerequisite for a discrete Grönwall inequality.

648 **Lemma 6.6.** *Under Assumptions 1 to 4 and for sufficiently small $h > 0$,*

$$649 \quad (6.45) \quad \|\varepsilon((n+1)h)\|_h \leq \begin{cases} Kh^2 + (1+Kh)\|\varepsilon^{(0)}(nh)\|, & \text{if } p \geq 1, \\ K(T)h + (1+Kh)\|\varepsilon^{(0)}(nh)\|, & \text{if } p < 1. \end{cases}$$

651 *Proof.* We observe, by the triangle inequality for $\|\cdot\|_h$, that

$$652 \quad \|\varepsilon((n+1)h)\|_h = \left\| \Psi_{P(nh),h}(\mathbf{m}(nh)) - \Phi_h(x^{(0)}(nh)) \right\|_h$$

$$653 \quad \leq \left\| \Psi_{P(nh),h}(\mathbf{m}(nh)) - \Phi_h(m^{(0)}(nh)) \right\|_h + \left\| \Phi_h(m^{(0)}(nh)) - \Phi_h(x^{(0)}(nh)) \right\|_h.$$

655 The proof is concluded by applying Lemma 6.5 to the first and Lemma 3.3 to the second
656 summand of this bound (as well as recalling that $\|\varepsilon^{(0)}(nh)\| = \|m^{(0)}(nh) - x^{(0)}(nh)\|$). ■

657 6.4. Global convergence rates.

658 **Theorem 6.7 (Global truncation error).** *Under Assumptions 1 to 4 and for sufficiently*
659 *small $h > 0$,*

$$660 \quad (6.46) \quad \max_{n \in [T/h+1]} \|\varepsilon^{(0)}(nh)\| \leq \max_{n \in [T/h+1]} \|\varepsilon(nh)\|_h \leq \begin{cases} K(T)h, & \text{if } p \geq 1, \\ K(T), & \text{if } p < 1. \end{cases}$$

662 **Remark 6.8.** Theorem 6.7 not only implies that the truncation error $\|\varepsilon^{(0)}(nh)\|$ on the
663 solution of (1.1) has global order h , but also (by (3.5)) that the truncation error $\|\varepsilon^{(1)}(nh)\|$
664 on the derivative is uniformly bounded by a constant K independent of h . The rates of this
665 theorem are sharp—in the sense that it cannot be improved over all f satisfying Assump-
666 tion 1—since, for $p \geq 1$, it is (in light of the sharp Theorem 5.2) the best global rate possible
667 and, for $p < 1$, the bound cannot be improved because the state misalignment is not globally
668 convergent (see (6.31)), which yields a maximum local added error of order Kh in Lemma 6.5
669 and hence only a global error of order K in Theorem 6.7.

670 *Proof.* Using $\|\varepsilon^{(0)}(nh)\| \leq \|\varepsilon(nh)\|_h$ (due to (3.5)), the bound (6.45), a telescoping sum,

671 and $\|\varepsilon(0)\|_h \leq Kh^2$ (by [Assumption 3](#)), we obtain, for sufficiently small $h > 0$, that

$$\begin{aligned}
 672 \quad & \|\varepsilon((n+1)h)\|_h - \|\varepsilon(nh)\|_h \stackrel{(3.5)}{\leq} \|\varepsilon((n+1)h)\|_h - \|\varepsilon^{(0)}(nh)\| \\
 673 \quad & \stackrel{(6.45)}{\leq} \begin{cases} Kh^2 + Kh\|\varepsilon^{(0)}(nh)\| \stackrel{(3.5)}{\leq} Kh^2 + Kh\|\varepsilon(nh)\|_h, & \text{if } p \geq 1, \\ K(T)h + Kh\|\varepsilon^{(0)}(nh)\| \stackrel{(3.5)}{\leq} K(T)h + Kh\|\varepsilon(nh)\|_h, & \text{if } p < 1, \end{cases} \\
 674 \quad & \stackrel{(\text{tel. sum})}{=} \begin{cases} Kh^2 + Kh \sum_{l=0}^{n-1} (\|\varepsilon((l+1)h)\|_h - \|\varepsilon(lh)\|_h) + \|\varepsilon(0)\|_h, & \text{if } p \geq 1, \\ K(T)h + Kh \sum_{l=0}^{n-1} (\|\varepsilon((l+1)h)\|_h - \|\varepsilon(lh)\|_h) + \|\varepsilon(0)\|_h, & \text{if } p < 1, \end{cases} \\
 675 \quad (6.47) \quad & \stackrel{(\|\varepsilon(0)\|_h \leq Kh^2)}{\leq} \begin{cases} Kh^2 + Kh \sum_{l=0}^{n-1} (\|\varepsilon((l+1)h)\|_h - \|\varepsilon(lh)\|_h), & \text{if } p \geq 1, \\ K(T)h + Kh \sum_{l=0}^{n-1} (\|\varepsilon((l+1)h)\|_h - \|\varepsilon(lh)\|_h), & \text{if } p < 1. \end{cases} \\
 676 \quad &
 \end{aligned}$$

677 Now, recall that—by a special version of the discrete Grönwall inequality [\[4\]](#)—if z_n and g_n are
 678 sequences of real numbers (with $g_n \geq 0$), $c \geq 0$ is a nonnegative constant, and if

$$\begin{aligned}
 679 \quad (6.48) \quad & z_n \leq c + \sum_{l=0}^{n-1} g_l z_l, \quad \text{for all } n \in \mathbb{N}, \\
 680 \quad &
 \end{aligned}$$

681 then

$$\begin{aligned}
 682 \quad (6.49) \quad & z_n \leq c \prod_{l=0}^{n-1} (1 + g_l) \leq c \exp\left(\sum_{l=0}^{n-1} g_l\right), \quad \text{for all } n \in \mathbb{N}. \\
 683 \quad &
 \end{aligned}$$

684 Application of this inequality to [\(6.47\)](#) with $z_n := \|\varepsilon((n+1)h)\|_h - \|\varepsilon(nh)\|_h$, $g_n := Kh$, and
 685 $c := Kh^2$, if $p \geq 1$ (or $c := K(T)h$, if otherwise $p < 1$), yields

$$\begin{aligned}
 686 \quad (6.50) \quad & \|\varepsilon((n+1)h)\|_h - \|\varepsilon(nh)\|_h \leq \begin{cases} K(T)h^2 \exp(nKh) \stackrel{n \leq T/h}{\leq} K(T)h^2, & \text{if } p \geq 1, \\ K(T)h \exp(nKh) \stackrel{n \leq T/h}{\leq} K(T)h, & \text{if } p < 1. \end{cases} \\
 687 \quad &
 \end{aligned}$$

688 By another telescoping sum argument and $\|\varepsilon(0)\|_h \leq Kh^2$, we obtain

$$\begin{aligned}
 689 \quad & \|\varepsilon(nh)\|_h \stackrel{(\text{tel. sum})}{=} \sum_{l=0}^{n-1} (\|\varepsilon((l+1)h)\|_h - \|\varepsilon(lh)\|_h) + \|\varepsilon(0)\|_h \\
 690 \quad (6.51) \quad & \stackrel{(6.50)}{\leq} \begin{cases} nK(T)h^2 + Kh^2 \stackrel{n \leq T/h}{\leq} K(T)h + Kh^2 \leq K(T)h, & \text{if } p \geq 1, \\ nK(T)h + Kh^2 \stackrel{n \leq T/h}{\leq} K(T) + Kh^2 \leq K(T), & \text{if } p < 1, \end{cases} \\
 691 \quad &
 \end{aligned}$$

692 for sufficiently small $h > 0$. Recalling that, by [\(3.5\)](#), $\|\varepsilon^{(0)}(nh)\| \leq \|\varepsilon(nh)\|_h$ concludes the
 693 proof. ■

7. Calibration of confidence intervals. In PN, one way to judge calibration of a Gaussian output $\mathcal{N}(m, V)$ is to check whether the implied 0.95 confidence interval $[m - 2\sqrt{V}, m + 2\sqrt{V}]$ contracts at the same rate as the convergence rate of the posterior mean to the true quantity of interest (even if this amounts to matching a worst-case error to a probabilistic confidence interval). For the filter, this would mean that the rate of contraction of $\max_n \sqrt{P_{00}(nh)}$ should contract at half the rate of the (sharp) bounds on $\max_{n \in [T/h+1]} \|\varepsilon^{(0)}(nh)\|$ from [Theorem 6.7](#). Otherwise, for a higher or lower rate of the interval it would eventually be under- or overconfident, as $h \rightarrow 0$. The following proposition shows—in light of the sharp bound (6.46) on the global error—that the confidence intervals are well-calibrated in this way if and only if $p \geq 1$.

Proposition 7.1. *Under [Assumptions 2 and 4](#) and for sufficiently small $h > 0$,*

$$(7.1) \quad \max_{n \in [T/h+1]} P_{00}^-(nh) \leq K(T)h^{(p+1)\wedge 2}, \quad \text{and}$$

$$(7.2) \quad \max_{n \in [T/h+1]} P_{00}(nh) \leq K(T)h^{(p+1)\wedge 2}.$$

Proof. Again, w.l.o.g. $d = 1$. We first show that the bounds (7.1) and (7.2) hold and then argue that they are sharp. The recursion for $P_{00}^-(nh)$ is given by

$$(7.3) \quad \begin{aligned} P_{00}^-((n+1)h) &\stackrel{(2.8),(2.3)}{=} P_{00}(nh) + 2hP_{01}(nh) + h^2P_{11}(nh) + \frac{\sigma^2}{3}h^3 \\ &= P_{00}^-(nh) - \beta^{(0)}(nh)P_{01}^-(nh) + 2hR\beta^{(0)}(nh) + h^2R\beta^{(1)}(nh) + \frac{\sigma^2}{3}h^3, \end{aligned}$$

where we used $P_{00}(nh) = P_{00}^-(nh) - \beta^{(0)}(nh)P_{01}^-(nh)$ and $P_{11}(nh) = R\beta^{(1)}(nh)$ (both due to (2.12) and (2.9)), as well as $P_{01}(nh) = R\beta^{(0)}(nh)$ (see (6.25)), for the last equality in (7.3). By $P_{01}^-(nh) \leq P_{01}(nh)$ and $|\beta^{(1)}| \leq 1$ (due to (2.9)), application of the triangle inequality to (7.3) yields

$$(7.4) \quad P_{00}^-((n+1)h) \leq P_{00}^-(nh) + |\beta^{(0)}(nh)| |P_{01}(nh)| + 2hR|\beta^{(0)}(nh)| + h^2R + \frac{\sigma^2}{3}h^3,$$

which, by (6.11) and (6.12), implies

$$(7.5) \quad P_{00}^-((n+1)h) \leq P_{00}^-(nh) + Kh^{(p+2)\wedge 3}.$$

This, by $N = T/h$, implies (7.1). Since $P_{00}(nh) \leq P_{00}^-(nh)$, this bound is also valid for P_{00} , i.e. (7.2) holds. The bound (7.1) is sharp, since, e.g. when the covariance matrices are in the steady state, the covariance matrix keeps growing by a rate of $Kh^{(p+2)\wedge 3}$ for sufficiently small $h > 0$, since the only negative summand in (7.3) is given by

$$(7.6) \quad \begin{aligned} &\beta^{\infty,(0)}P_{01}^\infty = \\ &\underbrace{\frac{1}{2}h^2}_{\in \Theta(h^2)} \underbrace{\sqrt{(\sigma^2h)^2 + 4(\sigma^2h)R}}_{\in \Theta(h^{1 \wedge \frac{p+1}{2}})} \underbrace{\left((\sigma^2h)^2 + 4(\sigma^2h)R + ((\sigma^2h) + 2R)\sqrt{(\sigma^2h)^2 + 4(\sigma^2h)R} \right)}_{\in \Theta(h^{2 \wedge (p+1)})}, \end{aligned}$$

due $R \equiv Kh^p$ under [Assumption 4](#). Hence, the sole negative summand $-\beta^{\infty,(0)}P_{01}^\infty$ of [\(7.3\)](#) is in $\Theta(h^{5 \wedge \frac{3p+7}{2}})$ and thereby of higher order than the remaining sum of positive summands:

$$(7.7) \quad \underbrace{2hR}_{\in \Theta(h^{p+1})} \underbrace{\beta^{\infty,(0)}(nh)}_{\in \Theta(h), \text{ by (6.7)}} + \underbrace{h^2R}_{\in \Theta(h^{p+2})} \underbrace{\beta^{\infty,(1)}(nh)}_{\in \Theta(1), \text{ by (6.8)}} + \underbrace{\frac{\sigma^2}{3}h^3}_{\in \Theta(h^3)} \in \Theta(h^{3 \wedge (p+2)}).$$

that—by $R \equiv Kh^p$ as well as [\(6.7\)](#) and [\(6.8\)](#)—is, for all $p \geq 0$, in $\Theta(h^{3 \wedge (p+2)})$. Hence, for sufficiently small $h > 0$, it still holds in the steady state that $P_{00}^-(n+1)h - P_{00}^-(nh) \geq Kh^{(p+2) \wedge 3}$, and therefore [\(7.1\)](#) is sharp. The sharpness of [\(7.1\)](#) is inherited by [\(7.2\)](#) since, in the steady state, by [\(2.12\)](#) and [\(2.9\)](#), $P_{00}(nh) = P_{00}^-(nh) - \beta^{(0),\infty}P_{01}^{-,\infty}$ and the subtracted quantity $\beta^{(0),\infty}P_{01}^{-,\infty}$ is—as shown above—only of order $\Theta(h^{5 \wedge \frac{3p+7}{2}})$. ■

8. Discussion. We presented a worst-case convergence rate analysis of the Gaussian ODE filter, comprising both local and global convergence rates. While local convergence rates of h^{q+1} were shown to hold for all $q \in \mathbb{N}$, IBM and IOUP prior as well as any noise model $R \geq 0$, our global convergence results is restricted to the case of $q = 1$, IBM prior and fixed noise model $R \equiv Kh^p$ with $p \geq 1$. While the restriction of the noise model seems inevitable (due to the non-vanishing state misalignment for $p < 1$, see [Lemma 6.4](#)), we believe that the other two restrictions can be lifted: In light of [Theorem 5.2](#), global convergence rates for the IOUP prior might only require an additional assumption that ensures that all possible data sequences $\{y(nh); n = 1, \dots, T/h\}$ (and thereby all possible q^{th} -state sequences $\{m^{(q)}(nh); n = 0, \dots, T/h\}$) remain uniformly bounded. For the case of $q \geq 2$, it seems plausible that a proof analogous to the presented one would already yield global convergence rates of order h^q .⁶

For the use of the filter in inverse problem, the calibration of confidence intervals might turn out to be significant: Interestingly, all of the studied noise models, which turned out to be globally convergent, give well-calibrated confidence intervals in the sense discussed in [Section 7](#) (maybe because probabilistic state-space models jointly treat uncertain quantities in a self-consistent way). We seek to exploit this property in inverse problems to compensate for over-confidence of numerical forward solutions; [\[6, Section 3.2.\]](#) and [\[1, Section 7\]](#), e.g., discuss how another class of PN forward solvers (reviewed above in [Subsection 1.2](#)) can achieve this. The orders of the predictive confidence intervals can also help to intuitively explain the threshold of $p = 1$ below which many properties of the filter break: According to [\[13, Equation \(20\)\]](#), the ‘true’ (push-forward) variance on $y(t)$ given the predictive distribution $\mathcal{N}(m^-(t), P^-(t))$ is equal to the integral of ff^\top with respect to $\mathcal{N}(m^-(t), P^-(t))$, whose maximum over all time steps, by [\(7.1\)](#), has order $\mathcal{O}(h^{\frac{p+1}{2} \wedge 1})$ if ff^\top is globally Lipschitz—since $P^-(t)$ enters the argument of the integrand ff^\top , after a change of variable, only under a square root. Hence, the added ‘statistical’ noise R on the evaluation of f is of lower order than the accumulated ‘numerical’ variance $P^-(t)$ (thereby preventing numerical convergence) if and only if $p < 1$. Maybe this, in the spirit of [\[9, Subsection 3\(d\)\]](#), can serve as a criterion for vector fields f that are too roughly approximated for a numerical solver to output a trustworthy result, even

⁶According to [\[16\]](#), the filter might, however, suffer from numerical instability for high choices of q . (See [\[28, Section 3.1.\]](#) for an explanation of how such results on spline-based methods concern the ODE filter.)

as $h \rightarrow 0$.

While the rates of h^{q+1} (local) and h^q (global) are optimal given that q derivatives are modeled, we believe that future research should seek to go beyond worst-case analyses. The competitive practical performance of the filter, as demonstrated in [28, Section 5], might only become completely captured by an average-case analysis in the sense of [24]. To see this, recall that Kalman filtering is optimal in the (average) L^2 -sense for time series analysis in linear dynamical systems with Gaussian initial distribution (as defined by our prior (2.1)), i.e. when the data is not evaluations of f but real i.i.d. measurements, and in the special case of $\dot{x}(t) = f(t)$, i.e. when the IVP simplifies to a quadrature problem—see [26, Section 4.3.] and [19, Section 2.2.] respectively. In fact, the entire purpose of the update step is to correct the prediction in the (on average) right direction, while a worst-case analysis must assume that it corrects in the worst possible direction in every step (which we execute by the application of the triangle inequality in (4.2) resulting in a worst-case upper bound that is the sum of the worst-case errors from prediction and update step). An analysis of the probabilities of ‘good’ vs. ‘bad’ updates might therefore pave the way for such an average-case analysis in the setting of this paper. Since, in practice, truncation errors of ODE solvers tend to be significantly smaller than the worst case, such an analysis might be very useful for applications. Lastly, in light of our byproduct of closed-form expressions for the steady states, more equivalences with classical methods for $R > 0$ could now be within close reach (via a reformulation as a Nordsieck method; as carried out in [28] for $R \equiv 0$).

Appendix A. Derivation of A and Q . As derived in [25, section 2.2.6.] the solution of the SDE (2.1), i.e.

$$(A.1) \quad d\mathbf{X}_{j,t} = \begin{pmatrix} dX_{j,t}^{(0)} \\ \vdots \\ dX_{j,t}^{(q-1)} \\ dX_{j,t}^{(q)} \end{pmatrix} = \underbrace{\begin{pmatrix} 0 & 1 & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & & 0 \\ \vdots & & \ddots & & 1 \\ a_0 & \dots & \dots & \dots & a_q \end{pmatrix}}_{=:F} \underbrace{\begin{pmatrix} X_t^{(0)} \\ \vdots \\ X_t^{(q-1)} \\ X_t^{(q)} \end{pmatrix}}_{=:X_t} dt + \underbrace{\begin{pmatrix} 0 \\ \vdots \\ 0 \\ \sigma \end{pmatrix}}_{=:L} dB_t,$$

is a Gauss-Markov process with mean $m(t)$ and covariance matrix $P(t)$ given by

$$(A.2) \quad m(t) = A(t)m(0), \quad P(t) = A(t)P(0)A(t)^\top + Q(t, \sigma),$$

where the matrices $A(h)$, $Q(h, \sigma) \in \mathbb{R}^{(q+1) \times (q+1)}$ are explicitly defined by

$$(A.3) \quad A(h) = \exp(hF),$$

$$(A.4) \quad Q(h, \sigma) := \int_0^h \exp(F(t - \tau)) L \sigma^2 L^\top \exp(F(t - \tau))^\top d\tau.$$

A.1. Integrated Brownian motion. The result of this calculation is also given in [13, 28]. If we choose $a_0, \dots, a_q = 0$ in (A.1) the unique strong solution of the SDE is a q -times IBM. By (A.3) and

$$(A.5) \quad ((hF)^k)_{i,j} = h^k \mathbb{1}_{j-i=k},$$

it follows that

$$(A.6) \quad A^{\text{IBM}}(h)_{i,j} = \left(\sum_{k=0}^{\infty} \frac{(hF)^k}{k!} \right)_{i,j} = \left(\mathbb{1}_{i \leq j} \frac{h^{j-i}}{(j-i)!} \right)_{i,j}.$$

Analogously, it follows that

$$(A.7) \quad \exp(F(t - \tau)) = \left(\mathbb{1}_{i \leq j} \frac{(t-\tau)^{j-i}}{(j-i)!} \right)_{i,j}.$$

If we insert this into the integrand in (A.4) and do the matrix multiplication, we obtain, by the sparsity of L , that

$$(A.8) \quad Q(h)_{i,j} = \frac{\sigma^2}{(q-i)!(q-j)!} \int_0^h (h-\tau)^{2q-i-j} d\tau = \sigma^2 \frac{h^{2q+1-i-j}}{(2q+1-i-j)(q-i)!(q-j)!}.$$

A.2. Integrated Ornstein-Uhlenbeck process. Parts of this calculation can be found in [17]. If we choose $a_0, \dots, a_{q-1} = 0$ and $a_q = -\theta$ (for an arbitrary $\theta > 0$) in (A.1) the unique strong solution of the SDE is a q -times IOUP; see e.g. [12, Example 6.8.]. By (A.3) and

$$(A.9) \quad ((Fh)^k)_{i,j} = h^k \left[\mathbb{1}_{j-i=k} + (-\theta)^{k+i-q} \mathbb{1}_{\{j=q, i+k \geq q\}} \right].$$

Hence,

$$(A.10) \quad A(h)_{ij} := \exp(Fh)_{ij} = \begin{cases} \mathbb{1}_{j \geq i} \frac{h^{j-i}}{(j-i)!}, & \text{if } j \neq q, \\ \frac{1}{(-\theta)^{q-i}} \sum_{k=q-i}^{\infty} \frac{(-\theta h)^k}{k!}, & \text{if } j = q. \end{cases}$$

If we insert (A.10) into (A.4), we obtain, by the sparsity of L , that

$$(A.11) \quad Q(h)_{ij} = \frac{\sigma^2}{\theta^{2q-i-j}} \int_0^h \left(\sum_{k=q-i}^{\infty} \frac{(-\theta \tau)^k}{k!} \right) \left(\sum_{l=q-j}^{\infty} \frac{(-\theta \tau)^l}{l!} \right) d\tau,$$

and, by the dominated convergence theorem (with dominating function $\tau \mapsto e^{2\theta\tau}$), we obtain

$$(A.12) \quad Q(h)_{ij} = \frac{\sigma^2}{\theta^{2q-i-j}} \sum_{k=q-i}^{\infty} \sum_{l=q-j}^{\infty} \int_0^h \frac{(-\theta \tau)^{k+l}}{k!l!} d\tau = \frac{\sigma^2}{\theta^{2q-i-j}} \sum_{k=q-i}^{\infty} \sum_{l=q-j}^{\infty} \theta^{k+l} \frac{h^{k+l+1}}{(k+1+l)k!l!}.$$

Now, by recalling from (A.8) that the summand for $(k, l) = (q-i, q-j)$ in the above double series is equal to $Q(h)_{ij}$ for IBM by (A.8) and the rest of the series is in $\Theta(h^{2q+2-i-j})$, it follows that

$$(A.12) \quad Q(h)_{ij} = Q^{\text{IBM}}(h)_{ij} + \Theta(h^{2q+2-i-j}).$$

Acknowledgments. The authors are grateful to Han Cheng Lie for discussions and feedback to early versions of what is now Sections 3 and 4 of this work, as well as Subsection 6.4. The authors also thank Michael Schober for valuable discussions and helpful comments on the manuscript.

TJS is partially supported by the Freie Universität Berlin within the Excellence Initiative of the German Research Foundation (DFG), by the DFG through grant CRC 1114 “Scaling Cascades in Complex Systems”, and by the National Science Foundation under grant DMS-1127914 to the Statistical and Applied Mathematical Sciences Institute (SAMSI) and SAMSI’s QMC Working Group II “Probabilistic Numerics”. Any opinions, findings, and conclusions or recommendations expressed in this article are those of the authors and do not necessarily reflect the views of the above-named institutions and agencies.

REFERENCES

- [1] A. ABDULLE AND G. GAREGNANI, *Random time step probabilistic methods for uncertainty quantification in chaotic and geometric numerical integration*, arXiv:1703.03680 [math.NA], (2018).
- [2] B. ANDERSON AND J. MOORE, *Optimal Filtering*, Prentice-Hall, 1979.
- [3] O. A. CHKREBTHI, D. A. CAMPBELL, B. CALDERHEAD, AND M. A. GIROLAMI, *Bayesian solution uncertainty quantification for differential equations*, Bayesian Analysis, 11 (2016), pp. 1239–1267.
- [4] D. S. CLARK, *Short proof of a discrete Gronwall inequality*, Discrete Applied Mathematics, 8 (1987), pp. 279–281.
- [5] J. COCKAYNE, C. OATES, T. SULLIVAN, AND M. GIROLAMI, *Bayesian probabilistic numerical methods*, arXiv:1702.03673 [stat.ME], (2017).
- [6] P. R. CONRAD, M. GIROLAMI, S. SÄRKKÄ, A. STUART, AND K. ZYGALAKIS, *Statistical analysis of differential equations: introducing probability measures on numerical solutions*, Statistics and Computing, 27 (2017), pp. 1065–1082.
- [7] P. DIACONIS, *Bayesian numerical analysis*, Statistical decision theory and related topics, IV (1988), pp. 163–175.
- [8] J. HARTIKAINEN AND S. SÄRKKÄ, *Kalman filtering and smoothing solutions to temporal Gaussian process regression models*, in IEEE International Workshop on Machine Learning for Signal Processing (MLSP), 2010, pp. 379–384.
- [9] P. HENNIG, M. OSBORNE, AND M. GIROLAMI, *Probabilistic numerics and uncertainty in computations*, Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences, 471 (2015).
- [10] B. L. HULME, *Piecewise polynomial Taylor methods for initial value problems*, Numerische Mathematik, 17 (1971), pp. 367–381.
- [11] A. JAZWINSKI, *Stochastic Processes and Filtering Theory*, Academic Press, 1970.
- [12] I. KARATZAS AND S. SHREVE, *Brownian Motion and Stochastic Calculus*, Springer, 1991.
- [13] H. KERSTING AND P. HENNIG, *Active uncertainty calibration in Bayesian ODE solvers*, in Proceedings of the Thirty-Second Conference on Uncertainty in Artificial Intelligence (UAI), A. Ihler and D. Janzing, eds., AUAI Press, 2016, pp. 309–318.
- [14] P. LANCASTER AND L. RODMAN, *Algebraic Riccati Equations*, Oxford Science Publications, 1995.
- [15] H. LIE, A. STUART, AND T. SULLIVAN, *Strong convergence rates of probabilistic integrators for ordinary differential equations*, arXiv:1703.03680 [math.NA], (2017).
- [16] F. LOSCALZO AND T. TALBOT, *Spline function approximations for solutions of ordinary differential equations*, SIAM Journal on Numerical Analysis, 4 (1967).
- [17] E. MAGNANI, H. KERSTING, M. SCHÖBER, AND P. HENNIG, *Bayesian filtering for ODEs with bounded derivatives*, arXiv:1709.08471 [cs.NA], (2017).
- [18] A. NORDSIECK, *On numerical integration of ordinary differential equations*, Mathematics of Computation, 16 (1962), pp. 22–49.
- [19] A. O’HAGAN, *Bayes–Hermite quadrature*, Journal of Statistical Planning and Inference, 29 (1991), pp. 245–260.
- [20] A. O’HAGAN, *Some Bayesian numerical analysis*, Bayesian Statistics, 4 (1992), pp. 345–363.
- [21] B. ØKSENDAL, *Stochastic Differential Equations: An Introduction with Applications*, Springer, 6 ed., 2003.
- [22] H. POINCARÉ, *Calcul des probabilités*, Gauthier-Villars, Paris, 1896.

- 884 [23] C. RASMUSSEN AND C. WILLIAMS, *Gaussian Processes for Machine Learning*, The MIT Press, 2006.
- 885 [24] K. RITTER, *Average-Case Analysis of Numerical Problems*, Springer, 2000.
- 886 [25] S. SÄRKKÄ, *Recursive Bayesian Inference on Stochastic Differential Equations*, PhD thesis, Helsinki
- 887 University of Technology, 2006.
- 888 [26] S. SÄRKKÄ, *Bayesian Filtering and Smoothing*, Cambridge University Press, 2013.
- 889 [27] M. SCHÖBER, D. DUVENAUD, AND P. HENNIG, *Probabilistic ODE solvers with Runge-Kutta means*, in Ad-
- 890 vances in Neural Information Processing Systems (NIPS) 27, Z. Ghahramani, M. Welling, C. Cortes,
- 891 N. D. Lawrence, and K. Q. Weinberger, eds., Curran Associates, Inc., 2014, pp. 739–747.
- 892 [28] M. SCHÖBER, S. SÄRKKÄ, AND P. HENNIG, *A probabilistic model for the numerical solution of initial*
- 893 *value problems*, Statistics and Computing, (2018).
- 894 [29] J. SKILLING, *Bayesian solution of ordinary differential equations*, Maximum Entropy and Bayesian Meth-
- 895 ods, Seattle, (1991).
- 896 [30] G. TESCHL, *Ordinary Differential Equations and Dynamical Systems*, American Mathematical Society,
- 897 2012.
- 898 [31] O. TEYMUR, K. ZYGALAKIS, AND B. CALDERHEAD, *Probabilistic linear multistep methods*, in Advances in
- 899 Neural Information Processing Systems (NIPS) 29, D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon,
- 900 and R. Garnett, eds., Curran Associates, Inc., 2016, pp. 4314–4321.
- 901 [32] J. WANG, J. COCKAYNE, AND C. OATES, *On the Bayesian solution of differential equations*,
- 902 arXiv:1703.03680 [math.NA], (2018).