
A Probabilistic State Space Model for Joint Inference from Differential Equations and Data

Jonathan Schmidt
University of Tübingen
Tübingen, Germany

jonathan.schmidt@uni-tuebingen.de

Nicholas Krämer
University of Tübingen
Tübingen, Germany

nicholas.kraemer@uni-tuebingen.de

Philipp Hennig
University of Tübingen
Max Planck Institute for Intelligent Systems
Tübingen, Germany
philipp.hennig@uni-tuebingen.de

Abstract

Mechanistic models with differential equations are a key component of scientific applications of machine learning. Inference in such models is usually computationally demanding, because it involves repeatedly solving the differential equation. The main problem here is that the numerical solver is hard to combine with standard inference techniques. Recent work in probabilistic numerics has developed a new class of solvers for ordinary differential equations (ODEs) that phrase the solution process directly in terms of Bayesian filtering. We here show that this allows such methods to be combined very directly, with conceptual and numerical ease, with latent force models in the ODE itself. It then becomes possible to perform approximate Bayesian inference on the latent force as well as the ODE solution in a single, linear complexity pass of an extended Kalman filter / smoother — that is, at the cost of computing a single ODE solution. We demonstrate the expressiveness and performance of the algorithm by training, among others, a non-parametric SIRD model on data from the COVID-19 outbreak.

1 Introduction

Mechanistic models based on ordinary differential equations (ODEs) are popular across a wide range of scientific disciplines. To increase the descriptive power of such models, it is common to consider parametrized versions of ODEs and find a set of parameters such that the dynamics reproduce empirical observations as accurately as possible. Algorithms for this purpose typically involve repeated forward simulations in the context of, e.g., Markov-chain Monte Carlo or optimization. The need for iterated computation of ODE solutions may demand simplifications in the model to meet limits in the computational budget.

This work describes an algorithm that merges mechanistic knowledge in the form of an ODE with a non-parametric model over the parameters controlling the ODE – a *latent force* that represents quantities of interest. The algorithm then infers a trajectory that is informed

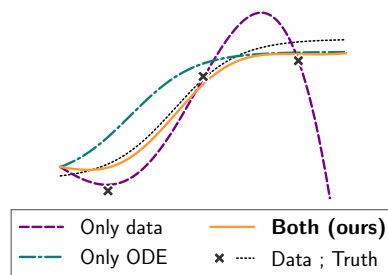


Figure 1: Inferring an unknown function with a Gaussian Process and differential sources of information.

by the observations but also follows sensible dynamics, as defined by the ODE, in the absence of observations (Figure 1). The main insight enabling this approach is that if probabilistic ODE solvers use the language of (extended) Kalman filters, conditioning on observations and solving the ODE itself is possible in one and the same process of Bayesian filtering and smoothing. Instead of iterated computation of ODE solutions, a posterior distribution arises from *a single* forward simulation, which has complexity equivalent to numerically computing an ODE solution, once, with a filtering-based, probabilistic ODE solver [39]. Intuitively, one can think of this as opening up the black box ODE solver and acknowledging that each task – solving the ODE and discovering a latent force – is probabilistic inference in a state-space model.

The main contribution of this work is formalizing this intuition. Several experiments empirically prove the efficiency and the expressivity of the resulting algorithm. In particular, a practical model for the dynamics of the COVID-19 pandemic is considered, in which a non-parametric latent force captures the effect of policy measures that continuously change the contact rate among the population.

2 Problem setting

Let $x : [t_0, t_{\max}] \rightarrow \mathbb{R}^d$ be a process that is observed at a discrete set of points $\mathcal{T}_N^{\text{OBS}} := (t_0^{\text{OBS}}, \dots, t_N^{\text{OBS}})$ through a sequence of measurements $y_{0:N} := (y_0, \dots, y_N) \in \mathbb{R}^{(N+1) \times k}$. Assume that these measurements are subject to additive i.i.d. Gaussian noise, according to the observation model

$$y_n = Hx(t_n) + \epsilon_n, \quad \epsilon_n \sim \mathcal{N}(0, R), \quad (1)$$

for $n = 0, \dots, N$ and matrices $H \in \mathbb{R}^{k \times d}$ and $R \in \mathbb{R}^{k \times k}$. Further suppose that $x(t)$ solves the ODE

$$\dot{x}(t) = f(x(t); u(t)), \quad (2)$$

and satisfies the initial condition $x(t_0) = x_0 \in \mathbb{R}^d$. The vector field $f : \mathbb{R}^d \times \mathbb{R}^\ell \rightarrow \mathbb{R}^d$ is assumed to be autonomous, which is no loss of generality (e.g. [22]) but simplifies the notation. The *latent force* $u : [t_0, t_{\max}] \rightarrow \mathbb{R}^\ell$ parametrizes f and shall be unknown.

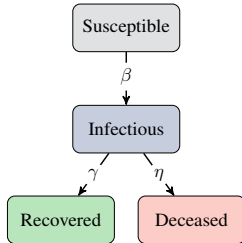


Figure 2: SIRD dynamics.

SIR-type models (e.g. [9]) are a common choice to describe the evolution of the COVID-19 pandemic. In SIR-type models, a population partitions into a discrete set of compartments. The differential equation then describes the transition of counts of individuals between these compartments. For example, the SIRD model [13] formulates the transitions between susceptible, infectious, recovered, and deceased people as

$$\begin{aligned} \dot{S}(t) &= -\beta(t)S(t)I(t)/P, & \dot{R}(t) &= \gamma I(t), \\ \dot{I}(t) &= \beta(t)S(t)I(t)/P - \gamma I(t) - \eta I(t), & \dot{D}(t) &= \eta I(t), \end{aligned} \quad (3)$$

governed by contact rate $\beta(t) : [t_0, t_{\max}] \rightarrow [0, 1]$, recovery rate $\gamma \in [0, 1]$, and mortality rate $\eta \in [0, 1]$ (Figure 2). S , I , R , and D evolve over time, but the total population P (as the sum of the compartments) is assumed to remain constant. In this context, the contact rate $\beta(t)$ is the latent force and varies over time (in the notation from Eq. (2), β is u). A time-varying contact rate provides a model for the impact of governmental measures on the dynamics of the pandemic. The experiments in Section 5 isolate the impact of the contact rate on the course of the infection counts, by assuming that γ and η are fixed and known. The method is by no means restricted to inference over a single latent force, as will also be shown in Section 5.1. In this SIRD setting, the goal is to infer an (approximate) joint posterior over $\beta(t)$ and the dynamics of $S(t)$, $I(t)$, $R(t)$, and $D(t)$ as well as to use the reconstructed dynamics to extrapolate into the future. Section 3 explains the conceptual details, Section 4 distinguishes the method from related work, and Section 5 evaluates the performance.

3 Method

This section explains how to infer the unknown process $u(t)$ and the ODE solution $x(t)$ in a single forward solve. Section 3.1 defines the prior model, Section 3.2 describes the probabilistic numerical ODE inference setup, and Section 3.3 describes approximate Gaussian filtering and smoothing in this context. Section 3.4 summarizes the resulting algorithm. The exposition of classic concepts here is necessarily compact. In-depth introductions can be found, e.g., in the book by Särkkä and Solin [34].

3.1 Gauss–Markov prior

Let $\nu \in \mathbb{N}$. Define two independent Gauss–Markov processes $U : [t_0, t_{\max}] \rightarrow \mathbb{R}^\ell$ and $X : [t_0, t_{\max}] \rightarrow \mathbb{R}^{d(\nu+1)}$ that solve the linear, time-invariant stochastic differential equations [29],

$$dU(t) = F_U U(t) dt + L_U dW_U(t), \quad dX(t) = F_X X(t) dt + L_X dW_X(t), \quad (4)$$

with drift matrices $F_U \in \mathbb{R}^{\ell \times \ell}$ and $F_X \in \mathbb{R}^{d(\nu+1) \times d(\nu+1)}$, as well as dispersion matrices $L_U \in \mathbb{R}^{\ell \times s}$ and $L_X \in \mathbb{R}^{d(\nu+1) \times d}$. $W_U : [t_0, t_{\max}] \rightarrow \mathbb{R}^s$ and $W_X : [t_0, t_{\max}] \rightarrow \mathbb{R}^d$ are Wiener processes. U and X satisfy the Gaussian initial conditions,

$$U(t_0) \sim \mathcal{N}(m_U, P_U), \quad X(t_0) \sim \mathcal{N}(m_X, P_X), \quad (5)$$

defined by $m_U \in \mathbb{R}^\ell$, $P_U \in \mathbb{R}^{\ell \times \ell}$, $m_X \in \mathbb{R}^{d(\nu+1)}$, and $P_X \in \mathbb{R}^{d(\nu+1) \times d(\nu+1)}$. $U(t)$ models the unknown function $u(t)$ and can be any Gauss–Markov process that admits a representation as the solution of a linear SDE with Gaussian initial conditions. $X(t) = (X^{(0)}(t), \dots, X^{(\nu)}(t)) \in \mathbb{R}^{d(\nu+1)}$ models the ODE dynamics, in light of which we require $X^{(i)}(t) = \frac{d^i}{dt^i} X^{(0)}(t) \in \mathbb{R}^d$, $i = 0, \dots, \nu$. In other words, the first element in $X(t)$ is an estimate for $x(t)$, the second element is an estimate for $\frac{d}{dt}x(t)$, et cetera. Encoding that the state X consists of a model for $x(t)$ as well as its first ν derivatives imposes structure on F_X and L_X (see e.g. [21]). Examples include the Matérn, integrated Ornstein-Uhlenbeck, and integrated Wiener processes; the canonical choice for probabilistic ODE solvers would be integrated Wiener processes [35, 39, 5, 22].

The class of Gauss–Markov priors inherits its wide generalizability from Gaussian process models; recall that Gauss–Markov processes like U and X are Gaussian processes with the Markov property. While not every Gaussian process with one-dimensional input space is Markovian, a large number of descriptions of Gauss–Markov processes emerge by translating a covariance function into an (approximate) SDE representation [34, Chapter 12]. For example, this applies to (quasi-)periodic, squared-exponential, or rational quadratic kernels; in particular, sums and products of Gauss–Markov processes admit a state-space representation [36, 34]. Recent research has considered approximate SDE representations of general Gaussian processes in one dimension [24]. With these tools, prior knowledge over U or X can be encoded straightforwardly into the model.

3.2 Two likelihoods: for observations and for the ordinary differential equation

A functional relationship between the processes $U(t)$, $X(t)$ and the data $y_{0:N}$ emerges by combining two likelihood functions: one for the observations $y_{0:N}$ (recall Equation (1)), and one for the ordinary differential equation. The present section formalizes both. Let $\mathcal{T} = \mathcal{T}_N^{\text{OBS}} \cup \mathcal{T}_M^{\text{ODE}}$ be the union of the observation-grid $\mathcal{T}_N^{\text{OBS}}$, which has been introduced in Section 2, and an ODE-grid $\mathcal{T}_M^{\text{ODE}} := (t_0^{\text{ODE}}, \dots, t_M^{\text{ODE}})$. The name ‘‘ODE-grid’’ expresses that this grid contains the locations on which the ODE information will enter the inference problem, as described below.

$\mathcal{T}_N^{\text{OBS}}$ contains the locations of $y_{0:N}$, in light of which the first of two observation models is

$$Y_n | X(t_n^{\text{OBS}}) \sim \mathcal{N}\left(HX^{(0)}(t_n^{\text{OBS}}), R\right), \quad (6)$$

for $n = 0, \dots, N$. This is a reformulation of the relationship between process x and observations $y_{0:N}$ in Eq. (1) in terms of X (instead of x , which is modeled by $X^{(0)}$). Including this first measurement model ensures that the inferred solution remains close to the data points. $\mathcal{T}_M^{\text{ODE}}$ contains the locations on which $U(t)$ connects to $X(t)$ through the ODE. Specifically, the set of random variables $Z_{0:M} \in \mathbb{R}^{(M+1) \times d}$, defined as

$$Z_m | X(t_m^{\text{ODE}}), U(t_m^{\text{ODE}}) \sim \delta\left(X^{(1)}(t_m^{\text{ODE}}) - f\left(X^{(0)}(t_m^{\text{ODE}}); U(t_m^{\text{ODE}})\right)\right), \quad (7)$$

where δ is the Dirac delta, describes the discrepancy between the current estimate of the derivative of the ODE solution (i.e. $X^{(1)}$) and its desired value (i.e. $f(X^{(0)}; U)$), as prescribed by the vector field f . If the random variables $Z_{0:M}$ realize small values everywhere, $X^{(0)}$ solves the ODE as parametrized by U . This motivates introducing artificial data points $z_{0:M} \in \mathbb{R}^{(M+1) \times d}$ that are equal to zero, $z_m = 0 \in \mathbb{R}^d$, $m = 0, \dots, M$. Conditioning the stochastic processes X and U on attaining this (artificial) zero data ensures that the inferred solution follows ODE dynamics throughout the domain. Figure 3 shows the discretized state-space model.

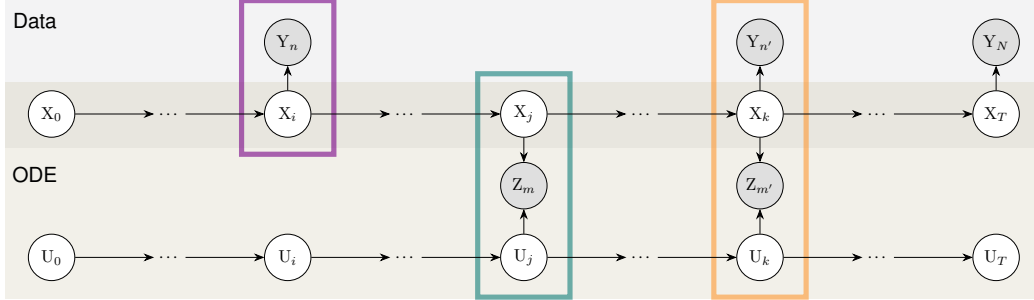


Figure 3: **Instance of the described state-space model, visualized as a directed graphical model.** Shaded variables are observed. Either *only data*, *only mechanistic knowledge*, or *both sources of information* can be conditioned on during inference (recall Figure 1).

3.3 Approximate inference with an extended Kalman filter

Both X and U enter the likelihood in Eq. (7) through a possibly non-linear vector field f . Therefore, the posterior distribution (recall $z_{0:M} = 0$)

$$p(U(t), X(t) \mid Z_{0:M} = z_{0:M}, Y_{0:N} = y_{0:N}) \quad (8)$$

is intractable, but can be approximated efficiently. Even though the problem is discretized, the posterior distribution is continuous [34, Chapter 10]. There are mainly two approaches to computing a tractable approximation of the intractable posterior distribution in Eq. (8): approximate Gaussian filtering and smoothing [33], which computes a cheap, Gaussian approximation of this posterior, and sequential Monte Carlo methods [28], whose approximate posterior may be more descriptive, but also more expensive to compute. Like the literature on probabilistic ODE solvers [39, 5], this work uses approximate Gaussian filtering and smoothing techniques for their low computational complexity.

The continuous-discrete state-space model inherits its non-linearity from the ODE vector field f . Linearizing f with a first-order Taylor series expansion creates a tractable inference problem; more specifically, it gives rise to the extended Kalman filter (EKF) [16, 26]. Loosely speaking, if the random variable Z is large in magnitude, then X and U are poor estimates for the ODE and its parameter. An EKF update, based on the first-order linearization of f , approximately corrects this misalignment. If sufficiently many ODE measurements $z_{0:M}$ are available, a sequence of such updates preserves sensible ODE dynamics over time. An alternative to a Taylor-series linearization is the unscented transform, which yields the unscented Kalman filter [41, 18]. The computational complexity of both algorithms is linear in the number of grid points and cubic in the dimension of the state-space. Detailed implementation schemes can be found, for instance, in the book by Särkkä [33].

The EKF approximates the filtering distribution

$$p(U(t), X(t) \mid Z_{0:m} = z_{0:m}, Y_{0:n} = y_{0:n}, \text{ such that } t_m^{\text{ODE}}, t_n^{\text{OBS}} \leq t). \quad (9)$$

It describes the current state of the system given all the previous measurements and allows updates in an online fashion as soon as new observations emerge. If desired, the Rauch-Tung-Striebel smoother turns the filtering distribution into an approximation of the full (smoothing) posterior (in Eq. (8)). In doing so, all observations – that is, measurements according to both Eq. (6) and Eq. (7) – are taken into account for the posterior distribution at each location t . As special cases, this setup recovers: (i) a Kalman filter/Rauch-Tung-Striebel smoother [19] if the ODE likelihood (Eq. (7)) is omitted; (ii) an ODE solver [39], if the data likelihood (Eq. (6)) is omitted. In the present setting, however, both likelihoods play an important role.

3.4 Algorithm and implementation

The procedure is summarized in Algorithm 1. The prediction step is determined by the prior and is available in closed-form (Appendix A.2). At times at which data is observed according to the linear Gaussian measurement model in Eq. (6), the update step follows the rules of the standard Kalman filter. Before updating on pseudo-observations according to the ODE likelihood (Eq. (7)), the non-linear measurement model is linearized at the predicted mean. More details are provided

Algorithm 1 Compute the filtering distribution by conditioning on both $y_{0:N}$ and $z_{0:M}$.

Input: data $y_{0:N}$, time grid $\mathcal{T} = \mathcal{T}_N^{\text{OBS}} \cup \mathcal{T}_M^{\text{ODE}}$, vector field f , m_X, P_X, m_U, P_U [Eq. (9)]
Output: Filtering distribution [Eq. (9)]
Initialize $X_0 = \mathcal{N}(m_X, P_X)$ and $U_0 = \mathcal{N}(m_U, P_U)$ [Eq. (5)]
for $t_j \in \mathcal{T}$ **do**
 Predict X_j from X_{j-1} and predict U_j from U_{j-1}
 if $t_j \in \mathcal{T}_N^{\text{OBS}}$ **then** update X_j on y_j **end if** [Eq. (6)]
 if $t_j \in \mathcal{T}_M^{\text{ODE}}$ **then** linearize measurement model and update X_j and U_j on z_j **end if** [Eq. (7)]
end for

in Appendix A. The filtering distribution can be turned into a smoothing posterior by running a backwards-pass with a Rauch-Tung-Striebel smoother (e.g. [33]).

The computational cost of obtaining either, the filtering or the smoothing posterior, are both linear in the number of grid points and cubic in the dimension of the state-space, i.e. $\mathcal{O}((N + M)(d^3\nu^3 + \ell^3))$. Only a single forward-backward pass is required. If desired, the approximate Gaussian posterior can be refined iteratively by means of posterior linearization and iterated Gaussian filtering and smoothing, which yields the maximum-a-posteriori (MAP) estimate [3, 38]. The experiments presented in Section 5 show how a single forward-backward pass already approximates the MAP estimate accurately.

4 Related work

Latent forces and ODE solvers: The explained method closely relates to probabilistic ODE solvers and latent force models [44], especially the kind of latent force model that exploits the state-space formulation of the prior [12]. The difference is that, in the spirit of probabilistic numerical algorithms, the mechanistic knowledge in the form of an ODE is injected through the likelihood function instead of the prior. A similar approach of linking observations to mechanistic constraints has previously been used in the literature on constrained Gaussian processes [17] and gradient matching [6, 43]. Probabilistic ODE solvers have been used by Kersting et al. [20] for efficient ODE inverse problem algorithms, but their approach is different to the present algorithm, in which the need for iterated optimization or sampling is avoided altogether.

Monte Carlo methods: (Markov-chain) Monte Carlo methods are also able to infer a time-dependent ODE latent force from a set of state observations. Options that are compatible with a setup similar to the present work would include sequential Monte Carlo techniques [28], elliptical slice sampling [27], or Hamiltonian Monte Carlo [4] (for instance realized as the No-U-Turn sampler [15]). The shared disadvantage of Monte Carlo methods applied to the resulting ODE inverse problem is that the complexity of obtaining *a single* Monte Carlo sample is of the same order of magnitude as computing the *full* Gaussian approximation of the posterior distribution. In Appendix B we show results from a parametric version of the SIRD-latent force model (using the No-U-Turn sampler as provided by NumPyro [30]). This sampler requires *thousands* of numerical ODE solutions, compared to the single solve of our method. This fact is also reflected in the wall-clock time needed for both types of inference. While the MCMC experiment in Appendix B takes in the order of hours, each experiment with our approach takes under one minute to complete. In other words, the algorithm in the present work poses an efficient yet expressive alternative to Monte Carlo methods for approximate inference with dynamical systems.

5 Experiments

This section describes three blocks of experiments. The implementation is based on ProbNum [42] and all experiments use a conventional, consumer-level CPU. First, a range of artificial datasets is generated by sampling ODE parameters from a prior state-space model and simulating a solution of the corresponding ODE. Inference in such a controlled environment allows comparing to the ground truth, thereby assessing the quality of the approximate inference. We consider three ODE models to this end. Second, a COVID-19 dataset will probe the predictive performance of the probabilistic model and the resulting approximate posterior distribution. Third, some changes to the model from the COVID-19 experiments, for instance, ensuring that the number of case counts must be positive,

will improve the interpretability (for example, of the credible intervals). Controlling the range of values that the prior state-space can realize introduces additional non-linearity into the model – which can also be locally approximated by the EKF – and makes the solution more physically meaningful.

5.1 Simulated environments

As a first test for the capabilities of the proposed method, we consider three simulated environments. To this end, the training data is generated as follows. The starting point is always an initial value problem with dynamics defined by a vector field f and a Gauss–Markov prior over the dynamics x and the unknown parameters u of the vector field. Then, (i) we sample the time-varying parameter trajectories from the Gauss–Markov prior; (ii) we solve the ODE, as parametrized by the sampled trajectories from (i), using LSODA [14] with adaptive step sizes using SciPy [40]; (iii) we subsample the ground-truth solution on a uniform grid (which will become $\mathcal{T}_N^{\text{OBS}}$) to generate artificial state observations $y_{0:N}$; (iv) we add Gaussian i.i.d. noise to the observations.

The procedure described above generates both a ground truth to compare to and a noisy, artificially observed data set. Given such a set of observations, Algorithm 1 computes a posterior distribution over the true trajectories under appropriate model assumptions. In this posterior, we look for the proximity of the mean estimate to the underlying ground truth; the closer, the better. We measure this proximity in the root-mean-square error. Furthermore, the width of the posterior (expressed by the posterior covariance) should deliver an appropriate quantification of the mismatch. We report the χ^2 -statistic [2], which suggests that the posterior distribution is well-calibrated if the χ^2 -statistic is close to the dimension d of the ground truth. Three mechanistic models serve as examples.

Van-der-Pol: The first of three test problems is the van-der-Pol oscillator [11]. It has one parameter μ (sometimes referred to as a stiffness constant, because for large μ , the van-der-Pol system is stiff). As a prior over the dynamics we choose a twice-integrated Wiener process with diffusion intensity $\sigma_X^2 = 300$. The stiffness parameter μ is modeled as a Matérn- $\frac{3}{2}$ process with lengthscale $\ell_U = 10$ and diffusion intensity $\sigma_U^2 = 0.3$. The posterior is computed on a grid from $t_0 = 0$ to $t_{\max} = 25$ units of time with step size $\Delta t = 0.025$.

Lotka-Volterra: The Lotka-Volterra equations [25] describe the change in the size of two populations, predators and prey. There are four parameters, which we call a , b , c , and d , which describe the interaction and death/reproduction rates of the populations. As a prior over the dynamics we choose a twice-integrated Wiener process with diffusion intensity $\sigma_X^2 = 10$. The four parameters are modeled as Matérn- $\frac{3}{2}$ processes with lengthscales $\ell_{U_a} = \ell_{U_b} = \ell_{U_c} = \ell_{U_d} = 40$. The diffusion intensities are $\sigma_{U_a}^2 = \sigma_{U_c}^2 = 0.01$ and $\sigma_{U_b}^2 = \sigma_{U_d}^2 = 0.001$. The posterior is computed on a grid from $t_0 = 0$ to $t_{\max} = 60$ units of time with step size $\Delta t = 0.1$.

SIRD: As detailed in Section 2, the SIRD model is governed by a contact rate $\beta(t)$. Recall that we assume a time-dependent β to account for governmental measures in reaction to the spread of COVID-19. The recovery rate γ and fatality rate η are fixed at $\gamma = 0.06$ and $\eta = 0.002$, like they will be in the experiments with real data in Sections 5.2 and 5.3 below. As a prior over the dynamics we choose a twice-integrated Wiener process with diffusion intensity $\sigma_X^2 = 50$. The contact rate β is modeled as a Matérn- $\frac{3}{2}$ process with lengthscale $\ell_U = 14$ and diffusion intensity $\sigma_U^2 = 0.1$. The posterior is computed on a grid from $t_0 = 0$ to $t_{\max} = 100$ units of time with step size $\Delta t = 0.1$.

The model allows for straightforward restriction of parameter values by using link functions. The natural support for the SIRD-contact rate is the interval $[0, 1]$, but $U(t)$, as a Gauss–Markov process, takes values on the real line. A change in the basis of $\beta(t)$ with a logistic sigmoid function ϑ before it enters the likelihood fixes this misspecification. Similarly, the Lotka-Volterra parameters are inferred in log-space to ensure positivity. It is an appealing aspect of the EKF that these non-linear transformations do not require significant adaptation of the algorithm. Instead, the EKF treats it as merely another level of linearization of Eq. (7). Section 5.3 extends this to the state dynamics.

The results are shown in Figure 4. On all test problems, the algorithm recovers the true states and the true latent force accurately. The recovery is not exact, which shows how the Gaussian posterior is only an approximation of the true posterior. The χ^2 -statistic for the van-der-Pol stiffness parameter μ is 1.11, which lies in $(0.0039, 3.8415)$, the 90% confidence interval of the χ^2 distribution with 1 degree of freedom. The root-mean-square error (RMSE) to the truth is 0.14. The χ^2 -statistic for the Lotka-Volterra parameters is 8.06, which lies in $(0.7107, 9.4877)$, the 90% confidence interval of the χ^2 distribution with 4 degrees of freedom. The RMSE to the truth is 0.04 in log space and 0.018 in

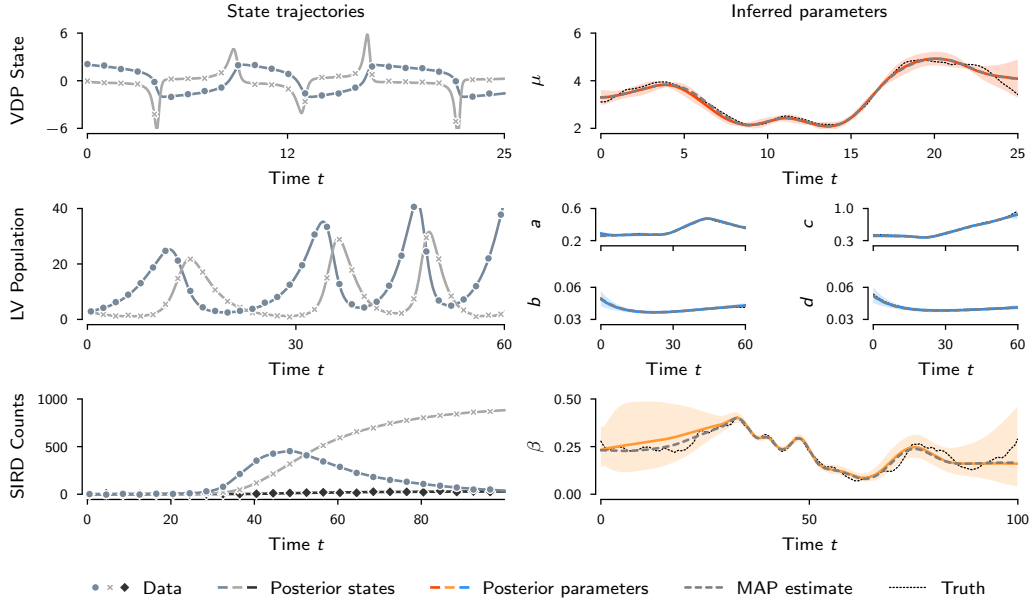


Figure 4: **State recovery in simulated environments.** The stiffness parameter of the van-der-Pol oscillator (top row) and the Lotka-Volterra parameters (middle row) are inferred accurately with appropriately high certainty. For the SIRD experiment (bottom row), the uncertainty is high, where low case counts provide little information about the latent contact rate. With more fluctuations in the observed counts, the approximated contact rate displays more certainty.

linear space. The χ^2 -statistic for the contact rate β is 0.91, which lies in $(0.0039, 3.8415)$, the 90% confidence interval of the χ^2 distribution with 1 degree of freedom. The RMSE to the truth is 0.2 in logit space and 0.033 in linear space.

5.2 COVID-19 data

We continue with the SIRD model introduced in Eq. (3), now using data collected in Germany over the period from January 22, 2020, to May 27, 2021. Throughout the pandemic, the German government has imposed mitigation measures of varying severity. Together with seasonal effects, summer vacations, etc., they caused a continual change in the contact rate. The next experiments aim to recover said contact rate (and the SIRD counts) from the German dataset.

The Center for Systems Science and Engineering at the Johns Hopkins University publishes daily counts of confirmed ($y_n^{\text{confirmed}}$), recovered ($y_n^{\text{recovered}}$), and deceased (y_n^{deceased}) individuals [7]. One can transform this data to suit the SIRD model

$$I_n := y_n^{\text{confirmed}} - R_n - D_n, \quad R_n := y_n^{\text{recovered}}, \quad D_n := y_n^{\text{deceased}}. \quad (10)$$

The counts I_n , R_n , and D_n are available for each day, starting with January 22, 2020. Assuming a constant population over time, the numbers of susceptible individuals S_n are always evident from the other quantities, thus left out of the visualizations. We fix the population at $P = 83\,783\,945$, based on public record. We rescale the data to cases per one thousand people (CPT).

As a prior over $X(t)$, due to its popularity in constructing probabilistic ODE solvers [39], we assume a twice-integrated Wiener process. $\beta(t)$ is modelled as a Matérn-3/2 process with length scale $\ell_q = 75$ and diffusion intensity $\sigma_q^2 = 0.05$. The state-space model is straightforwardly extendable to sums and products of (more) processes [36, 34]. Inferring parameters that are constant over time, however, is not straightforward due to potentially singular transition models [33, Section 12.3.1].

As described in Section 5.1, the contact rate is inferred in logit space. We shift the logistic sigmoid function such that it fulfills $\vartheta(0) = 0.1$ in which case the stationary mean $\bar{U} = 0$ translates to a stationary mean $\vartheta(\bar{U}) = \bar{\beta} = 0.1$ of the Matérn process that models the contact rate. The recovery

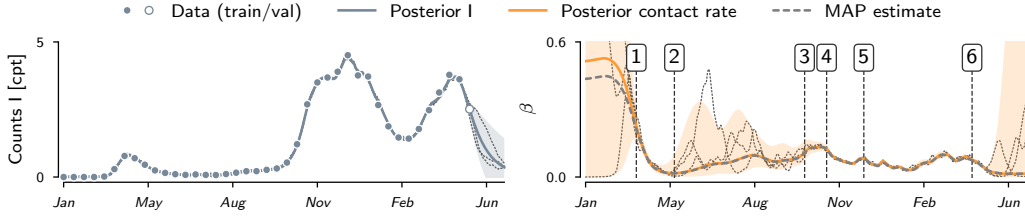


Figure 5: **Estimated counts of infectious cases and contact rate based on real COVID-19 data.** The case counts of infectious people are scaled to cases per thousand (cpt). The uncertainty over the contact rate increases when the case counts are low. After a single forward solve, the inferred mean is already close to the MAP estimate. The shaded areas show the 95 % credible interval and the dotted black lines are samples from the posterior.

Table 1: List of selected governmental measures imposed in Germany with the aim to contain the spread of COVID-19. These events are depicted in Figures 5 and 6 (see column ‘Mark’). Links to the sources are provided in Appendix C.

| Mark | Governmental Measures |
|-------|--|
| 1 | Stringent contact restrictions, partial shutdown of public life |
| 2 - 3 | Continual relaxations of measures |
| 4 | Partial shutdown of public life (‘lockdown light’) |
| 5 | Hard lockdown, stringent contact restrictions |
| 6 | First nationwide decree of restrictions, increased intensification of measures |

rate and mortality rate are considered known and fixed at $\gamma = 0.06$ and $\eta = 0.002$ to isolate the effect of the inference procedure on recovering the evolution of the contact rate $U(t) = \beta(t)$. We set the mean of the Gaussian initial conditions to the first data point that is available. The diffusion intensity of the prior process $X(t)$ is set to $\sigma_X^2 = 10$. The latent process U and all derivatives are initialized at zero. Note that due to the logistic sigmoid transform, an initial value $U_0 = 0$ amounts to an initial contact rate $\beta_0 = 0.1$.

In the present scenario, we cannot take the SIRD model as an accurate description of the underlying data but merely as a tool that aids the inference engine in recovering physically meaningful states and forces. In order to account for this model mismatch, the Dirac likelihood from Eq. (7) is relaxed towards a Gaussian likelihood with measurement noise $\lambda^2 = 0.01$. This equals the data observation noise and thus balances the respective impact of either (misspecified) source of information. Intuitively, adding ODE measurement noise reduces how strictly the vector field dynamics are enforced during inference and therefore avoids overconfident estimates of $\beta(t)$.

The mesh-size of the ODE is $\Delta t = 1/24$ days, i.e. ODE updates are computed on an hourly basis. The final 14 observations are excluded from the training set to serve as validation data for evaluating the extrapolation behavior of the proposed method. Figure 5 shows the results. The mean of the state X estimates the case counts accurately in both interpolation and extrapolation tasks. The estimated contact rate rapidly decreases around late March, remains low until fall, increases momentarily, and is dampened again soon after. This aligns with a set of political measures imposed by the government (compare Figure 5 to Table 1). The uncertainty over the estimated contact rate is high in the early beginning when the case counts are still low. It then increases again in summer and with the beginning of the extrapolation phase.

If the experiment is taken as-is, the credibility intervals of the posterior over $X(t)$ include negative numbers (mostly where the case counts are low and the uncertainty high, and when extrapolating). Of course, in a system that models counts of people in different stages of a disease, negative numbers should be excluded altogether. The proposed method provides straightforward means to address this issue. Section 5.3 explains the details.

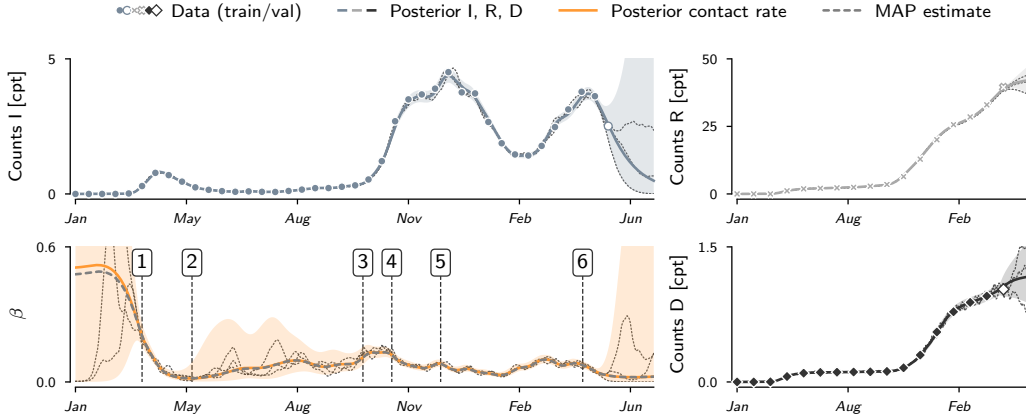


Figure 6: **Estimated case counts and contact rate**, inferred in the logarithmic basis on real COVID-19 and vaccination data. The case counts of infectious people are scaled to cases per thousand (cpt). Again, the uncertainty of the contact rate increases where the case counts are low. Now, the posterior credible interval is restricted to the positive reals. The shaded areas show the 95 % credible interval and the dotted black lines are samples from the posterior.

5.3 Non-negative state estimates

The following experiment evaluates how the proposed method performs in combination with a state-space model that constrains the support of the dynamics. Concretely, let $X(t)$ model the logarithm of the SIRD dynamics and the respective derivatives. With a slight abuse of notation, we will continue writing “ X ” even though it lives in a different space than in the previous sections. The structure of the dynamic model is the same. The diffusion intensity of the prior process $X(t)$ is $\sigma_X^2 = 0.05$. The diffusion is not comparable to the value in the previous section because the state dynamics moved to log-space. Using $\frac{d}{dt} \exp(x(t)) = \exp(x(t))\dot{x}(t)$, the ODE likelihood becomes

$$Z_m | X_m^{\text{ODE}}, U_m^{\text{ODE}} \sim \mathcal{N}(\zeta_1 - f(\zeta_2; \zeta_3), \lambda^2 I_d), \quad (11)$$

with auxiliary quantities (recall the logistic sigmoid ϑ)

$$\zeta_1 := \exp\left(X^{(0)}(t_m^{\text{ODE}})\right) X^{(1)}(t_m^{\text{ODE}}), \quad \zeta_2 := \exp\left(X^{(0)}(t_m^{\text{ODE}})\right), \quad \zeta_3 := \vartheta(U(t_m^{\text{ODE}})). \quad (12)$$

The exponential function introduces an additional non-linearity into the state-space model, which necessitates smaller step-sizes for the ODE measurements (see below).

The observed case count data $y_{0:N}$ is transformed into the log-space, too, in which we assume additive, i.i.d. Gaussian noise. On the one hand, transforming the measurements into log-space implies that the measurement model for the counts remains linear; on the other hand, it imposes a log-normal noise model (if viewed back in “linear space”). Log-normal noise underlines how the estimated states cannot be negative. Again, we scale the counts to cases per thousand.

As depicted in Figure 6, the reconstruction of the driving processes in this setting yields results that at first glance, look similar to the previous experiment. The states match the data points well. However, the extrapolation is more realistic in that the credible intervals encode that negative values are impossible (which is due to the log-transform). The mean of the recovered contact rate closely resembles the estimate of the previous experiment. Again, upon implementation of strict governmental measures, the uncertainty decreases, whereas in the context of relaxations, the uncertainty is high.

6 Statement on Societal Impact

This work performs methods research to develop an efficient numerical algorithm to infer latent forces governing ordinary differential equations. As a testbed, we use data from the COVID-19 pandemic. We do so to motivate and visualize the practical value of our methods. The results of this

algorithm, however, should not be taken as policy advice. The model used in the paper is deliberately simplistic. The presented work therefore should not be misunderstood as epidemiological research. The machine learning community has, over time, frequently used data of contemporary societal concern to motivate and test new algorithmic concepts (well-known examples from the UCI collection include the Wisconsin Breast Cancer Dataset, the mushroom classification dataset, and the German credit data set). Our work follows in this line. Of course, if this algorithm, or any competitor, is used to derive policy advice, the underlying differential equation and latent states must be carefully considered by domain experts, which we are not.

7 Conclusion

By coupling mechanistic and data-driven inference so directly, the algorithm builds on the core premise of probabilistic numerics – that computation itself is a data source that does not differ, formally, from observational data. Information from observations and mechanistic knowledge (in the form of an ODE) can thus be described in the same language of Bayesian filtering and smoothing. This removes the need for an outer loop over multiple forward solves and thus drastically reduces the computational cost. Our experimental evaluation corroborates that the resulting approximate posterior is close to the ground truth and drastically reduces computational cost over Monte Carlo alternatives. It faithfully captures multiple sources of uncertainty from the data, numerical (discretization) error, and epistemic uncertainty about the mechanism. We hope this framework helps empower practitioners, not just by reducing computational burden but also by providing a more flexible modelling platform.

Acknowledgements

The authors gratefully acknowledge financial support by the European Research Council through ERC StG Action 757275 / PANAMA; the DFG Cluster of Excellence “Machine Learning - New Perspectives for Science”, EXC 2064/1, project number 390727645; the German Federal Ministry of Education and Research (BMBF) through the Tübingen AI Center (FKZ: 01IS18039A); and funds from the Ministry of Science, Research and Arts of the State of Baden-Württemberg. The authors thank the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for supporting N. Krämer. Moreover, the authors thank Nathanael Bosch and Marvin Pförtner for valuable discussions.

References

- [1] P. Axelsson and F. Gustafsson. Discrete-time solutions to the continuous-time differential Lyapunov equation with applications to Kalman filtering. *IEEE Transactions on Automatic Control*, 60(3):632–643, 2015.
- [2] Y. Bar-Shalom, X. R. Li, and T. Kirubarajan. *Estimation With Applications to Tracking and Navigation: Theory Algorithms and Software*. John Wiley & Sons, 2004.
- [3] B. M. Bell. The iterated Kalman smoother as a Gauss–Newton method. *SIAM Journal on Optimization*, 4(3):626–636, 1994.
- [4] M. Betancourt. A conceptual introduction to Hamiltonian Monte Carlo. *arXiv:1701.02434*, 2017.
- [5] N. Bosch, P. Hennig, and F. Tronarp. Calibrated adaptive probabilistic ODE solvers. In *International Conference on Artificial Intelligence and Statistics*, 2021.
- [6] B. Calderhead, M. Girolami, and N. Lawrence. Accelerating Bayesian inference over nonlinear differential equations with Gaussian processes. In *Advances in Neural Information Processing Systems*, 2009.
- [7] E. Dong, H. Du, and L. Gardner. An interactive web-based dashboard to track COVID-19 in real time. *The Lancet Infectious Diseases*, 20(5):533–534, 2020.
- [8] J. Dormand and P. Prince. A family of embedded Runge-Kutta formulae. *Journal of Computational and Applied Mathematics*, 6(1):19–26, 1980.

- [9] G. Giordano, F. Blanchini, R. Bruno, P. Colaneri, A. Di Filippo, A. Di Matteo, and M. Colaneri. Modelling the COVID-19 epidemic and implementation of population-wide interventions in Italy. *Nature Medicine*, pages 1–6, 2020.
- [10] M. Grewal and A. Andrews. *Kalman Filtering: Theory and Practice Using MATLAB*. Wiley, 2011.
- [11] J. Guckenheimer. Dynamics of the van der Pol equation. *IEEE Transactions on Circuits and Systems*, 27(11):983–989, 1980.
- [12] J. Hartikainen, M. Seppänen, and S. Särkkä. State-space inference for non-linear latent force models with application to satellite orbit prediction. In *International Conference on Machine Learning*, 2012.
- [13] H. W. Hethcote. The mathematics of infectious diseases. *SIAM Review*, 42(4):599–653, 2000.
- [14] A. Hindmarsh and L. Petzold. LSODA, ordinary differential equation solver for stiff or non-stiff system. 2005.
- [15] M. D. Hoffman and A. Gelman. The No-U-Turn sampler: adaptively setting path lengths in Hamiltonian Monte Carlo. *Journal of Machine Learning Research*, 15(1):1593–1623, 2014.
- [16] A. H. Jazwinski. *Stochastic Processes and Filtering Theory*. Academic Press, 1970.
- [17] C. Jidling, N. Wahlström, A. Wills, and T. B. Schön. Linearly constrained gaussian processes. In *Advances in Neural Information Processing Systems*, 2017.
- [18] S. J. Julier and J. K. Uhlmann. Unscented filtering and nonlinear estimation. *Proceedings of the IEEE*, 92(3):401–422, 2004.
- [19] R. E. Kalman. A new approach to linear filtering and prediction problems. *Transactions of the ASME—Journal of Basic Engineering*, 82(Series D):35–45, 1960.
- [20] H. Kersting, N. Krämer, M. Schiegg, C. Daniel, M. Tiemann, and P. Hennig. Differentiable likelihoods for fast inversion of ‘likelihood-free’ dynamical systems. In *International Conference on Machine Learning*, 2020.
- [21] H. Kersting, T. J. Sullivan, and P. Hennig. Convergence rates of Gaussian ODE filters. *Statistics and Computing*, 30(6):1791–1816, 2020.
- [22] N. Krämer and P. Hennig. Stable implementation of probabilistic ODE solvers. *arXiv:2012.10106*, 2020.
- [23] D. C. Liu and J. Nocedal. On the limited memory BFGS method for large scale optimization. *Mathematical Programming*, 45(1-3):503–528, 1989.
- [24] J. Loper, D. M. Blei, J. P. Cunningham, and L. Paninski. General linear-time inference for Gaussian processes on one dimension. *arXiv:2003.05554*, 2020.
- [25] A. J. Lotka. The growth of mixed populations: two species competing for a common food supply. In *The Golden Age of Theoretical Ecology: 1923–1940*, pages 274–286. Springer, 1978.
- [26] P. S. Maybeck. *Stochastic Models, Estimation, and Control*. Academic Press, 1982.
- [27] I. Murray, R. Adams, and D. MacKay. Elliptical slice sampling. In *International Conference on Artificial Intelligence and Statistics*, 2010.
- [28] C. A. Naesseth, F. Lindsten, and T. B. Schön. Elements of sequential Monte Carlo. *Foundations and Trends® in Machine Learning*, 12(3):307–392, 2019.
- [29] B. Øksendal. *Stochastic Differential Equations*. Springer, 2003.
- [30] D. Phan, N. Pradhan, and M. Jankowiak. Composable effects for flexible and accelerated probabilistic programming in NumPyro. *arXiv:1912.11554*, 2019.

- [31] A. Rahimi and B. Recht. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems*, 2008.
- [32] C. Rasmussen and C. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- [33] S. Särkkä. *Bayesian Filtering and Smoothing*. Cambridge University Press, 2013.
- [34] S. Särkkä and A. Solin. *Applied Stochastic Differential Equations*. Cambridge University Press, 2019.
- [35] M. Schober, S. Särkkä, and P. Hennig. A probabilistic model for the numerical solution of initial value problems. *Statistics and Computing*, 29:99–122, 2019.
- [36] A. Solin and S. Särkkä. Explicit link between periodic covariance functions and state space models. In *International Conference on Artificial Intelligence and Statistics*, 2014.
- [37] R. Stengel. *Optimal Control and Estimation*. Dover Publications, 1994.
- [38] F. Tronarp, Á. F. García-Fernández, and S. Särkkä. Iterative filtering and smoothing in nonlinear and non-Gaussian systems using conditional moments. *IEEE Signal Processing Letters*, 25(3): 408–412, 2018.
- [39] F. Tronarp, H. Kersting, S. Särkkä, and P. Hennig. Probabilistic solutions to ordinary differential equations as nonlinear Bayesian filtering: a new perspective. *Statistics and Computing*, 29(6): 1297–1315, 2019.
- [40] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, et al. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nature Methods*, 17(3):261–272, 2020.
- [41] E. A. Wan and R. Van Der Merwe. The unscented Kalman filter for nonlinear estimation. In *IEEE 2000 Adaptive Systems for Signal Processing, Communications, and Control Symposium*, pages 153–158, 2000.
- [42] J. Wenger, N. Krämer, M. Pförtner, J. Schmidt, N. Bosch, N. Effenberger, J. Zenn, A. Gessner, T. Karvonen, F.-X. Briol, et al. ProbNum: Probabilistic numerics in Python. *arXiv preprint arXiv:2112.02100*, 2021.
- [43] P. Wenk, G. Abbati, M. A. Osborne, B. Schölkopf, A. Krause, and S. Bauer. ODIN: ODE-informed regression for parameter and state inference in time-continuous dynamical systems. *AAAI Conference on Artificial Intelligence*, 34(04):6364–6371, 2020.
- [44] M. Álvarez, D. Luengo, and N. D. Lawrence. Latent force models. In *International Conference on Artificial Intelligence and Statistics*, 2009.

A Implementation details

This section provides detailed information about the state-space model and approximate Gaussian inference therein. Appendix A.1 defines the augmented state-space model that formalizes the dynamics of the Gauss–Markov processes introduced in Section 3.1. Appendix A.2 provides the equations for prediction and update steps of the extended Kalman filter in such a setup, which is described in Section 3.4 (in particular, Algorithm 1).

A.1 Augmented state-space model

Section 3 describes the joint inference of both a latent process $u(t) : [t_0, t_{\max}] \rightarrow \mathbb{R}^l$ that parametrizes an ODE and $x(t) : [t_0, t_{\max}] \rightarrow \mathbb{R}^d$, the solution of said ODE. The dynamics of the processes are modeled by the stochastic differential equation

$$d \begin{pmatrix} \mathbf{U}(t) \\ \mathbf{X}(t) \end{pmatrix} = \underbrace{\begin{pmatrix} F_U & 0 \\ 0 & F_X \end{pmatrix}}_{=:F} \begin{pmatrix} \mathbf{U}(t) \\ \mathbf{X}(t) \end{pmatrix} dt + \underbrace{\begin{pmatrix} L_U & 0 \\ 0 & L_X \end{pmatrix}}_{=:L} d \begin{pmatrix} \mathbf{W}_U(t) \\ \mathbf{W}_X(t) \end{pmatrix}, \quad (\text{A.1})$$

with Gaussian initial conditions

$$\begin{pmatrix} \mathbf{U}(t_0) \\ \mathbf{X}(t_0) \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} m_U(t_0) \\ m_X(t_0) \end{pmatrix}, \begin{pmatrix} P_U(t_0) & 0 \\ 0 & P_X(t_0) \end{pmatrix} \right). \quad (\text{A.2})$$

The block-diagonal structure is due to the independent dynamics of the prior processes. The *drift matrices* F_U and F_X , as well as the *dispersion matrices* L_U and L_X depend on the choice of the respective processes U and X. The measurement models are given in Eq. (6) (for observed data) and in Eq. (7) (for ODE measurements).

In the experiments presented in Sections 5.2 and 5.3 we model the latent contact rate $\beta(t)$ as a Matérn- $3/2$ process with characteristic length scale ℓ_q . Hence,

$$d\mathbf{U}(t) = \underbrace{\begin{pmatrix} 0 & 1 \\ -(\sqrt{3}/\ell_q)^2 & -2\sqrt{3}/\ell_q \end{pmatrix}}_{F_U} \mathbf{U}(t) dt + \underbrace{\begin{pmatrix} 0 \\ 1 \end{pmatrix}}_{L_U} d\mathbf{W}_U(t). \quad (\text{A.3})$$

The SIRD counts are modeled as the twice-integrated Wiener process

$$d\mathbf{X}(t) = \underbrace{\begin{pmatrix} 0 & I_d & 0 \\ 0 & 0 & I_d \\ 0 & 0 & 0 \end{pmatrix}}_{F_X} \mathbf{X}(t) dt + \underbrace{\begin{pmatrix} 0 \\ 0 \\ I_d \end{pmatrix}}_{L_X} d\mathbf{W}_X(t), \quad (\text{A.4})$$

such that $\mathbf{X} = (\mathbf{X}^{(0)}, \mathbf{X}^{(1)}, \mathbf{X}^{(2)})^\top$ models the SIRD counts and the first two derivatives. Notice that $F_X \in \mathbb{R}^{d(\nu+1) \times d(\nu+1)}$ and $L_X \in \mathbb{R}^{d(\nu+1) \times d}$ are block matrices. I_d denotes the $d \times d$ identity matrix. In the context of the experiments, $d = 4$ (S, I, R, and D) and $\nu = 2$ (*twice-integrated* Wiener process). More details on the use of integrated Wiener processes in probabilistic ODE solvers can be found in, for instance, the work by Kersting et al. [21].

A.2 Kalman filter equations

This section is concerned with the exact steps that make up the algorithm summarized in Section 3.4. The stochastic differential equation defined in Eq. (A.1) formalizes the dynamics of the processes $\mathbf{U}(t)$ and $\mathbf{X}(t)$ that model $u(t)$ and $x(t)$, respectively. Define $\Delta t := t_j - t_{j-1} > 0$ for all $t_j = t_1, \dots, t_{\max}$. The *transition densities* of U and X are [10]

$$\mathbf{U}(t + \Delta t) \mid \mathbf{U}(t) \sim \mathcal{N}(\Phi_U(\Delta t)\mathbf{U}(t), Q_U(\Delta t)), \quad (\text{A.5a})$$

$$\mathbf{X}(t + \Delta t) \mid \mathbf{X}(t) \sim \mathcal{N}(\Phi_X(\Delta t)\mathbf{X}(t), Q_X(\Delta t)), \quad (\text{A.5b})$$

where transition matrices $\Phi_U(\Delta t) \in \mathbb{R}^{\ell \times \ell}$ and $\Phi_X(\Delta t) \in \mathbb{R}^{d(\nu+1) \times d(\nu+1)}$, as well as the process noise covariances $Q_U(\Delta t) \in \mathbb{R}^{\ell \times \ell}$ and $Q_X(\Delta t) \in \mathbb{R}^{d(\nu+1) \times d(\nu+1)}$ are available in closed form and can be computed, for instance, with matrix fraction decomposition [37, 1].

Define the transition matrix and process noise covariance of the process in Eq. (A.1) as

$$\Phi(\Delta t) := \begin{pmatrix} \Phi_U(\Delta t) & 0 \\ 0 & \Phi_U(\Delta t) \end{pmatrix}, \quad Q(\Delta t) := \begin{pmatrix} Q_U(\Delta t) & 0 \\ 0 & Q_U(\Delta t) \end{pmatrix}. \quad (\text{A.6})$$

Further, let

$$\begin{pmatrix} U(t_j) \\ X(t_j) \end{pmatrix} \sim \mathcal{N}(m_j, P_j), \quad (\text{A.7})$$

for time points $t_j \in \mathcal{T} = \mathcal{T}^{\text{OBS}} \cup \mathcal{T}^{\text{ODE}}$. The predicted mean and covariance m_j^- and P_j^- are

$$m_j^- = \Phi(\Delta t) m_{j-1}, \quad (\text{A.8})$$

$$P_j^- = \Phi(\Delta t) P_{j-1} \Phi(\Delta t)^\top + Q(\Delta t), \quad (\text{A.9})$$

for given initial conditions m_0, P_0 . The prediction step is the same, for both $t_j \in \mathcal{T}^{\text{OBS}}$ and $t_j \in \mathcal{T}^{\text{ODE}}$.

As detailed in Section 3, two different update steps are defined for two kinds of observations. When observing data $y_{0:N}$, i.e. $t_n \in \mathcal{T}^{\text{OBS}}$, the update step follows the rules of a standard Kalman filter. The updated mean m_n and covariance P_n at time t_n are computed as

$$v_n = y_n - H m_n^-, \quad (\text{A.10})$$

$$S_n = H P_n^- H^\top + R, \quad (\text{A.11})$$

$$K_n = P_n^- H^\top S_n^{-1}, \quad (\text{A.12})$$

$$m_n = m_n^- + K_n v_n, \quad (\text{A.13})$$

$$P_n = P_n^- - K_n S_n K_n^\top. \quad (\text{A.14})$$

The matrices H and R are defined as in Eq. (6) in the paper.

Recall the ODE measurement model from Eq. (7), which we here denote as h , as

$$h \left(\begin{pmatrix} U(t) \\ X(t) \end{pmatrix} \right) = X^{(1)} - f \left(X^{(0)}; U(t) \right). \quad (\text{A.15})$$

At locations $t_m \in \mathcal{T}^{\text{ODE}}$, we condition on the ODE measurements $z_{0:M}$. Recall that these pseudo-observations are all zero. According to Eq. (10.79) in the book by Särkkä and Solin [34],

$$v_m = z_m - h(m_m^-), \quad (\text{A.16})$$

$$S_m = [Dh(m_m^-)] P_m^- [Dh(m_m^-)]^\top + \lambda^2 I_d, \quad (\text{A.17})$$

$$K_m = P_m^- [Dh(m_m^-)]^\top S_m^{-1}, \quad (\text{A.18})$$

$$m_m = m_m^- + K_m v_m, \quad (\text{A.19})$$

$$P_m = P_m^- - K_m S_m K_m^\top, \quad (\text{A.20})$$

where $[Dh(m_m^-)]$ denotes the Jacobian of h at m_m^- . In the case of a Dirac likelihood (see Eq. (7)), $\lambda^2 = 0$ holds. For numerical stability (especially for $\lambda^2 = 0$) one can instead implement square-root filtering (see, e.g., [10, 22]). All experiments in Section 5 use square-root filtering.

B Parametric model for MCMC sampling

This section first introduces a functional form for $\beta(t)$ that connects to the non-parametric model introduced in Section 3. Then, a generative model for Markov-chain Monte Carlo (MCMC) inference over the unknown parameters of $\beta(t)$ is set up.

We establish a parametric model for the latent, time-varying contact rate in an SIRD model in terms of Fourier features. In light of Mercer's theorem and the fact that stationary covariance functions have complex-exponential eigenfunctions [32, Chapter 4.3], this closely connects to the Matérn- $3/2$ process used in Sections 5.2 and 5.3 (see also [31]).

Concretely, we proceed as follows. Let \mathbb{T} denote a dense time grid. First, (i) compute the kernel Gram matrix K on \mathbb{T} , such that $(K)_{ij} = k(x_i, x_j)$ with $x_i, x_j \in \mathbb{T}$. k is the Matérn- $3/2$ covariance function.

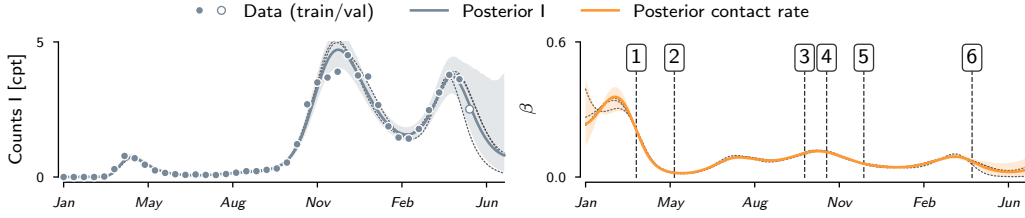


Figure 7: **Estimated counts of infectious cases and contact rate.** The estimates are obtained from MCMC sampling in an SIRD model with a parametric function for the contact rate $\beta(t)$. The case counts of infectious people are scaled to cases per thousand (cpt). The shaded areas show the 95% credible interval and the dotted black lines are samples from the posterior. Compared to the non-parametric approach presented in the paper, the estimate over $\beta(t)$ is very confident in general. The posterior mean closely resembles the results obtained in Sections 5.2 and 5.3. The numbered markers in the right plot are explained in Table 1 in the paper.

As in the experiments before, we set the characteristic lengthscale to $\ell = 75$. Then, (ii) compute the eigendecomposition of K . In order to keep the dimensionality of the inference problem feasible, select $r \ll |\mathbb{T}|$ eigenvectors that correspond to the r largest eigenvalues of K . In this experiment, we choose $r = 25$. (iii) For each eigenvector, the strongest frequency component ω is determined by the discrete Fourier decomposition. This yields a set of frequencies $\{\omega_i : i = 1, \dots, r\}$. Finally, the parametric model is defined as the sum of parametrized Fourier features of the form

$$\beta(t) = \vartheta \left(\sum_{i=1}^r a_i \cos(2\pi\omega_i t) + b_i \sin(2\pi\omega_i t) \right), \quad (\text{B.1})$$

where ϑ is the logistic sigmoid function as described in Section 5. We aim to compute a posterior contact rate $\beta(t)$ by MCMC inference over the coefficients a_i and b_i , $i = 1, \dots, r$. To this end, we define a prior over the parameter vector $\theta := (a_1, b_1, \dots, a_r, b_r)^\top$ and a likelihood for the COVID-19 case counts $y_{0:N}$ with respect to θ .

In order to ensure non-negative case counts, as in Section 5.3, we assume log-normally distributed measurements with i.i.d. noise

$$p(y_{0:N} | \theta) = \prod_{n=0}^N \text{LogNormal} \left(y_n; \log \left(x^{(\theta)}(t_n) \right), \sigma^2 I_{2r} \right), \quad (\text{B.2})$$

where σ^2 is inferred from the data along with θ . $x^{(\theta)}(t_n)$ denotes the solution of the SIRD system at time t_n , parametrized by the vector of coefficients θ through the contact rate from Eq. (B.1). Notably, each evaluation of the likelihood involves numerically integrating the SIRD system, which significantly increases the computational cost entailed by the inference algorithm. This is done by NumPyro's DOPRI-5 implementation [30, 8].

The prior distributions over the Fourier-feature coefficients and over σ^2 are chosen as

$$p(\theta) = \mathcal{N}(\theta; \mu_\theta, \Sigma_\theta), \quad p(\sigma^2) = \text{HalfCauchy}(\sigma^2; 0.01). \quad (\text{B.3})$$

The mean μ_θ of the prior over θ is set to a maximum-likelihood estimate by minimizing the negative logarithm of Eq. (B.2) with SciPy's L-BFGS optimization algorithm [40, 23]. The covariance is chosen as $\Sigma_\theta = 0.1 \cdot I_{2r}$.

The goal of the experiment is to compute a posterior over the coefficients θ (and the measurement covariance σ^2) that is comparable to the results obtained in Sections 5.2 and 5.3. Like before, recovery rate and fatality rate are assumed fixed and known at $\gamma = 0.06$ and $\eta = 0.002$. We compute the posterior $p(\theta | y_{0:N})$ using NumPyro's implementation of the No-U-Turn sampler [15].

Figure 7 shows the estimated number of infectious people and the contact rate over time as inferred by the MCMC algorithm. The state estimate matches the data points well and the uncertainty increases when extrapolating. Like in the experiments in Sections 5.2 and 5.3, the final 14 observations serve

as a validation set and the model extrapolates 31 days into the future. The posterior mean closely resembles the results obtained from our method. However, the uncertainty is lower in general, especially in the beginning and over the summer months.

C Sources for governmental measures in Germany

This section provides the sources used to list the governmental measures in Table 1. In order to provide reliable sources, we refer to the official press releases, as published by the German government. For each policy change, we provide a very brief idea of the imposed measures and official sources by the German government (only available in German language).

C.1 March 22, 2020 (Mark 1)

Citizens are urged to restrict social contacts as much as possible and the formation of groups is sanctioned in public spaces as well as at home.

<https://www.bundesregierung.de/breg-de/themen/coronavirus/besprechung-der-bundeskanzlerin-mit-den-regierungschefinnen-und-regierungschefs-der-laender-vom-22-03-2020-1733248>

<https://www.bundesregierung.de/resource/blob/975226/1733246/e6d6ae0e89a7f1ea1ebf6f32cf472736/2020-03-22-mpk-data.pdf?download=1>

C.2 May 6, 2020 (Mark 2)

The government puts the federal states in charge of appropriately relaxing the imposed measures. Different states handle the situation differently, according to the respective incidences (*'hotspot strategy'*).

<https://www.bundesregierung.de/breg-de/aktuelles/pressekonferenzen/pressekonferenz-von-bundeskanzlerin-merkel-ministerpraesident-soeder-und-dem-ersten-buergermeister-tschentscher-im-anschluss-an-das-gespraech-mit-den-regierungschefinnen-und-regierungschefs-der-laender-1751050>

C.3 October 7, 2020 (Mark 3) and October 14, 2020

The population is again urged to restrict contacts if possible.

<https://www.bundeskanzlerin.de/bkin-de/aktuelles/telefonschaltkonferenz-des-chefs-des-bundeskanzleramts-mit-den-chefinnen-und-chefs-der-staats-und-senatskanzleien-der-laender-am-7-oktober-2020-1796770>

One week later, new light restrictions are imposed. The number of people allowed in social gatherings is limited, according to local incidences.

<https://www.bundesregierung.de/resource/blob/997532/1798920/9448da53f1fa442c24c37abc8b0b2048/2020-10-14-beschluss-mpk-data.pdf?download=1>

C.4 November 2, 2020 (Mark 4)

Partial shutdown of public life (*'lockdown light'*). Across the country, the number of people allowed in social gatherings is limited to ten, where the number of households present must not exceed two. Most of public services are closed or offered only virtually, if possible.

<https://www.bundesregierung.de/breg-de/aktuelles/videokonferenz-der-bundeskanzlerin-mit-den-regierungschefinnen-und-regierungschefs-der-laender-am-28-oktober-2020-1805248>

C.5 December 16, 2020 (Mark 5)

Across the country, the number of people allowed in social gatherings is limited to five, where the number of households present must not exceed two. Except for stores of systemic importance, the retail sector is mostly shut down.

<https://www.bundesregierung.de/resource/blob/997532/1827366/69441fb68435a7199b3d3a89bff2c0e6/2020-12-13-beschluss-mpk-data.pdf?download=1>

C.6 April 23, 2021 (Mark 6)

The aforementioned measures were mostly governed and implemented by the respective federal states. On April 22, 2021, the German government decides on a nationwide decree of measures to come into effect on the following day (April 23, 2021). Depending on the seven-day incidence, curfews, contact restrictions, and a shutdown of large parts of public life are imposed.

https://www.bundesgesundheitsministerium.de/fileadmin/Dateien/3_Downloads/Gesetze_und_Verordnungen/GuV/B/4_BevSchG_BGBL.pdf