

Journal Pre-proof

Incorporating physical constraints in a deep probabilistic machine learning framework for coarse-graining dynamical systems

Sebastian Kaltenbach, Phaedon-Stelios Koutsourelakis

PII: S0021-9991(20)30447-2
DOI: <https://doi.org/10.1016/j.jcp.2020.109673>
Reference: YJCPH 109673

To appear in: *Journal of Computational Physics*

Received date: 6 January 2020
Revised date: 9 May 2020
Accepted date: 16 June 2020

Please cite this article as: S. Kaltenbach, P.-S. Koutsourelakis, Incorporating physical constraints in a deep probabilistic machine learning framework for coarse-graining dynamical systems, *J. Comput. Phys.* (2020), 109673, doi: <https://doi.org/10.1016/j.jcp.2020.109673>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2020 Published by Elsevier.



Highlights

- A generative model for the automated discovery of CG dynamics.
- The target density is augmented by virtual observables which reflect physical constraints.
- The incorporation of physical constraints leads to a reduction of the training data.
- A probabilistic formulation that is capable of quantifying predictive uncertainty.
- Full reconstruction of futures of the entire FG state vector as well as any FG observable.

Incorporating physical constraints in a deep probabilistic machine learning framework for coarse-graining dynamical systems

Sebastian Kaltenbach^a, Phaedon-Stelios Koutsourelakis^{a,*}

^a*Professorship of Continuum Mechanics, Technical University of Munich*

Abstract

Data-based discovery of effective, coarse-grained (CG) models of high-dimensional dynamical systems presents a unique challenge in computational physics and particularly in the context of multiscale problems. The present paper offers a data-based, probabilistic perspective that enables the quantification of predictive uncertainties. One of the outstanding problems has been the introduction of physical constraints in the probabilistic machine learning objectives. The primary utility of such constraints stems from the undisputed physical laws such as conservation of mass, energy etc. that they represent. Furthermore and apart from leading to physically realistic predictions, they can significantly reduce the requisite amount of training data which for high-dimensional, multiscale systems are expensive to obtain (Small Data regime). We formulate the coarse-graining process by employing a probabilistic state-space model and account for the aforementioned equality constraints as virtual observables in the associated densities. We demonstrate how deep neural nets in combination with probabilistic inference tools can be employed to identify the coarse-grained variables and their evolution model without ever needing to define a fine-to-coarse (restriction) projection and without needing time-derivatives of state variables.

We advocate a sparse Bayesian learning perspective which avoids over-fitting and reveals the most salient features in the CG evolution law. The formulation adopted enables the quantification of a crucial, and often ne-

*Corresponding author

Email addresses: `sebastian.kaltenbach@tum.de` (Sebastian Kaltenbach),
`p.s.koutsourelakis@tum.de` (Phaedon-Stelios Koutsourelakis)

glected, component in the CG process, i.e. the predictive uncertainty due to information loss. Furthermore, it is capable of reconstructing the evolution of the full, fine-scale system and therefore the observables of interest need not be selected a priori. We demonstrate the efficacy of the proposed framework by applying it to systems of interacting particles and a series of images of a nonlinear pendulum. In both cases we identify the underlying coarse dynamics and can generate extrapolative predictions including the forming and propagation of a shock for the particle systems and a stable trajectory in the phase space for the pendulum.

Keywords: Bayesian machine learning, virtual observables, multiscale modeling, reduced order modeling, coarse graining

1. Introduction

2 High-dimensional, nonlinear dynamical systems are ubiquitous in applied
 3 physics and engineering. The computational resources needed for their so-
 4 lution can grow exponentially with the dimension of the state-space as well
 5 as with the smallest time-scale that needs to be resolved and which deter-
 6 mines the discretization time-step. Hence the ability to construct reduced,
 7 *coarse-grained* descriptions and models that are nevertheless predictive of
 8 various observables and at time-scales much larger than the inherent ones, is
 9 an important task (Givon et al., 2004).

10 One strategy for learning such coarse-grained (CG) models is based on
 11 data generated by simulations of the fine-grained (FG) system. This can
 12 yield an automated solution especially in cases where domain knowledge is
 13 limited or absent. The derivation of CG models from data is also partic-
 14 ularly relevant in domains where FG models are not available, such as in
 15 social sciences or biophysics, but data abound (Bialek, 2012; Alber et al.,
 16 2019). Data-based methodologies have also been fueled by recent advances
 17 in statistical- (Ghahramani, 2015) or machine-learning (LeCun et al., 2015)
 18 which, in large part, have been enabled by large datasets (and the compu-
 19 tational means to leverage them). We note nevertheless that coarse-graining
 20 tasks based on FG simulation data exhibit some fundamental differences
 21 (Koutsourelakis et al., 2016). Firstly, the acquisition of FG simulation data
 22 is by definition expensive and the reduction of the required FG simulations is
 23 one of the objectives of CG model development. Secondly, in physical appli-
 24 cations, significant information about the underlying physical/mathematical

structure of the problem, and of the CG model in particular, is available. This
 26 information might come in the form of constraints that reflect e.g. undisputed
 physical principles such as conservation laws (e.g. mass, momentum, energy).
 28 Injecting this prior information into the CG models in combination with FG
 data in an automated fashion represents a significant challenge (Marcus and
 30 Davis, 2019), especially in the context of *probabilistic* models (Stinis et al.,
 2019). Such a capability would be instrumental not only in reducing the
 32 required amount of FG data, but more importantly, in enabling predictions
 under *extrapolative* settings as those arising e.g. when the initial conditions
 34 of the FG system are different from the ones in the training data.

In this paper, we propose a generative, probabilistic (Bayesian) machine
 36 learning framework (Koutsourelakis and Bilonis, 2011) which employs FG
 simulation data augmented by *virtual observables* to account for constraints.
 38 The latter concept which we elucidate in the sequel, enables the incorpora-
 tion of domain knowledge in probabilistic models and represents, in our
 40 opinion the most novel contribution of this paper. Furthermore and within
 the Bayesian framework advocated, it allows us to introduce appropriate pri-
 42 ors that promote the discovery of *slow-varying* CG state-variables which is a
 highly-desirable feature for multiscale systems (Kevrekidis et al., 2003). In
 44 contrast to most existing techniques which consider the problems of CG state
 variable discovery and CG model construction in two or more steps (Schmid,
 46 2010; Williams et al., 2015; Wu and Noé, 2017; Froyland et al., 2014), we
 address both simultaneously (Felsberger and Koutsourelakis, 2019). The
 48 framework proposed consists of two building blocks: a probabilistic coarse-
 to-fine map (Schöberl et al., 2017) and an evolution law for the CG dynamics.
 50 The former can be endowed with great flexibility in discovering appropriate
 CG variables when combined with deep neural nets (Raissi et al., 2017, 2019;
 52 Yang and Perdikaris, 2019), which is especially challenging if the number of
 training data is small¹. We demonstrate nevertheless the efficacy of such an
 54 approach when physical information is incorporated a-priori into the model.
 The CG variables identified are not restricted to indicator functions of sub-
 56 domains of the state-space as in other generative models (Mardt et al., 2018;
 Wu and Noé, 2017; Wu et al., 2018) and which are difficult to learn when the
 58 simulation data is limited and has not sufficiently populated all important

¹In the dynamical systems investigated the size of the dataset depends on the length
 of the FG time-sequences as well as the number of such sequences employed for training.

regions of the state-space.

60 The second component of the proposed framework pertains to the discovery of the CG evolution law which is learned by employing a large vocabulary
 62 of feature functions and sparsity-inducing priors. This leads to interpretable solutions (Duncker et al., 2019), even in the *Small Data* regime that avoid
 64 overfitting and reveal salient characteristics of the CG system (Grigo and Koutsourelakis, 2019). The premise of sparsity (Pantazis and Tsamardinos, 2019)
 66 has been employed in the past for the discovery of the CG dynamics as e.g. in the SINDy method (Brunton et al., 2016; Kaiser et al., 2018; Champion et al., 2019).
 68 This however requires the availability of time-derivatives of the CG variables and does not directly lead to a posterior on the model parameters that can reflect inferential uncertainties. Nonparametric models
 70 for the CG dynamics have also been proposed (Ohkubo, 2011) but have been restricted to low dimensions. The learned CG dynamics are in general non-linear
 72 in contrast to efforts based on transfer operators (Klus et al., 2018) and particularly the Koopman operator (Koopman, 1931; Mezić, 2005; Brunton et al., 2016).
 74 While the associated theory guarantees the existence of a linear operator, this is possible in the infinite dimensional space of observables, it does not specify how many should be used to obtain a good approximation,
 76 and more importantly, how one can predict future FG states given predictions on the evolution of those observables i.e. the reconstruction step.

80 The latter constitutes the main difference of the proposed model with non-generative ones based e.g. on information-theoretic concepts (Katsoulakis and Plecháč, 2013; Harmandaris et al., 2016; Katsoulakis and Vilanova, 2019)
 82 or on the Mori-Zwanzig (MZ) formalism (Mori, 1965; Zwanzig, 1973; Chorin and Stinis, 2007). Apart from the difficulties in approximating the right-hand-side of the MZ-prescribed CG dynamics, and particularly the memory
 84 term (Lei et al., 2016; Zhu et al., 2018), this can only guarantee correct predictions of the CG variables' evolution. If observables not depending on CG variables are of interest, then a reconstruction operator would need to
 86 be added. In contrast, in the proposed model this reconstruction operator is represented by the probabilistic coarse-to-fine map which is simultaneously learned from the data and can quantify predictive uncertainties associated
 88 with the information loss that unavoidably takes place in any CG process as well as due to the fact that finite (and preferably, small) data has been used for training.

96 The enabling computational technology for training the proposed model

is based on probabilistic inference. In order to resolve the intractable posterior on latent variables and model parameters in our Bayesian framework, we make use of Stochastic Variational Inference (Hoffman et al., 2013) as MCMC is cumbersome in high dimensions. We operate on the discretized time domain (Archambeau and Oppé, 2011) and demonstrate how amortized (Krishnan et al., 2017; Fortuin et al., 2019) and non-amortized approximations can be employed.

The remainder of the paper is structured as follows: In Section 2 we present the general methodological framework with special attention on the two building blocks of the state-space model proposed i.e. the transition law for the CG dynamics and the incorporation of virtual observables (section 2.2), as well as the emission law which provides the link between CG and FG description through a probabilistic *coarse-to-fine* map (section 2.3). Computational aspects related to inference and prediction are discussed in sections 2.4 and 2.5 respectively. Section 3 contains illustrative applications involving coarse-graining of high-dimensional systems of interacting particles (section 3.1) as well as learning the dynamics of a nonlinear pendulum (section 3.2) from a sequence of images. We conclude in section 4 which also contains a discussion on possible extensions.

2. Methodology

In general, we use the subscript f or lower-case letters to denote variables associated with the (high-dimensional) fine-grained(FG)/full-order model and the subscript c or upper-case letters for quantities of the (lower-dimensional) coarse-grained(CG)/reduced-order description. We also use a circumflex $\hat{\cdot}$ to denote observed/known variables. We begin with the presentation of the FG and the CG model and subsequently explain the essential ingredients of the proposed formulation.

2.1. The FG and CG models

We consider a, generally high-dimensional, FG system with state variables \mathbf{x} of dimension d_f ($d_f \gg 1$) such that $\mathbf{x} \in \mathcal{X}_f \subset \mathbf{R}^{d_f}$. The dynamics of the FG system are dictated by system of deterministic or stochastic ODEs i.e.,

$$\dot{\mathbf{x}}_t = \mathbf{f}(\mathbf{x}_t, t), \quad t > 0 \quad (1)$$

130 The initial condition \mathbf{x}_0 might be deterministic or drawn from a specified
 132 distribution. In the following we do not make explicit use of the FG dynamics
 but rely purely on FG data i.e. time sequences simulated from Equation (1)
 with a time-step, say δt . That is, our observables consists of n data sequences
 134 over $T + 1$ FG time-steps δt i.e.,

$$\mathcal{D}_{T,n} = \{\hat{\mathbf{x}}_{0:T\delta t}^{(1:n)}\} \quad (2)$$

We denote the (unknown) CG state variables by \mathbf{X} and assume $\mathbf{X} \in \mathcal{X}_c \subset$
 136 \mathbf{R}^{d_c} , where d_c is the dimension of the CG system. We presuppose Markovian
 dynamics² for the CG system of the form:

$$\dot{\mathbf{X}}_t = \mathbf{F}(\mathbf{X}_t, t) \quad (3)$$

138 which we discretize using a linear multistep method and a CG time step Δt :

$$\mathbf{R}_l(\mathbf{X}) = \sum_{k=0}^K (\alpha_k \mathbf{X}_{(l-k)\Delta t} + \Delta t \beta_k \mathbf{F}(\mathbf{X}_{(l-k)\Delta t})) = 0, \quad l = K, K+1, \dots \quad (4)$$

where α_k, β_k are the parameters of the discretization scheme and \mathbf{R}_l the cor-
 140 responding residual at time step l (Butcher, 2016). We note that depending
 on the values of the parameters K, α_k, β_k , several of the well-known, ex-
 142 plicit/implicit, numerical time-integration schemes can be recovered. In this
 work, our goal is two-fold:

- 144 a) to identify the CG state-variables \mathbf{X} and their relation with the FG
 description \mathbf{x} ,
- 146 b) to identify the right-hand side of Equation (3),

in view of enabling predictions of the FG system over longer time horizons.
 148 Traditionally, the aforementioned tasks are *not* considered simultaneously.
 Usually the CG state variables are specified a priori using domain-knowledge
 150 (physical insight) or based on the observables of interest (Harmandaris et al.,
 2016). In other efforts, linear or non-linear dimensionality reduction proce-
 152 dures are first employed in order to identify such a lower-dimensional set of
 collective variables \mathbf{X} (e.g. (Coifman et al., 2008)). In both of these cases,

²As discussed in section 3, this assumption can be relaxed.

154 \mathbf{X} are defined using a *fine-to-coarse*, projection map e.g. $\mathbf{X} = \Pi(\mathbf{x})$ where
 155 $\Pi : \mathcal{X}_f \subset \mathbb{R}^{d_f} \rightarrow \mathcal{X}_c \subset \mathbb{R}^{d_c}$. Irrespective of whether this map is prescribed
 156 from the physics or learned from data, it is generally a many-to-one function
 157 that does not have an inverse i.e. if the CG states \mathbf{X} are known one cannot
 158 readily reconstruct \mathbf{x} (Trashorras and Tsagkarogiannis, 2010).

We note that that this has nothing to do with the quality of the CG
 160 evolution law (problem b) above). Even if the Mori-Zwanzig (MZ) formal-
 161 ism were employed, which in principle provides an exact, closed system of
 162 evolution equations for any observable of the FG states and therefore for
 163 $\mathbf{X} = \Pi(\mathbf{x})$, even if all the terms in the right-hand side were available, one
 164 would simply be able to predict the future evolution of \mathbf{X} but not \mathbf{x} . This
 165 might be sufficient for a lot of problems of practical interest where the CG
 166 variables (or observables thereof) are of sole interest. Our goal however is
 167 a bit more ambitious, i.e. we seek to find a \mathbf{X} that would allow us to re-
 168 construct as accurately as possible the whole FG vector \mathbf{x} into the future.
 169 As with any coarse-graining process, we recognize that this would unavoid-
 170 ably imply some information loss which in turn will give rise to predictive
 171 uncertainty (Katsoulakis and Trashorras, 2006). In this work, we advocate
 172 a probabilistic framework that quantifies this uncertainty.

With regards to problem b) above, we note that its solution hinges upon
 174 the CG variables \mathbf{X} employed (problem a)). Irrespective of the breadth of
 175 the model forms considered (i.e. functions \mathbf{F} in Equation (3)), the evolution
 176 of some \mathbf{X} might fall outside this realm. For example, it is known from MZ
 177 theory that memory terms can become significant for certain observables. It
 178 is well-known that such memory terms can be substituted or approximated by
 179 additional variables (Kondrashov et al., 2015) which would in turn imply an
 180 augmented CG description \mathbf{X} in Equation (3) that contains these auxiliary
 181 internal state variables (Coleman and Gurtin, 1967).

182 We address problems a) and b) in the coarse-graining process *simultane-*
 183 *ously* by employing a probabilistic state-space model. This consists of two
 184 densities i.e.

- 185 • the transition law which dictates the evolution of the CG variables \mathbf{X}
 186 (section 2.2). Special attention is paid to the definition of *virtual ob-*
 187 *servables* with which the CG states and their dynamics can be injected
 188 with physical information.
- 189 • the emission law which provides the link between CG and FG descrip-
 190 tion through a probabilistic *coarse-to-fine* map (section 2.3, (Felsberger

and Koutsourelakis, 2019)).

We emphasize that in our formulation, the CG state-variables \mathbf{X} are implicitly defined as latent generators of the FG description \mathbf{x} . As discussed in detail in the sequel, this enables a straightforward, *probabilistic* reconstruction of \mathbf{x} when \mathbf{X} is known. The inverse map (analogous to Π above) arises naturally through probabilistic inference as explained in section 2.4. An overview of the essential elements of the proposed model can be seen in the probabilistic graphical model of Figure 1.

2.2. Transition Law: CG dynamics and virtual observables

Typical state-space models (Cappe et al., 2005; Ghahramani, 2004; Durstewitz, 2017; Krishnan et al., 2017) postulate Markovian, stochastic dynamics for the hidden variables \mathbf{X} , in the form of a diffusion process, which are subsequently discretized *explicitly* using e.g. a Euler-Maruyama scheme with time step Δt . This gives rise to a, generally Gaussian, conditional density $p(\mathbf{X}_{(l+1)\Delta t}|\mathbf{X}_{l\Delta t})$ which can be stacked over multiple time-instants in order to formulate a generalized prior on the CG-space.

When the CG state-variables \mathbf{X} are given (in part or in whole) physical meaning (e.g. as thermodynamic state variables), then some of the equations for their evolution are prescribed by associated physical principles e.g. conservation of mass, momentum, energy. These can be reflected in the residuals \mathbf{R}_l of the governing equations as in Equation (4) or alternatively as equality constraints of the form:

$$\mathbf{c}_l(\mathbf{X}_{l\Delta t}) = \mathbf{0}, \quad l = 0, 1, \dots \quad (5)$$

which must hold at each time-step. The function $\mathbf{c}_l : \mathcal{X}_c \subset \mathbb{R}^{d_c} \rightarrow \mathbb{R}^{M_c}$ enforces these known constraints at each time-step (see specific examples in section 3) and the only requirement we will impose is that of differentiability of \mathbf{c}_l (see section 2.4). In order to account for the aforementioned constraints in the transition law of the CG state variables, we employ the novel (to the best of our knowledge) concept of *virtual observables*. In particular for each of the residuals \mathbf{R}_l in Equation (4), we define a new variable/vector $\hat{\mathbf{R}}_l$ which relates to \mathbf{R}_l as follows:

$$\hat{\mathbf{R}}_l = \mathbf{R}_l(\mathbf{X}) + \sigma_R \boldsymbol{\epsilon}_R, \quad \boldsymbol{\epsilon}_R \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (6)$$

We further assume that $\hat{\mathbf{R}}_l$ have been *virtually observed* and $\hat{\mathbf{R}}_l = 0$ leading to an augmented version of the data in Equation (2), by a set of virtual

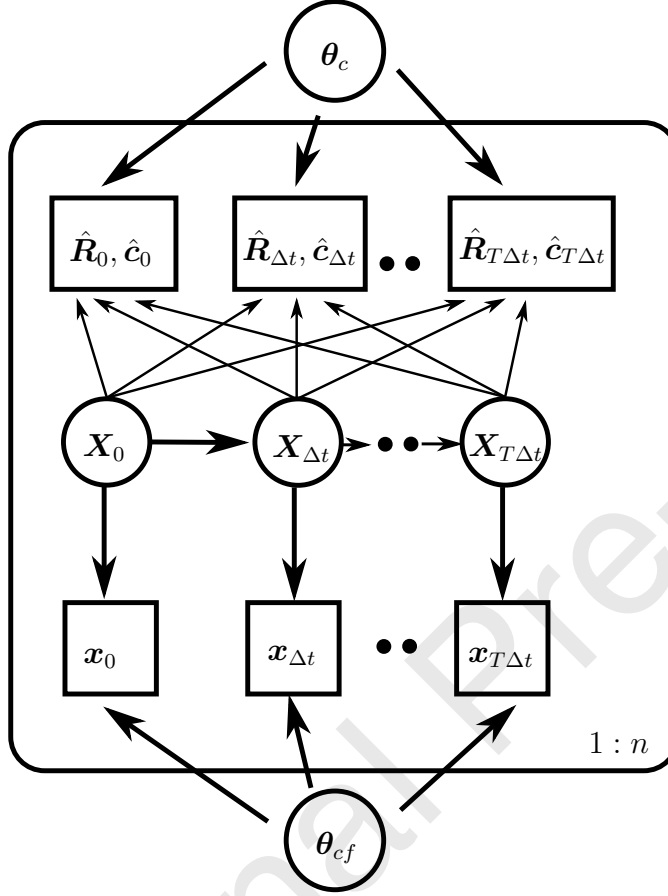


Figure 1: Proposed probabilistic graphical model. The CG variables \mathbf{X} are latent and are inferred together with the parameters θ_c and θ_{cf} . Apart from the the FG states \mathbf{x} , the observables are augmented by *virtual observables* $\hat{\mathbf{R}}, \hat{\mathbf{c}}$ (see section 2.2). These virtual observables can depend on all CG variables but more often this dependence is restricted to only a few of them.

observations and therefore virtual likelihoods of the type:

$$p(\hat{\mathbf{R}}_l = \mathbf{0} \mid \mathbf{X}, \sigma_R) = \mathcal{N}(\mathbf{0} \mid \mathbf{R}_l(\mathbf{X}), \sigma_R^2 \mathbf{I}) \quad (7)$$

224 The “noise” parameter σ_R determines the intensity of the enforcement of the
virtual observations and is analogous to the tolerance parameter with which
226 residuals are enforced in a deterministic solution of the dynamics. Similarly,
for constraints of the form of Equation (5), additional variables and virtual
228 observables of the type:

$$\mathbf{0} = \hat{\mathbf{c}}_l = \mathbf{c}_l(\mathbf{X}_{l\Delta t}) + \sigma_c \boldsymbol{\epsilon}_c, \quad \boldsymbol{\epsilon}_c \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (8)$$

can be defined which would lead to an augmented (virtual) likelihood with
230 terms of the type:

$$p(\hat{\mathbf{c}}_l = \mathbf{0} \mid \mathbf{X}_{l\Delta t}, \sigma_c) = \mathcal{N}(\mathbf{0} \mid \mathbf{c}_l(\mathbf{X}_{l\Delta t}), \sigma_c^2 \mathbf{I}) \quad (9)$$

where the role of σ_c^2 is analogous to σ_R^2 above.

232 Since the goal is to identify the right-hand side of the evolution laws in
Equation (3), we denote by $\boldsymbol{\theta}_c$ the parameters appearing in \mathbf{F} i.e. $\mathbf{F}(\mathbf{X}_t, t; \boldsymbol{\theta}_c)$.
234 Accordingly, the virtual observations in Equation (6) or Equation (8) would
depend on $\boldsymbol{\theta}_c$. We defer until section 3 a detailed discussion on the form,
236 the parametrization as well as the prior specifications in the Bayesian set-
ting adopted. The latter plays an important role as with sparsity-inducing
238 priors we can avoid overfitting and obtain a parsimonious and physically-
interpretable solution for \mathbf{F} . We finally remark that physical information
240 taking the form of equalities can also be available for the FG states \mathbf{x} . While
this can be incorporated using appropriate virtual observables as above, the
242 inference framework would exhibit significant differences (in brief, FG states
would need to be inferred as well) and in order to avoid confusion we do not
244 discuss such cases here.

2.3. Emission law: Coarse-to-Fine map

246 We make use of a probabilistic *generative* model in the definition of the
CG state-variables through a *coarse-to-fine* map (Felsberger and Koutsourelakis,
248 2019) as opposed to traditional, many-to-one maps from the FG de-
scription to the CG one. We denote the associated (conditional) density
250 by:

$$p_{cf}(\mathbf{x}_t \mid \mathbf{X}_t; \boldsymbol{\theta}_{cf}) \quad (10)$$

Observables \mathcal{D}	$\hat{\mathbf{x}}_{0:T\Delta t}^{(1:n)}$	FG simulation Data
	$\hat{\mathbf{R}}_{0:T}^{(1:n)}$	Virtual Observables corresponding to CG model residuals
	$\hat{\mathbf{c}}_{0:T}^{(1:n)}$	Virtual Observables corresponding to CG constraints
Latent variables	$\mathbf{X}_{0:T\Delta t}^{(1:n)}$	CG state variable
Model parameters $\boldsymbol{\theta}$	$\boldsymbol{\theta}_{cf}$	parameters in the coarse-to-fine mapping
	$\boldsymbol{\theta}_c$	parameters in the CG evolution law

Table 1: Data, latent variables and model parameters

where $\boldsymbol{\theta}_{cf}$ denote the (unknown) parameters that will be learned from the data. The form of p_{cf} can be adapted to the particulars of the problem and can be endowed with various levels of domain knowledge. In section 3, we provide various examples, from particle-systems where p_{cf} is fully determined by the physics, to a more abstract case where deep neural networks are employed in order to learn the full p_{cf} . We note finally that a (probabilistic) fine-to-coarse map can still be learned in the current setting, and would correspond to the *posterior* of \mathbf{X}_t given \mathbf{x}_t . We discuss this as well as all aspects pertaining to inference and learning in the next section.

2.4. Inference and Learning

We start this section by summarizing the main elements of the model presented (i.e. data, latent variables and parameters - see also Table 1) and subsequently describe a fully Bayesian inference scheme based on Stochastic Variational Inference (SVI, (Hoffman et al., 2013)) tools.

We adopt an enlarged definition of *data* which we cumulatively denote by \mathcal{D} and which encompasses:

- FG simulation data as in Equation (2) consisting of n sequences of the FG state-variables. As the likelihood model implied by the p_{cf} in Equation (10) involves only the observables at each *coarse* time-step we denote those by $\{\hat{\mathbf{x}}_{0:T\Delta t}^{(1:n)}\}$. We assume that the number of observations in each sequence is the same although this is not necessary. In fact, the length of each time-sequence and the number of time-sequences needed could be the subject of an active learning scheme. This would be particularly important in cases where very expensive, high-dimensional FG simulators are employed. The generative, proposed formulation can account for any type of (in)direct or (in)complete/partial, experimental or computational observations relating to FG states which we omit

here for simplicity of the presentation. We nevertheless illustrate this capability of the model in the example of section 3.2.

- Virtual observables relating to the CG states \mathbf{X} at each time-step l consisting of residuals $\hat{\mathbf{R}}_l^{(1:n)}$ as in Equation (6) and/or constraints $\hat{\mathbf{c}}_l^{(1:n)}$ as in Equation (8) (the superscript pertains to the time sequence $i = 1, \dots, n$). Assuming they pertain to all time-steps, we denote them by $\{\hat{\mathbf{R}}_{0:T}^{(1:n)}, \hat{\mathbf{c}}_{0:T}^{(1:n)}\}$.

The latent (unobserved) variables of the model are represented by the CG state-variables $\{\mathbf{X}_{0:T\Delta t}^{(1:n)}\}$ which relate to the FG data through the p_{cf} (in Equation (10)) and to the virtual observables through Equation (7) or Equation (9).

Finally, the (unknown) parameters of the model which we denote cumulatively by $\boldsymbol{\theta}$ consist of³:

- $\boldsymbol{\theta}_c$ which parametrize the right-hand-side of the CG evolution law (see end of section 2.2),
- $\boldsymbol{\theta}_{cf}$ which parametrize the probabilistic coarse-to-fine map (Equation (10)),
- σ_R, σ_c involved in the enforcement of virtual observables in Equation (6) and Equation (8) respectively, and,
- *hyperparameters* associated with the priors employed on the latent variables or the previous parameters.

Following a fully-Bayesian formulation, we can express the posterior of the unknowns (i.e. latent variables and parameters) as follows:

$$p(\mathbf{X}_{0:T\Delta t}^{(1:n)}, \boldsymbol{\theta} \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \mathbf{X}_{0:T\Delta t}^{(1:n)}, \boldsymbol{\theta}) p(\mathbf{X}_{0:T\Delta t}^{(1:n)}, \boldsymbol{\theta})}{p(\mathcal{D})} \quad (11)$$

where $p(\mathbf{X}_{0:T\Delta t}^{(1:n)}, \boldsymbol{\theta})$ denotes the prior on the latent variables and parameters.

³If any of the parameters in this list are prescribed, then they are omitted from $\boldsymbol{\theta}$.

302 We discuss first the likelihood term $p(\mathcal{D}|\mathbf{X}_{0:T\Delta t}^{(1:n)}, \boldsymbol{\theta})$ which can be decom-
 304 posed into the product of three (conditionally) independent terms, one for
 each data-type, i.e.:

$$p(\mathcal{D} | \mathbf{X}_{0:T\Delta t}^{(1:n)}, \boldsymbol{\theta}) = p(\hat{\mathbf{x}}_{0:T\Delta t}^{(1:n)} | \mathbf{X}_{0:T\Delta t}^{(1:n)}, \boldsymbol{\theta}) p(\hat{\mathbf{R}}_{0:T}^{(1:n)} | \mathbf{X}_{0:T\Delta t}^{(1:n)}, \boldsymbol{\theta}) p(\hat{\mathbf{c}}_{0:T}^{(1:n)} | \mathbf{X}_{0:T\Delta t}^{(1:n)}, \boldsymbol{\theta}) \quad (12)$$

We further note that (from Equation (10)):

$$p(\hat{\mathbf{x}}_{0:T\Delta t}^{(1:n)} | \mathbf{X}_{0:T\Delta t}^{(1:n)}, \boldsymbol{\theta}) = \prod_{i=1}^n \prod_{l=0}^T p_{cf}(\mathbf{x}_l^{(i)} | \mathbf{X}_l^{(i)}, \boldsymbol{\theta}_{cf}) \quad (13)$$

306 and (from Equation (7)):

$$\begin{aligned} p(\hat{\mathbf{R}}_{0:T}^{(1:n)} | \mathbf{X}_{0:T\Delta t}^{(1:n)}, \boldsymbol{\theta}) &= \prod_{i=1}^n \prod_{l=0}^T \mathcal{N}(\mathbf{0} | \mathbf{R}_l(\mathbf{X}^{(i)}), \sigma_R^2 \mathbf{I}) \\ &\propto \prod_{i=1}^n \prod_{l=0}^T \frac{1}{\sigma_R^{dim(\mathbf{R})}} \exp \left\{ -\frac{1}{2\sigma_R^2} |\mathbf{R}_l(\mathbf{X}^{(i)})|^2 \right\} \end{aligned} \quad (14)$$

and (from Equation (9)):

$$\begin{aligned} p(\hat{\mathbf{c}}_{0:T}^{(1:n)} | \mathbf{X}_{0:T\Delta t}^{(1:n)}, \boldsymbol{\theta}) &= \prod_{i=1}^n \prod_{l=0}^T \mathcal{N}(\mathbf{0} | \mathbf{c}_l(\mathbf{X}_l^{(i)}), \sigma_c^2 \mathbf{I}) \\ &\propto \prod_{i=1}^n \prod_{l=0}^T \frac{1}{\sigma_c^{dim(\mathbf{c})}} \exp \left\{ -\frac{1}{2\sigma_c^2} |\mathbf{c}_l(\mathbf{X}_l^{(i)})|^2 \right\} \end{aligned} \quad (15)$$

308 While the complexity of the expressions involved imply a non-analytic solu-
 tion for the posterior, we emphasize that the terms above encode actual and
 310 virtual observables (constraints) and they are differentiable, a property that
 is crucial for carrying out Variational Inference.

312 Before presenting the inference procedure, we mention an interesting pos-
 sibility for encoding prior information for the latent CG states $\mathbf{X}_{0:T\Delta t}^{(1:n)}$ through
 314 the prior term $p(\mathbf{X}_{0:T\Delta t}^{(1:n)})$. A desirable property of the CG state-variables is
 that of *slowness* i.e. that they should capture features of the system that
 316 evolve over (much) larger time-scales (Kevrekidis et al., 2003). The discovery
 of such features has been the goal of several statistical analysis procedures
 318 (e.g. Slow Feature Analysis (Wiskott and Sejnowski, 2002)) as well as in
 physics/chemistry literature (see a recent review in (Klus et al., 2018)). In
 320 this work we promote the discovery of such slow features by appropriate prior
 selection, and in particular by penalizing the jumps between two successive

time-instants, i.e.:

$$\begin{aligned}
p(\mathbf{X}_{0:T\Delta t}^{(1:n)}) &= \prod_{i=1}^n p_{c,0}(\mathbf{X}_0^{(i)}) \prod_{l=0}^{T-1} p(\mathbf{X}_{(l+1)\Delta t}^{(i)} | \mathbf{X}_{l\Delta t}^{(i)}, \sigma_X^2 \mathbf{I}) \\
&= \prod_{i=1}^n p_{c,0}(\mathbf{X}_0^{(i)}) \prod_{l=0}^{T-1} \mathcal{N}(\mathbf{X}_{(l+1)\Delta t}^{(i)} | \mathbf{X}_{l\Delta t}^{(i)}, \sigma_X^2 \mathbf{I}) \\
&\propto \prod_{i=1}^n p_{c,0}(\mathbf{X}_0^{(i)}) \prod_{l=0}^{T-1} \frac{1}{\sigma_X^{d_c}} \exp \left\{ -\frac{1}{\sigma_X^2} \left| \mathbf{X}_{(l+1)\Delta t}^{(i)} - \mathbf{X}_{l\Delta t}^{(i)} \right|^2 \right\}
\end{aligned} \tag{16}$$

where $p_{c,0}$ is a prior density for the initial CG state. We observe that the strength of the penalty is inversely proportional to the hyperparameter σ_X^2 and in the limit $\sigma_X^2 \rightarrow 0$ it implies a constant time history of \mathbf{X}_t . As the appropriate value for σ_X^2 depends on the problem, we include this in the parameter vector $\boldsymbol{\theta}$ that is inferred/learned from the data.

Given the intractability of the actual posterior, we advocate in this work Variational Inference. This operates on a parameterized family of densities, say $q_\phi(\mathbf{X}_{0:T\Delta t}^{(1:n)} | \boldsymbol{\theta})$ and attempts to find the one (i.e. the value of ϕ) that most closely approximates the posterior by minimizing their Kullback-Leibler divergence. It can be readily shown (Bishop, 2006), that the optimal q_ϕ , maximizes the Evidence Lower Bound (ELBO) $\mathcal{F}(q_\phi(\mathbf{X}_{0:T\Delta t}^{(1:n)} | \boldsymbol{\theta}))$ below:

$$\begin{aligned}
\log p(\mathcal{D}) &= \log \int p(\mathcal{D}, \mathbf{X}_{0:T\Delta t}^{(1:n)} | \boldsymbol{\theta}) d\mathbf{X}_{0:T\Delta t}^{(1:n)} d\boldsymbol{\theta} \\
&= \log \int \frac{p(\mathcal{D} | \mathbf{X}_{0:T\Delta t}^{(1:n)} | \boldsymbol{\theta}) p(\mathbf{X}_{0:T\Delta t}^{(1:n)} | \boldsymbol{\theta})}{q_\phi(\mathbf{X}_{0:T\Delta t}^{(1:n)} | \boldsymbol{\theta})} q_\phi(\mathbf{X}_{0:T\Delta t}^{(1:n)} | \boldsymbol{\theta}) d\mathbf{X}_{0:T\Delta t}^{(1:n)} d\boldsymbol{\theta} \\
&\geq \int \log \frac{p(\mathcal{D} | \mathbf{X}_{0:T\Delta t}^{(1:n)} | \boldsymbol{\theta}) p(\mathbf{X}_{0:T\Delta t}^{(1:n)} | \boldsymbol{\theta})}{q_\phi(\mathbf{X}_{0:T\Delta t}^{(1:n)} | \boldsymbol{\theta})} q_\phi(\mathbf{X}_{0:T\Delta t}^{(1:n)} | \boldsymbol{\theta}) d\mathbf{X}_{0:T\Delta t}^{(1:n)} d\boldsymbol{\theta} \\
&= \mathcal{F}(q_\phi(\mathbf{X}_{0:T\Delta t}^{(1:n)} | \boldsymbol{\theta}))
\end{aligned} \tag{17}$$

In the examples analyzed we decompose the approximate posterior as:

$$\begin{aligned}
q_\phi(\mathbf{X}_{0:T\Delta t}^{(1:n)} | \boldsymbol{\theta}) &= q_\phi(\mathbf{X}_{0:T\Delta t}^{(1:n)}) q_\phi(\boldsymbol{\theta}) \\
&= \left[\prod_{i=0}^n q_\phi(\mathbf{X}_{0:T\Delta t}^{(i)}) \right] q_\phi(\boldsymbol{\theta})
\end{aligned} \tag{18}$$

where the first line is the so-called mean-field approximation and the second is a direct consequence of the (conditional) independence of the time sequences in the likelihood. We note that evaluations of the ELBO \mathcal{F} involve

338 expectations with respect to q_ϕ i.e.:

$$\mathcal{F}\left(q_\phi(\mathbf{X}_{0:T\Delta t}^{(1:n)}, \boldsymbol{\theta})\right) = \mathbb{E}_{q_\phi}\left[\log p(\mathcal{D} | \mathbf{X}_{0:T\Delta t}^{(1:n)}, \boldsymbol{\theta})\right] + \mathbb{E}_{q_\phi}\left[\log \frac{p(\mathbf{X}_{0:T\Delta t}^{(1:n)}, \boldsymbol{\theta})}{q_\phi(\mathbf{X}_{0:T\Delta t}^{(1:n)}, \boldsymbol{\theta})}\right] \quad (19)$$

and in order to maximize it (with respect to ϕ), gradients of those are needed.
 340 Given the intractability of these expectations and their derivatives, we make
 use of Monte Carlo estimates in combination with stochastic gradient ascent
 342 for the ϕ -updates. In order to reduce the Monte Carlo error in these es-
 timates, we make use of the reparametrization trick (Kingma and Welling,
 344 2014), for which the differentiability of the residuals/constraints is necessary.
 We specify the particulars of the algorithm more precisely in the numerical
 346 illustration section (see e.g. Algorithm 3 or 4).

We note that maximum likelihood or maximum-a-posteriori (MAP) point
 348 estimates for any of the parameters involved can be obtained as a special case
 of the aforementioned scheme by employing a q_ϕ that is equal to a Dirac-delta
 350 function. Furthermore, amortized versions of the approximate posterior q_ϕ
 i.e. forms that explicitly account on the dependence on the data values, can
 352 be employed in part or in whole. These have the capability of being able
 to transfer information across data points and are necessary in the realm of
 354 Big Data. We note though that we operate in the *Small Data* regime, i.e.
 the number of time sequences n (and time-steps T) is not particularly large.
 356 Hybrid versions between amortized and non-amortized posteriors could also
 be employed (Kim et al., 2018).

We note finally that while the ELBO \mathcal{F} is used purely as the objective
 358 function for the determination of the approximate posterior, its role can be
 quite significant in model validation and refinement. In particular since \mathcal{F}
 360 approximates the model evidence (denominator of Equation (11)), once eval-
 362 uated, it can be used to comparatively assess different models. These could
 have different CG states \mathbf{X} (in type and/or number) or different parametriza-
 364 tions $\boldsymbol{\theta}$. In this regard, the ELBO \mathcal{F} could serve as the primary driver for
 the adaptive refinement of the CG model (Grigo and Koutsourelakis, 2019)
 366 in order to better explain the observables and lead to superior predictions.

2.5. Prediction

368 An essential feature of the proposed modeling framework is the abil-
 ity to produce *probabilistic* predictive estimates. These encompass the
 370 information-loss due to the coarse-graining process as well as the epistemic

uncertainty arising from finite (and small) datasets. We distinguish between
 372 two settings:

- 374 a) the "*interpolative*" i.e. predictions into the future of a sequence i observed up to time-step T i.e. $\hat{\mathbf{x}}_{0:T\Delta t}^{(i)}$ which was used in the training phase - see section 3, or
- 376 b) the "*extrapolative*" i.e. predictions for a completely new initial condition $\hat{\mathbf{x}}_0$ - see section 3.

378 We note that any predictions should account for the domain knowledge incorporated in the training through the residuals \mathbf{R}_l or constraints \mathbf{c}_l . Formally that is, one should enlarge the posterior density defined in Equation
 380 (11), in order to account for the residuals and/or constraints at future time-steps. This would in turn imply, that future (FG or CG) states should be
 382 inferred from such an augmented posterior i.e. prediction would imply an enlarged inference process. In the examples presented we adopt a simpler
 384 procedure that retains the essential features (i.e. probabilistic nature) but is more computationally expedient. In particular, for case a) above and if
 386 $q_\phi(\mathbf{X}_{T\Delta t}^{(i)})$ is the (marginal) posterior of the last, hidden CG state and $q(\boldsymbol{\theta})$ the posterior of the model parameters, then we (see also Algorithm 1):
 388

- sample from $q(\mathbf{X}_{T\Delta t}^{(i)}), q(\boldsymbol{\theta})$
- 390 • for each sample, we propagate the CG dynamics dynamics of Equation (3) (e.g. by solving the corresponding residual Equations (4)) in order
 392 to obtain $\mathbf{X}_{(T+1)\Delta t}^{(i)}, \mathbf{X}_{(T+2)\Delta t}^{(i)}, \dots$, and,
- we sample $\mathbf{x}_{(T+1)\Delta t}^{(i)}$ from $p_{cf}(\mathbf{x}_{(T+1)\Delta t}^{(i)} | \mathbf{X}_{(T+1)\Delta t}^{(i)}, \boldsymbol{\theta}_{cf})$, $\mathbf{x}_{(T+2)\Delta t}^{(i)}$ from
 394 $p_{cf}(\mathbf{x}_{(T+2)\Delta t}^{(i)} | \mathbf{X}_{(T+2)\Delta t}^{(i)}, \boldsymbol{\theta}_{cf})$ etc.

We note that this procedure does not necessarily ensure enforcement of the
 396 constraints by future CG states. Nevertheless it gives rise to samples of the full FG state evolution from which any observable of interest as well as the
 398 predictive uncertainty can be computed.

For the *extrapolative* setting above, i.e. for a new FG initial condition
 400 $\hat{\mathbf{x}}_0$, the evolution equations of the CG states as well as the emission density p_{cf} can be employed as long as the initial state \mathbf{X}_0 is specified or better

Algorithm 1: Prediction - Algorithm for *interpolative* setting

Result: Sample of $\mathbf{x}_{(T+P)\Delta t}^{(i)}$ **Data:** $q_\phi(\mathbf{X}_{T\Delta t}), q_\phi(\boldsymbol{\theta})$

- 1 Sample from $q_\phi(\mathbf{X}_{T\Delta t}^{(i)})$ and $q_\phi(\boldsymbol{\theta})$;
 - 2 **while** *Time-step $(T + P)\Delta t$ of interest not reached* **do**
 - 3 | Apply the CG evolution law as described in Equation (4);
 - 4 **end**
 - 5 Sample from $p_{cf}(\mathbf{x}_{(T+P)\Delta t} | \mathbf{X}_{(T+P)\Delta t}, \boldsymbol{\theta})$
-

Algorithm 2: Prediction - Algorithm for *extrapolative* setting

Result: Sample of $\mathbf{x}_{P\Delta t}$ **Data:** $p_\phi(\hat{\mathbf{x}}_0), q_\phi(\boldsymbol{\theta})$

- 1 Apply Bayesian Inference as described in Equation (20) to infer $p(\mathbf{X}_0 | \hat{\mathbf{x}}_0)$;
 - 2 Sample from $p(\mathbf{X}_0 | \hat{\mathbf{x}}_0)$ and $q(\boldsymbol{\theta})$;
 - 3 **while** *Time-step $P\Delta t$ of interest not reached* **do**
 - 4 | Apply the CG evolution law as described in Equation (4);
 - 5 **end**
 - 6 Sample from $p_{cf}(\mathbf{x}_{P\Delta t} | \mathbf{X}_{P\Delta t}, \boldsymbol{\theta})$
-

402 yet inferred. For that purpose, the posterior $p(\mathbf{X}_0 | \hat{\mathbf{x}}_0)$ of \mathbf{X}_0 needs to be determined which according to Bayes rule will be proportional to:

$$p(\mathbf{X}_0 | \hat{\mathbf{x}}_0) \propto p_{cf}(\hat{\mathbf{x}}_0 | \mathbf{X}_0, \boldsymbol{\theta}_{cf}) p_{c,0}(\mathbf{X}_0) \quad (20)$$

404 where $p_{c,0}(\mathbf{X}_0)$ is the initial state's prior (see also Equation (16)). For each sample of $\boldsymbol{\theta}_{cf}$ from the (approximate) posterior $q_\phi(\boldsymbol{\theta}_{cf})$, samples of \mathbf{X}_0 must
 406 be drawn from $p(\mathbf{X}_0 | \hat{\mathbf{x}}_0)$ and subsequently propagated as in the 3 steps above in order to obtain predictive samples of the full FG state vector (see Algo-
 408 rithm 2).

2.6. Computational considerations

410 We note that in multiscale dynamical systems of physical interest, the computational cost stems primarily from the simulation of the FG system
 412 due to its generally very high-dimensional state-vector \mathbf{x} and very small time-step δt . Hence, one of the main objectives of this work is to enable the

414 learning of the CG dynamics with the fewest possible and shortest possible
 415 FG time-sequences.

416 We note that once such FG simulation (or experimental) data have been
 417 obtained, neither the training phase of the CG model (section 2.4) nor the
 418 prediction phase (section 2.5) require any additional FG simulations. The
 419 cost of training depends on the dimension of the CG states \mathbf{X} as well as the
 420 number of parameters θ_c (for the CG dynamics), θ_{cf} (for the coarse-to-fine
 421 map) and ϕ (for the approximate posterior).

422 We emphasize that this is a one-time, offline cost i.e. once the CG model
 423 has been trained, it can be used to produce probabilistic predictive estimates
 424 of the whole FG state-vector into the future without any further recourse to
 425 the FG model. One needs only to simulate in such case the CG dynamics
 426 which due to the lower-dimensional state-vector \mathbf{X} and the much larger CG
 427 time-step Δt are much less cumbersome than the FG system.

428 Finally, if more FG data (e.g. longer or new sequences) become available
 429 at a later stage, the SVI algorithm can be re-initialized from the previous
 430 values and incorporate the new likelihood terms. If a modest amount of data
 431 is introduced, one would expect small (or even no changes for faraway states)
 432 changes and therefore rapid convergence. Naturally the introduction of ob-
 433 servables at new time instants would introduce additional latent variables for
 434 the corresponding CG states.

3. Numerical Illustrations

436 We demonstrate the capabilities of the proposed framework in discovering
 437 predictive, coarse-grained evolution laws as well as effective coarse-grained
 438 descriptions, on three examples. Two of those involve very high-dimensional
 439 systems of stochastically interacting particles (section 3.1, (Felsberger and
 440 Koutsourelakis, 2019)) and the third, a nonlinear pendulum, the dynamics
 441 of which we attempt to identify simply from sequences of images (section
 442 3.2, (Champion et al., 2019)). In the sequel, we specify the elements of the
 443 proposed model that were presented generically in the previous sections and
 444 concretize parametrizations and their meaning. The goals of the numerical
 445 illustrations are:

- 446 • to assess the predictive performance of the model under “interpolative”
 447 and “extrapolative” conditions (see section 2.5). By “interpolative” we
 448 mean the ability to predict the evolution of an FG states-sequence when
 data from this sequence has been used for training. By “extrapolative”,

we mean the ability to predict the full FG state evolution from new initial conditions that were *not* used in training.

- to examine the effect of the number n and length T of the data sequences and assess the model's ability to learn the correct structure with small n, T and partial observations.
- to examine the enforcement of the residuals/constraints (e.g. conservation of mass) in the inferred and predicted states.
- to examine the ability of the model to identify *sparse*, interpretable solutions for the CG dynamics.
- to assess the magnitude and time evolution of the predictive uncertainty estimates.
- to assess the ability of the model to learn effective CG state variables and accurate coarse-to-fine maps.

Some of the simulation results as well as the corresponding code will be made available at the following github repository⁴ upon publication.

3.1. Particle systems

3.1.1. FG model

The FG model consists of d_f identical particles which can move in the bounded one-dimensional domain $[-1, 1]$ (under periodic boundary conditions). The FG variables \mathbf{x}_t consist therefore of the coordinates of the particles at each time instant t and the dimension of the system d_f is equal to the number of particles. We consider two types of stochastic dynamics that correspond to an advection-diffusion-type (section 3.1.5) and an inviscid-Burgers-type behavior (section 3.1.6). The particulars of the microscopic dynamics are described in the corresponding sections. In the following, we discuss common aspects of both problems that pertain to the CG description, the CG evolution law and the inference procedures.

⁴https://github.com/SebastianKaltenbach/PhysicalConstraints_ProbabilisticCG.git

3.1.2. CG variables and coarse-to-fine mapping

For the CG representation, we employ the normalized particle density $\rho(s, t)$, $s \in [-1, 1]$ (Li et al., 2007) which we discretize in d_c bins. The state vector $\mathbf{X}_t = \{X_{t,j}\}_{j=1}^{d_c}$ contains the particle density values in each of the bins j , i.e. $\sum_{j=1}^{d_c} X_{t,j} = 1$ and $X_{t,j} \geq 0 \forall t, j$. We emphasize that CG and FG variables are of a different nature (i.e. proportion of particles in each bin vs. coordinates of particles) and, more importantly for practical purposes, of very different dimension.

The nature of the CG variables \mathbf{X}_t suggests a *multinomial* for the coarse-to-fine density p_{cf} (section 2.3) i.e.:

$$p_{cf}(\mathbf{x}_t | \mathbf{X}_t) = \frac{d_f!}{m_1(\mathbf{x}_t)! m_2(\mathbf{x}_t)! \dots m_{d_c}(\mathbf{x}_t)!} \prod_{j=1}^{d_c} X_{t,j}^{m_j(\mathbf{x}_t)}, \quad (21)$$

where $m_j(\mathbf{x}_t)$ is the number of particles in bin j . The underlying assumption is that, given the CG state \mathbf{X}_t , the coordinates of the particles \mathbf{x}_t are *conditionally* independent. This does *not* imply that they move independently nor that they cannot exhibit coherent behavior (Felsberger and Koutsourelakis, 2019). The practical consequence of Equation (21) is that no parameters need to be learned for p_{cf} (in contrast to section 3.2).

3.1.3. The CG evolution law and the virtual observables

With regards to the evolution law of the CG states (Equation (3)), we postulate a right-hand side $\mathbf{F}(\mathbf{X}_t; \boldsymbol{\theta}_c) = \{F_j(\mathbf{X}_t; \boldsymbol{\theta}_c)\}_{j=1}^{d_c}$ of the form:

$$\begin{aligned} F_j(\mathbf{X}_t, \boldsymbol{\theta}_c) &= \sum_{m=1}^M \theta_{c,m} \psi_m^{(j)}(\mathbf{X}_t) \\ &= \underbrace{\sum_{h=-H}^H \theta_{c,h}^{(1)} X_{t,j+h}}_{1^{st} order} + \underbrace{\sum_{h_1=-H}^H \sum_{h_2 \geq h_1}^H \theta_{c, (h_1, h_2)}^{(2)} X_{t,j+h_1} X_{t,j+h_2}}_{2^{nd} order} \end{aligned} \quad (22)$$

which consists of first- and second-order interactions over a window of size H with $\boldsymbol{\theta}_c^{(1)}$ and $\boldsymbol{\theta}_c^{(2)}$ denoting the vectors of the corresponding unknown coefficients. In this case, the total number of unknown coefficients $\boldsymbol{\theta}_c$, is $M = \dim(\boldsymbol{\theta}_c) = (2H + 1) + (H + 1)(2H + 1)$ and grows quadratically with the neighborhood-size H . Since each of the CG variables $X_{t,j}$ refers to the particle density at bin j (and at time t), the neighborhood size H corresponds to the number of bins to the left or to the right of bin j that affect its

evolution in time The feature functions that we generically denote with $\psi_m^{(j)}$ in Equation (22) can also involve higher-order interactions or be of non-polynomial type. Non-Markovian models could be accommodated as well by accounting for memory terms. It is obviously impossible to know a priori which feature functions are relevant in the evolution of the CG states or what types of interactions are essential (e.g. first, second-order etc). At the same time, and especially in the Small Data regime considered, employing a large vocabulary of feature functions can lead to *overfitting*, lack of interpretability and poor predictions, particularly under “extrapolative” conditions. This highly-important *model selection* issue has been of concern in several coarse-graining studies (Noid, 2013). We propose of automatically addressing this within the Bayesian framework advocated by employing appropriate sparsity-inducing priors for θ_c (Felsberger and Koutsourelakis, 2019). In particular, we make use of the Automatic Relevance Determination (ARD, (Mackay, 1995)) model according to which

$$p(\theta_{c,m} \mid \tau_m) = \mathcal{N}(\theta_{c,m} \mid 0, \tau_m^{-1}), \quad m = 1, 2, \dots, M = \dim(\theta_c). \quad (23)$$

The following hyperprior for the precision hyperparameters $\tau = \{\tau_m\}_{m=1}^M$ was used:

$$p(\tau_k \mid \gamma_0, \delta_0) = \text{Gamma}(\tau_k \mid \gamma_0, \delta_0) \quad (24)$$

The hyperparameters γ_0 and δ_0 are set to very small values 10^{-9} in all ensuing studies (Bishop and Tipping, 2000). As we demonstrate in the sequel, the hyperprior proposed can give rise to parsimonious solutions for the CG dynamics even in the Small Data setting considered.

A discretized version of the CG evolution law (Equation (3) and Equation (22)) with time step Δt is considered by employing a forward Euler scheme⁵ which implies the following residual vector \mathbf{R}_l at each time-step l (Equation (4)):

$$\mathbf{R}_l(\mathbf{X}) = \mathbf{X}_{(l+1)\Delta t,j} - \mathbf{X}_{l\Delta t,j} - \Delta t \mathbf{F}(\mathbf{X}_{l\Delta t,j}, \theta_c), \quad \forall l \quad (25)$$

and the corresponding virtual observables $\hat{\mathbf{R}}_l$ (Equation (6)).

More importantly, the nature of the CG variables suggests a *conservation of mass* constraint that has to be fulfilled at each time step l . In view of

⁵This corresponds to a multistep method in Equation (4) with $K = 1$, $a_0 = 1$, $a_1 = -1$, $\beta_0 = 0$ and $\beta_1 = -1$.

the discussion of section 2.2, this suggests the scalar constraint function as in Equation (5):

$$c_l(\mathbf{X}_{l\Delta t}) = \sum_{j=1}^{d_c} X_{l\Delta t,j} - 1 = 0, \quad \forall l \quad (26)$$

and the corresponding virtual observables \hat{c}_l (Equation (8)).

3.1.4. Inference and Learning

Given the multinomial p_{cf} in Equation (21), we employed the following procedure for generating training data which consists of n numerical experiments in which the FG model is randomly initialized and propagated for one coarse time-step Δt i.e. for $T = \frac{\Delta t}{\delta t}$ microscopic time-steps. In particular:

- For $i = 1, \dots, n$, we:
 - sample CG initial state $\hat{\mathbf{X}}_0^{(i)}$ from a density $p_{c,0}(\hat{\mathbf{X}}_0^{(i)})$.
 - sample FG initial state $\hat{\mathbf{x}}_0^{(i)}$ from $p_{cf}(\hat{\mathbf{x}}_0^{(i)} | \mathbf{X}_0^{(i)})$.
 - solve the (discretized) FG model for $\frac{\Delta t}{\delta t}$ microscopic time-steps and record final state $\hat{\mathbf{x}}_{\Delta t}^{(i)}$

The generated FG data $\{\hat{\mathbf{x}}_{\Delta t}^{(i)}\}_{i=1}^n$ over a *single* CG time-step are used subsequently to draw inferences on the CG model states and parameters (section 2.4). We note that longer time sequences could readily be generated (albeit at an increased cost). The number of samples n is also something that can be selected adaptively since inferences and predictions can be updated as soon as more data become available. The density $p_{c,0}(\mathbf{X}_0^{(i)})$ from which initial CG states are drawn, can be selected quite flexibly and some indicative samples are shown in Figure 2 for the advection-diffusion case, and in Figure 12 for the inviscid-Burgers' case. In summary, the data \mathcal{D} employed, apart from $\{\hat{\mathbf{x}}_{\Delta t}^{(i)}\}_{i=1}^n$ above consists of the virtual observables $\{\hat{\mathbf{R}}_0^{(1:n)}, \hat{\mathbf{c}}_1^{(1:n)}\}$.

As a result of the data employed and the parametrization adopted, we have $\mathbf{X}_{\Delta t}^{(1:n)}$ as the sole latent vector and $\boldsymbol{\theta}_c, \boldsymbol{\tau}$ as the unknown (hyper)parameters. Since only a single CG time-step was considered, we omitted the slowness prior (see Equation (16)). Hence we sought an approximate posterior $q_\phi(\mathbf{X}_{\Delta t}, \boldsymbol{\theta}_c, \boldsymbol{\tau})$ (Equation (17)) which we factorized as in Equation (18) as follows:

$$q_\phi(\mathbf{X}_{\Delta t}^{(1:n)}, \boldsymbol{\theta}_c, \boldsymbol{\tau}) = \left[\prod_{i=1}^n q_\phi(\mathbf{X}_{\Delta t}^{(i)}) \right] q(\boldsymbol{\theta}_c) q(\boldsymbol{\tau}) \quad (27)$$

Upon substitution in Equation (19), this yields the following ELBO:

$$\begin{aligned} \mathcal{F}(q_\phi(\mathbf{X}_{\Delta t}^{(1:n)}, \boldsymbol{\theta}_c, \boldsymbol{\tau})) &= \mathbb{E}_{q_\phi} \left[\log p(\mathcal{D} | \mathbf{X}_{\Delta t}^{(1:n)}, \boldsymbol{\theta}_c) \right] + \mathbb{E}_{q_\phi} [\log p(\boldsymbol{\theta}_c | \boldsymbol{\tau})] \\ &\quad + \mathbb{E}_{q_\phi} [\log p(\boldsymbol{\tau})] - \mathbb{E}_{q_\phi} [\log q_\phi] \end{aligned} \quad (28)$$

564 where:

$$p(\mathcal{D} | \mathbf{X}_{\Delta t}^{(1:n)}, \boldsymbol{\theta}_c) = p(\hat{\mathbf{x}}_{\Delta t}^{(1:n)} | \mathbf{X}_{\Delta t}^{(1:n)}) p(\hat{\mathbf{R}}_0^{(1:n)} | \mathbf{X}_{\Delta t}^{(1:n)}, \boldsymbol{\theta}_c) p(\hat{\mathbf{c}}_1^{(1:n)} | \mathbf{X}_{\Delta t}^{(1:n)}) \quad (29)$$

Based on Equation (28) the optimal variational posterior densities can be
566 obtained as:

$$\log q^{opt}(\boldsymbol{\theta}_c) = \mathbb{E}_{q_\phi(\mathbf{X}_{\Delta t}^{(1:n)})} \left[\log p(\hat{\mathbf{R}}_0^{(1:n)} | \mathbf{X}_{0:1\Delta t}^{(1:n)}, \boldsymbol{\theta}_c) \right] + \mathbb{E}_{q(\boldsymbol{\tau})} [\log p(\boldsymbol{\theta}_c | \boldsymbol{\tau})] \quad (30)$$

$$\log q^{opt}(\boldsymbol{\tau}) = \mathbb{E}_{q_\phi(\boldsymbol{\theta}_c)} [\log p(\boldsymbol{\theta}_c | \boldsymbol{\tau})] + \log p(\boldsymbol{\tau}) \quad (31)$$

568

$$\begin{aligned} \log q_\phi^{opt}(X_{\Delta t}^{(i)}) &= \log p_{cf}(\mathbf{x}_{\Delta t}^i | \mathbf{X}_{\Delta t}^i) + \mathbb{E}_{q_\phi(\boldsymbol{\theta}_c)} \left[\log p(\hat{\mathbf{R}}_0^{(i)} | \mathbf{X}_{0:1\Delta t}^{(i)}, \boldsymbol{\theta}_c) \right] \\ &\quad + \log p(\hat{\mathbf{c}}_1^{(i)} | \mathbf{X}_{\Delta t}^{(i)}) \end{aligned} \quad (32)$$

The equations above are coupled and a closed-form solution can be ob-
570 tained only for the first two. In particular, the optimal posterior approxima-
tion for $\boldsymbol{\theta}_c$ is a multivariate normal with mean $\boldsymbol{\mu}_{\boldsymbol{\theta}_c}$ and covariance $\mathbf{S}_{\boldsymbol{\theta}_c}$.

572

$$\mathbf{S}_{\boldsymbol{\theta}_c}^{-1} = \sigma_R^{-2} \sum_{i=1}^n \sum_{j=1}^{d_c} \mathbb{E}_{q_\phi(\mathbf{X}_{\Delta t}^{(i)})} \left[\boldsymbol{\psi}^{(j)}(\mathbf{X}_{\Delta t}^{(i)}) \left(\boldsymbol{\psi}^{(j)}(\mathbf{X}_{\Delta t}^{(i)}) \right)^T \right] + \mathbb{E}_{q_\phi(\boldsymbol{\tau})} [\text{diag}(\boldsymbol{\tau})] \quad (33)$$

$$\mathbf{S}_{\boldsymbol{\theta}_c}^{-1} \boldsymbol{\mu}_{\boldsymbol{\theta}_c} = \sigma_R^{-2} \sum_{i=1}^n \sum_{j=1}^{d_c} \mathbb{E}_{q_\phi(\mathbf{X}_{\Delta t}^{(i)})} \left[\boldsymbol{\psi}^{(j)}(\mathbf{X}_{\Delta t}^{(i)}) \right] \quad (34)$$

574 where the vector $\boldsymbol{\psi}^{(j)}$ consists of the M feature functions $\psi_m^{(j)}$ in Equation
(22). The optimal posterior approximation for the vector $\boldsymbol{\tau}$ of the hyper-
576 parameters $\{\tau_m\}_{m=1}^M$ reduces to a product of independent Gamma-densities
(Bishop and Tipping, 2000) with parameters γ_m and δ_m which are given by:

$$\gamma_m = \gamma_0 + 0.5, \quad \delta_m = \delta_0 + \frac{1}{2} (\boldsymbol{\mu}_{\boldsymbol{\theta}_c, m} + S_{\boldsymbol{\theta}_c, (m, m)}), \quad m = 0, 1, \dots, M = \dim(\boldsymbol{\theta}_c) \quad (35)$$

Algorithm 3: Inference algorithm for particle systems

Result: $\{q_\phi(\mathbf{X}_{\Delta t}^{(i)})\}_{i=1}^n, q(\boldsymbol{\theta}_c), q(\boldsymbol{\tau})$
Data: $\{\mathbf{X}_0^{(i)}, \hat{\mathbf{x}}_{\Delta t}^{(i)}\}_{i=1}^n$

- 1 Initialize the parameters for the variational distributions;
- 2 Set iteration counter w to zero;
- 3 Set convergence limit ϵ ;
- 4 **while** $\|parameters_w - parameters_{w-1}\|^2 > \epsilon$ **do**
- 5 **for** $i \leftarrow 1$ **to** n **do**
- 6 Update $q_\phi(\mathbf{X}_{\Delta t}^{(i)})$ by maximizing the ELBO (see Equation (28))
- 7 **end**
- 8 update $q(\boldsymbol{\theta}_c)$ according to Equation (33) and Equation (34) ;
- 9 update $q(\boldsymbol{\tau})$ according to Equation (35) ;
- 10 update the iteration counter by one ;
- 11 **end**

578 Finally and since closed-form updates for the optimal posterior $q_\phi^{opt}(\mathbf{X}_{\Delta t}^{(i)})$
are impossible, we employed Stochastic Variational Inference (SVI) as de-
580 tailed in section 2.4 by assuming a multivariate lognormal (in order to en-
sure positivity of $X_{\Delta t,j}$) with parameters $\boldsymbol{\phi} = \{\boldsymbol{\mu}_i, \mathbf{S}_i\}_{i=1}^n$ ⁶. Noisy gradients
582 with respect to the parameters $\boldsymbol{\phi}$ were estimated with Monte Carlo and the
reparametrization trick (Kingma and Welling, 2014) and $\boldsymbol{\phi}$ were updated
584 using stochastic gradient ascent (the ADAM algorithm of (Kingma and Ba,
2014) in particular). The inference steps are summarized in Algorithm 3.

	$d_f = \dim(\mathbf{x})$	$d_c = \dim(\mathbf{X})$	FG time-step δt	CG time-step Δt
Advection-Diffusion	250×10^3	≤ 64	2.5×10^{-3}	2
inviscid Burgers	250×10^3	≤ 128	2.5×10^{-3}	4

Table 2: FG/CG state-space dimensions and FG/CG time-steps for particle systems in-
vestigated.

⁶Diagonal covariances \mathbf{S}_i were employed.

586 3.1.5. Advection-Diffusion system

For the simulations presented in this section $d_f = 250 \times 10^3$ particles were used, which, at each microscopic time step $\delta t = 2.5 \times 10^{-3}$ performed random, non-interacting, jumps of size $\delta s = \frac{1}{640}$, either to the left with probability $p_{left} = 0.1875$ or to the right with probability $p_{right} = 0.2125$. The positions were restricted in $[-1, 1]$ with periodic boundary conditions. It is well-known (Cottet and Koumoutsakos, 2000) that in the limit (i.e. $d_f \rightarrow \infty$) the particle density $\rho(s, t)$ can be modeled with an advection-diffusion PDE with diffusion constant $D = (p_{left} + p_{right}) \frac{\delta s^2}{2\delta t}$ and velocity $v = (p_{right} - p_{left}) \frac{\delta s}{\delta t}$:

$$\frac{\partial \rho}{\partial t} + v \frac{\partial \rho}{\partial s} = D \frac{\partial^2 \rho}{\partial s^2}, \quad s \in (-1, 1).. \quad (36)$$

For the CG description, 64 bins were employed i.e. $d_c = 64$ and a time step $\Delta t = 2$ (see Table 2). Furthermore we employed first- and second-order feature function as in Equation (22) with a neighborhood size $H = 5$ which implies a total of $M = 77$ unknown parameters θ_c . We incorporate virtual observables pertaining to the residuals $\hat{\mathbf{R}}_0$ with $\sigma_R^2 = 10^{-6}$ (Equation (7))⁷ and the virtual observables $\hat{\mathbf{c}}_1$ pertaining to the conservation-of-mass constraint with $\sigma_c^2 = 10^{-10}$ (Equation (9)).

We employed $n = 32$ and $n = 64$ time sequences for training that were generated as detailed in section 3.1.4 with initial conditions $\{\mathbf{X}_0^{(i)}\}_{i=1}^n$ such as the ones seen in Figure 2. The initial conditions were generated by sampling the amplitude of a *sine* function, which was shifted up to ensure all values are positive and then normalized.

Figure 3 provides a histogram of the function values of the conservation-of-mass constraint $\{c_1(\mathbf{X}_{\Delta t}^{(i)})\}_{i=1}^n$ upon convergence. The small values suggest that this has been softly incorporated in the CG states. A similar histogram for the norm of the residuals $\{\mathbf{R}_0(\mathbf{X}^{(i)})\}_{i=1}^n$ is depicted in Figure 4 which also suggests enforcement of the CG evolution with the parameters θ_c learned from the data. The evolution of the posterior mean μ_{θ_c} (Equation (34)) of (a subset of) these parameters over the iterations of the SVI is depicted in Figure 5. Therein, and more clearly in Figure 6, one can observe the

⁷A very interesting possibility which is not explored here would be to learn σ_R^2 i.e. the strength of the enforcement of the CG evolution law from the data. This would increase the flexibility of the model in cases where the vocabulary of the feature functions selected in the right-hand side of the CG dynamics is not rich enough.

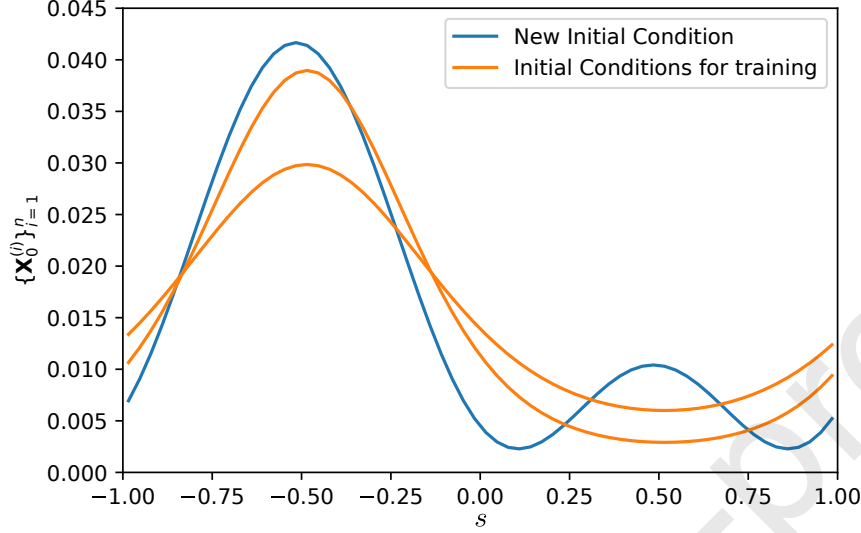


Figure 2: Sample initial conditions $\{\mathbf{X}_0^{(i)}\}_{i=1}^n$ for the Advection-Diffusion problem (orange) and an initial condition (blue) used for “extrapolative” predictions.

ability of the ARD prior to deactivate the vast majority of the right-hand-side feature functions and reveal a small subset of non-zero, salient terms. Both with $n = 32$ and $n = 64$ training data sequences, only parameters θ_c associated with first-order-interactions (Equation (22)) are activated. In particular, these are $\theta_{c,-3}^{(1)}$ and $\theta_{c,1}^{(1)}$ which are associated with the feature functions $X_{t,j-3}$ and $X_{t,j+1}$ respectively in Equation (22). This shares similarities with a finite-difference discretization scheme for the advection-diffusion and could be considered as an upwind scheme. The two identified coefficients do not form a centered difference operator but the center of the operator is shifted to the left and therefore takes into account the direction of the particle movement. As the value of the coefficients is not exactly the same the diffusive part is also captured.

Figure 7 depicts one of the inferred CG states $\mathbf{X}_{\Delta t}^{(i)}$ as well as the associated posterior uncertainty. Once the CG evolution law is learned, this state can be propagated into the future as detailed in section 2.5 in order to generate predictions. Indicative predictions (under “interpolative” conditions)

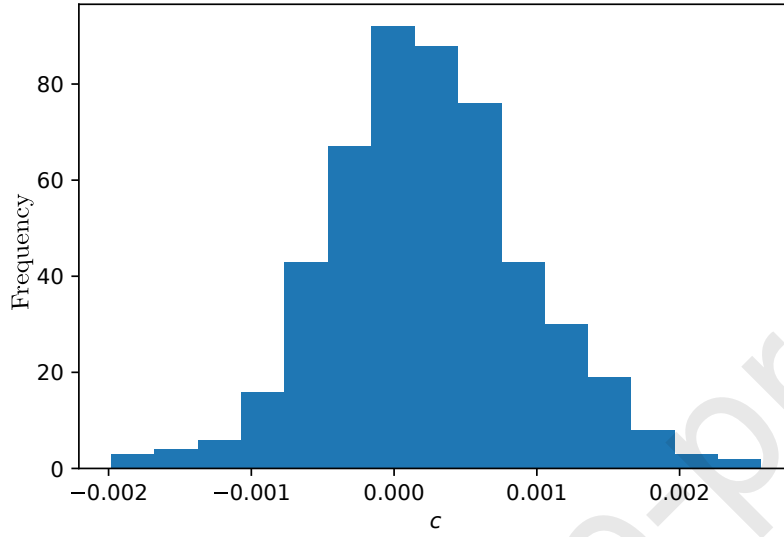


Figure 3: Histogram of the mass constraint c_1

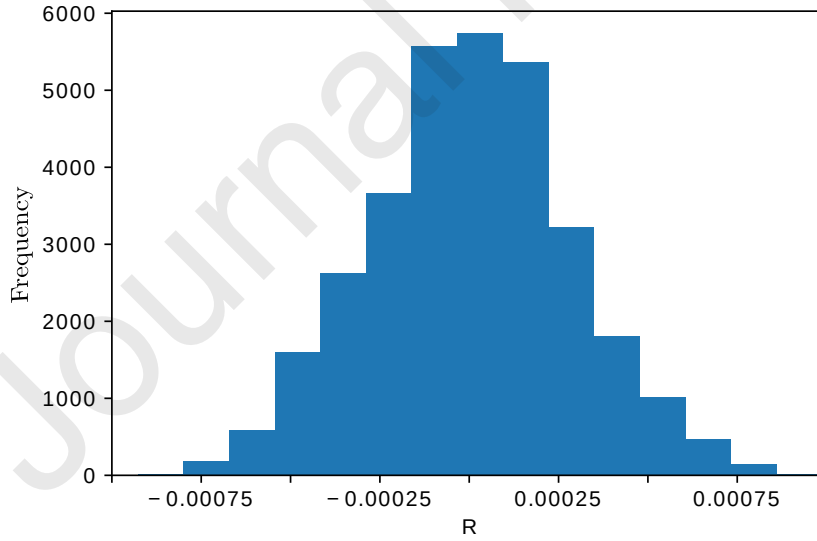


Figure 4: Histogram of the norm of the residual R_0

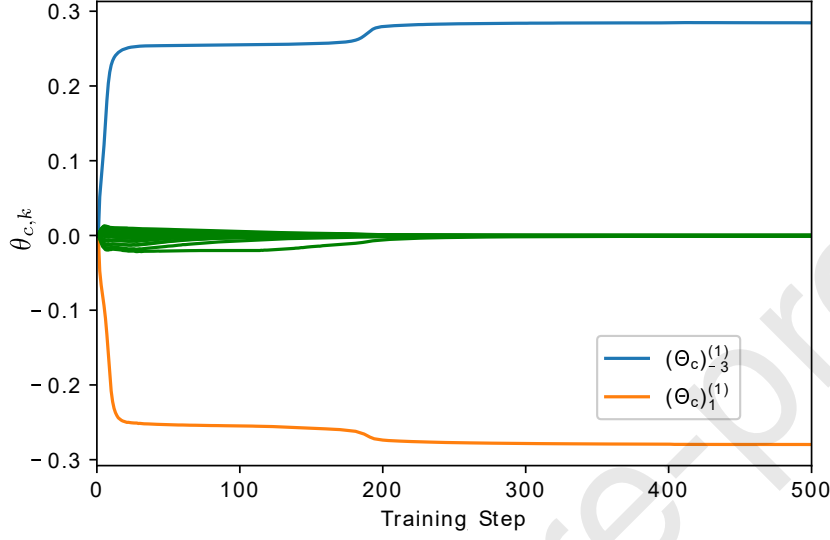


Figure 5: Evolution of a subset of θ_c parameters with respect to the iterations of the SVI for $n = 64$.

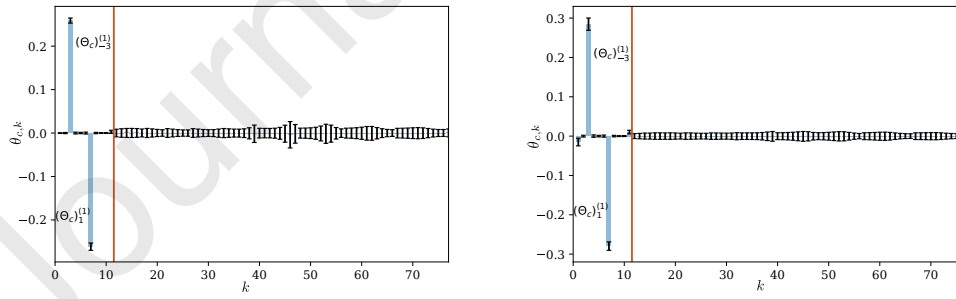


Figure 6: Comparison of the inferred parameters θ_c for $n = 32$ (left) and $n = 64$ (right) training data sequences. The black bars indicate ± 1 standard deviation. The red vertical line separates first- from second-order coefficients.

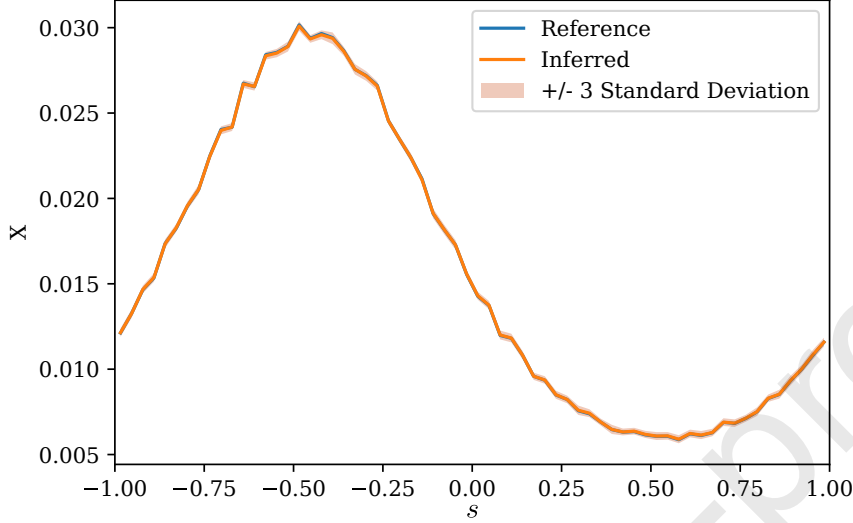


Figure 7: Inferred CG state $\mathbf{X}_{\Delta t}^{(i)}$ for a data sequence i . Reference is obtained by sorting the particles into bins according to their position.

can be seen in Figure 8 where the particle density $\rho_x(t, s)$ up to $25\Delta t$ into the future is drawn. The latter as well as the associated uncertainty bounds are estimated directly from the reconstructed FG states. As one would expect, the predictive uncertainty grows, the further into the future one tries to predict. Figure 9 compares the predictive performance as a function of the training data used i.e. $n = 32$ or $n = 64$. In both cases, the ground truth is envelopped and as one would expect, more training data lead to smaller uncertainty bounds.

640

We also tested the trained model (on $n = 64$) under “extrapolative” conditions i.e. for a different initial condition than the ones included in the training data (Figure 2). The predictive estimates in Figure 10 show very good agreement with the reference solution. It is important to point out that the model can correctly advect and diffuse the particle-bump initially introduced around $s = 0.5$ which suggests that the CG dynamics learned reflect the most important features of the problem.

648

Finally, in Figure 11, the evolution of the mass constraint into the future is depicted and good agreement with the target value is observed.

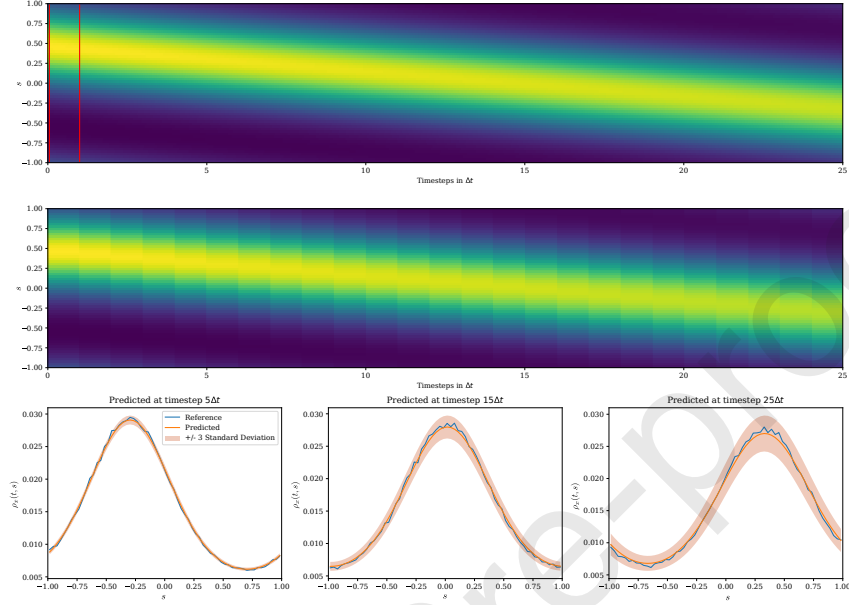


Figure 8: Prediction based on an initial condition contained in the training data. Top: Reference data (the vertical lines indicate the time instances with given data), Middle: Predictive posterior mean, Bottom: snapshots at three different time instances.

3.1.6. Burgers' system

The second test-case involved an FG system of $d_f = 500 \times 10^3$ particles which perform *interactive* random walks i.e. the jump performed at each fine-scale time-step $\delta t = 2.5 \times 10^{-3}$ depends on the positions of the other walkers. In particular we adopted interactions as described in Roberts (1989); Chertock and Levy (2001); Li et al. (2007) so as, in the limit (i.e. when $d_f \rightarrow \infty, \delta t \rightarrow 0, \delta s \rightarrow 0$), the particle density $\rho(s, t)$ follows the inviscid Burgers' equation:

$$\frac{\partial \rho}{\partial t} + \frac{1}{2} \frac{\partial \rho^2}{\partial s} = 0, \quad s \in (-1, 1). \quad (37)$$

For the CG description, 128 bins were employed i.e. $d_c = 128$ and a time step $\Delta t = 4$ (see Table 2). As compared with the previous case, we

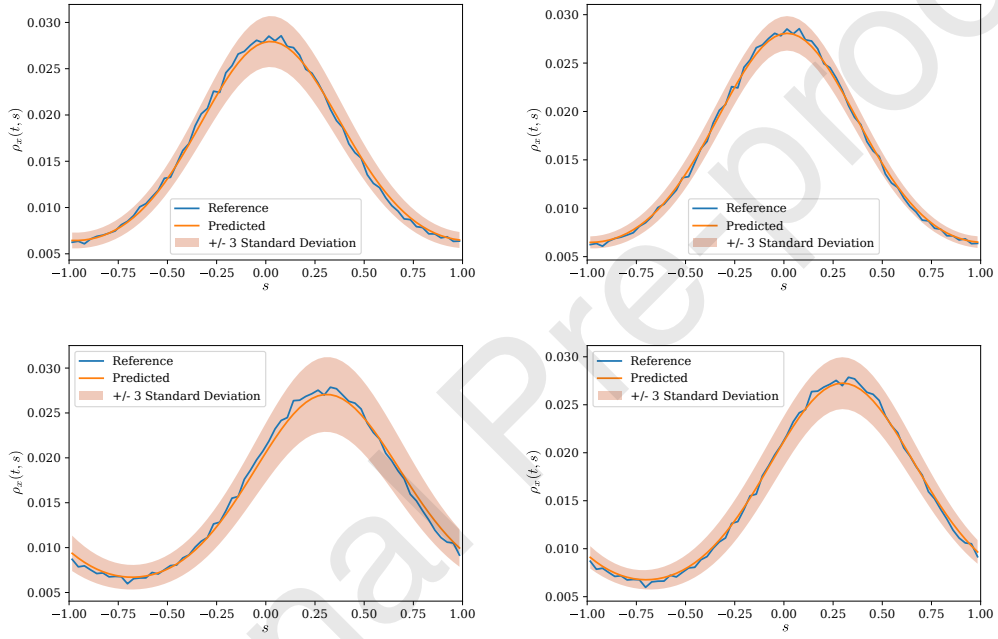


Figure 9: Comparison of the predictions for $n = 32$ (left) and $n = 64$ (right) at $15\Delta t$ (top) and $25\Delta t$ (bottom).

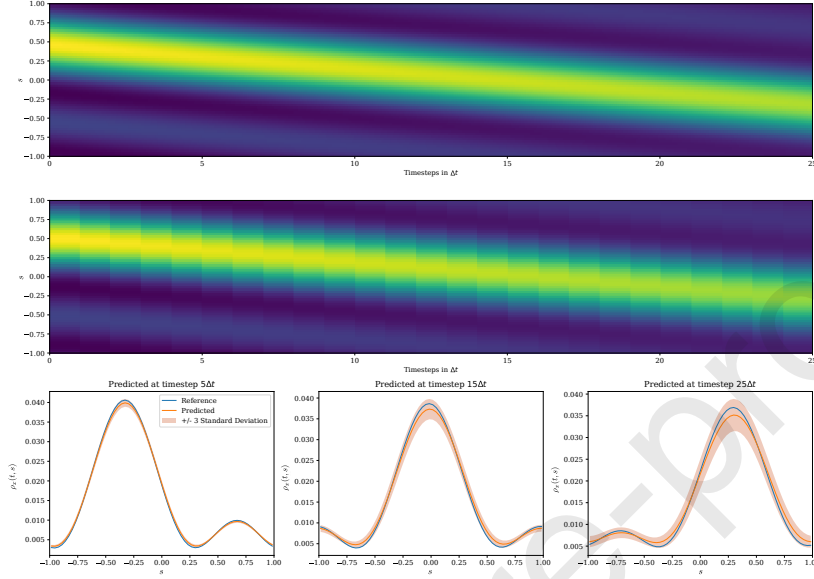


Figure 10: Prediction based on an initial condition NOT contained in the training data. Top: Reference data, Middle: Predictive posterior mean, Bottom: snapshots at three different time instances

660 enlarged the neighborhood size H in the first- and second-order interactions
to $H = 8$, which yielded $M = 170$ right-hand-side terms in Equation (22). We
662 incorporate virtual observables pertaining to the residuals $\hat{\mathbf{R}}_0$ with $\sigma_R^2 = 10^{-7}$
(Equation (7)) and the virtual observables $\hat{\mathbf{c}}_1$ pertaining to conservation-of-
664 mass constraint with $\sigma_c^2 = 10^{-10}$ (Equation (9)).

We employed $n = 32$, $n = 64$ and $n = 128$ time sequences for training that
666 were generated as detailed in section 3.1.4 with initial conditions $\{\mathbf{X}_0^{(i)}\}_{i=1}^n$
such as the ones seen in Figure 12. They were generated by randomizing the
668 width and height of a triangular profile.

Figure 13 provides a histogram of the function values of the conservation-
670 of-mass constraint $\{c_1(\mathbf{X}_{\Delta t}^{(i)})\}_{i=1}^n$ upon convergence. The small values suggest
that this has been softly incorporated in the CG states. A similar histogram
672 for the norm of the residuals $\{\mathbf{R}_0(\mathbf{X}^{(i)})\}_{i=1}^n$ is depicted in Figure 14 which

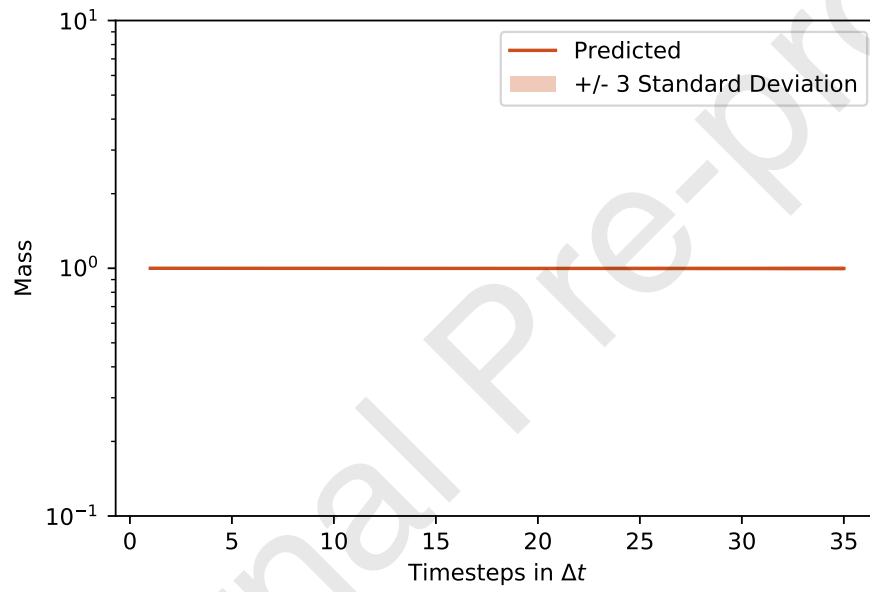


Figure 11: Evolution of the mass constraint (target value is 1) in time including future time-instants. "Predicted" corresponds to the posterior mean.

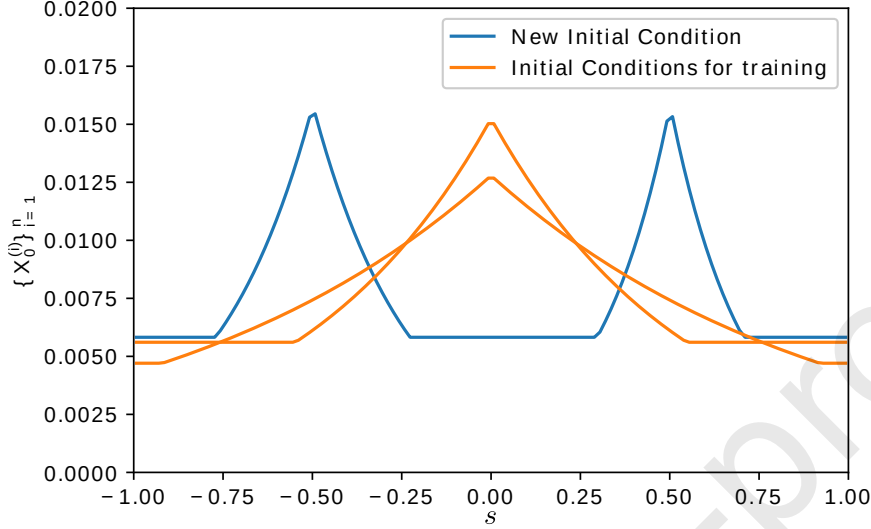


Figure 12: Sample initial conditions $\{\mathbf{X}_0^{(i)}\}_{i=1}^n$ for the Burgers' problem (orange) and an initial condition (blue) used for “extrapolative” predictions.

also suggests enforcement of the CG evolution with the parameters $\boldsymbol{\theta}_c$ learned from the data. The evolution of the posterior mean $\boldsymbol{\mu}_{\boldsymbol{\theta}_c}$ (Equation (34)) of (a subset of) these parameters over the iterations of the SVI is depicted in Figure 15. As in the previous example, in Figure 16 one can observe the ability of the ARD prior model to yield sparse solutions for the right-hand side of the CG evolution law. For all three training datasets with $n = 32, 64, 128$ time-sequences, only parameters $\boldsymbol{\theta}_c$ associated with second-order-interactions (Equation (22)) are activated. In particular, these are the negative coefficient $\boldsymbol{\theta}_{c,(0,0)}^{(2)}$ (in all three cases) as well as different second-order coefficients. In the cases of $n = 32$ and $n = 64$ two coefficients are found with positive mean and high posterior uncertainty, but they also have negative posterior correlation (correlation coefficient of -0.88). As all activated coefficients pertain to feature-functions involving the actual bin or bins to the left, the learned evolution law could be interpreted as an upwind scheme, which takes the direction of the Burgers' flow into account. Such schemes are considered advantageous for numerical simulations of fluid flows.

Figure 17 depicts one of the inferred CG states $\mathbf{X}_{\Delta t}^{(i)}$ as well as the asso-

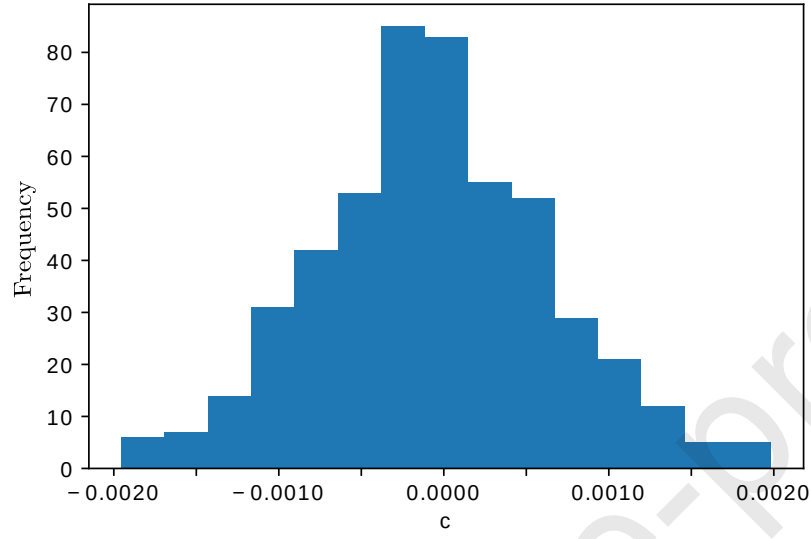


Figure 13: Histogram of the mass constraint c_1

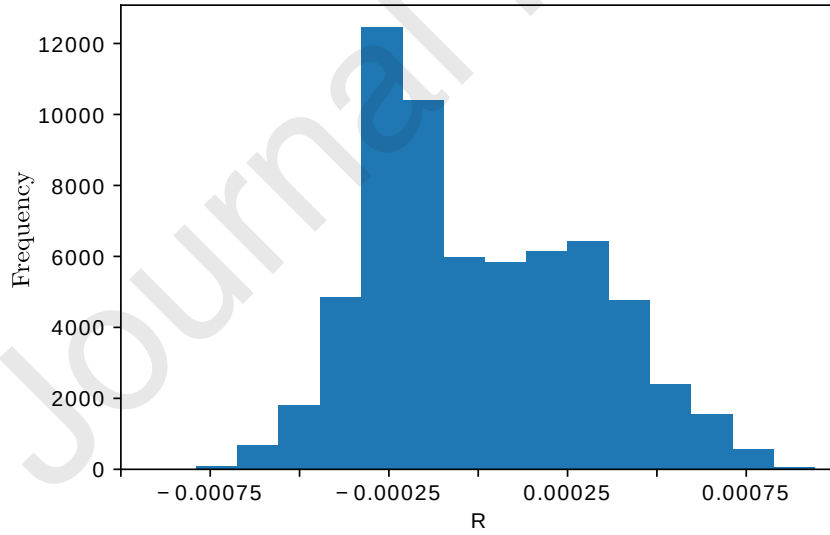


Figure 14: Histogram of the norm of the residual constraint R_0

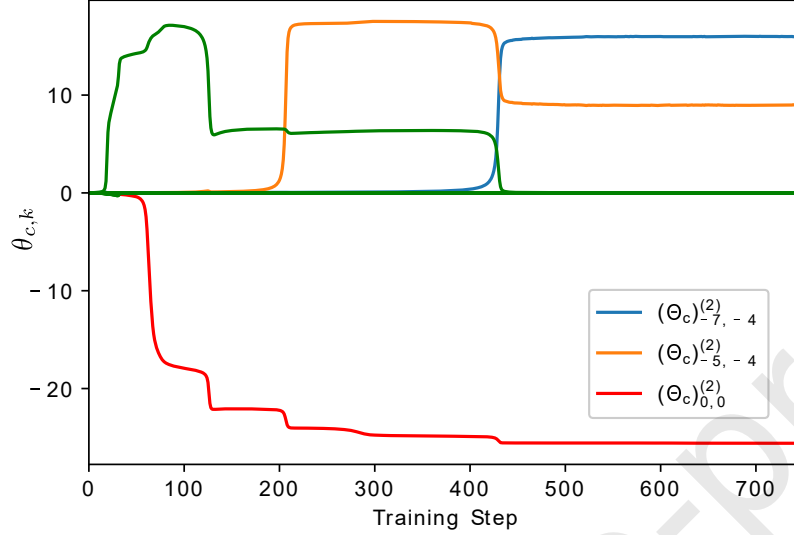


Figure 15: Evolution of a subset of θ_c parameters with respect to the iterations of the SVI for $n = 64$.

ciated posterior uncertainty. Given the learned CG dynamics, this state can be propagated into the future as detailed in section 2.5 in order to generate predictions. Indicative predictions (under “interpolative” conditions) can be seen in Figure 18 where the particle density up to $25\Delta t$ into the future is drawn. The latter as well as the associated uncertainty bounds are estimated directly from the reconstructed FG states. As in the previous example, the predictive uncertainty grows, the further into the future one tries to predict. Figure 19 compares the predictive performance as a function of the training data used i.e. $n = 32$ or $n = 64$. The increase in data leads for this example to a better fit of the posterior mean to the reference, which captures the location of the shock more precisely. The predictive uncertainty bounds are particularly large at the location of the shock which is the most challenging component in such systems.

We also test the trained model (on $n = 64$) under “extrapolative” conditions i.e. for a “bimodal” initial condition which was quite different from the ones included in the training data (Figure 12). The predictive estimates in Figure 20 show very good agreement with the reference solution. We want

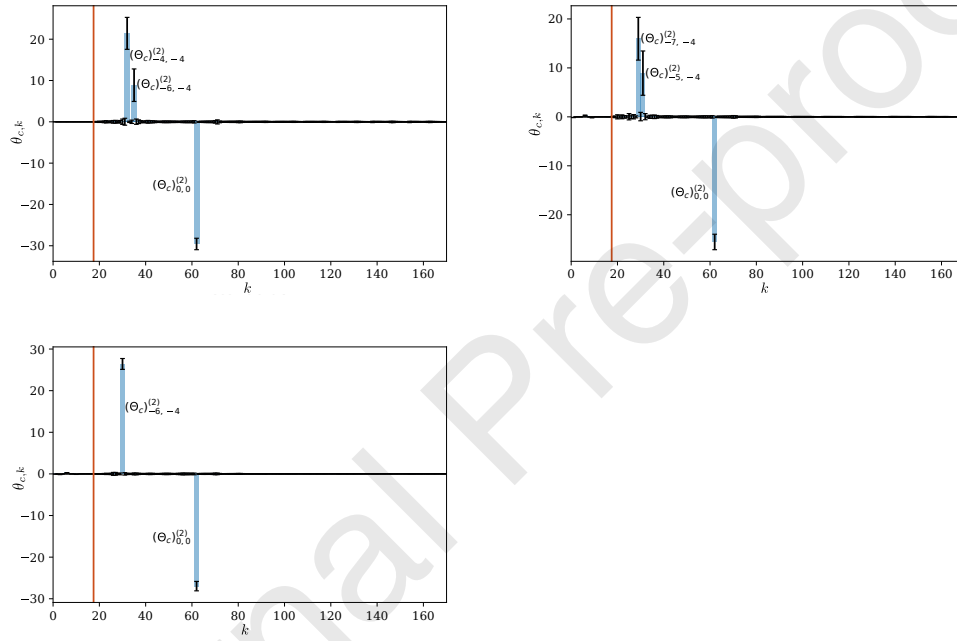


Figure 16: Comparison of the inferred parameters θ_c for $n = 32$ (top-left), $n = 64$ (top-right) and $n = 128$ (bottom-left) training data. The black bars indicate ± 1 standard deviation. The red vertical line separates first- from second-order coefficients.

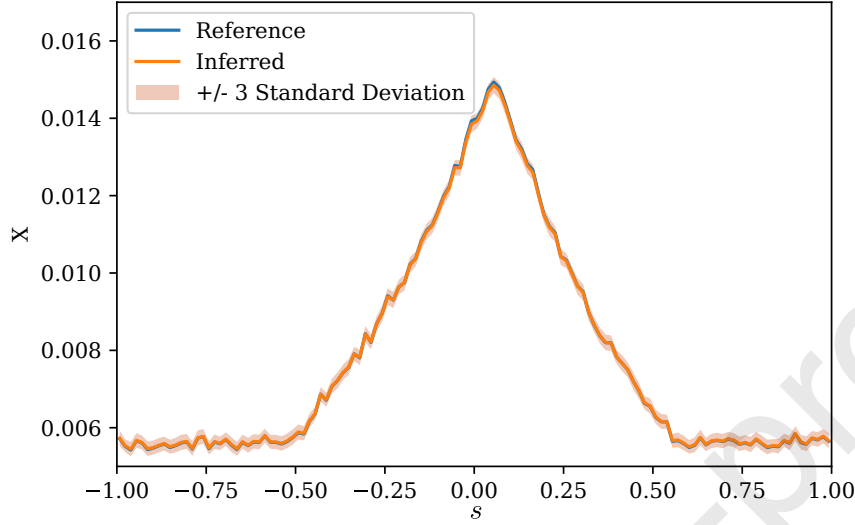


Figure 17: Example of inferred CG state $\mathbf{X}_{\Delta t}^{(i)}$ for data sequence i .

708 to point out that the trained model is capable of capturing the development,
 the position as well as the propagation of a shock front. Finally, in Figure
 710 21, the evolution of the mass constraint into the future is depicted and good
 agreement with the target value is observed.

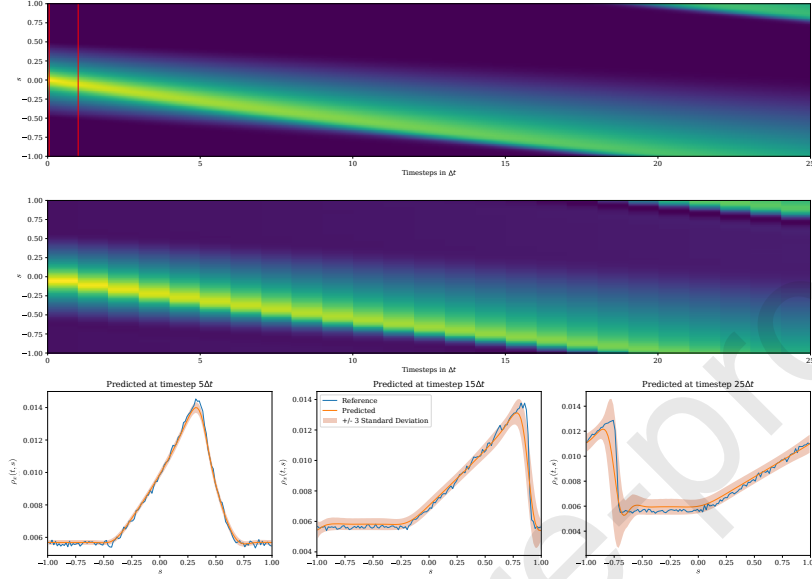


Figure 18: Prediction based on an initial condition contained in the training data. Top: Reference data (the vertical lines indicate the time instances with given data), Middle: Predictive posterior mean, Bottom: snapshots at three different time instances

712 3.2. Nonlinear Pendulum

In this final example we consider time sequences of images of a nonlinear
714 pendulum in two dimensions as in (Champion et al., 2019).

716 3.2.1. FG model

For the FG data we generate a series of black-and-white images of a mov-
ing disc tied on a string and forming a pendulum (see Figure 31). Each image
718 consists of 29×29 pixels each and each pixel's value was either 1 (occupied)
or -1 (unoccupied). Hence \mathbf{x}_t was a $d_f = 29^2 = 581$ -dimensional vector
720 of binary variables. The dynamics of the pendulum can be fully described
by the rotation angle y_t which follows a nonlinear, second-order ODE of the
722 form:

$$\ddot{y}_t + \sin(y_t) = 0 \quad (38)$$

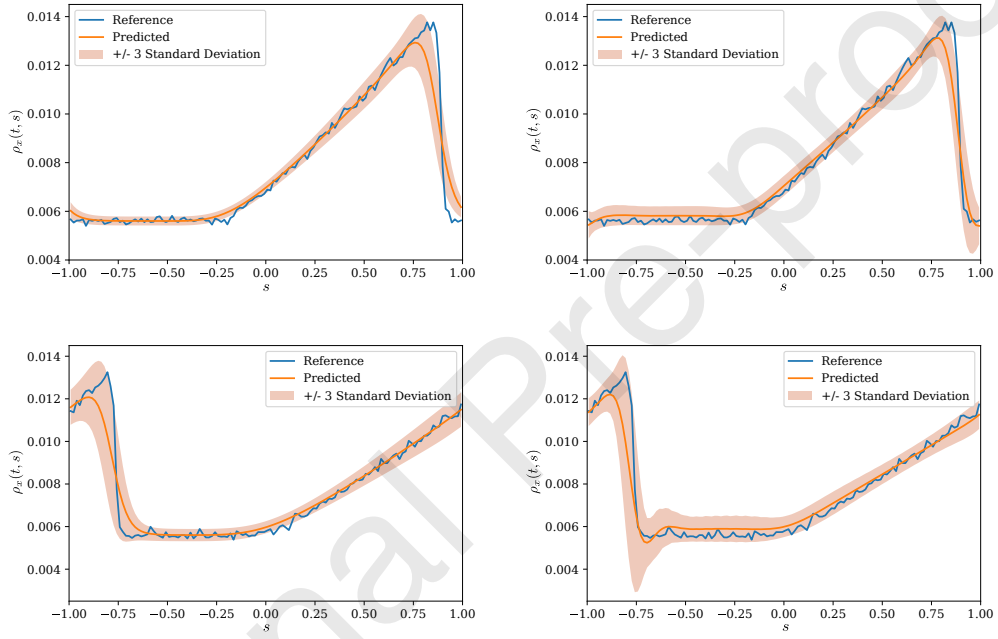


Figure 19: Comparison of the predictions for $n = 32$ (left) and $n = 64$ (right) training data at $15\Delta t$ (top) and $25\Delta t$ (bottom).

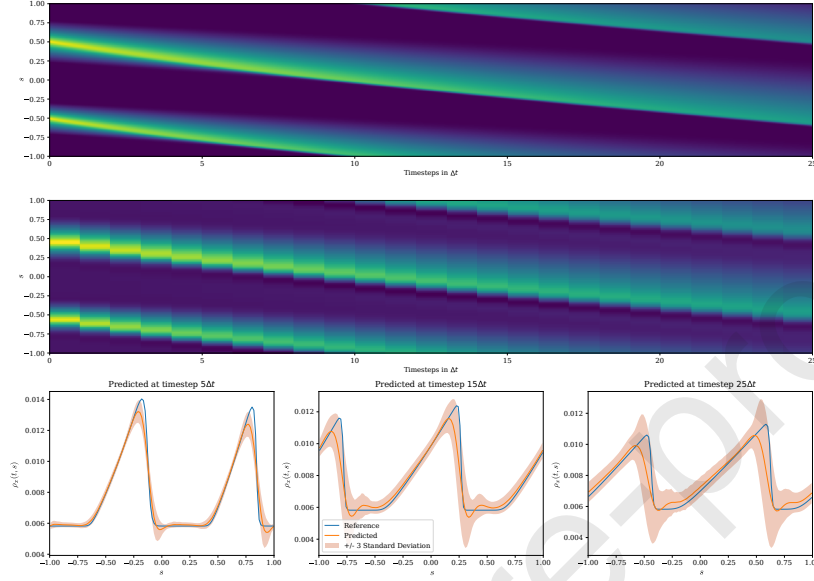


Figure 20: Prediction based on an initial condition NOT contained in the training data. Top: Reference data, Middle: Predictive posterior mean, Bottom: snapshots at three different time instances

The primary goal is to identify the right CG variables as well as CG dynamics solely from image data i.e. binary vectors $\{\hat{\mathbf{x}}_{0:T\Delta t}^{(i)}\}_{i=1}^n$ collected over T time-steps as the pendulum is initialized from n states/positions. The length of time sequences in the following numerical results was $T = 74$ and the CG time-step $\Delta t = 0.05^8$. We also considered the effect of missing data i.e. only observing a subset of the $T + 1$ values in each sequence and present respective results in Section 3.2.5.

3.2.2. CG variables and coarse-to-fine mapping

The only knowledge introduced a priori with regards to the CG variables \mathbf{X}_t is that $\dim(\mathbf{X}) = d_c = 2$. We intend to investigate procedures that can

⁸For the generation of images a *microscopic* time-step $\delta t = 0.01$ for the integration of Equation (38) was used.

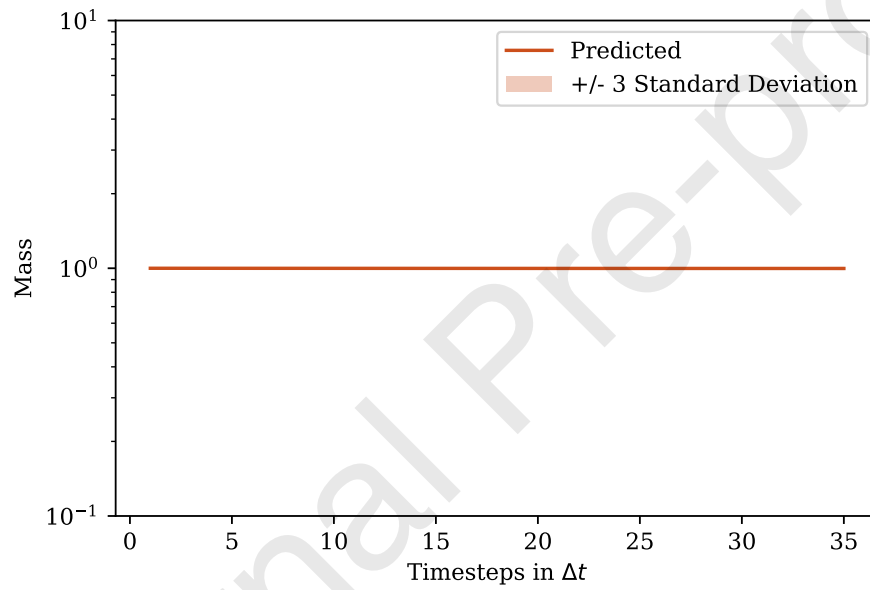


Figure 21: Evolution of the mass constraint (target value is 1) in time including future time-instants. "Predicted" corresponds to the posterior mean.

734 automatically identify d_c i.e. the number of CG variables. We note at this stage that such efforts could be guided by the ELBO \mathcal{F} (e.g. Equation (19)) which approximates the model evidence and therefore provides a natural
736 Bayesian score for comparing models with different numbers of CG variables.

The other pertinent model component is the coarse-to-fine map which
738 is enabled by the $p_{cf}(\mathbf{x}_t|\mathbf{X}_t)$ (section 2.3). To that end, we employed the following logistic model⁹:

$$p_{cf}(\mathbf{x}|\mathbf{X}) = \prod_{s=1}^{d_f} p_{cf}(x_s|\mathbf{X}) \quad (39)$$

740 with

$$p_{cf}(x_s|\mathbf{X}) = \begin{cases} \frac{1}{1 + \exp(-G_s(\mathbf{X}; \boldsymbol{\theta}_{cf}))} & \text{for } x_s = 1 \\ \frac{1}{1 + \exp(+G_s(\mathbf{X}; \boldsymbol{\theta}_{cf}))} & \text{for } x_s = 0 \end{cases} \quad (40)$$

where x_s is the value (1,0) of each of the pixels $s = 1, \dots, d_f$. For the
742 link functions $\{G_s\}_{s=1}^{d_f}$, we employed a deep neural net with weights $\boldsymbol{\theta}_{cf}$, the details of which are shown in Figure 22. One fully connected layer followed
744 by two transposed convolutional layers were found to be flexible enough to accurately represent the functions G_s . The CNNs were specifically chosen
746 because of their ability to extract/map features from/to images.

3.2.3. The CG evolution law and the virtual observables

748 With regards to the evolution law of the CG states $\mathbf{X}_t = \{X_{t,1}, X_{t,2}\}$, we postulate the following form:

$$\begin{aligned} \dot{X}_{t,1} &= F_1(\mathbf{X}_t, \boldsymbol{\theta}_c) = X_{t,2} \\ \dot{X}_{t,2} &= F_2(\mathbf{X}_t, \boldsymbol{\theta}_c) = \boldsymbol{\theta}_c^T \boldsymbol{\psi}(X_{t,1}) = \sum_{m=0}^M \theta_{c,m} \psi_m(X_{t,1}) \end{aligned} \quad (41)$$

750 where $\boldsymbol{\theta}_c$ denote the associated parameters. In total we employed $M = 101$ feature functions of the following type:

$$\psi_m(X) = \begin{cases} 1, & m = 0 \\ \sin(mX), & m = 1, \dots, M/2 = 50 \\ \cos((m - 50)X), & m = 51, \dots, M = 100 \end{cases} \quad (42)$$

⁹We omit the time-index t for clarity.

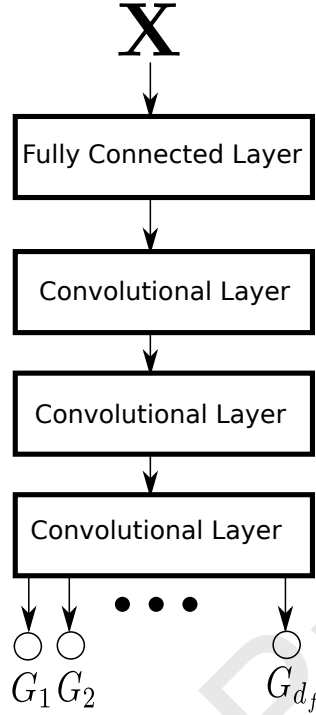


Figure 22: Deep neural net employed for the link functions G_s (Equation (39)). After one dense layer which $32 \cdot 7 \cdot 7$ nodes and rectified linear unit activation function (ReLU), two two-dimensional transposed convolutional layers with 32 filters and a kernel size of 3 as well as a ReLU activation function are applied followed by one-last two-dimensional transposed convolutional layers with one filter, kernel size 3 and without activation to generate the functions G_s .

752 The form of Equation (41) implies a second-order ODE where the second CG
 754 variable plays the role of the velocity. With regards to the parameters θ_c ,
 the sparsity-inducing ARD prior detailed in section 3.1.2 was employed.

756 To enforce the associated dynamics, we made use of the symplectic Euler
 time-discretization scheme, which is a first-order integrator, that is explicit
 in the first variable ($X_{t,1}$) and *implicit* in the other ($X_{t,2}$)¹⁰. The associated

¹⁰This corresponds to a multistep method in Equation (4) with $K = 1$, $a_0 = 1, a_1 = -1, \beta_0 = 0$ and $\beta_1 = -1$ for the explicit part and $K = 1$, $a_0 = 1, a_1 = -1, \beta_0 = -1$ and $\beta_1 = 0$ for the implicit part.

virtual observables (see Equation (6)) were enforced with $\sigma_R^2 = 10^{-5}$.

3.2.4. Inference and Learning

As in the previous examples (Equation (27)), the approximate posterior was factorized as:

$$q_\phi(\mathbf{X}_{0:T\Delta t}^{(1:n)}, \boldsymbol{\theta}_c, \boldsymbol{\tau}) = \left[\prod_{i=1}^n q_\phi(\mathbf{X}_{0:T\Delta t}^{(i)}) \right] q(\boldsymbol{\theta}_c) q(\boldsymbol{\tau}) \quad (43)$$

and closed-form updates were used for $q(\boldsymbol{\theta}_c)$ (see Equations (33) and (34)) and $q(\boldsymbol{\tau})$ (see Equation (35)).

SVI was applied for the posterior densities $q_\phi(\mathbf{X}_{0:T\Delta t}^{(i)})$ on the vector of the latent CG states $\mathbf{X}_{0:T\Delta t}^{(i)}$ which we approximated with multivariate Gaussians. Since the posterior reveals the fine-to-coarse map which apart from insight can be used for predictive purposes as well, we employed an *amortized* version of SVI ((Kingma and Welling, 2014)) i.e. explicitly accounted for the dependence of each $q_\phi(\mathbf{X}_{0:T\Delta t}^{(i)})$ on the corresponding FG observables $\hat{\mathbf{x}}_{0:T\Delta t}^{(i)}$ i.e.:

$$q_\phi(\mathbf{X}_{0:T\Delta t}^{(i)}) = \mathcal{N}(\boldsymbol{\mu}_\phi(\hat{\mathbf{x}}_{0:T\Delta t}^{(i)}), \mathbf{S}_\phi(\hat{\mathbf{x}}_{0:T\Delta t}^{(i)})) \quad (44)$$

The parameters ϕ were the weights of a deep convolutional neural net, the architecture of which is shown in Figure 23. This was chosen because it mirrors the DNN architecture employed for the coarse-to-fine map in Figure 22.

Finally it should be mentioned that the "slowness" prior was employed on the hidden states $\mathbf{X}_{0:T\Delta t}^{(1:n)}$ as described in Equation (16)¹¹. Maximum-likelihood estimates for the hyperparameter $\sigma_{\mathbf{X}}^2$ were employed which readily arise by differentiating the ELBO \mathcal{F} and which yield the following update equation:

$$\sigma_{\mathbf{X}}^2 = \frac{1}{n T d_c} \sum_{i=1}^n \sum_{l=0}^{T-1} \mathbb{E}_{q_\phi(\mathbf{X}_{0:T\Delta t}^{(i)})} \left[\left| \mathbf{X}_{(l+1)\Delta t}^{(i)} - \mathbf{X}_{l\Delta t}^{(i)} \right|^2 \right] \quad (45)$$

Maximum likelihood estimates were also obtained for the parameters $\boldsymbol{\theta}_{cf}$ (Equation (39)) by numerically differentiating the ELBO \mathcal{F} and performing Stochastic Gradient Ascent (SGA).

¹¹For the prior distribution $p_{c,0}(\mathbf{X}_0^{(i)})$ a Gaussian mixture distribution with means +1.5 and -1.5 and standard deviation 1.5 was used.

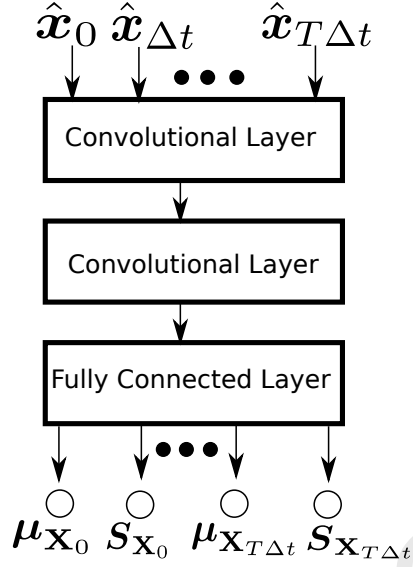


Figure 23: DNN architecture for approximate posterior q_ϕ . The input consists of a time series of pictures of the pendulum and can therefore be considered to be three-dimensional, where the first and second dimension are the number of pixels and the third dimension is the number of time steps available for training. This input is given to a three-dimensional convolutional layer with kernel size $(3, 3, 2)$, 32 filters and a ReLU activation followed by another three-dimensional convolutional layer with kernel size 2 in each dimension, 64 filters and a ReLU activation. The last layer is a fully connected layer with $2d_c \cdot T$ nodes and without activation to generate the mean and variance values for each time step of the inferred \mathbf{X} coordinates.

784 A general summary of the steps involved for the inference procedure is-
can be found in Algorithm 4. For the implementation we made use of the
Tensorflow framework (Abadi et al., 2016).

Algorithm 4: Algorithm for the Pendulum system

Result: $\phi, q(\theta_c), q(\tau), \theta_{cf}, \sigma_{\mathbf{X}}$
Data: $\hat{\mathbf{x}}_{0:T\Delta t}^{(1:n)}$

- 1 Initialize all required parameters;
- 2 Set iteration counter w to zero;
- 3 **while** $\|ELBO_w - ELBO_{w-1}\|^2 > \epsilon$ **do**
- 786 4 Update the parameters θ_{cf} and ϕ by SGA of the ELBO (Equation (19)) ;
- 5 update $q(\theta_c)$ according to Equation (33) and Equation (34) ;
- 6 update $q(\tau)$ according to Equation (35) ;
- 7 update the parameter $\sigma_{\mathbf{X}}$ according to Equation (45);
- 8 update the iteration counter by one;
- 9 **end**

3.2.5. Results

Each data sequence $\hat{\mathbf{x}}_{0:T\Delta t}^{(i)}$ used consisted of 75 images, i.e. $T = 74$, generated with a time-step $\Delta t = 0.05$ (Figure 24). We investigated two cases for the number of data sequences i.e. $n = 16$ and $n = 64$. The data generation involved sampling uniformly the initial angle $y_0 \in [-\pi, \pi]$ and assuming zero initial velocity i.e. $\dot{y}_0 = 0$. We emphasize that none of the data sequences contained a complete oscillation of the pendulum i.e. always partial trajectories were observed.

Figure 25 indicates the posterior means of the inferred θ_c that parametrize the CG evolution law (Equation (41)) for $n = 16$ and $n = 64$. Of the 101 possible terms, only 2 are activated due the ARD prior.

Figure 26 illustrates trajectories in the two-dimensional CG state-space obtained with various initial conditions for the CG model identified with $n = 16$ and $n = 64$ data sequences. The blue curves correspond to “interpolative” settings i.e. to the CG states of an observed sequence of images, whereas the orange curves to “extrapolative settings” i.e. to the CG states inferred by initializing the pendulum from an arbitrary position *not* contained in the training data. In Figure 27 the predicted evolution in time of both coarse-grained variables is shown. The periodic nature of the CG dynamics is obvious, even though the CG state variables implicitly identified do not correspond to the natural ones i.e. y_t and \dot{y}_t .

This can be seen in Figure 28 where for data-sequences $\mathbf{x}_{0:T\Delta t}^{(i)}$ (corresponding to the pendulum at various positions i.e. angles $y_{0:T\Delta t}$), we compute

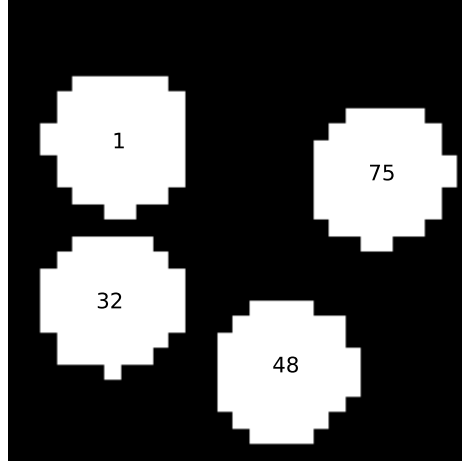


Figure 24: Indicative positions of the pendulum in a data sequence $\hat{\mathbf{x}}_{0:T\Delta t}^{(i)}$. The number indicates the corresponding time-step.

from the approximate posterior $q_\phi(\mathbf{X}_{0:T\Delta t}^{(i)}|\mathbf{x}_{0:T\Delta t}^{(i)})$ (Equation (44)) the mean of the corresponding CG states $\mathbf{X}_{0:T\Delta t}^{(i)}$ as well as the (in this case negligible) standard deviation. For each time instant $l = 0, 1, \dots, T$, we plot the pairs of $y_{l\Delta t}$ and (the mean of) $X_{l\Delta t,1}$ (i.e. the first of the CG variables identified) to show the relation between the two variables. While it is obvious from the scales that the first CG variable identified is *not* the angle, it appears to be isomorphic to y . The latter property persists for $n = 64$ even though the sign of the relation has been reversed. The difference between the first CG variable identified and the natural angle y explains the difference between the CG evolution law identified (Figure 25) and the reference one Equation (38).

Figure 29 provides predictive estimates of the position of the center of mass in time. These were obtained by propagating the CG variables in time and for each time instant, sampling p_{cf} for corresponding images \mathbf{x} . From the latter, the center of mass was computed from the activated pixels i.e. the pixels with value 1. Naturally, predictive uncertainty arises due the stochasticity in the initial conditions of \mathbf{X} as well as in p_{cf} . The latter is quantified by the standard deviation and plotted in Figure 29. As in the

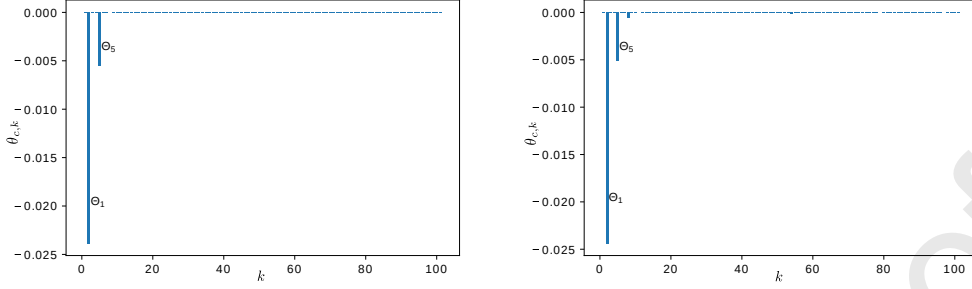


Figure 25: Posterior means of the inferred θ_c that parametrize the CG evolution law (Equation (41)) for $n = 16$ (left) and $n = 64$ (right) training data.

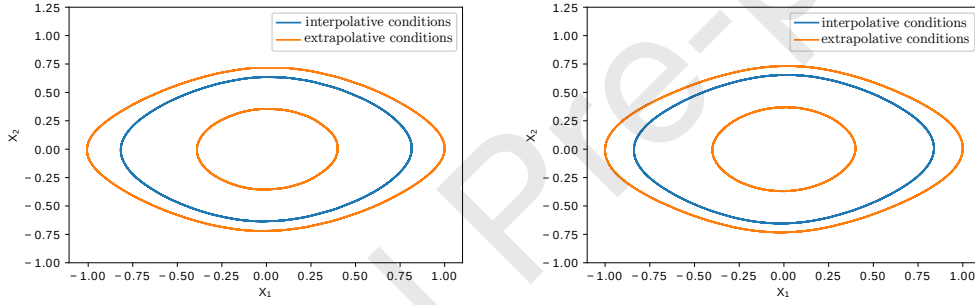


Figure 26: Comparison of trajectories in state space \mathbf{X} of the CG dynamics learned for $n = 16$ (left) and $n = 64$ (right) training data.

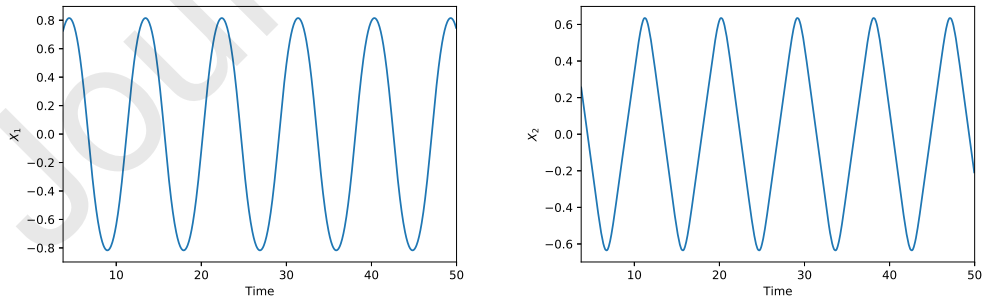


Figure 27: Predicted posterior mean of CG state variables \mathbf{X}_t

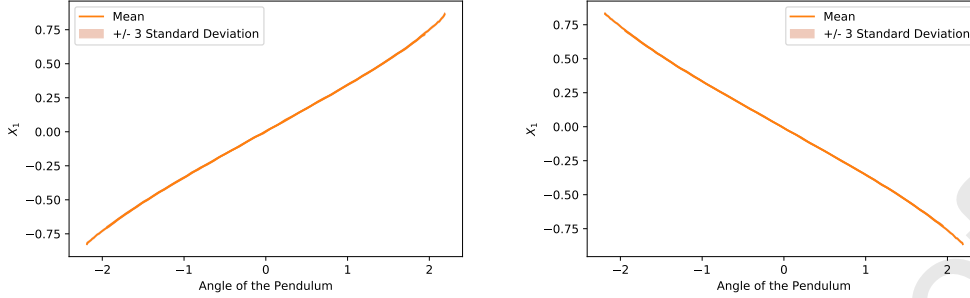


Figure 28: Mapping between the angle of the pendulum and the coarse-grained coordinates for 32 training data and 64 (right) training data.

previous examples, the predictive uncertainty grows, albeit modestly, with time.

Figure 30 depicts predictions in time for two pixels in the image. One can clearly distinguish the change-points i.e. when the pendulum crosses the pixel and its value is changed from 0 to 1 as well as the predictive uncertainty which is concentrated at those change-points. This demonstrates one of the strengths of our approach as due to the coarse-to-fine mapping the whole FG state is reconstructed and every observable can be computed together with the associated predictive uncertainty.

Finally, Figure 31 compares actual images obtained by the reference dynamics of the pendulum with the predictive posterior mean obtained by the CG model and p_{cf} trained on the data. Even though these extend up to 875 time-steps i.e. more than 11 times longer than the time-window over which observations were available, they match the reference quite accurately, a strong indication that the right CG variables and CG dynamics have been identified. An animation containing all frames can be found by following this [link](#).

3.2.6. Missing data

The generative nature of the proposed model makes it highly suitable for handling missing FG data either in the form of partial observations of the FG state vector \mathbf{x}_t or observations over a portion/subset of the time-sequence considered. We investigate the latter case in this section but note that in both situations the only modification required is removing the likelihood terms corresponding to the missing data from Equation (13).

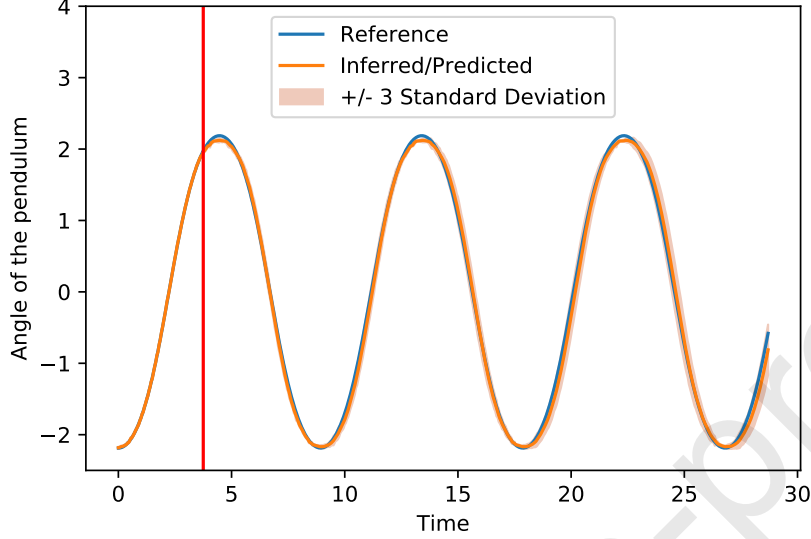


Figure 29: Inferred/Predicted evolution of the center of mass of the pendulum. The vertical line separates the inferred states from the predictions

In particular, we investigated the performance of the model when every second FG state \mathbf{x}_t in the training sequences was *not* observed i.e. the FG observables consisted of $\{\mathbf{x}_0^{(i)}, \mathbf{x}_{2\Delta t}^{(i)}, \mathbf{x}_{4\Delta t}^{(i)}, \dots, \mathbf{x}_{T\Delta t}^{(i)}\}$ for each data sequence i (where $T = 74$ as before). As one would expect, fewer observations lead to higher inferential uncertainties as seen when comparing Figure 28 (fully observed case) with Figure 32 (partially observed case). More importantly, fewer observations lead to higher predictive uncertainty as seen when comparing the predictions for the center of pendulum in Figure 29 (fully observed case) with Figure 33 (partially observed case).

4. Conclusions

We proposed a probabilistic generative model for the automated discovery of coarse-grained variables and dynamics based on fine-grained simulation data. The FG simulation data are augmented in a fully Bayesian fashion by virtual observables that enable the incorporation of physical constraints at the CG level that appear in the form of equalities. These could be residuals of the CG evolution law or more importantly conservation laws that are available when CG variables have physical meaning. This is particu-

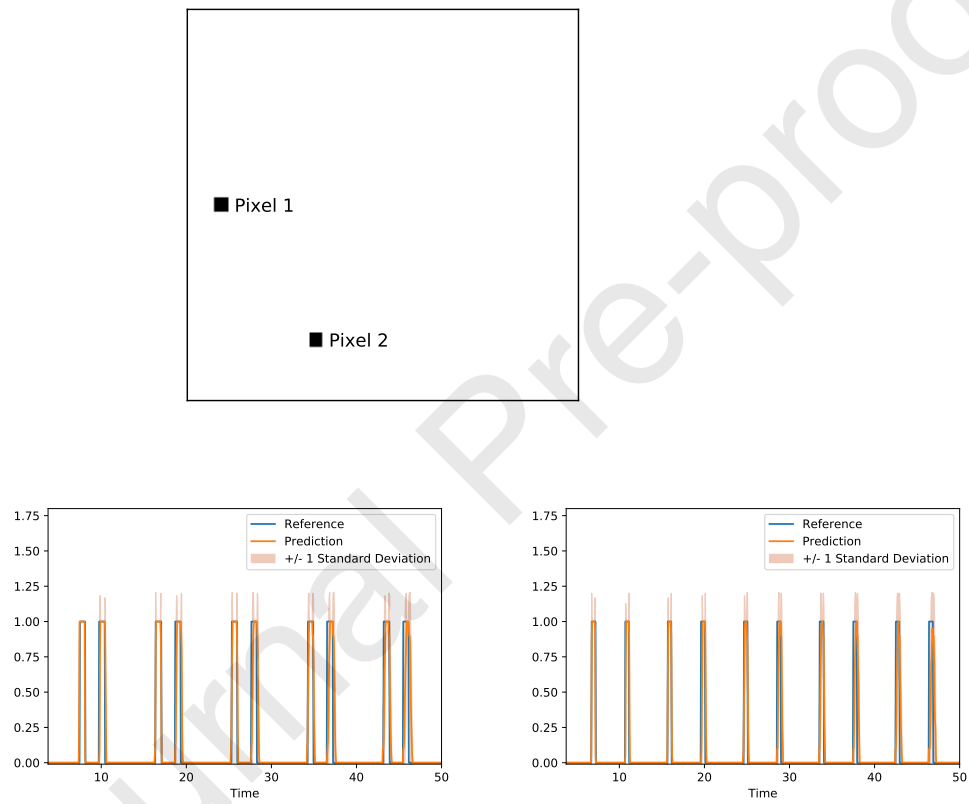


Figure 30: Predicted time history of a single pixel: Pixel 1 (left) and Pixel 2 (right)

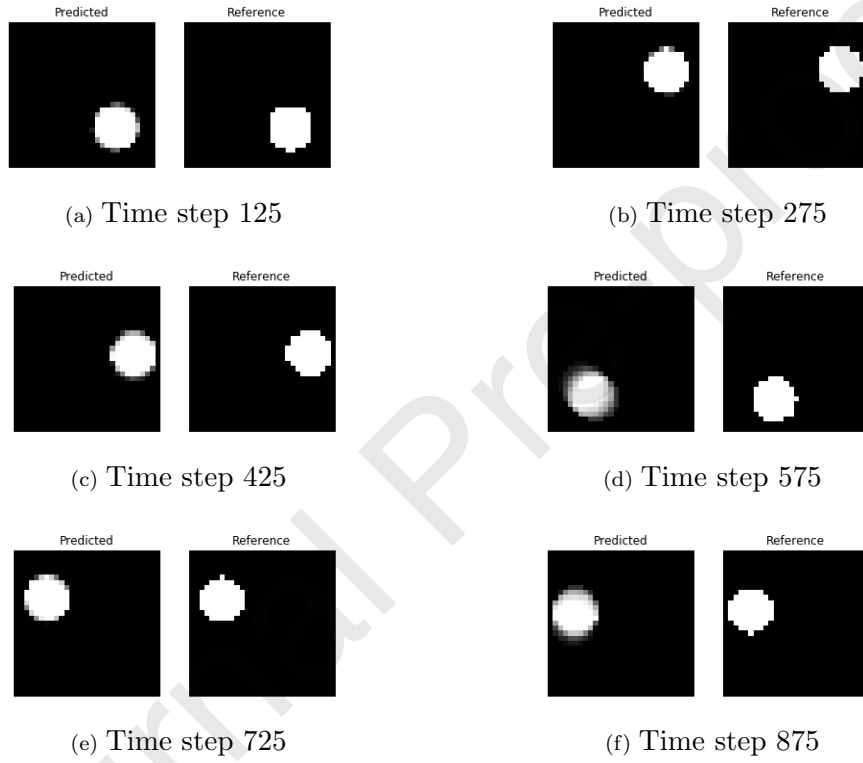


Figure 31: Predictive posterior means of images of the pendulum compared to the reference data

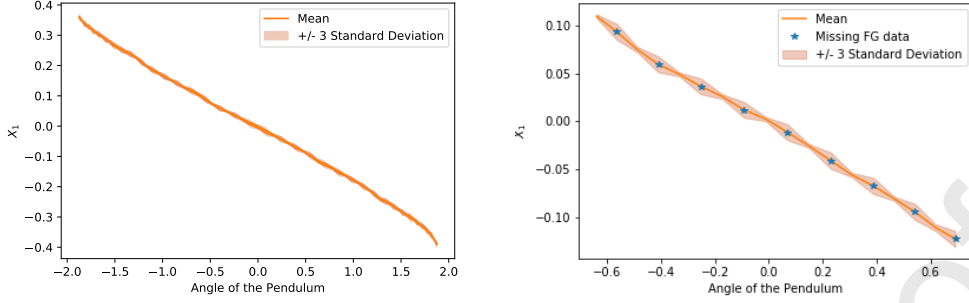


Figure 32: Effect of missing data on the CG variables. The figure on the right is zoomed-in to show the higher uncertainty associated with CG states with missing data

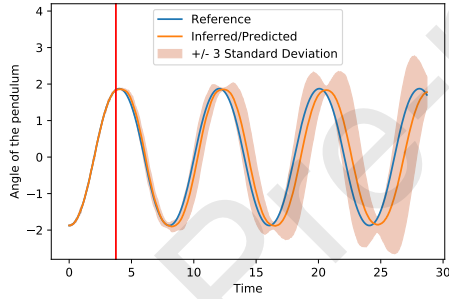


Figure 33: Inferred/Predicted evolution of the center of mass of the pendulum for the missing data case. The vertical line separates the inferred states from the predictions

larily important in the context of physical modeling as in many cases such
 870 domain knowledge is a priori available and its inclusion can, not only re-
 duce the amount of training data, but endow the CG model learned with the
 872 necessary features that would allow it to provide accurate predictions in out-
 of-distribution settings. Our approach learns simultaneously a coarse-to-fine
 874 mapping and an evolution law for the coarse-grained dynamics by employing
 probabilistic inference tools for the latent variables and model parameters.
 876 The use of deep neural nets for the former component can endow great expres-
 siveness and flexibility. The concept of sparsity, which is invoked in learning
 878 CG dynamics from a large vocabulary of right-hand-side terms, is readily
 incorporated using sparsity-inducing Bayesian priors without any hyperpa-
 880 rameter tuning. Furthermore, appropriate priors can promote the discovery
 of slow-varying CG variables which better capture the macroscopic features

of the system. As a result of the aforementioned characteristics, the framework can learn from *Small Data* (i.e. shorter and fewer FG time-sequences) which is a crucial advantage in multiscale models where the simulation of the FG dynamics is expensive and slow in exploring the state-space. The model proposed was successfully tested on coarse-graining tasks from different areas. In all three examples, the method performed well under interpolative, and more importantly under extrapolative settings i.e. in cases where initial conditions different from the ones seen during training, are prescribed. Partial or incomplete FG observations can readily be handled due to its generative nature. Moreover, as it is able to reconstruct the entire FG state vector at any future time instant, it is capable of producing predictions of any FG observable of interest as well as quantify the associated predictive uncertainty.

There exists various possibilities to extend the proposed framework, both methodologically as well as in terms of applications. In the latter case and apart from using it for predictive purposes, the CG model learned could also be employed in optimization and control applications. On the methodological front an obvious extension would be to account for the virtual observables at future time-instants as well. This would ensure their enforcement by future CG states but would unavoidably complicate their simulation as a probabilistic inference scheme would need to be employed in order to draw samples.

Another important question pertains to the stability of the CG dynamics identified (Pan and Duraisamy, 2020). This is not currently guaranteed in the discretized nor in the continuous version. This could potentially be achieved by an a-priori parametrization of the CG dynamics in a way that guarantees stability which could in turn reduce the expressivity of the model. Finally, we note that, in our opinion, the most difficult question in coarse-graining multiscale systems, is finding the number of CG state variables that are needed. In physics problems, very often one has an idea of which variables would be suitable either based on the analysis-objectives and/or physical insight. Almost never though does one have a guarantee that these variables are sufficient. Assuming they are, the problem then reduces to finding the appropriate closures (i.e. right-hand sides in the CG dynamics) which is the problem we try to address in this paper. The discovery of additional, potentially non-physical CG state variables, would require additional advances for which we believe the ELBO, i.e. the (approximate) model evidence, could serve as the guiding objective.

References

- D. Givon, R. Kupferman, A. Stuart, Extracting Macroscopic Dynamics: Model Problems and Algorithms, *Nonlinearity* (2004).
- W. Bialek, *Biophysics: Searching for Principles*, Princeton University Press, 2012.
- M. Alber, A. Buganza Tepole, W. R. Cannon, S. De, S. Dura-Bernal, K. Garikipati, G. Karniadakis, W. W. Lytton, P. Perdikaris, L. Petzold, E. Kuhl, Integrating machine learning and multiscale modeling—perspectives, challenges, and opportunities in the biological, biomedical, and behavioral sciences, *NPJ digital medicine* 2 (2019) 115. doi:10.1038/s41746-019-0193-y.
- Z. Ghahramani, Probabilistic machine learning and artificial intelligence, *Nature* 521 (2015) 452–459. URL: <http://www.nature.com/nature/journal/v521/n7553/full/nature14541.html>. doi:10.1038/nature14541.
- Y. LeCun, Y. Bengio, G. Hinton, Deep learning, *nature* 521 (2015) 436–444.
- P.-S. Koutsourelakis, N. Zabaras, M. Girolami, Big data and predictive computational modeling, *Journal of Computational Physics* 321 (2016) 1252–1254.
- G. Marcus, E. Davis, *Rebooting AI: Building Artificial Intelligence We Can Trust*, Pantheon, 2019.
- P. Stinis, T. Hagge, A. M. Tartakovsky, E. Yeung, Enforcing constraints for interpolation and extrapolation in generative adversarial networks, *Journal of Computational Physics* 397 (2019) 108844.
- P.-S. Koutsourelakis, E. Bilionis, Scalable Bayesian Reduced-Order Models for Simulating High-Dimensional Multiscale Dynamical Systems, *Multiscale Modeling & Simulation* 9 (2011) 449–485. doi:10.1137/100783790.
- I. Kevrekidis, C. Gear, J. Hyman, P. Kevrekidis, O. Runborg, K. Theodoropoulos, Equation-free multiscale computation: enabling microscopic simulators to perform system-level tasks, *Communications in Mathematical Sciences* 1 (2003) 715–762.

- P. J. Schmid, Dynamic mode decomposition of numerical and experimental data, *Journal of Fluid Mechanics* 656 (2010) 5–28. URL: <https://www.cambridge.org/core/journals/journal-of-fluid-mechanics/article/dynamic-mode-decomposition-of-numerical-and-experimental-data/AA4C763B525515AD4521A6CC5E10BD4>. doi:10.1017/S0022112010001217.
- M. O. Williams, I. G. Kevrekidis, C. W. Rowley, A Data-Driven Approximation of the Koopman Operator: Extending Dynamic Mode Decomposition, *Journal of Nonlinear Science* 25 (2015) 1307–1346. doi:10.1007/s00332-015-9258-5.
- H. Wu, F. Noé, Variational approach for learning Markov processes from time series data, arXiv:1707.04659 [math, stat] (2017). URL: <http://arxiv.org/abs/1707.04659>.
- G. Froyland, G. A. Gottwald, A. Hammerlindl, A Computational Method to Extract Macroscopic Variables and Their Dynamics in Multiscale Systems, *SIAM Journal on Applied Dynamical Systems* 13 (2014) 1816–1846. URL: <https://epubs.siam.org/doi/abs/10.1137/130943637>. doi:10.1137/130943637.
- L. Felsberger, P. Koutsourelakis, Physics-constrained, data-driven discovery of coarse-grained dynamics, *Communications in Computational Physics* 25 (2019) 1259–1301. doi:10.4208/cicp.0A-2018-0174.
- M. Schöberl, N. Zabaras, P.-S. Koutsourelakis, Predictive coarse-graining, *Journal of Computational Physics* 333 (2017) 49–77.
- M. Raissi, P. Perdikaris, G. E. Karniadakis, Physics informed deep learning (part i): Data-driven solutions of nonlinear partial differential equations, arXiv preprint arXiv:1711.10561 (2017).
- M. Raissi, P. Perdikaris, G. E. Karniadakis, Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations, *Journal of Computational Physics* 378 (2019) 686–707.

- 982 Y. Yang, P. Perdikaris, Conditional deep surrogate models for stochastic,
high-dimensional, and multi-fidelity systems, *Computational Mechanics*
984 (2019) 1–18.
- A. Mardt, L. Pasquali, H. Wu, F. Noé, VAMPnets for deep learn-
986 ing of molecular kinetics, *Nature Communications* 9 (2018). URL:
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5750224/>. doi:10.
988 1038/s41467-017-02388-1.
- H. Wu, A. Mardt, L. Pasquali, F. Noe, Deep generative markov state models,
990 in: *Advances in Neural Information Processing Systems*, 2018, pp. 3975–
3984.
- 992 L. Duncker, G. Böhner, J. Boussard, M. Sahani, Learning interpretable
continuous-time models of latent stochastic dynamical systems, *arXiv*
994 preprint arXiv:1902.04420 (2019).
- C. Grigo, P.-S. Koutsourelakis, A physics-aware, probabilistic machine learn-
996 ing framework for coarse-graining high-dimensional systems in the small
data regime, *arXiv preprint arXiv:1902.03968* (2019).
- 998 Y. Pantazis, I. Tsamardinos, A unified approach for sparse dynamical system
inference from temporal measurements, *Bioinformatics* 35 (2019) 3387–
1000 3396. URL: [https://academic.oup.com/bioinformatics/article/35/](https://academic.oup.com/bioinformatics/article/35/18/3387/5305020)
18/3387/5305020. doi:10.1093/bioinformatics/btz065.
- 1002 S. L. Brunton, J. L. Proctor, J. N. Kutz, Discovering governing equations
from data by sparse identification of nonlinear dynamical systems, *Pro-*
1004 *ceedings of the National Academy of Sciences* 113 (2016) 3932–3937.
- E. Kaiser, J. N. Kutz, S. L. Brunton, Sparse identification of nonlinear
1006 dynamics for model predictive control in the low-data limit, *Proceedings*
of the Royal Society A 474 (2018) 20180335.
- 1008 K. Champion, B. Lusch, J. N. Kutz, S. L. Brunton, Data-driven discovery
of coordinates and governing equations, *arXiv preprint arXiv:1904.02107*
1010 (2019).
- J. Ohkubo, Nonparametric model reconstruction for stochastic differen-
1012 tial equations from discretely observed time-series data, *Physical Re-*
view E 84 (2011) 066702. URL: [https://link.aps.org/doi/10.1103/](https://link.aps.org/doi/10.1103/PhysRevE.84.066702)
1014 *PhysRevE*.84.066702. doi:10.1103/PhysRevE.84.066702.

- 1016 S. Klus, F. Nüske, P. Koltai, H. Wu, I. Kevrekidis, C. Schütte, F. Noé, Data-
Driven Model Reduction and Transfer Operator Approximation, *Journal of*
1018 *Nonlinear Science* 28 (2018) 985–1010. URL: <https://doi.org/10.1007/s00332-017-9437-7>. doi:10.1007/s00332-017-9437-7.
- 1020 B. O. Koopman, Hamiltonian Systems and Transformations in Hilbert Space,
Proceedings of the National Academy of Sciences of the United States
1022 of America 17 (1931) 315–318. URL: <https://www.jstor.org/stable/86114>.
- 1024 I. Mezić, Spectral properties of dynamical systems, model reduction and
decompositions, *Nonlinear Dynamics* 41 (2005) 309–325.
- 1026 S. L. Brunton, B. W. Brunton, J. L. Proctor, J. N. Kutz, Koopman invariant
subspaces and finite linear representations of nonlinear dynamical systems
for control, *PloS one* 11 (2016) e0150171.
- 1028 M. A. Katsoulakis, P. Plecháč, Information-theoretic tools for parametrized
coarse-graining of non-equilibrium extended systems, *The Journal of chem-*
1030 *ical physics* 139 (2013) 074115.
- 1032 V. Harmandaris, E. Kalligiannaki, M. Katsoulakis, P. Plecháč, Path-
space variational inference for non-equilibrium coarse-grained systems,
Journal of Computational Physics 314 (2016) 355–383. URL: [http://](http://www.sciencedirect.com/science/article/pii/S002199911600173X)
1034 www.sciencedirect.com/science/article/pii/S002199911600173X.
doi:10.1016/j.jcp.2016.03.021.
- 1036 M. A. Katsoulakis, P. Vilanova, Data-driven, variational model reduction
of high-dimensional reaction networks, *Journal of Computational Physics*
1038 (2019) 108997.
- 1040 H. Mori, Transport, collective motion, and brownian motion, *Progress of*
theoretical physics 33 (1965) 423–455.
- 1042 R. Zwanzig, Nonlinear generalized langevin equations, *Journal of Statistical*
Physics 9 (1973) 215–220.
- 1044 A. Chorin, P. Stinis, Problem reduction, renormalization, and memory, *Com-*
munications in Applied Mathematics and Computational Science 1 (2007)
1–27.

- 1046 H. Lei, N. A. Baker, X. Li, Data-driven parameterization of the generalized
 1048 langevin equation, *Proceedings of the National Academy of Sciences* 113
 (2016) 14183–14188.
- Y. Zhu, J. M. Dominy, D. Venturi, On the estimation of the Mori-Zwanzig
 1050 memory integral, *Journal of Mathematical Physics* 59 (2018) 103501. URL:
<https://aip.scitation.org/doi/10.1063/1.5003467>. doi:10.1063/1.
 1052 5003467.
- M. D. Hoffman, D. M. Blei, C. Wang, J. Paisley, Stochastic variational
 1054 inference, *The Journal of Machine Learning Research* 14 (2013) 1303–
 1347.
- 1056 C. Archambeau, M. Opper, Approximate inference for continuous-time
 Markov processes, *Bayesian Time Series Models* (2011) 125–140.
- 1058 R. G. Krishnan, U. Shalit, D. Sontag, Structured inference networks for non-
 linear state space models, in: *Thirty-First AAAI Conference on Artificial*
 1060 *Intelligence*, 2017.
- V. Fortuin, G. Rätsch, S. Mandt, Multivariate time series imputation with
 1062 variational autoencoders, *arXiv preprint arXiv:1907.04155* (2019).
- J. C. Butcher, *Numerical methods for ordinary differential equations*, John
 1064 Wiley & Sons, 2016.
- R. Coifman, I. Kevrekidis, S. Lafon, M. Maggioni, B. Nadler, Diffusion maps,
 1066 reduction coordinates and low dimensional representation of stochastic sys-
 tems, *Multiscale Modeling & Simulation* 7 (2008) 842 – 864.
- 1068 J. Trashorras, D. Tsagkarogiannis, From Mesoscale Back to Microscale:
 Reconstruction Schemes for Coarse-Grained Stochastic Lattice Sys-
 1070 tems, *SIAM Journal on Numerical Analysis* 48 (2010) 1647–1677.
 URL: <http://epubs.siam.org/doi/abs/10.1137/080722382>. doi:10.
 1072 1137/080722382.
- M. A. Katsoulakis, J. Trashorras, Information loss in coarse-graining
 1074 of stochastic particle dynamics, *Journal of statistical physics* 122
 (2006) 115–135. URL: [http://link.springer.com/article/10.1007/](http://link.springer.com/article/10.1007/s10955-005-8063-1)
 1076 [s10955-005-8063-1](http://link.springer.com/article/10.1007/s10955-005-8063-1).

- 1078 D. Kondrashov, M. D. Chekroun, M. Ghil, Data-driven non-Markovian
closure models, *Physica D: Nonlinear Phenomena* 297 (2015)
33–55. URL: [http://www.sciencedirect.com/science/article/pii/](http://www.sciencedirect.com/science/article/pii/S0167278914002413)
1080 [S0167278914002413](http://www.sciencedirect.com/science/article/pii/S0167278914002413). doi:10.1016/j.physd.2014.12.005.
- B. D. Coleman, M. E. Gurtin, Thermodynamics with Internal
1082 State Variables, *The Journal of Chemical Physics* 47 (1967) 597–
613. URL: [http://scitation.aip.org/content/aip/journal/jcp/47/](http://scitation.aip.org/content/aip/journal/jcp/47/2/10.1063/1.1711937)
1084 [2/10.1063/1.1711937](http://scitation.aip.org/content/aip/journal/jcp/47/2/10.1063/1.1711937). doi:10.1063/1.1711937.
- O. Cappe, E. Moulines, T. Ryden, *Inference in Hidden Markov Models*,
1086 Springer-Verlag, 2005.
- Z. Ghahramani, Unsupervised Learning, in: O. Bousquet, G. Raetsch,
1088 U. von Luxburg (Eds.), *Advanced Lectures on Machine Learning LNAI*
3176, Springer-Verlag, 2004.
- 1090 D. Durstewitz, A state space approach for piecewise-linear recurrent
neural networks for identifying computational dynamics from neural
1092 measurements, *PLOS Computational Biology* 13 (2017) e1005542. URL:
[http://journals.plos.org/ploscompbiol/article?id=10.1371/](http://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1005542)
1094 [journal.pcbi.1005542](http://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1005542). doi:10.1371/journal.pcbi.1005542.
- L. Wiskott, T. J. Sejnowski, Slow feature analysis: Unsupervised learn-
1096 ing of invariances, *Neural Computation* 14 (2002) 715–770. doi:10.1162/
089976602317318938.
- 1098 C. M. Bishop, *Pattern Recognition and Machine Learning*, Springer, 2006.
- D. P. Kingma, M. Welling, Auto-Encoding Variational Bayes, in: *The*
1100 *International Conference on Learning Representations (ICLR)*, volume
abs/1312.6114, Banff, Alberta, Canada, 2014. URL: [http://arxiv.org/](http://arxiv.org/abs/1312.6114)
1102 [abs/1312.6114](http://arxiv.org/abs/1312.6114).
- Y. Kim, S. Wiseman, A. C. Miller, D. Sontag, A. M. Rush, Semi-amortized
1104 variational autoencoders, *arXiv preprint arXiv:1802.02550* (2018).
- C. Grigo, P.-S. Koutsourelakis, Bayesian model and dimension reduction
1106 for uncertainty propagation: applications in random media, *SIAM/ASA*
Journal on Uncertainty Quantification 7 (2019) 292–323.

- 1108 J. Li, P. G. Kevrekidis, C. W. Gear, I. G. Kevrekidis, Deciding the
 1110 Nature of the Coarse Equation Through Microscopic Simulations: The
 Baby-Bathwater Scheme, *SIAM Rev.* 49 (2007) 469–487. URL: <http://dx.doi.org/10.1137/070692303>. doi:10.1137/070692303.
- 1112 W. G. Noid, Perspective: Coarse-grained models for biomolec-
 ular systems, *The Journal of Chemical Physics* 139 (2013).
 1114 URL: <http://scitation.aip.org/content/aip/journal/jcp/139/9/10.1063/1.4818908>. doi:<http://dx.doi.org/10.1063/1.4818908>.
- 1116 D. Mackay, Probable Networks and Plausible Predictions - a Re-
 view of Practical Bayesian Methods for Supervised Neural Networks,
 1118 *Network-Computation in Neural Systems* 6 (1995) 469–505. doi:10.1088/
 0954-898X/6/3/011.
- 1120 C. M. Bishop, M. E. Tipping, Variational relevance vector machines, in:
 Proceedings of the Sixteenth conference on Uncertainty in artificial intel-
 1122 ligence, Morgan Kaufmann Publishers Inc., 2000, pp. 46–53.
- D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, *arXiv*
 1124 preprint [arXiv:1412.6980](https://arxiv.org/abs/1412.6980) (2014).
- G.-H. Cottet, P. D. Koumoutsakos, *Vortex Methods: Theory and Practice*,
 1126 2 edition ed., Cambridge University Press, Cambridge ; New York, 2000.
- S. Roberts, Convergence of a Random Walk Method for the Burgers Equa-
 1128 tion, *Mathematics of Computation* 52 (1989) 647–673. URL: <http://www.jstor.org/stable/2008486>. doi:10.2307/2008486.
- 1130 A. Chertock, D. Levy, Particle Methods for Dispersive Equations, *Jour-
 nal of Computational Physics* 171 (2001) 708–730. URL: [http://](http://www.sciencedirect.com/science/article/pii/S0021999101968032)
 1132 www.sciencedirect.com/science/article/pii/S0021999101968032.
 doi:10.1006/jcph.2001.6803.
- 1134 M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S.
 Corrado, A. Davis, J. Dean, M. Devin, et al., Tensorflow: Large-scale
 1136 machine learning on heterogeneous distributed systems, *arXiv preprint*
[arXiv:1603.04467](https://arxiv.org/abs/1603.04467) (2016).
- 1138 S. Pan, K. Duraisamy, Physics-informed probabilistic learning of linear em-
 beddings of nonlinear dynamics with guaranteed stability, *SIAM Journal*
 1140 *on Applied Dynamical Systems* 19 (2020) 480–509.

Declaration of interests

☒ The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

☐ The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: