

# Learning Partially Known Stochastic Dynamics with Empirical PAC Bayes

...

Manuel Haussmann

Sebastian Gerwinn

Andreas Look

Barbara Rakitsch

Melih Kandemir

(Published at AISTATS 2021)



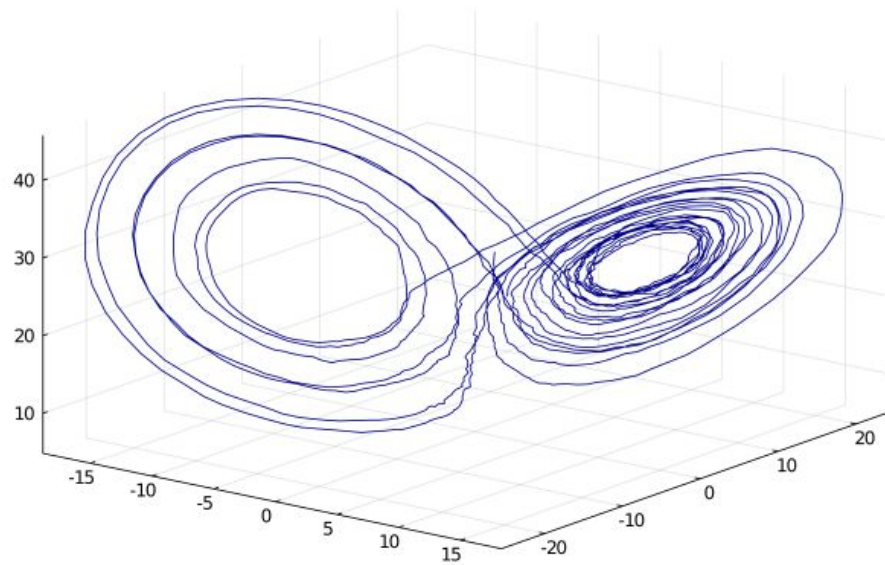
UNIVERSITÄT  
HEIDELBERG  
ZUKUNFT  
SEIT 1386



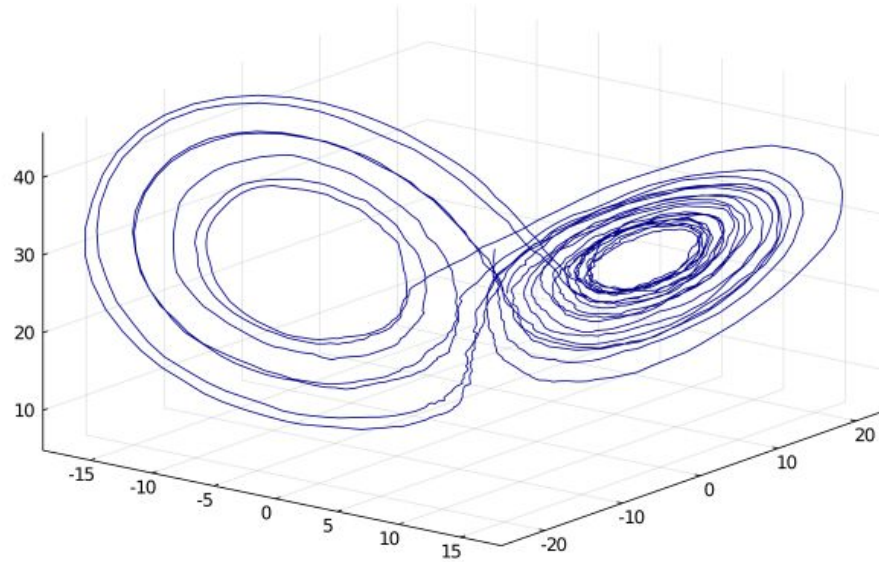
**BOSCH**

Invented for life

Given observed measurements...



**Given observed measurements...**



... we want to learn the underlying dynamical system

## Two complementary approaches

- The domain expert creates a system of differential equations
  - ✓ interpretable
  - ✓ can incorporate prior domain knowledge
  - ✓ few parameters
  - ✓ probably data efficient
  - ✗ not very flexible
  - ✗ what to do if only partial knowledge is available?

$$\text{WHITE BOX ODE}$$
$$d\mathbf{h}_t = r(\mathbf{h}_t, t)dt$$

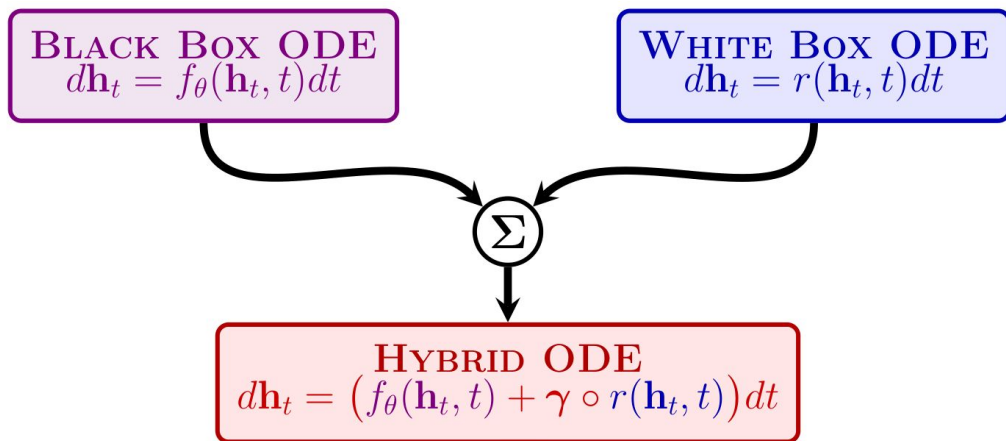
## Two complementary approaches

- The domain expert creates a system of differential equations
  - ✓ interpretable
  - ✓ can incorporate prior domain knowledge
  - ✓ few parameters
  - ✓ probably data efficient
  - ✗ not very flexible
  - ✗ what to do if only partial knowledge is available?
- The data scientist relies on deep learning and creates a neural ODE
  - ✓ very flexible
  - ✓ can adapt itself to complex relationship
  - ✗ lots of parameters
  - ✗ data hungry

$$\text{WHITE BOX ODE}$$
$$d\mathbf{h}_t = r(\mathbf{h}_t, t)dt$$

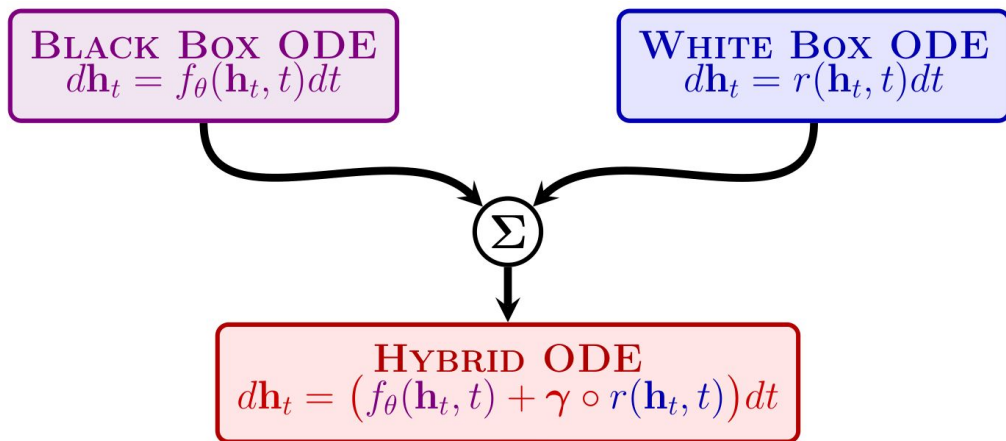
$$\text{BLACK BOX ODE}$$
$$d\mathbf{h}_t = f_\theta(\mathbf{h}_t, t)dt$$

## Can we combine the two?



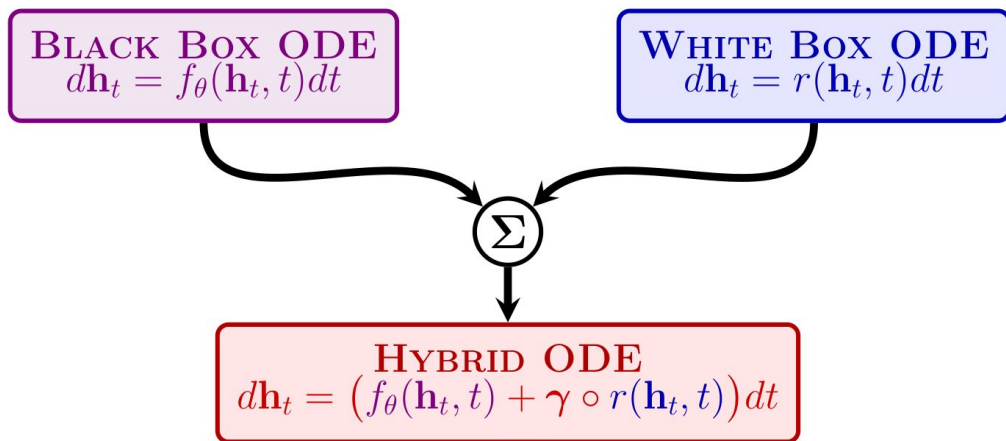
- Get the best of both worlds with in a hybrid model (with  $\gamma \in \{0, 1\}^D$ )

## Can we combine the two?



- Get the best of both worlds with in a hybrid model (with  $\gamma \in \{0, 1\}^D$ )
- Can we go one step further?

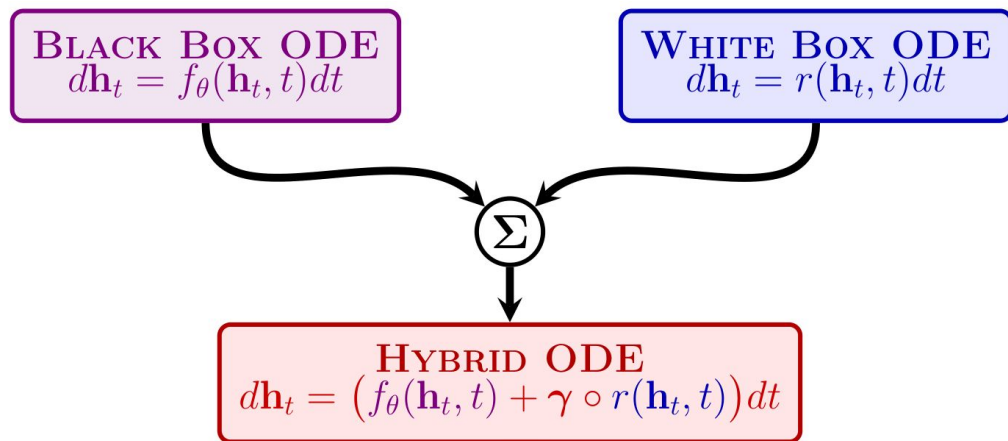
## Can we combine the two?



- Get the best of both worlds with in a hybrid model (with  $\gamma \in \{0, 1\}^D$ )
- Can we go one step further?
  - Model parameter (epistemic) uncertainty by switching to a BNN

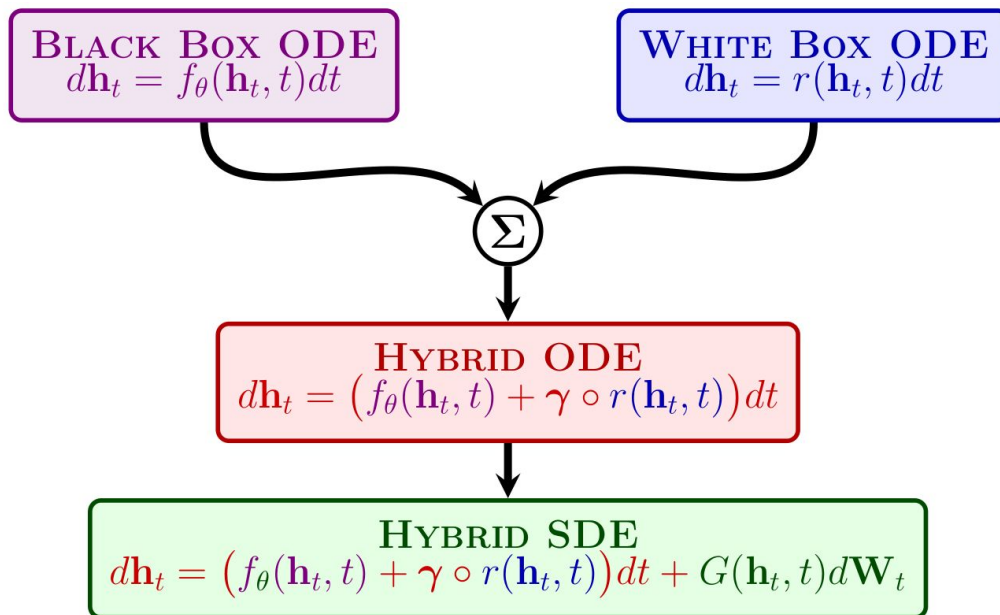


## Can we combine the two?



- Get the best of both worlds with in a hybrid model (with  $\gamma \in \{0, 1\}^D$ )
- Can we go one step further?
  - Model parameter (epistemic) uncertainty by switching to a BNN
  - Model (aleatoric) uncertainty by switching to an SDE

## The complete pipeline

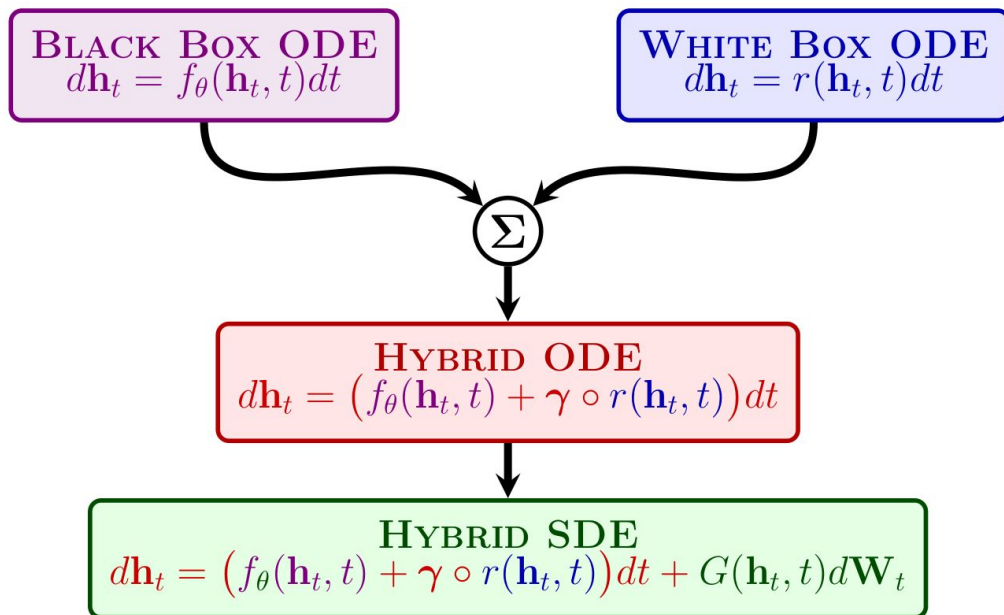


✓ flexibility & adaptability

✓ domain knowledge

✓ uncertainty

## The complete pipeline



✓ flexibility & adaptability

✓ domain knowledge

✓ uncertainty

Great, but... How do I actually learn such a thing?

## Step 1: Euler-Murayama discretization (see e.g. Särkkä and Solin, 2019)

Problem 1: How to deal with the SDE?

$$\text{HYBRID SDE}$$
$$d\mathbf{h}_t = (f_\theta(\mathbf{h}_t, t) + \gamma \circ r(\mathbf{h}_t, t))dt + G(\mathbf{h}_t, t)d\mathbf{W}_t$$

- $\mathbf{W}_t$  denotes a Wiener process,  $G$  gives the diffusion dynamics
- Discretization of the SDE into  $K$  steps gives us

$$\mathbf{h}_{t_{k+1}} = \mathbf{h}_{t_k} + (f_\theta(\mathbf{h}_{t_k}, t_k) + \gamma \circ r(\mathbf{h}_{t_k}, t_k))\Delta t_k + G(\mathbf{h}_{t_k}, t_k)\Delta W_k$$
$$\Delta W_k \sim \mathcal{N}(0, \Delta t_k \mathbf{1}_P) \quad \Delta t_k := t_{k+1} - t_k$$

## Step 1.5: Design a generative pipeline

- Our generative pipeline is then given as

$$\theta \sim p_\phi(\theta)$$

$$\mathbf{h}_0 \sim p(\mathbf{h}_0)$$

$$\mathbf{h}_{k+1} | \mathbf{h}_k, \theta \sim \mathcal{N}(\mathbf{h}_{k+1} | \mathbf{h}_k + d(\mathbf{h}_k, t_k) \Delta t_k, \Sigma_k), \quad k = 0, \dots, K - 1$$

$$\text{where } d(\mathbf{h}_k, t_k) := f_\theta(\mathbf{h}_k, t_k) + \gamma \circ r(\mathbf{h}_k, t_k)$$

$$\Sigma_k := G(\mathbf{h}_k, t_k) G(\mathbf{h}_k, t_k)^\top \Delta t_k$$

$$\Delta t_k := t_{k+1} - t_k$$

$$\mathbf{y}_k | \mathbf{h}_k \sim p(\mathbf{y}_k | \mathbf{h}_k), \quad k = 1, \dots, K$$

Where we have observed N trajectories  $\mathcal{D} = \{\mathbf{Y}_1, \dots, \mathbf{Y}_N\}$  with  $\mathbf{Y}_n = \{\mathbf{y}_1^n, \dots, \mathbf{y}_K^n\}$

- Classical approach: Infer a posterior over local and global parameters

$$p(\theta, \mathbf{H}_1, \dots, \mathbf{H}_N | \mathcal{D})$$

## Step 2: Empirical Bayes / Type II Maximum Likelihood

- Problem 2: Infer a posterior over local and global parameters

$$p(\theta, \mathbf{H}_1, \dots, \mathbf{H}_N | \mathcal{D})$$

- How?
  - MCMC? Becomes very quickly too expensive computationally
  - Variational Inference? Becomes too restrictive due to strong independence assumptions
- Instead, Type II Maximum Likelihood

$$\hat{\phi} = \arg \max_{\phi} \int p(\mathcal{D} | \mathbf{H}) p(\mathbf{H} | \theta) p_{\phi}(\theta) d\mathbf{H} d\theta$$

$$\approx \arg \max_{\phi} \log \left( \frac{1}{S} \sum_{s=1}^S p(\mathcal{D} | \mathbf{H}^s) \right) \quad \text{where } \mathbf{H}^s \sim p(\mathbf{H} | \theta^s), \theta^s \sim p_{\phi}(\theta)$$

### Step 3: PAC-based regularization

- Problem 3: We now have a tractable approach, but we lost a lot of regularization
- A very quick primer on PAC-Bayes (see e.g. McAllester, 1999,2003)

- A risk of a hypothesis  $h$  given a loss  $l$  is  $R(h) = \mathbb{E}_x [l(x, h(x))]$

- The empirical counterpart is given as  $R_{\mathcal{D}}(h) = \frac{1}{|\mathcal{D}|} \sum_{x \in \mathcal{D}} l(x, h(x))$

- Given a distribution over hypothesis  $Q$  and a prior distribution  $P$

$$\mathbb{P} \left( \forall Q : \mathbb{E}_{h \sim Q} [R(h)] \leq \mathbb{E}_{h \sim Q} [R_{\mathcal{D}}(h)] + \mathcal{C}(P, Q, \delta, N) \right) \geq 1 - \delta$$

### Step 3: PAC-based regularization

We place distributions over the hybrid and a prior process

$$Q_{0 \rightarrow T}(\mathbf{h}_{0 \rightarrow T}, \theta) = p_{\text{hyb}}(\mathbf{h}_{0 \rightarrow T} | \theta) p_{\phi}(\theta)$$

$$P_{0 \rightarrow T}(\mathbf{h}_{0 \rightarrow T}, \theta) = p_{\text{pri}}(\mathbf{h}_{0 \rightarrow T}) p_{\text{pri}}(\theta)$$

and define the risk as

$$R(H) \triangleq \mathbb{E}_{\mathbf{y}_k \sim \mathfrak{G}(t)} \left[ 1 - \frac{1}{\bar{B}_K} \prod_{k=1}^K p(\mathbf{y}_k | \mathbf{h}_k) \right]$$

$$\bar{B}_K \triangleq \max_{\mathbf{y}_k, \mathbf{h}_k} \prod_{k=1}^K p(\mathbf{y}_k | \mathbf{h}_k) \leq \left( \max_{\mathbf{y}_k, \mathbf{h}_k} p(\mathbf{y}_k | \mathbf{h}_k) \right)^K$$



### Step 3: PAC-based regularization

We place distributions over the hybrid and a prior process

$$Q_{0 \rightarrow T}(\mathbf{h}_{0 \rightarrow T}, \theta) = p_{\text{hyb}}(\mathbf{h}_{0 \rightarrow T} | \theta) p_{\phi}(\theta)$$

$$P_{0 \rightarrow T}(\mathbf{h}_{0 \rightarrow T}, \theta) = p_{\text{pri}}(\mathbf{h}_{0 \rightarrow T}) p_{\text{pri}}(\theta)$$

and can then show (Theorem 1) that

$$\text{with } \mathbb{P} \geq 1 - \delta \quad \mathbb{E}_{H \sim Q_{0 \rightarrow T}} [R(H)] \leq \mathbb{E}_{H \sim Q_{0 \rightarrow T}} [R_{\mathcal{D}}(H)] + \mathcal{C}_{\delta}(Q_{0 \rightarrow T}, P_{0 \rightarrow T})$$

Ultimate objective is then

$$\hat{\phi} := \arg \max_{\phi} \underbrace{\frac{1}{SN} \sum_{n=1}^N \sum_{s=1}^S \sum_{k=1}^K \log(p(\mathbf{y}_k^n | \mathbf{h}_k^{n,s}))}_{\text{approx marginal likelihood}} + \underbrace{\sqrt{\frac{\text{KL}(Q_{0 \rightarrow T} \| P_{0 \rightarrow T}) + \log(4\sqrt{N}/\delta)}{2N}}}_{\text{regularizer}}$$

## Theoretical Summary

- Step 0: Design a hybrid SDE model that allows for knowledge and ignorance
- Step 1: Turn into a probabilistic model via discretization
- Step 2: Empirical Bayes/Type II ML to get a principled, tractable objective (v1.0)
- Step 3: PAC-Bayes to get an improved principled objective (v2.0)

## Results: Lorenz Attractor as an SDE

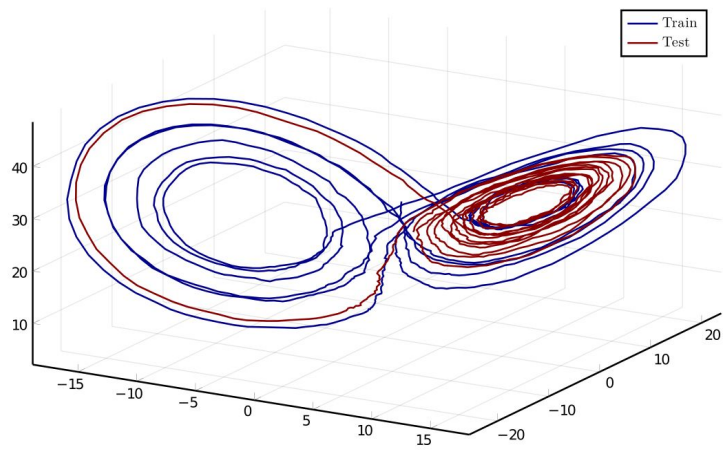
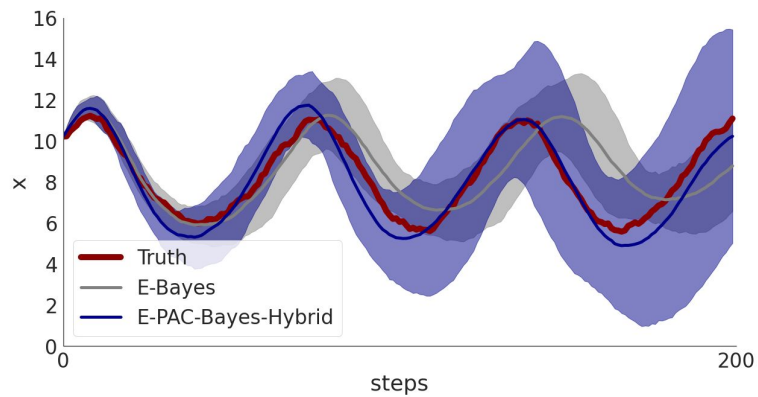
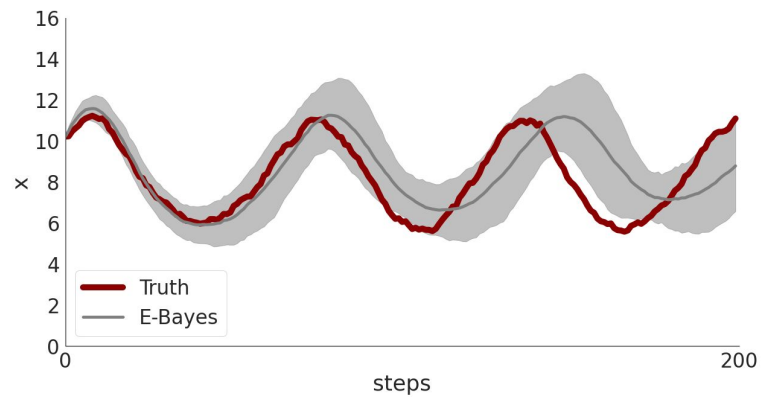
The underlying SDE

$$\begin{aligned}dx_t &= \zeta(y_t - x_t) dt + dW_t, \\dy_t &= (x_t(\kappa - z_t) - y_t) dt + dW_t, \\dz_t &= (x_t y_t - \rho z_t) dt + dW_t,\end{aligned}$$

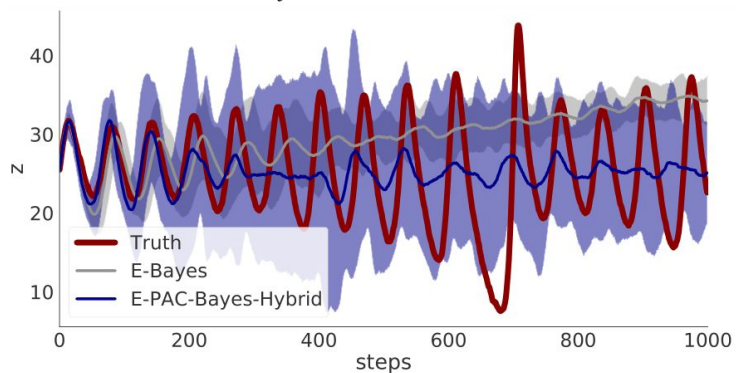
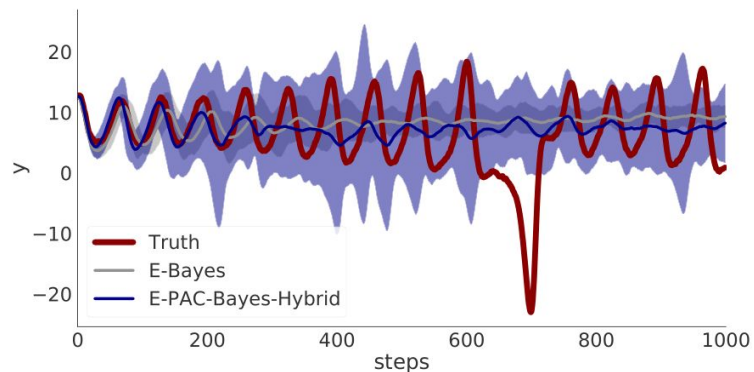
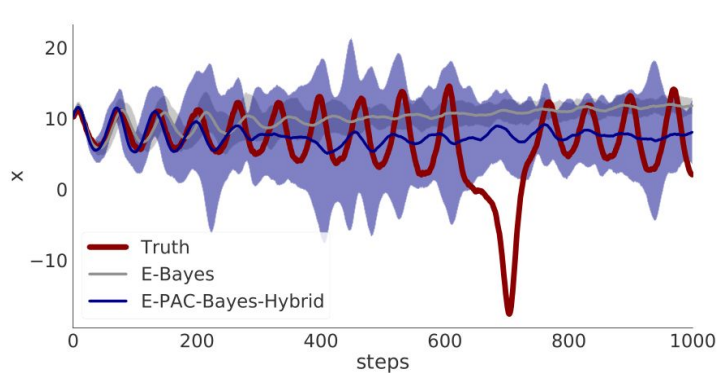
Prior Knowledge	Model	Test MSE
None	(i)	29.20 $\pm$ 0.19
	(ii)	29.05 $\pm$ 0.23
$\gamma = [1, 0, 0], \zeta \sim \mathcal{N}(10, 1)$	(iii)	27.58 $\pm$ 0.17
	(iv)	27.42 $\pm$ 0.16
$\gamma = [0, 1, 0], \kappa \sim \mathcal{N}(28, 1)$	(iii)	15.87 $\pm$ 0.46
	(iv)	15.06 $\pm$ 0.35
$\gamma = [0, 0, 1], \rho \sim \mathcal{N}(2.67, 1)$	(iii)	27.82 $\pm$ 0.26
	(iv)	28.37 $\pm$ 0.21
$\gamma = [1, 1, 1],$ $(\zeta, \kappa, \rho)^\top \sim \mathcal{N}((10, 28, 2.67)^\top, \mathbf{1}_3)$	(v)	16.40 $\pm$ 2.31

- (i) Empirical Bayes (EB) based without prior knowledge
- (ii) EB with PAC regularization
- (iii) EB with prior knowledge
- (iv) EB with prior knowledge and PAC regularization
- (v) A model with full prior knowledge over all three equations

## Results: Lorenz Attractor as an SDE

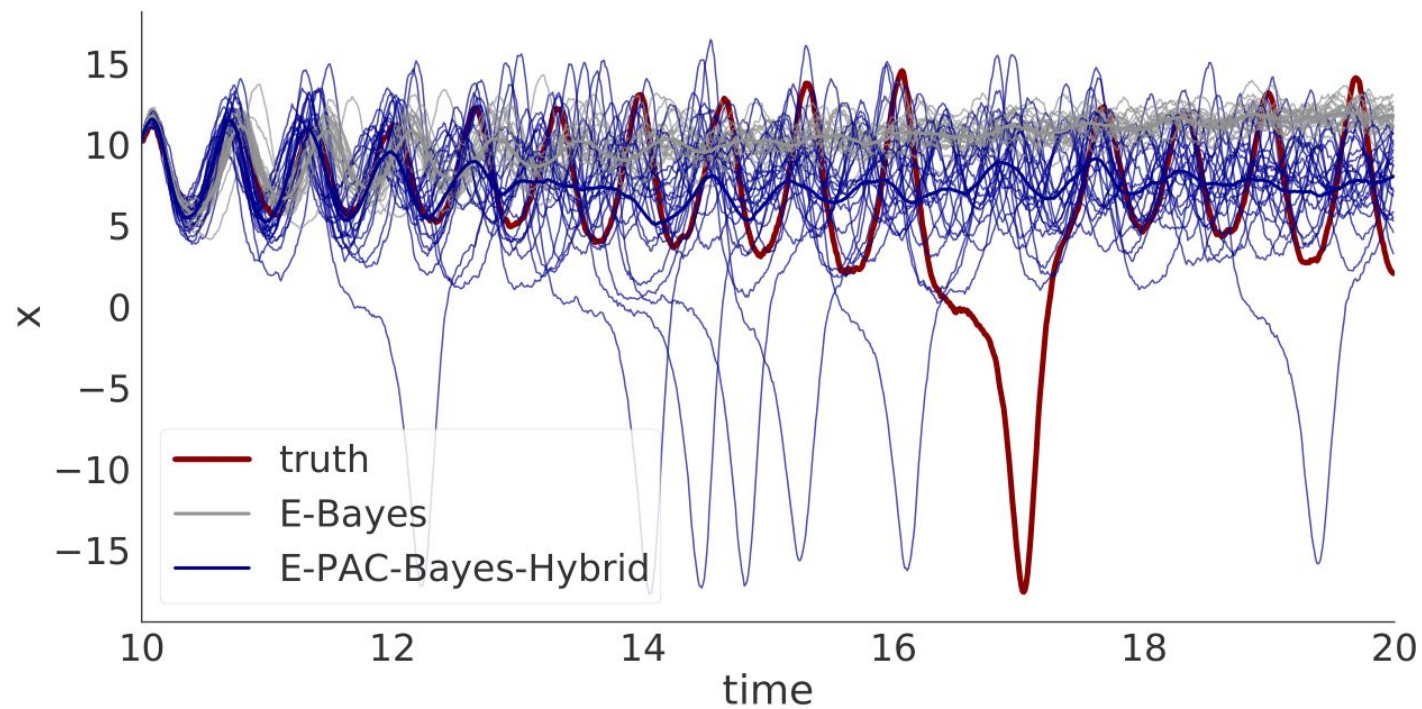


## Results: Lorenz Attractor as an SDE

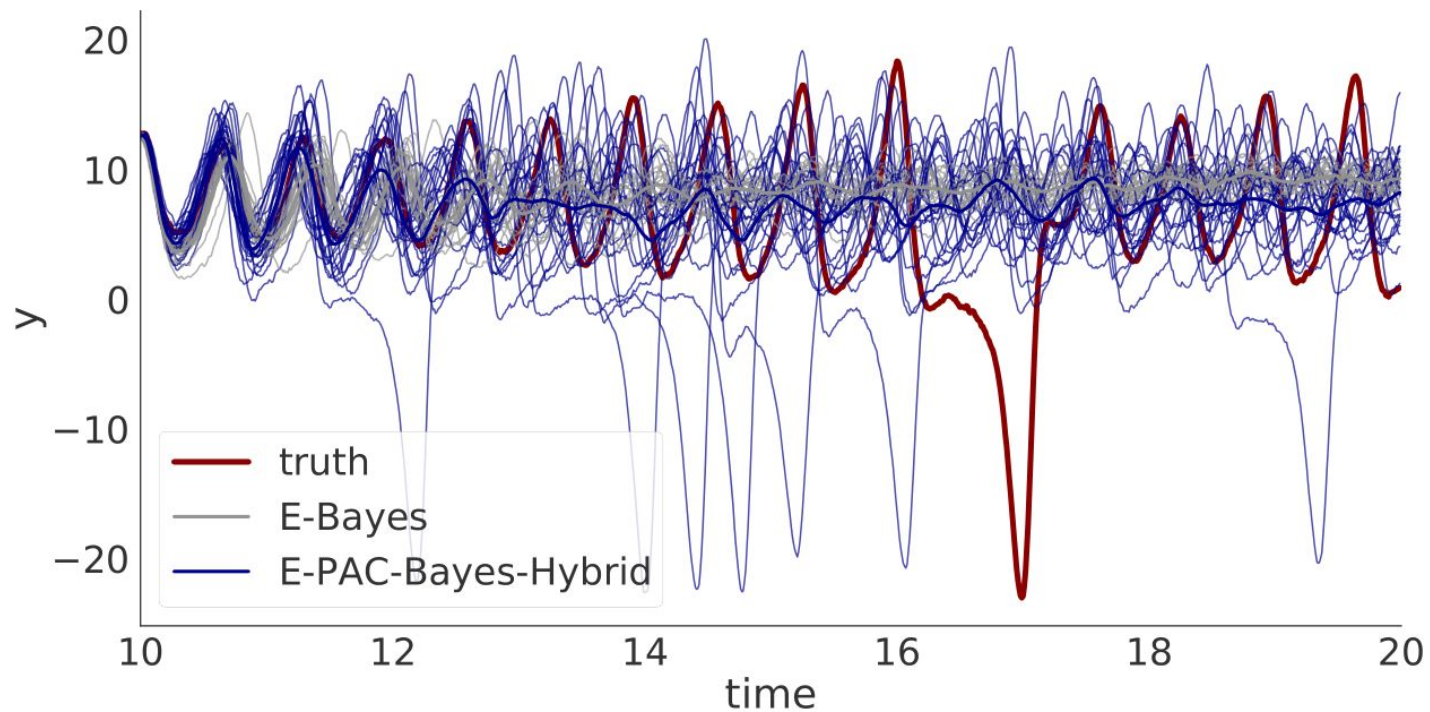


- Noisy knowledge over  $dZ$
- Std over 21 trajectories

## Results: Lorenz Attractor as an SDE

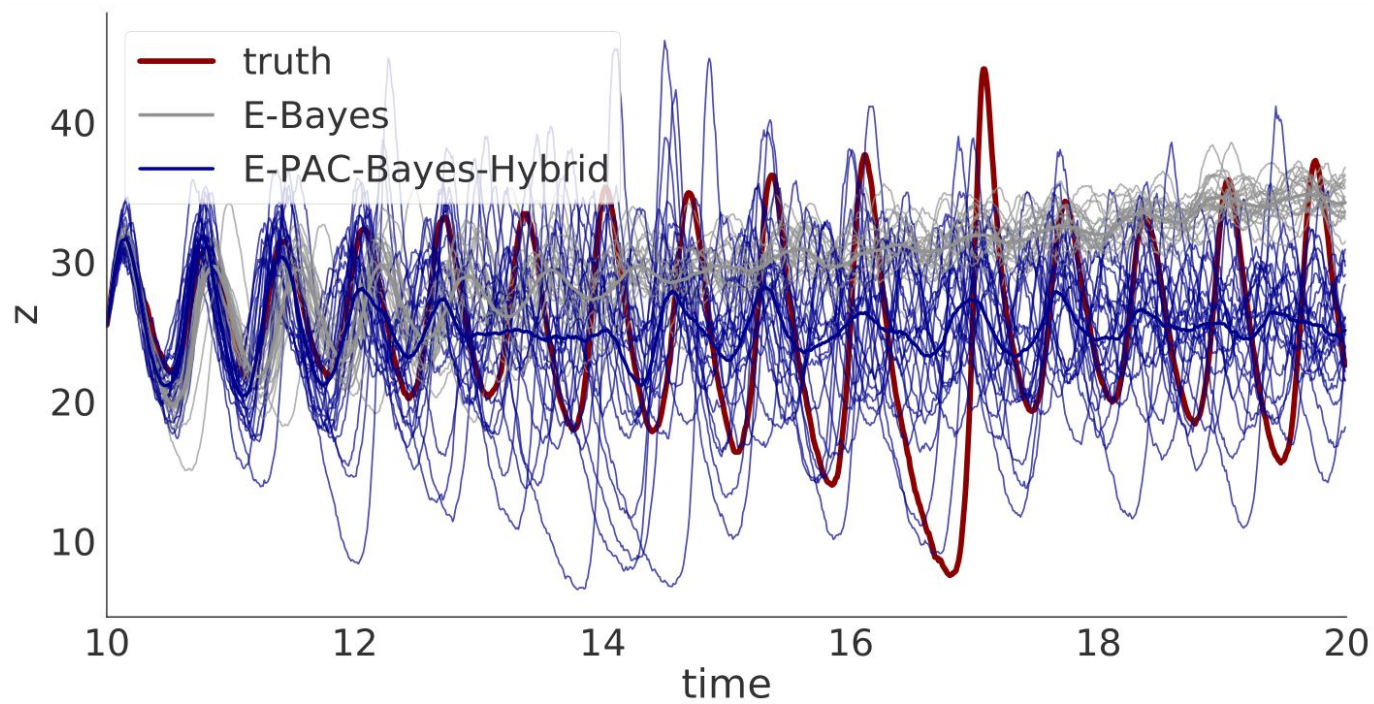


## Results: Lorenz Attractor as an SDE





## Results: Lorenz Attractor as an SDE





## Results: CMU Motion Capture Data (Sequences of human walking dynamics)

Method	Test MSE	Test NLL
DTSBN-S (Gan et al., 2015)	34.86 $\pm$ 0.02	Not Applicable
npODE (Heinonen et al., 2018)	22.96	Not Applicable
Neural-ODE (Chen et al., 2018a)	22.49 $\pm$ 0.88	Not Applicable
ODE <sup>2</sup> VAE (Yildiz et al., 2019)	10.06 $\pm$ 1.40	Not Reported
ODE <sup>2</sup> VAE-KL (Yildiz et al., 2019)	8.09 $\pm$ 1.95	Not Reported
D-BNN (SGLD) (Look and Kandemir, 2019)	13.89 $\pm$ 2.56	747.92 $\pm$ 58.49
D-BNN (VI) (Hegde et al., 2019)	9.05 $\pm$ 2.05	452.47 $\pm$ 102.59
(i) E-Bayes	8.68 $\pm$ 1.56	433.76 $\pm$ 77.78
(ii) E-PAC-Bayes	9.17 $\pm$ 1.20	489.82 $\pm$ 67.06
(iii) E-Bayes-Hybrid	9.25 $\pm$ 1.99	462.82 $\pm$ 99.61
(iv) E-PAC-Bayes-Hybrid	<b>7.84<math>\pm</math>1.41</b>	<b>415.38<math>\pm</math>80.37</b>

# Thank you for listening!

Please see [arxiv.org/abs/2006.09914](https://arxiv.org/abs/2006.09914) for further details  
Feel free to contact me at [manuel.haussmann@iwr.uni-heidelberg.de](mailto:manuel.haussmann@iwr.uni-heidelberg.de) with further questions