

Tree-structured Gaussian Process Approximations

Richard Turner and Thang Bui

May 21, 2015

Gaussian processes for time-series

$$x_t = \lambda x_{t-1} + \sigma_x \eta_t$$

$$\eta_t \sim \mathcal{N}(0, 1)$$

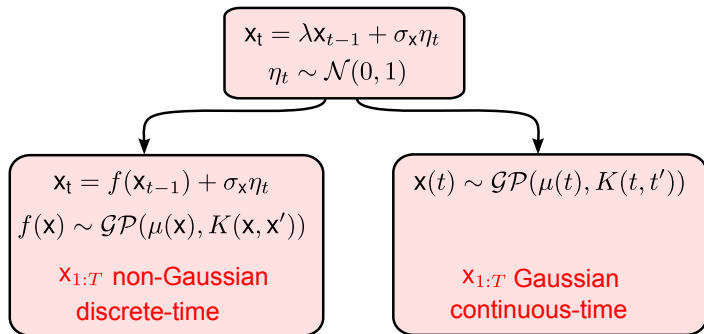
Gaussian processes for time-series

$$\begin{aligned}x_t &= \lambda x_{t-1} + \sigma_x \eta_t \\ \eta_t &\sim \mathcal{N}(0, 1)\end{aligned}$$

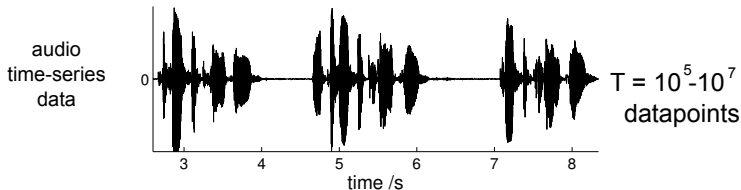
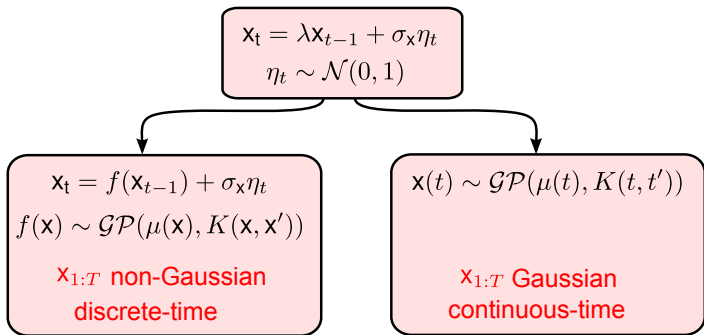
$$\begin{aligned}x_t &= f(x_{t-1}) + \sigma_x \eta_t \\ f(x) &\sim \mathcal{GP}(\mu(x), K(x, x'))\end{aligned}$$

**$x_{1:T}$ non-Gaussian
discrete-time**

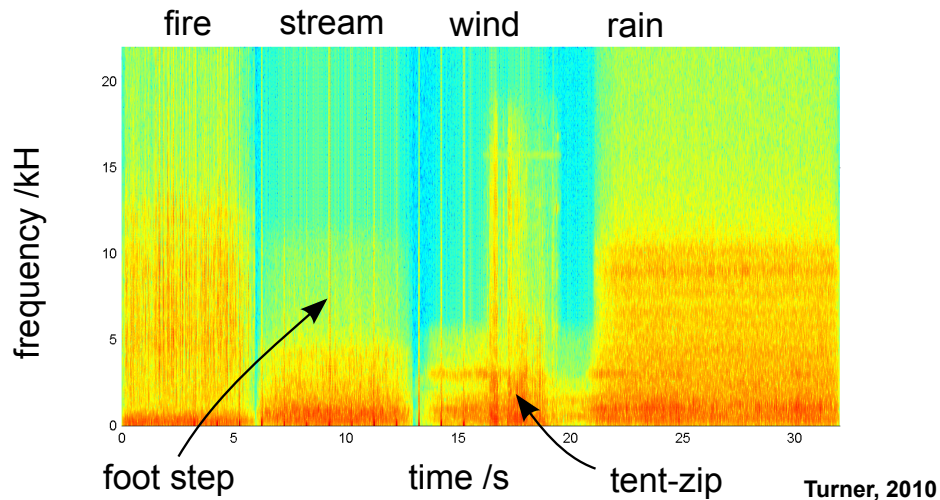
Gaussian processes for time-series



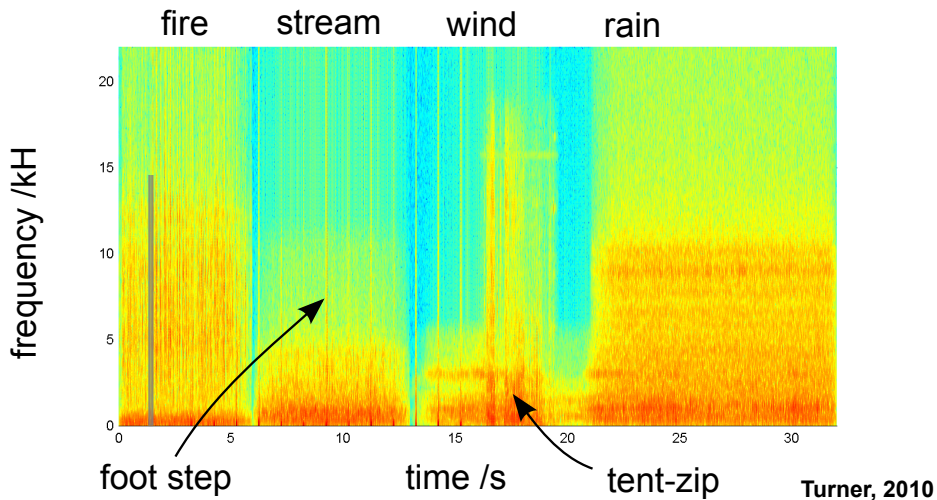
Gaussian processes for time-series



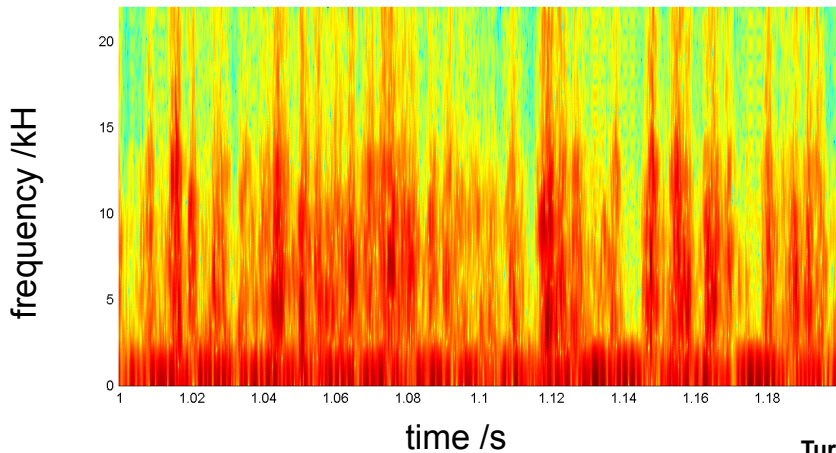
Audio texture modelling



Audio texture modelling

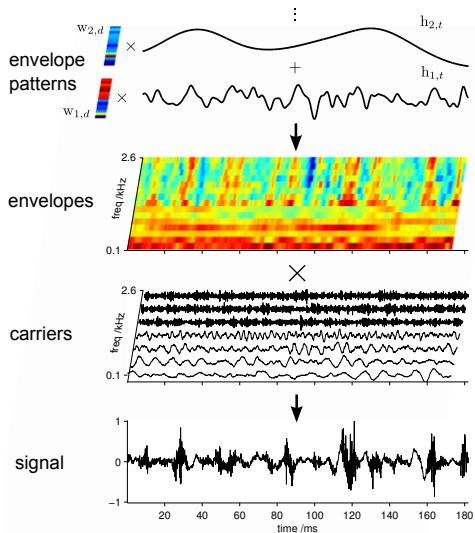


Audio texture modelling



Turner, 2010

Audio texture modelling



$$\log h_{l,t} \sim \text{GP} \left(\begin{array}{cc} \text{mean} & \text{spectrum} \\ \mu_k, & \frac{1}{0} f \end{array} \right)$$

= slow
Gaussian
process

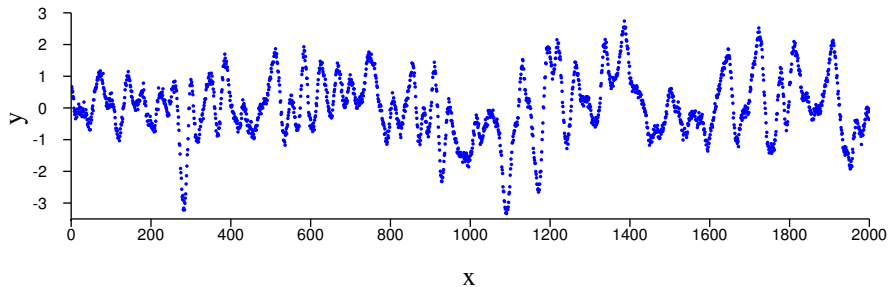
$$a_{d,t} = \sum_{l=1}^L h_{l,t} w_{l,d}$$

$$c_d(t) \sim \text{GP} \left(\begin{array}{cc} \text{mean} & \text{spectrum} \\ 0, & \frac{1}{0} f \end{array} \right)$$

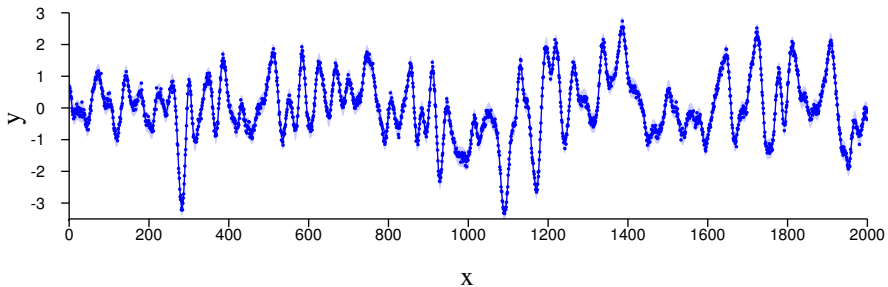
= bandpass
Gaussian
noise

$$y(t) = \sum_{d=1}^D \Re(x_d(t)) a_d(t)$$

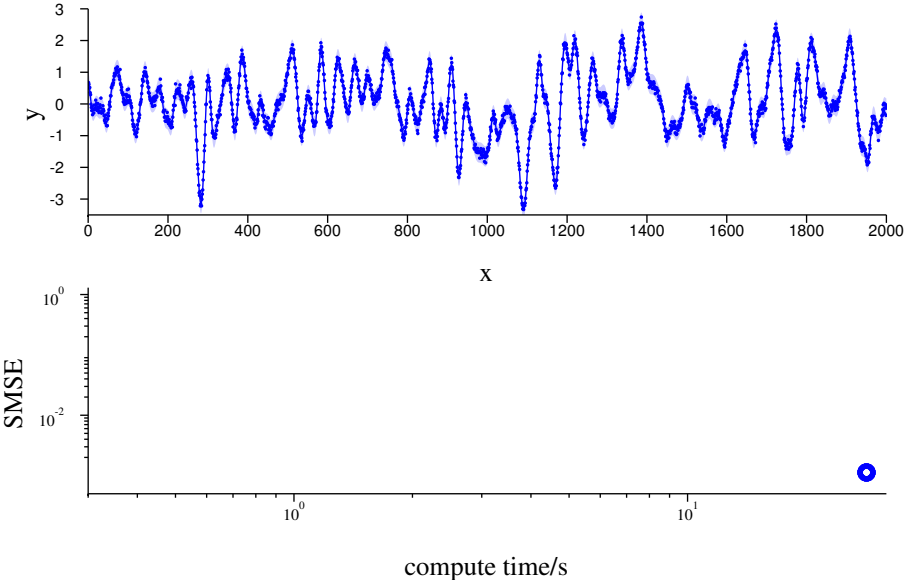
Many GP approximations are poor for time-series



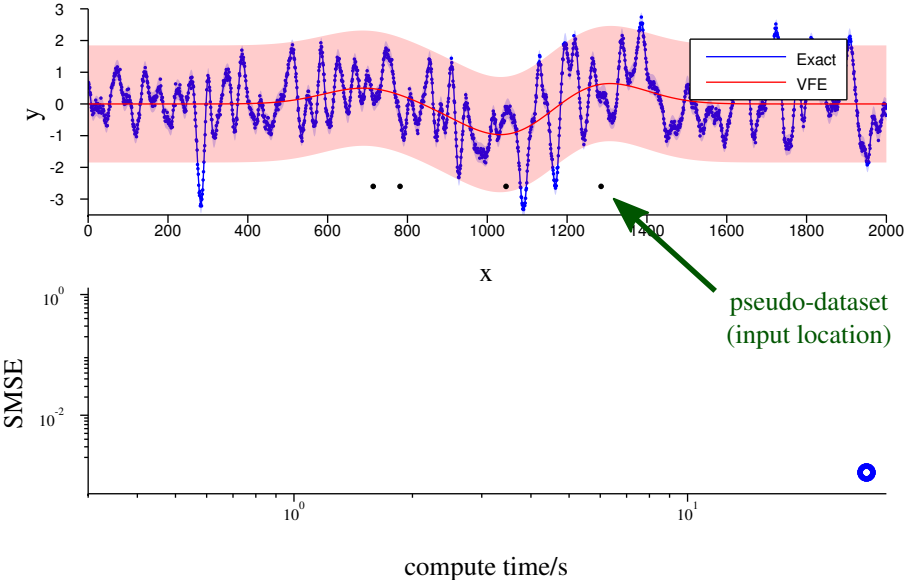
Many GP approximations are poor for time-series



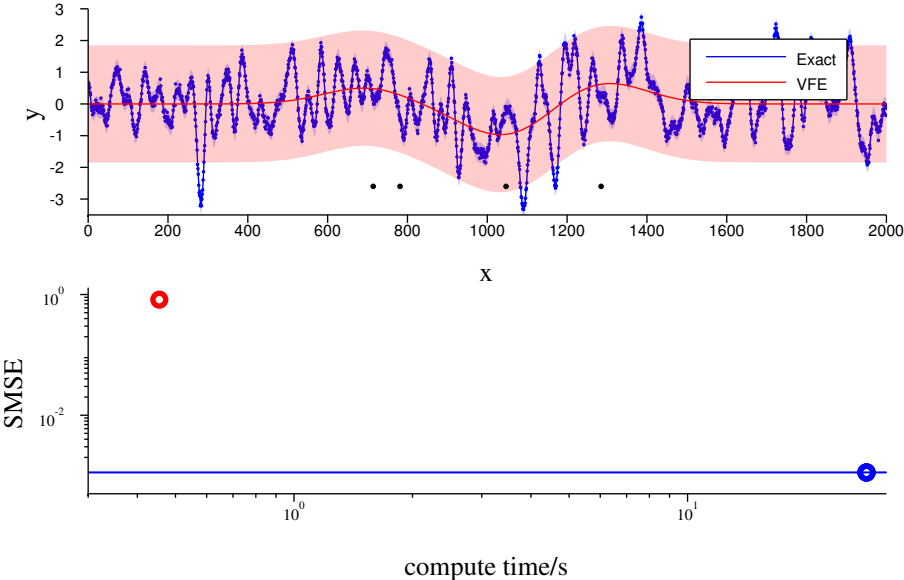
Many GP approximations are poor for time-series



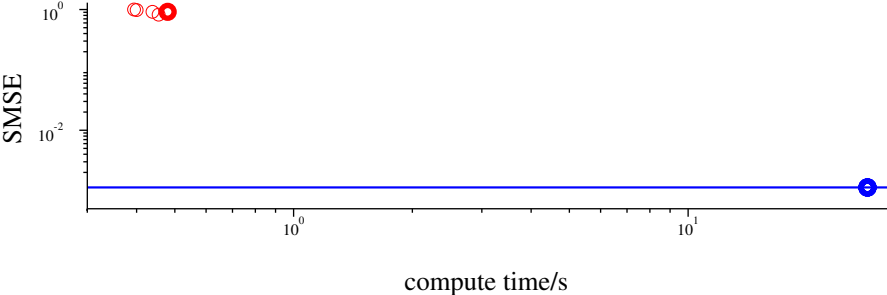
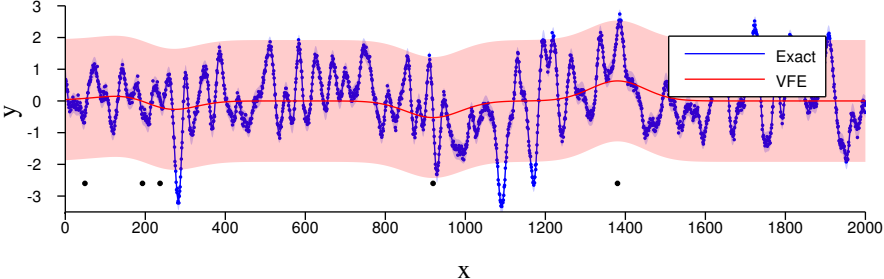
Many GP approximations are poor for time-series



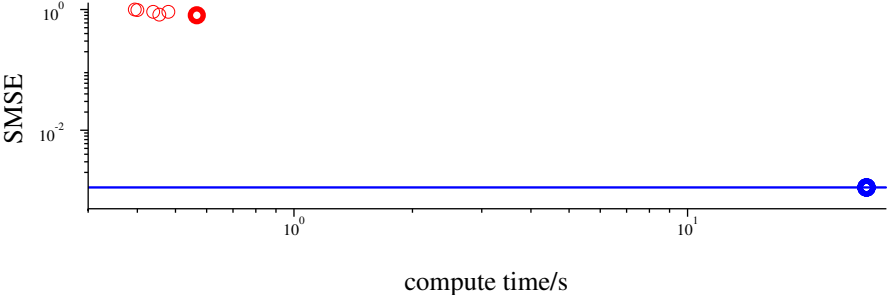
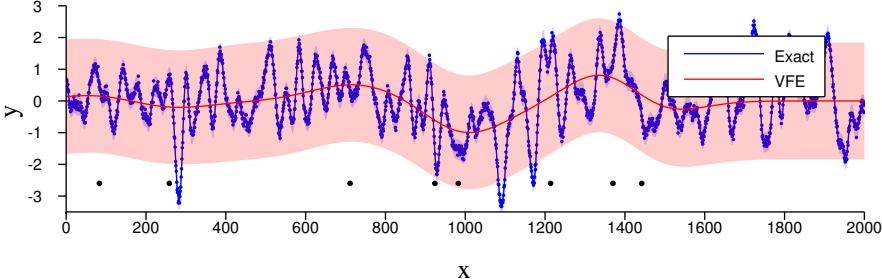
Many GP approximations are poor for time-series



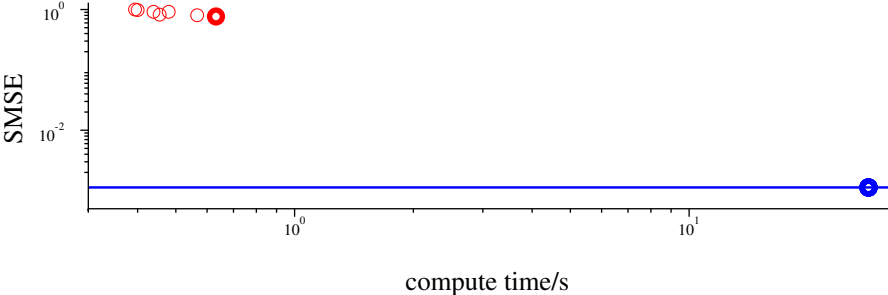
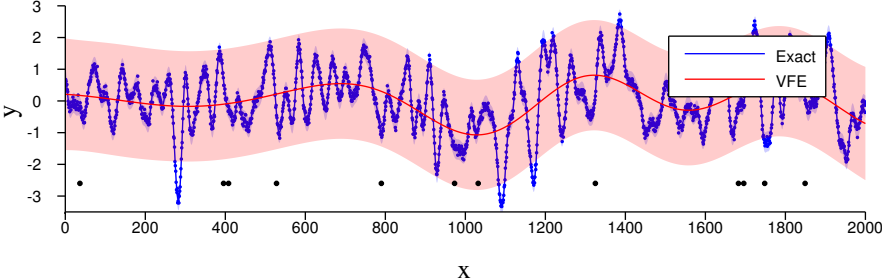
Many GP approximations are poor for time-series



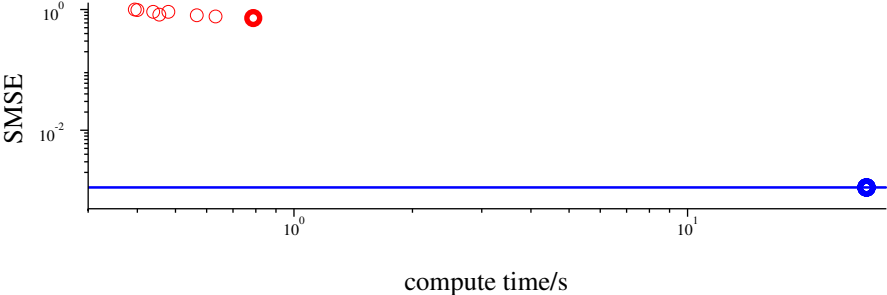
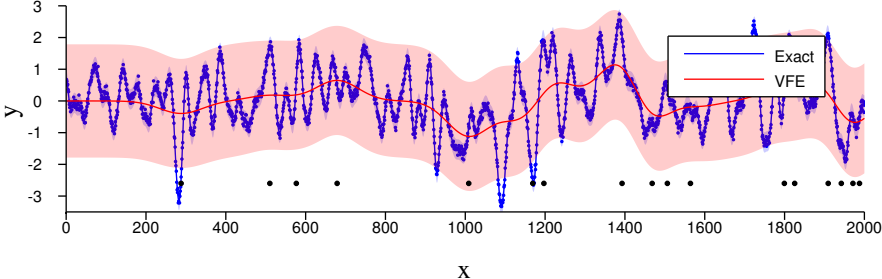
Many GP approximations are poor for time-series



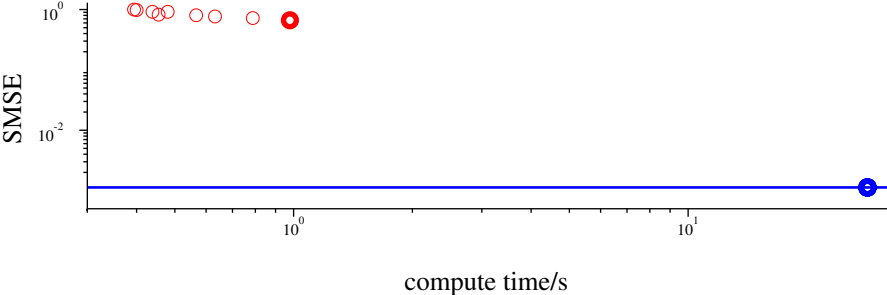
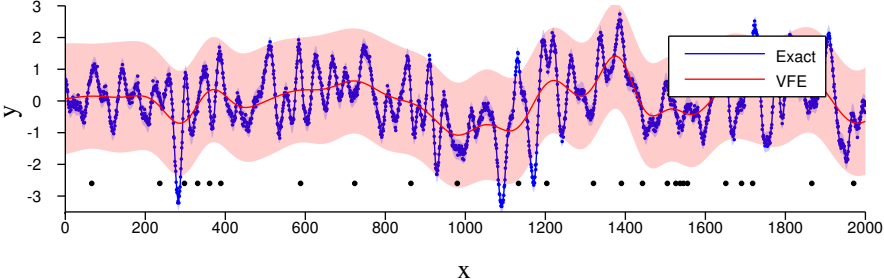
Many GP approximations are poor for time-series



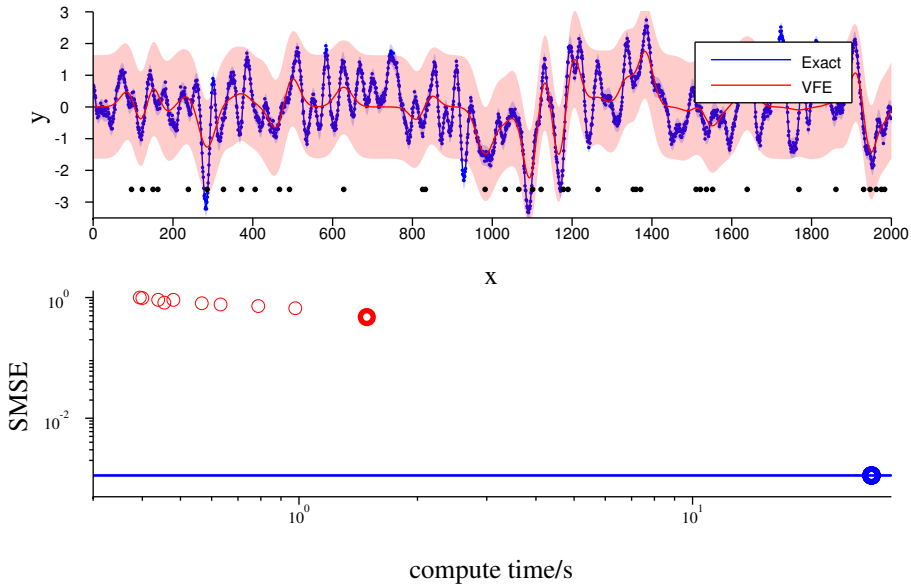
Many GP approximations are poor for time-series



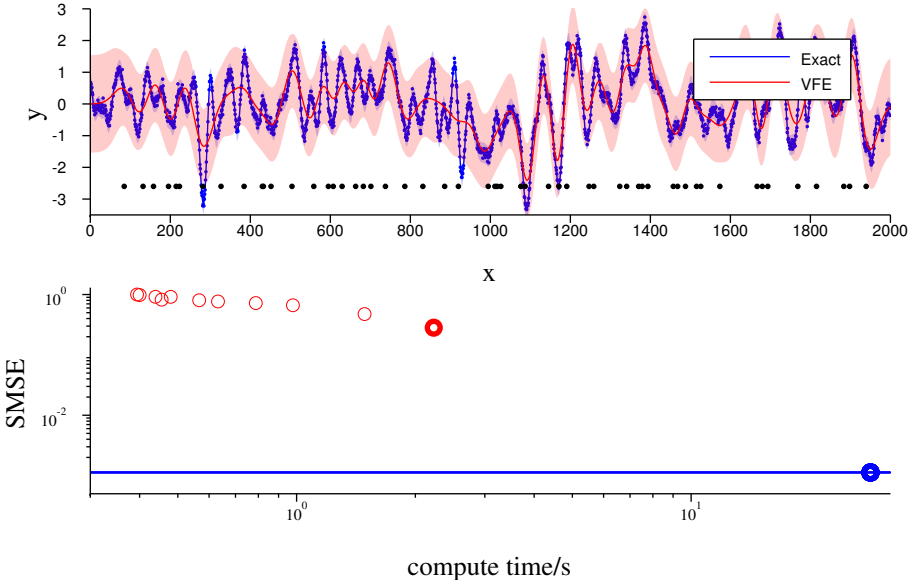
Many GP approximations are poor for time-series



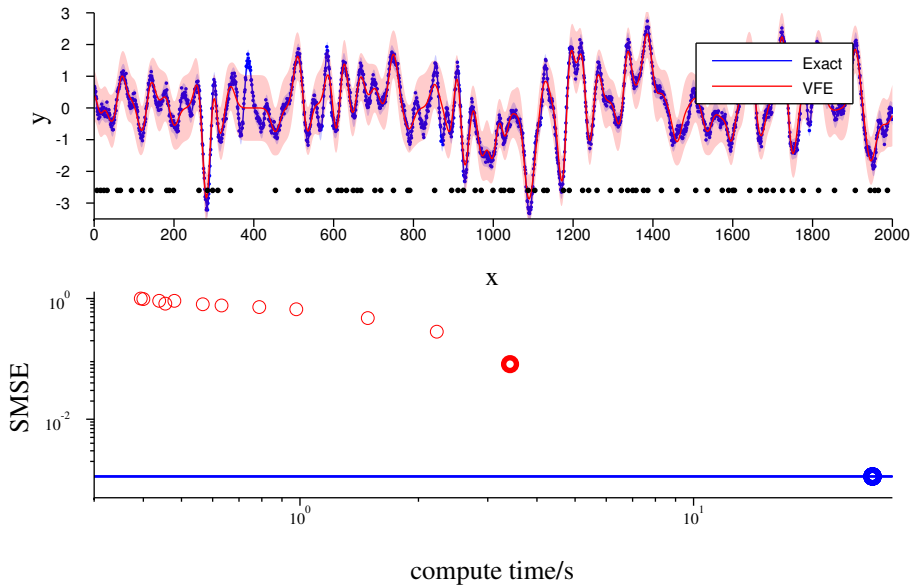
Many GP approximations are poor for time-series



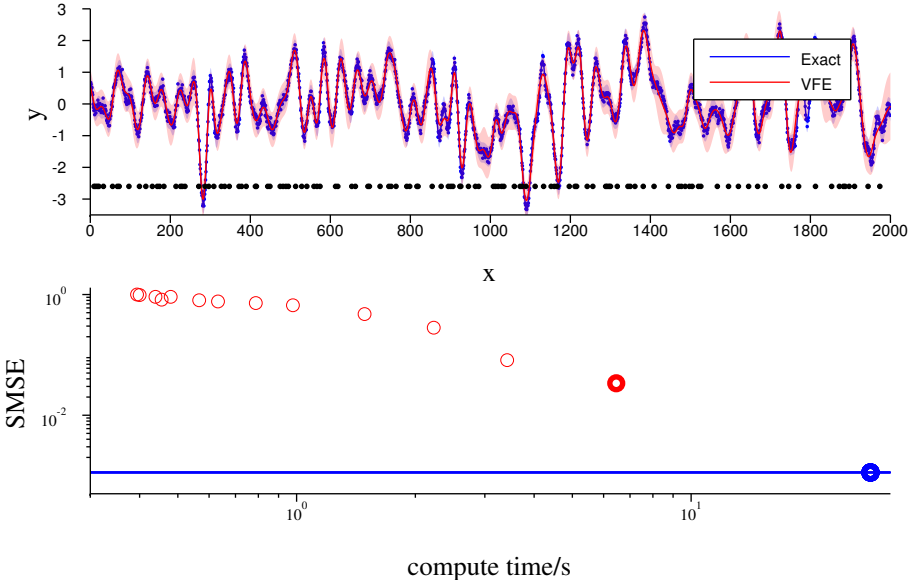
Many GP approximations are poor for time-series



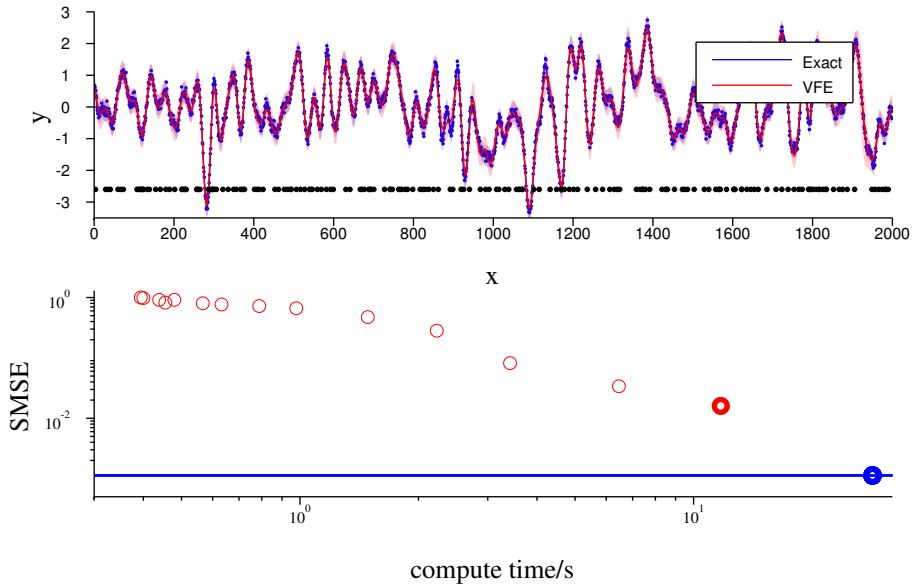
Many GP approximations are poor for time-series



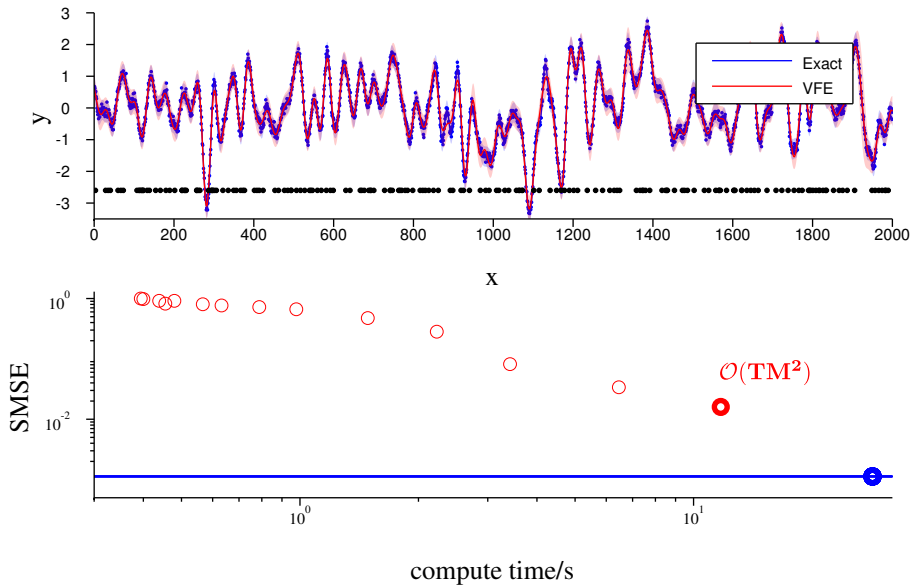
Many GP approximations are poor for time-series



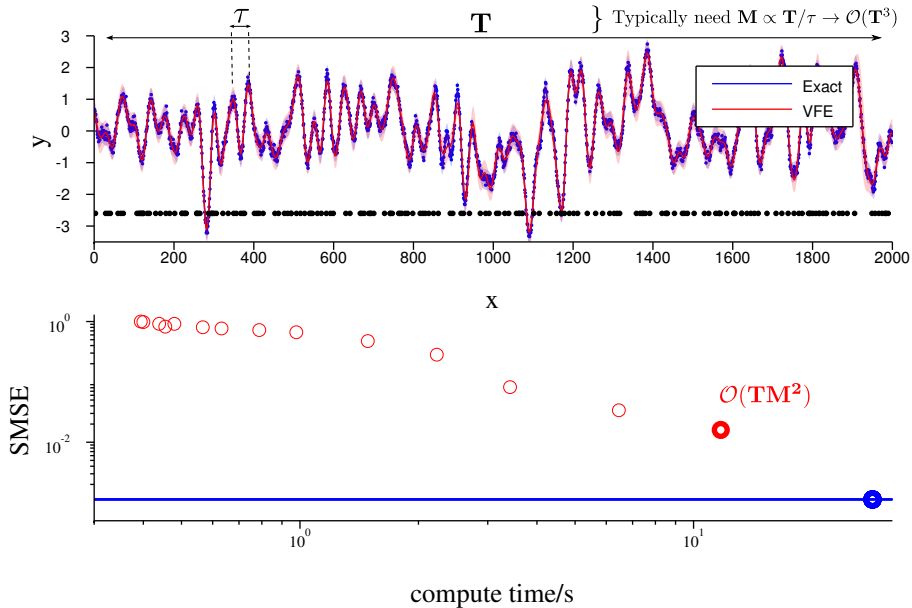
Many GP approximations are poor for time-series



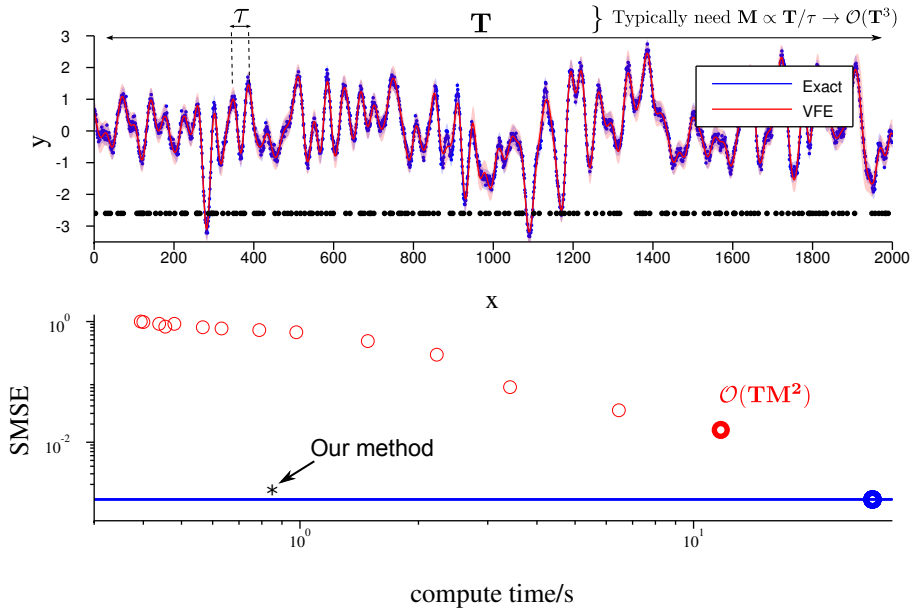
Many GP approximations are poor for time-series



Many GP approximations are poor for time-series

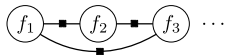


Many GP approximations are poor for time-series

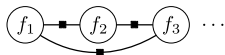


The Fully Independent Training Conditional (FITC) approximation (Snelson and Ghahramani 2006)

The Fully Independent Training Conditional (FITC) approximation (Snelson and Ghahramani 2006)

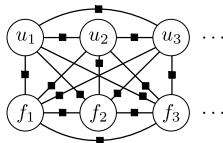


The Fully Independent Training Conditional (FITC) approximation (Snelson and Ghahramani 2006)



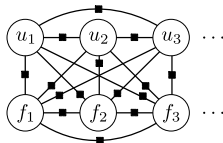
- 1 **Augment the model with inducing points** $\{x_m, u_m\}_{m=1}^M$

The Fully Independent Training Conditional (FITC) approximation (Snelson and Ghahramani 2006)



- 1 **Augment the model with inducing points** $\{x_m, u_m\}_{m=1}^M$

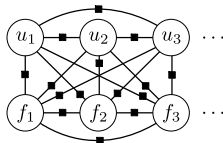
The Fully Independent Training Conditional (FITC) approximation (Snelson and Ghahramani 2006)



- 1 **Augment the model with inducing points** $\{x_m, u_m\}_{m=1}^M$

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{\mathbf{ff}} & \mathbf{K}_{\mathbf{fu}} \\ \mathbf{K}_{\mathbf{uf}} & \mathbf{K}_{\mathbf{uu}} \end{bmatrix} \right)$$

The Fully Independent Training Conditional (FITC) approximation (Snelson and Ghahramani 2006)

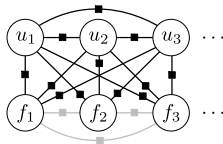


- 1 Augment the model with inducing points $\{x_m, u_m\}_{m=1}^M$

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{\mathbf{ff}} & \mathbf{K}_{\mathbf{fu}} \\ \mathbf{K}_{\mathbf{uf}} & \mathbf{K}_{\mathbf{uu}} \end{bmatrix} \right)$$

- 2 Remove *all* direct dependencies between function values \mathbf{f} ,
i.e. assume $f_i \perp\!\!\!\perp f_j | \mathbf{u}, \forall i, j$

The Fully Independent Training Conditional (FITC) approximation (Snelson and Ghahramani 2006)

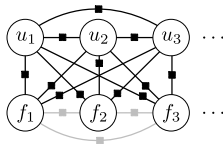


- 1 Augment the model with inducing points $\{x_m, u_m\}_{m=1}^M$

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{\mathbf{ff}} & \mathbf{K}_{\mathbf{fu}} \\ \mathbf{K}_{\mathbf{uf}} & \mathbf{K}_{\mathbf{uu}} \end{bmatrix} \right)$$

- 2 **Remove all direct dependencies between function values \mathbf{f} , i.e. assume $f_i \perp\!\!\!\perp f_j | \mathbf{u}, \forall i, j$**

The Fully Independent Training Conditional (FITC) approximation (Snelson and Ghahramani 2006)



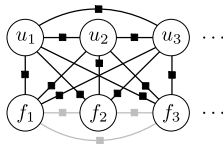
- 1 Augment the model with inducing points $\{x_m, u_m\}_{m=1}^M$

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{\mathbf{ff}} & \mathbf{K}_{\mathbf{fu}} \\ \mathbf{K}_{\mathbf{uf}} & \mathbf{K}_{\mathbf{uu}} \end{bmatrix} \right)$$

- 2 Remove *all* direct dependencies between function values \mathbf{f} , i.e. assume $f_i \perp\!\!\!\perp f_j | \mathbf{u}, \forall i, j$
- 3 **Calibrate model using a forward KL divergence**

$$\arg \min_{q(\mathbf{u}), \{q(f_i | \mathbf{u})\}_{i=1}^N} \text{KL}(p(\mathbf{f}, \mathbf{u}) || q(\mathbf{u}) \prod_{n=1}^N q(f_i | \mathbf{u}))$$

The Fully Independent Training Conditional (FITC) approximation (Snelson and Ghahramani 2006)



- 1 Augment the model with inducing points $\{x_m, u_m\}_{m=1}^M$

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{\mathbf{ff}} & \mathbf{K}_{\mathbf{fu}} \\ \mathbf{K}_{\mathbf{uf}} & \mathbf{K}_{\mathbf{uu}} \end{bmatrix} \right)$$

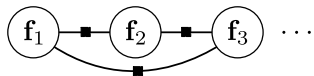
- 2 Remove *all* direct dependencies between function values \mathbf{f} , i.e. assume $f_i \perp\!\!\!\perp f_j | \mathbf{u}, \forall i, j$
- 3 **Calibrate model using a forward KL divergence**

$$\arg \min_{\mathbf{q}(\mathbf{u}), \{q(f_i|\mathbf{u})\}_{i=1}^N} \text{KL}(p(\mathbf{f}, \mathbf{u}) || q(\mathbf{u}) \prod_{n=1}^N q(f_i|\mathbf{u}))$$

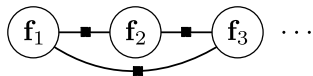
$$\Rightarrow q(\mathbf{u}) = p(\mathbf{u}) \quad , \quad q(f_i|\mathbf{u}) = p(f_i|\mathbf{u})$$

The chain-structured pseudo point approximation

The chain-structured pseudo point approximation

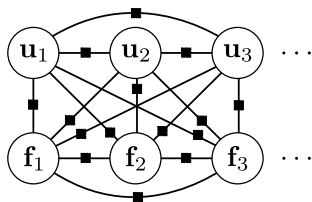


The chain-structured pseudo point approximation



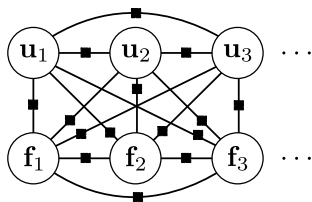
- 1 **Augment the model with inducing points** $\{x_m, u_m\}_{m=1}^M$

The chain-structured pseudo point approximation



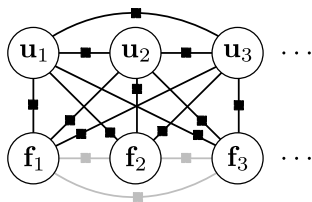
- 1 Augment the model with inducing points $\{x_m, u_m\}_{m=1}^M$

The chain-structured pseudo point approximation



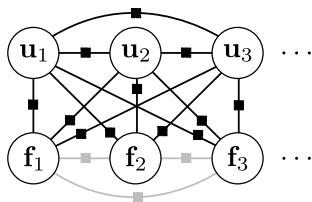
- 1 Augment the model with inducing points $\{x_m, u_m\}_{m=1}^M$
- 2 **Remove *some* direct dependencies between function values f ,**

The chain-structured pseudo point approximation



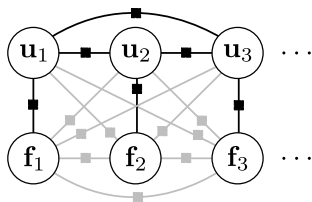
- 1 Augment the model with inducing points $\{x_m, u_m\}_{m=1}^M$
- 2 **Remove *some* direct dependencies between function values f ,**

The chain-structured pseudo point approximation



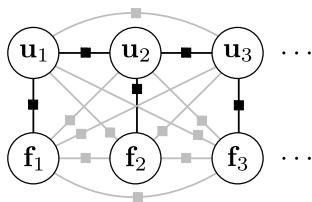
- 1 Augment the model with inducing points $\{x_m, u_m\}_{m=1}^M$
- 2 Remove *some* direct dependencies between function values f ,
and assume a chain structure on inducing points u

The chain-structured pseudo point approximation



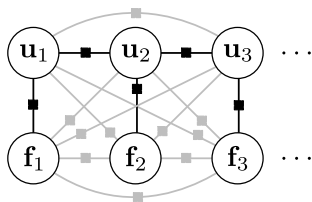
- 1 Augment the model with inducing points $\{x_m, u_m\}_{m=1}^M$
- 2 Remove *some* direct dependencies between function values f ,
and assume a chain structure on inducing points u

The chain-structured pseudo point approximation



- 1 Augment the model with inducing points $\{x_m, u_m\}_{m=1}^M$
- 2 Remove *some* direct dependencies between function values f ,
and assume a chain structure on inducing points u

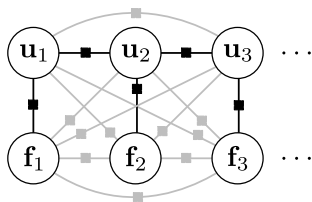
The chain-structured pseudo point approximation



- 1 Augment the model with inducing points $\{x_m, u_m\}_{m=1}^M$
- 2 Remove *some* direct dependencies between function values \mathbf{f} , and assume a chain structure on inducing points \mathbf{u}
- 3 **Calibrate model using a forward KL divergence**

$$\arg \min_{\{q(\mathbf{u}_k|\mathbf{u}_{k-1}), q(\mathbf{f}_k|\mathbf{u}_k)\}_{k=1}^K} \text{KL}(p(\mathbf{f}, \mathbf{u}) || \prod_k q(\mathbf{u}_k|\mathbf{u}_{k-1}) q(\mathbf{f}_k|\mathbf{u}_k))$$

The chain-structured pseudo point approximation

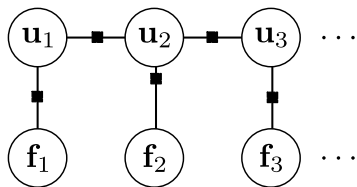


- 1 Augment the model with inducing points $\{x_m, u_m\}_{m=1}^M$
- 2 Remove *some* direct dependencies between function values \mathbf{f} , and assume a chain structure on inducing points \mathbf{u}
- 3 **Calibrate model using a forward KL divergence**

$$\arg \min_{\{q(\mathbf{u}_k|\mathbf{u}_{k-1}), q(\mathbf{f}_k|\mathbf{u}_k)\}_{k=1}^K} \text{KL}(p(\mathbf{f}, \mathbf{u}) || \prod_k q(\mathbf{u}_k|\mathbf{u}_{k-1}) q(\mathbf{f}_k|\mathbf{u}_k))$$

$$\Rightarrow q(\mathbf{u}_k|\mathbf{u}_{k-1}) = p(\mathbf{u}_k|\mathbf{u}_{k-1}) \quad , \quad q(\mathbf{f}_k|\mathbf{u}_k) = p(\mathbf{f}_k|\mathbf{u}_k)$$

The chain-structured pseudo point approximation



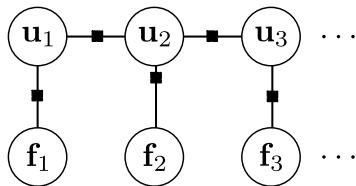
The chain-structured pseudo point approximation

New generative model:

$$q(\mathbf{u}) = \prod_{k=1}^K q(\mathbf{u}_k | \mathbf{u}_{k-1}),$$

$$q(\mathbf{f} | \mathbf{u}) = \prod_{k=1}^K q(\mathbf{f}_k | \mathbf{u}_k),$$

$$p(\mathbf{y} | \mathbf{f}) = \prod_{n=1}^N p(y_n; f_n, \sigma_n^2).$$



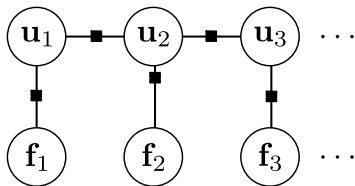
The chain-structured pseudo point approximation

New generative model:

$$q(\mathbf{u}) = \prod_{k=1}^K q(\mathbf{u}_k | \mathbf{u}_{k-1}),$$

$$q(\mathbf{f} | \mathbf{u}) = \prod_{k=1}^K q(\mathbf{f}_k | \mathbf{u}_k),$$

$$p(\mathbf{y} | \mathbf{f}) = \prod_{n=1}^N p(y_n; f_n, \sigma_n^2).$$



where

$$q(\mathbf{u}_{B_k} | \mathbf{u}_{B_{k-1}}) = p(\mathbf{u}_{B_k} | \mathbf{u}_{B_{k-1}}) = \mathcal{N}(\mathbf{u}_{B_k}; \mathbf{A}_k \mathbf{u}_{B_{k-1}}, \mathbf{Q}_k),$$

$$q(\mathbf{f}_{B_k} | \mathbf{u}_{B_k}) = p(\mathbf{f}_{B_k} | \mathbf{u}_{B_k}) = \mathcal{N}(\mathbf{f}_{B_k}; \mathbf{C}_k \mathbf{u}_{B_k}, \mathbf{R}_k).$$

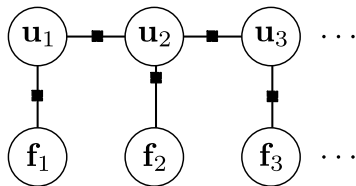
The chain-structured pseudo point approximation

New generative model:

$$q(\mathbf{u}) = \prod_{k=1}^K q(\mathbf{u}_k | \mathbf{u}_{k-1}),$$

$$q(\mathbf{f} | \mathbf{u}) = \prod_{k=1}^K q(\mathbf{f}_k | \mathbf{u}_k),$$

$$p(\mathbf{y} | \mathbf{f}) = \prod_{n=1}^N p(y_n; f_n, \sigma_n^2).$$



where

$$q(\mathbf{u}_{B_k} | \mathbf{u}_{B_{k-1}}) = p(\mathbf{u}_{B_k} | \mathbf{u}_{B_{k-1}}) = \mathcal{N}(\mathbf{u}_{B_k}; \mathbf{A}_k \mathbf{u}_{B_{k-1}}, \mathbf{Q}_k),$$

$$q(\mathbf{f}_{B_k} | \mathbf{u}_{B_k}) = p(\mathbf{f}_{B_k} | \mathbf{u}_{B_k}) = \mathcal{N}(\mathbf{f}_{B_k}; \mathbf{C}_k \mathbf{u}_{B_k}, \mathbf{R}_k).$$

- This is a Linear Dynamical System with a *strange* parameterisation!

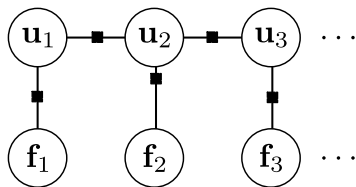
The chain-structured pseudo point approximation

New generative model:

$$q(\mathbf{u}) = \prod_{k=1}^K q(\mathbf{u}_k | \mathbf{u}_{k-1}),$$

$$q(\mathbf{f} | \mathbf{u}) = \prod_{k=1}^K q(\mathbf{f}_k | \mathbf{u}_k),$$

$$p(\mathbf{y} | \mathbf{f}) = \prod_{n=1}^N p(y_n; f_n, \sigma_n^2).$$



where

$$q(\mathbf{u}_{B_k} | \mathbf{u}_{B_{k-1}}) = p(\mathbf{u}_{B_k} | \mathbf{u}_{B_{k-1}}) = \mathcal{N}(\mathbf{u}_{B_k}; \mathbf{A}_k \mathbf{u}_{B_{k-1}}, \mathbf{Q}_k),$$

$$q(\mathbf{f}_{B_k} | \mathbf{u}_{B_k}) = p(\mathbf{f}_{B_k} | \mathbf{u}_{B_k}) = \mathcal{N}(\mathbf{f}_{B_k}; \mathbf{C}_k \mathbf{u}_{B_k}, \mathbf{R}_k).$$

- This is a Linear Dynamical System with a *strange* parameterisation!
- Inference using Kalman smoothing algorithm

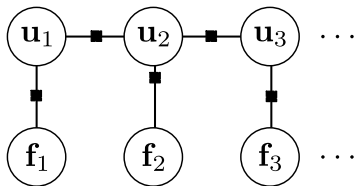
The chain-structured pseudo point approximation

New generative model:

$$q(\mathbf{u}) = \prod_{k=1}^K q(\mathbf{u}_k | \mathbf{u}_{k-1}),$$

$$q(\mathbf{f} | \mathbf{u}) = \prod_{k=1}^K q(\mathbf{f}_k | \mathbf{u}_k),$$

$$p(\mathbf{y} | \mathbf{f}) = \prod_{n=1}^N p(y_n; f_n, \sigma_n^2).$$



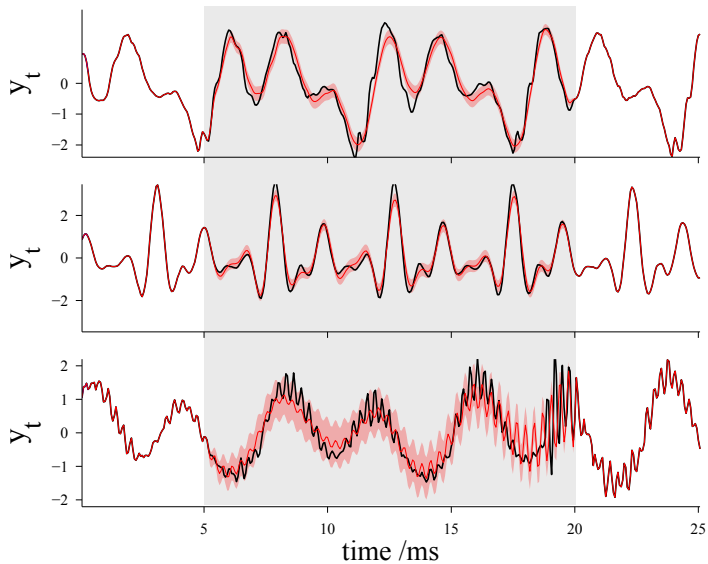
where

$$q(\mathbf{u}_{B_k} | \mathbf{u}_{B_{k-1}}) = p(\mathbf{u}_{B_k} | \mathbf{u}_{B_{k-1}}) = \mathcal{N}(\mathbf{u}_{B_k}; \mathbf{A}_k \mathbf{u}_{B_{k-1}}, \mathbf{Q}_k),$$

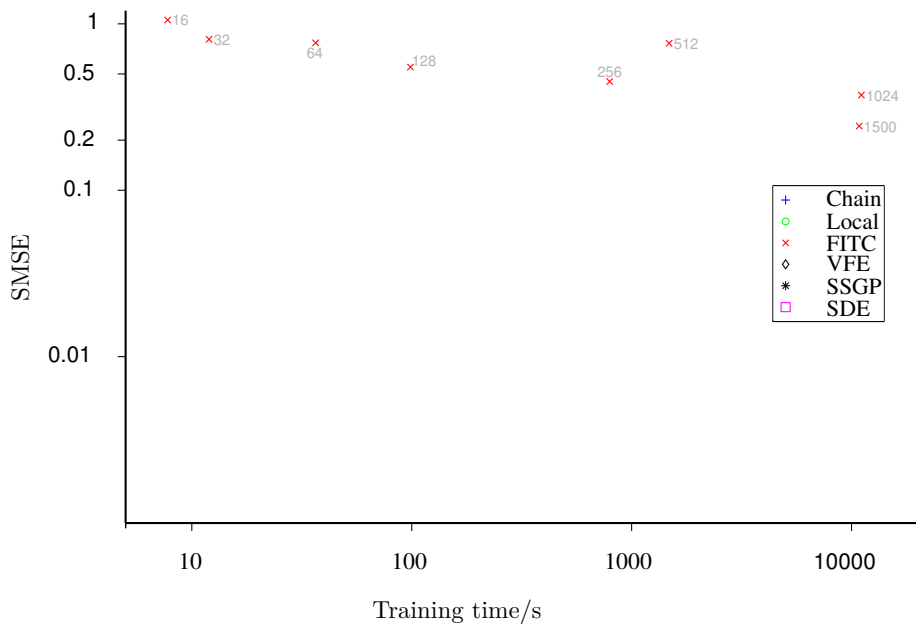
$$q(\mathbf{f}_{B_k} | \mathbf{u}_{B_k}) = p(\mathbf{f}_{B_k} | \mathbf{u}_{B_k}) = \mathcal{N}(\mathbf{f}_{B_k}; \mathbf{C}_k \mathbf{u}_{B_k}, \mathbf{R}_k).$$

- This is a Linear Dynamical System with a *strange* parameterisation!
- Inference using Kalman smoothing algorithm
- Complexity: $\mathcal{O}(TD^2)$, D : average number of observations per block

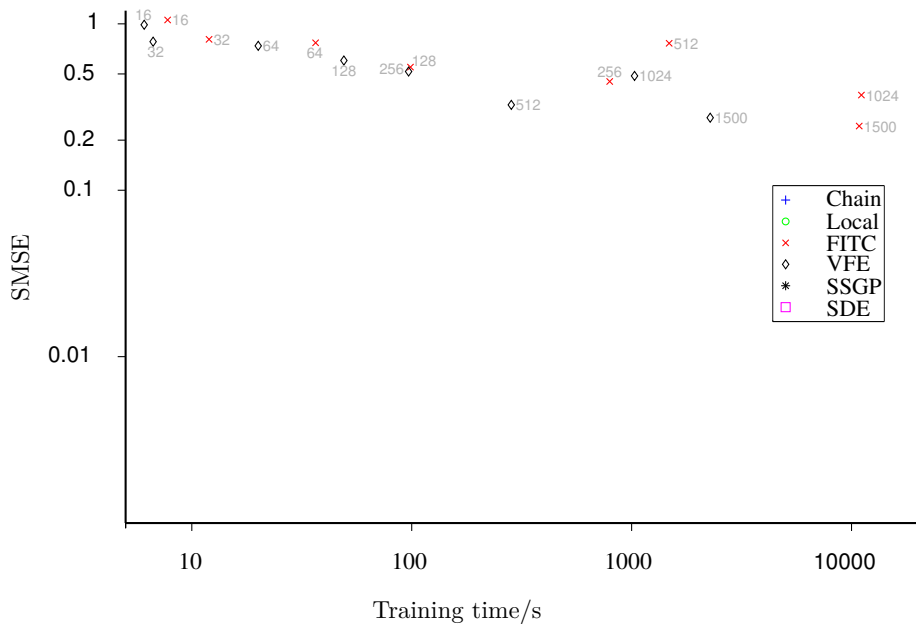
Results: Audio missing data imputation



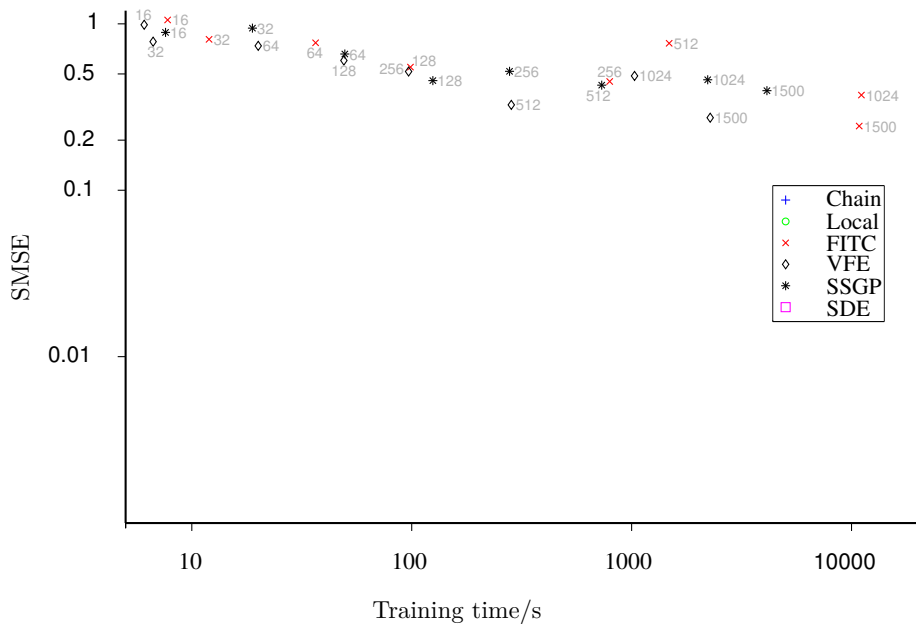
Results: Speed accuracy trade-off



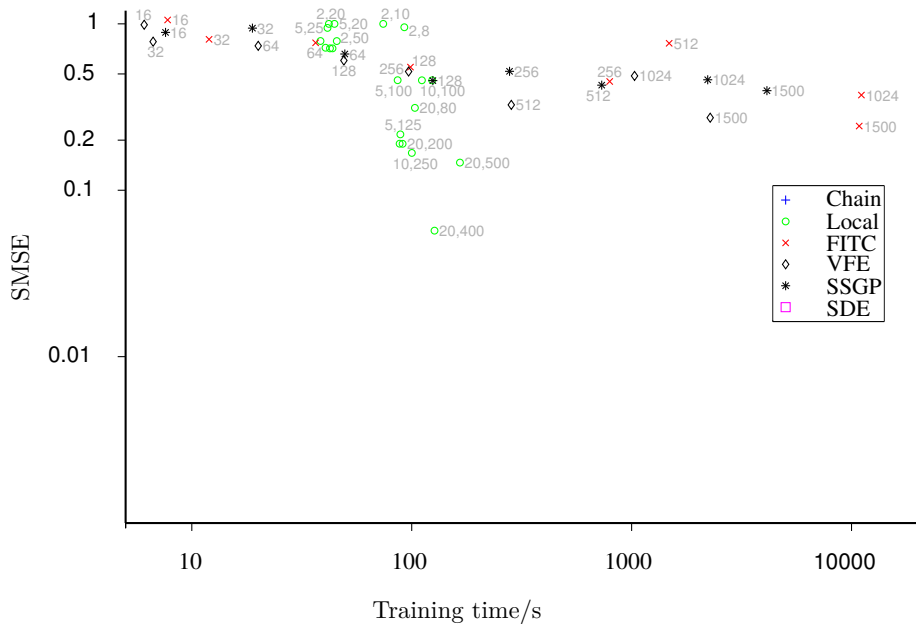
Results: Speed accuracy trade-off



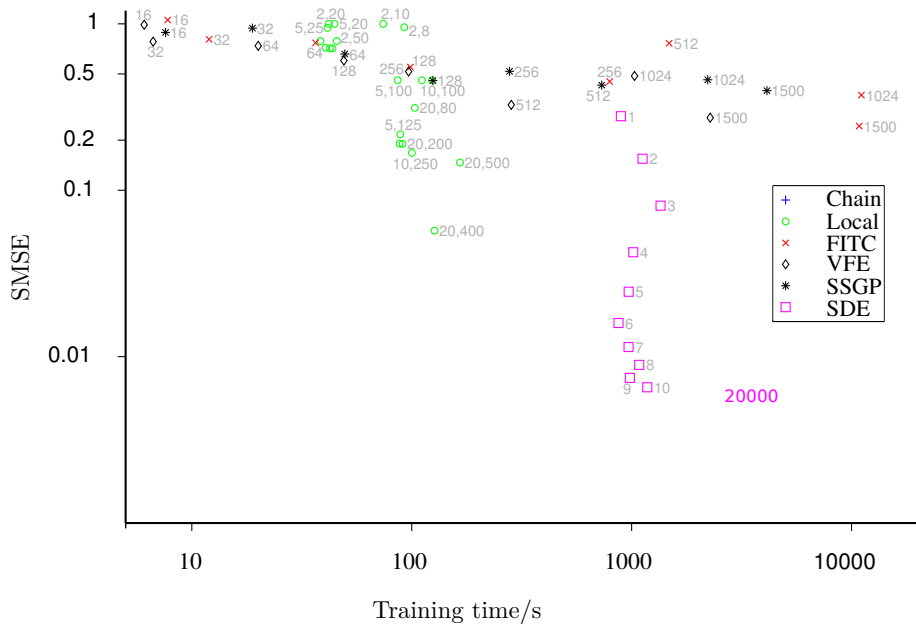
Results: Speed accuracy trade-off



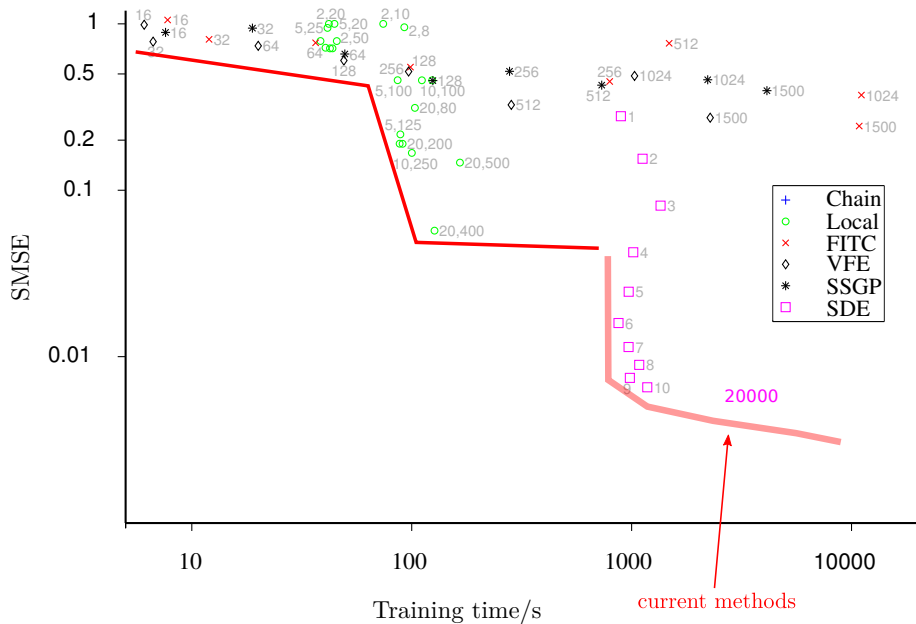
Results: Speed accuracy trade-off



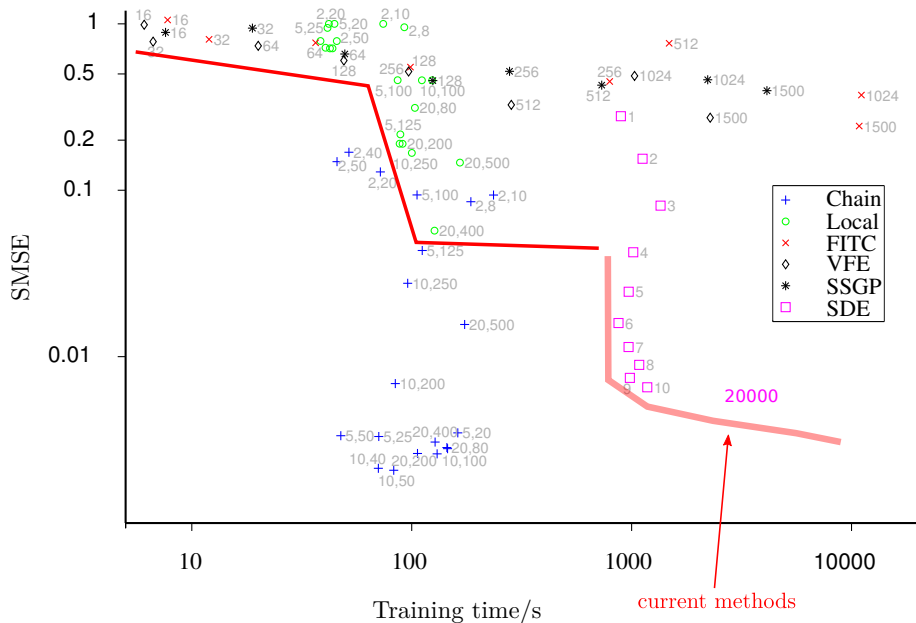
Results: Speed accuracy trade-off



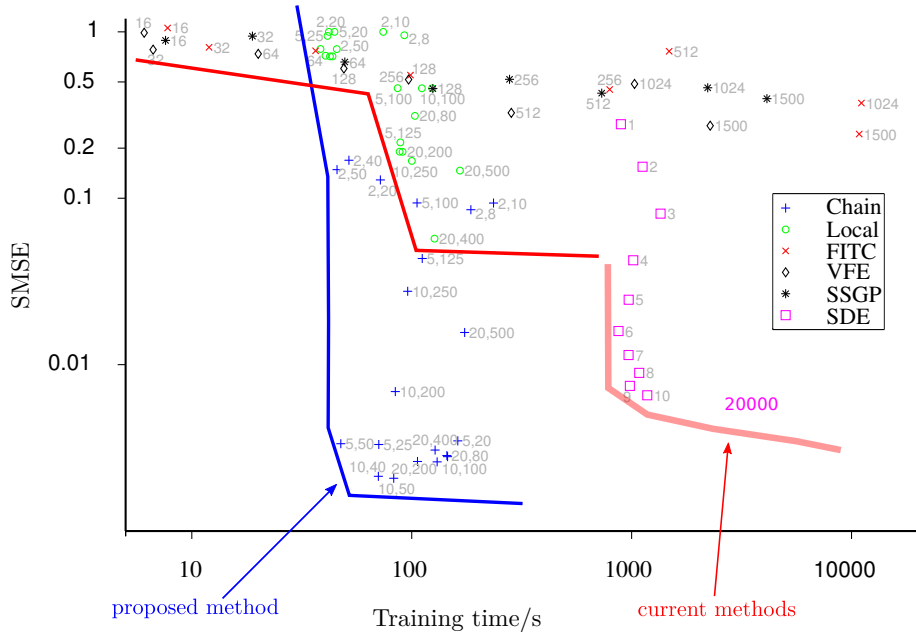
Results: Speed accuracy trade-off



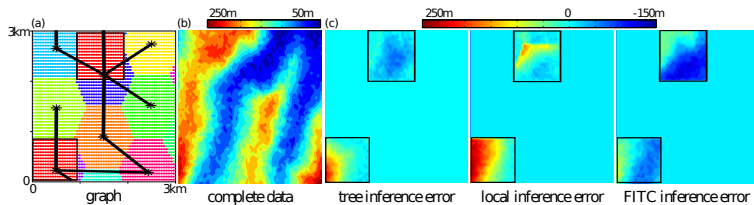
Results: Speed accuracy trade-off



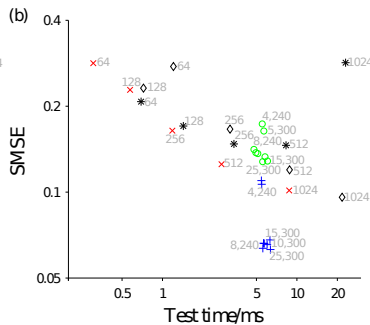
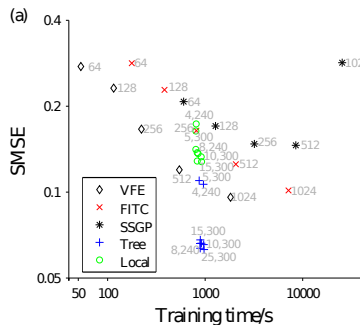
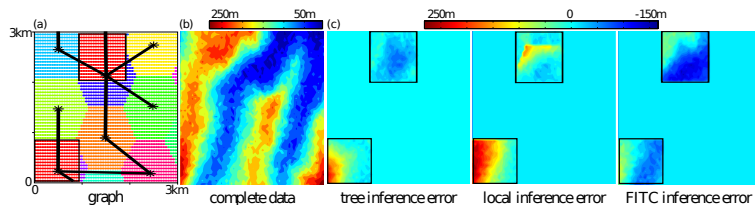
Results: Speed accuracy trade-off



Results: 2D spatial dataset (tree structured)



Results: 2D spatial dataset (tree structured)



Summary

- pseudo-dataset approximation methods **must grow in size with the length of the time-series**
- **simple extension to FITC** (or PITC) that imposes tree-structured conditional dependencies
- fast inference by the **up-down** algorithm

Open questions and current work

- indirect approximation method
 - ▶ involves exact inference in an approximate model
 - ▶ can we use similar ideas for direct approximation of the true posterior?
- connections between GPs and time-frequency analysis
 - ▶ multi-rate filters and striding as variational free-energy + FFT based approximations
 - ▶ rediscover Nyquist in the context of limits on GP approximation accuracy