# Supplementary Material:
# Sequential Inference for Deep Gaussian Process

**Yali Wang**
Chinese Academy of Sciences
yl.wang@siat.ac.cn

**Marcus Brubaker**
University of Toronto
mbrubake@cs.toronto.edu

**Brahim Chaib-draa**
Laval University
chaib@ift.ulaval.ca

**Raquel Urtasun**
University of Toronto
urtasun@cs.toronto.edu

## 1   Sparse Online Gaussian Process

In this section we briefly review sparse online GPs ($\text{GP}_{so}$) [1, 2]. The key idea is to learn GPs recursively by updating the posterior mean and covariance of the training set $\{(\mathbf{x}^n, y^n)\}_{n=1}^N$ in a sequential fashion. This online procedure is coupled with a sparsification mechanism in which a fixed-size subset of the training set (called the active set) is iteratively selected to avoid the unbounded computation growth of updating.

Specifically, the model parameters at the $(n-1)$-th step of $\text{GP}_{so}$ are [1, 2]

$$M^{n-1} = \{\tilde{\mathcal{X}}^{n-1}, \mu^{n-1}, \Sigma^{n-1}, Q^{n-1}\},$$

where $\tilde{\mathcal{X}}^{n-1}$ is the active set which is a training subset selected from the first $(n-1)$ training pairs, $\mathcal{N}(\mu^{n-1}, \Sigma^{n-1})$ is the posterior over $\tilde{\mathcal{X}}^{n-1}$, and $Q^{n-1}$ is the inverse covariance matrix of $\tilde{\mathcal{X}}^{n-1}$.

Once the $n$-th training pair $(\mathbf{x}^n, y^n)$ is available, the model parameters is updated from $M^{n-1}$ to $M^n$, based on the following strategy.

### 1.1   Update at the $n$-th step

Firstly, the following update is performed to take the information of $(\mathbf{x}^n, y^n)$ into account [2],

$$\gamma_n^2 = k(\mathbf{x}^n, \mathbf{x}^n) - \mathbf{k}_{n-1}^T(\mathbf{x}^n)\mathbf{q}_n, \tag{1}$$

$$\mu^n = \begin{bmatrix} \mu^{n-1} \\ \mathbf{q}_n^T \mu^{n-1} \end{bmatrix} + \frac{y^n - \mathbf{q}_n^T \mu^{n-1}}{\eta_n^2}\psi_n, \tag{2}$$

$$\Sigma^n = \begin{bmatrix} \Sigma^{n-1} & \Sigma^{n-1}\mathbf{q}_n \\ \mathbf{q}_n^T\Sigma^{n-1} & \gamma_n^2 + \mathbf{q}_n^T\Sigma^{n-1}\mathbf{q}_n \end{bmatrix} - \frac{\psi_n\psi_n^T}{\eta_n^2}, \tag{3}$$

where the relevant computation quantities are

$$\mathbf{q}_n = Q_{n-1}\mathbf{k}_{n-1}(\mathbf{x}^n),$$

$$\psi_n = \begin{bmatrix} \Sigma^{n-1}\mathbf{q}_n \\ \gamma_n^2 + \mathbf{q}_n^T\Sigma^{n-1}\mathbf{q}_n \end{bmatrix},$$

and the $i$-th entry of the vector $\mathbf{k}_{n-1}(\mathbf{x}^n)$ is computed by $k(\mathbf{x}^n, \mathbf{x}^i)$ and $\mathbf{x}^i$ is the $i$-th input in $\tilde{\mathcal{X}}^{n-1}$.

Next, depending on the value of $\gamma_n^2$ in Eq.(1), we decide if the $n$-th training point is added to the active set and how to further refine the model parameters $M^n$.

• When $\gamma_n^2 < \delta$ ($\delta = 10^{-6}$ in this work), the training pair $(\mathbf{x}^n, y^n)$ is not added to the active set. As a result, $\mu^n$ and $\Sigma^n$ in Eq.(2-3) are reduced to

$$\mu^n \leftarrow [\mu^n]_{-n}, \quad \Sigma^n \leftarrow [\Sigma^n]_{-n,-n}, \tag{4}$$

where $[\cdot]_{-n}$ deletes the $n$-th entry of a vector; $[\cdot]_{-n,-n}$ deletes the $n$-th row and column of a matrix. Additionally, since the active set does not change $\tilde{\mathcal{X}}^n = \tilde{\mathcal{X}}^{n-1}$, the inverse covariance matrix is $Q^n = Q^{n-1}$.

• When $\gamma_n^2 \geq \delta$, the training pair $(\mathbf{x}^n, y^n)$ is added to the active set, i.e., $\tilde{\mathcal{X}}^n = \tilde{\mathcal{X}}^{n-1} \cup (\mathbf{x}^n, y^n)$. The $\mu^n$ and $\Sigma^n$ are computed using Eq.(2-3). The inverse covariance matrix is then updated as follows

$$Q^n = \begin{bmatrix} Q^{n-1} & \mathbf{0} \\ \mathbf{0} & 0 \end{bmatrix} + \frac{1}{\gamma_n^2}\begin{bmatrix} \mathbf{q}_n \\ -1 \end{bmatrix}\begin{bmatrix} \mathbf{q}_n \\ -1 \end{bmatrix}^T. \tag{5}$$

In this case, we have to make sure that the size of the active set does not exceed our budget. If the size of $\tilde{\mathcal{X}}^n$ exceeds $N_{AC}$, we remove the data pair that has the lowest squared prediction error

$$k = \arg\min_k \left(\frac{[Q^n\mu^n]_k}{[Q^n]_{k,k}}\right)^2,$$

where $[\cdot]_k$ is the $k$-th entry of a vector, $[\cdot]_{k,k}$ is the $k$-th diagonal entry of a matrix. Consequentially, the active set $\tilde{\mathcal{X}}^n$ is reduced to $\tilde{\mathcal{X}}^n \leftarrow \tilde{\mathcal{X}}^n \setminus (\mathbf{x}^k, y^k)$ and $\mu^n$, $\Sigma^n$, $Q^n$ in Eq.(2),(3),(5) are updated to be

$$\mu^n \leftarrow [\mu^n]_{-k},$$
$$\Sigma^n \leftarrow [\Sigma^n]_{-k,-k},$$
$$Q^n \leftarrow [Q^n]_{-k,-k} - \frac{[Q^n]_{-k,k}[Q^n]_{-k,k}^T}{[Q^n]_{k,k}},$$

where $[\cdot]_{-k,k}$ is the $k$-th column of a matrix without the entry in the $k$-th row.

## 1.2 Prediction at the $n$-th step

Based on the model parameters $M^n$ at the $n$-th step, the predictive distribution at a given test input $\mathbf{x}^\star$ is Gaussian [1, 2]

$$p(y^\star|\mathbf{x}^\star, M^n) = \mathcal{N}(\mu_{sogp\star}, \sigma^2_{sogp\star}) \qquad (6)$$

with mean

$$\mu_{sogp\star} = \mathbf{q}_\star^T \mu^n, \quad \mathbf{q}_\star = Q^n \mathbf{k}_n(\mathbf{x}^\star)$$

and variance

$$\begin{aligned} \sigma^2_{sogp\star} = & \sigma^2 + k(\mathbf{x}^\star, \mathbf{x}^\star) + \mathbf{q}_\star^T \Sigma^n \mathbf{q}_\star \\ & - \mathbf{k}_n^T(\mathbf{x}^\star) Q^n \mathbf{k}_n(\mathbf{x}^\star) \end{aligned}$$

where $\sigma^2$ is the noise variance, the $i$-th entry of the vector $\mathbf{k}_n(\mathbf{x}^\star)$ is computed by $k(\mathbf{x}^\star, \mathbf{x}^i)$ and $\mathbf{x}^i$ is the $i$-th input in $\tilde{\mathcal{X}}^n$.

## 2 Experiments

In this section, we describe the full set of experiments we did to validate our approach. We mainly perform our sequential inference for Deep GP (our $\text{DGP}_{seq}$) and compare it to variational inference for Deep GP ($\text{DGP}_{var}$) [3] and sparse online GP ($\text{GP}_{so}$) [1]. Unless otherwise stated we keep the experimental settings (which we now describe) consistent. We normalize the output $y$ of all the data sets to $[0, 1]$ and then subtract the mean. We use five random train/test partitions for all the data sets. Based on these five partitions, we report average root mean squared error (RMSE) and mean negative log probability (MNLP) as a measure of predictive error, and report average training time (second) as a measure of computational complexity. We use PCA to initialize the latent layers of DGP, and choose an Automatic Relevance Determination (ARD) squared exponential kernel for all relevant approaches,

$$k(\mathbf{x}, \mathbf{x}') = \sigma^2_{ker} \exp[-0.5 \sum\nolimits_{i=1}^{d} c_i(\mathbf{x}_i - \mathbf{x}'_i)^2]$$

where $\sigma^2_{ker}$ is the amplitude and $c_1, \cdots, c_d$ are the ARD weights.

In the following we demonstrate the effectiveness of our approach for a number of important learning tasks in the GP community.

## 2.1 DGP for Unsupervised Learning

We start our experimental evaluation with two unsupervised learning applications: learning a deep dynamical prior for time series data and dimensionality reduction for image reconstruction.

### 2.1.1 DGP as Deep Dynamical Prior

DGP can be used as a flexible deep dynamical prior, where the input is the time step and the output is the multi-dimensional observation [4, 3]. We use three different datasets to illustrate the robustness of DGP: the *motion*, *flu* and *stock* datasets. The *motion* dataset[1] consists of 2465 time/pose pairs of five physical exercise activities (jumping jacks, two types of side twists, squats and jogging) from CMU Mocap database. Each pose is parameterized with a 62 dimensional vector containing the 3D rotations of all joints. The *flu* dataset[2] consists of 543 time/flu-activity-rate pairs from 2003-11-09 to 2014-03-30. Each flu-activity-rate is a 9 dimensional vector containing the flu rates of AB / BC / MB / NB / NL / NS / ON / SK / QC in Canada. The *stock* dataset[3] consists of 1500 time/Nasdaq index pairs collected from 1997-01-02 to 2014-11-18 (sampled every three working days). Each index is a 5 dimensional vector containing Nasdaq log-returns of Biotechnology, Composite, Industrial, Nasdaq100 and Telecommunications.

We randomly create five partitions of the data with training set sizes 1500/300/1000 for *motion*/*flu*/*stock* datasets respectively. We use $L = 2$, $D = 5$, $N_p = 100$ and $N_{AC} = 100/100/200$ as the basic settings for our $\text{DGP}_{seq}$. We run our sequential inference and hyperparameter learning twice (called two episodes) to get the hyperparameters into a reasonable region. Additionally, we run our sequential inference for both heteroscedastic GP and DGP, and denote them as our $\text{HGP}_{seq}$ and our $\text{HDGP}_{seq}$. We employ the same parameters of our $\text{DGP}_{seq}$ (whenever applicable) to our $\text{HGP}_{seq}$, our $\text{HDGP}_{seq}$, $\text{GP}_{so}$ [1] and $\text{DGP}_{var}$ [3].

**Complexity of the Deep Model**: We evaluate $\text{DGP}_{var}$ [3] and our $\text{DGP}_{seq}$ when varying the number of latent layers $L$ and dimensionality of each layer $D$. The results are shown in Figure 1 and 2. For the motion and flu datasets, our $\text{DGP}_{seq}$ outperforms the basic $\text{DGP}_{var}$[3] in terms of both prediction error and training efficiency. This illustrates that our sequential inference is more effective than the variational inference of [3]. For the stock dataset, our $\text{DGP}_{seq}$ (when changing $L$) and $\text{DGP}_{var}$[3] (when changing $D$) tend to achieve a poor performance because of the strong heteroscedasticity in this data. Additionally, in this data, the RMSE is similar for all approaches and settings. This might be because the stock observations are log-returns in which the fluctuation of the underlying function is flat. All the deep GP models can capture this fluctuation, and thus produce similar RMSEs.

---

[1]http://mocap.cs.cmu.edu/
[2]http://www.google.org/flutrends/ca/
[3]http://finance.yahoo.com/

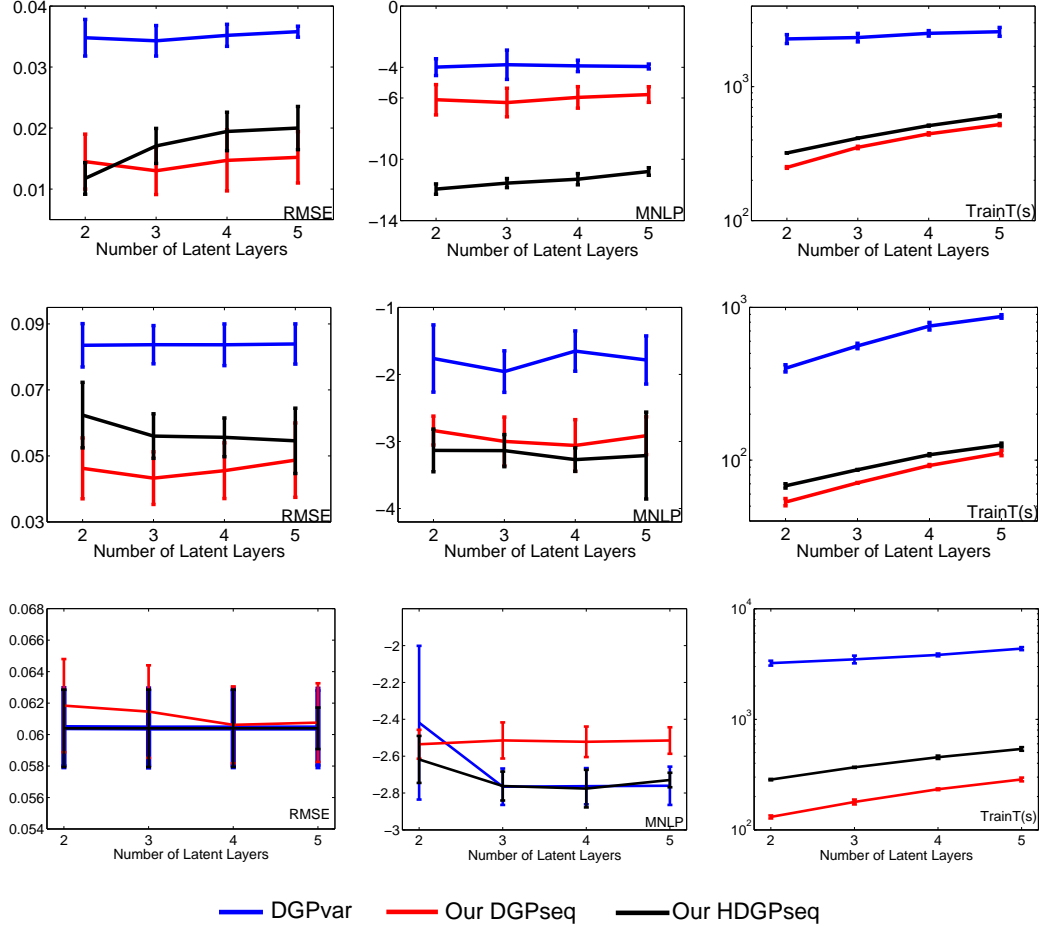**Yali Wang, Marcus Brubaker, Brahim Chaib-draa, Raquel Urtasun**



Figure 1: DGP as Deep Dynamical Prior (Robustness to Complexity of the Deep Model): RMSE/MNLP/Training Time (Column 1/2/3) as a function of the number of latent layers $L$. Row1/2/3: motion/flu/stock. The error bar is mean±standard deviation.

**Properties of Our Inference Framework**: We next evaluate the influence of the number of particles, $N_p$, and the size of the active set, $N_{AC}$ in our sequential inference framework. The results are shown in Figure 3 and 4. As expected, the performance of all our approaches tend to improve when $N_p$ and/or $N_{AC}$ increases. Additionally, our $DGP_{seq}$ outperforms $GP_{so}$ [1]. This illustrates that the predictive performance of deep GP models is generally better than shallow GP models.

**Heteroscedastic Learning**: We now evaluate the sequential inference procedure for the heteroscedastic DGP extension. As shown in Figure 1, 2, 3 and 4, our $HDGP_{seq}$ generally outperforms our $DGP_{seq}$ with a competitive RMSE but a lower MNLP. This is due to the fact that the heteroscedastic observation layer in our $HDGP_{seq}$ is able to achieve a lower prediction uncertainty in regions of the space which have less noise. Hence, compared to our $DGP_{seq}$, our $HDGP_{seq}$ can

further improve the prediction performance without significantly increasing computation.

**Missing Data in Multi-task Learning**: Since DGP can be used as a multi-output GP model, we assess the performance of our approaches for missing data cases. For *motion*, the observations of dimensions 27-33/49-55 (right arm / left leg) are missing during the time steps 700-720/1200-1220. For *flu*, the observations of provinces 1-5/6-9 are missing during the time steps 80-100/30-50. For *stock*, the observations of Nasdaq Biotechnology / Nasdaq Composite / Nasdaq Industrial / Nasdaq100 / Nasdaq Telecommunications are missing during the time steps 150-200/350-400/550-600/750-800/950-1000. We compare our $DGP_{seq}$ and $HDGP_{seq}$ to the baseline $GP_{so}$ [1]. Table 1 shows that our $DGP_{seq}$ and $HDGP_{seq}$ outperform $GP_{so}$ with a lower reconstruction error of missing dimensions. This is primarily due to the fact that deep structures can capture correlations between different outputs via
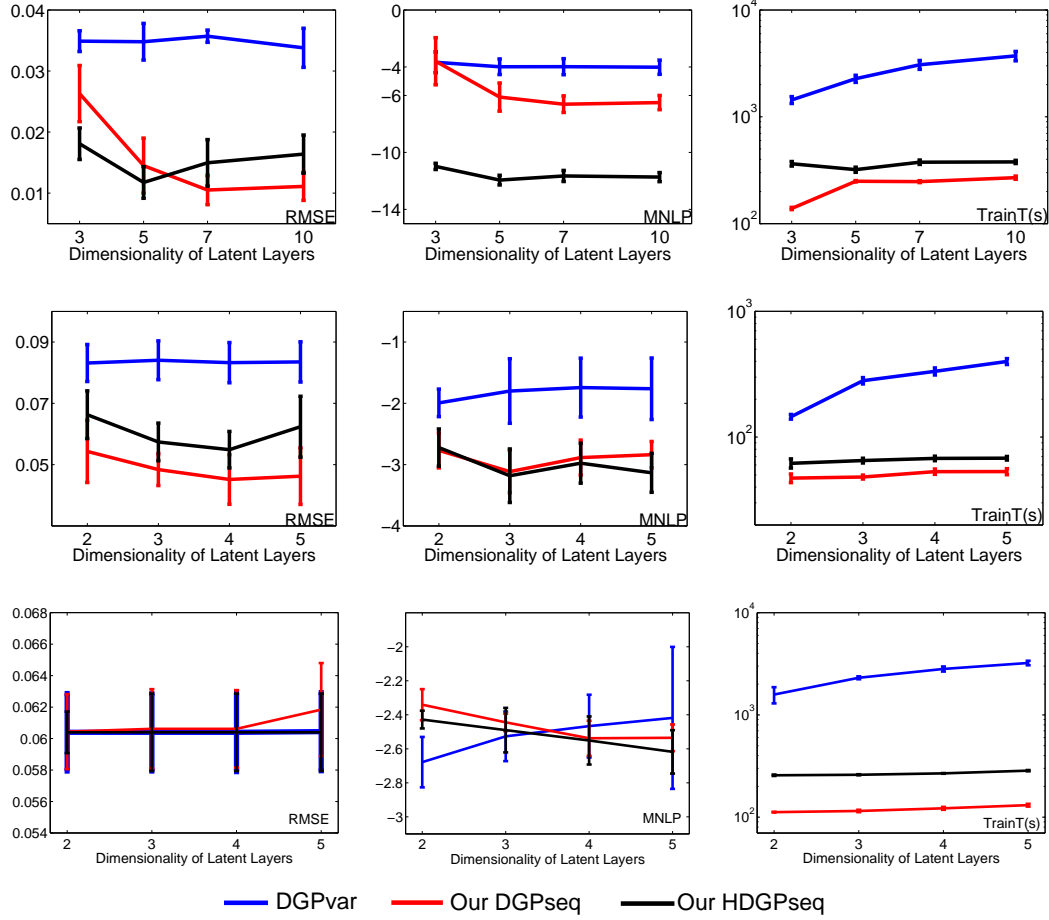
Figure 2: DGP as Deep Dynamical Prior (Robustness to Complexity of the Deep Model): RMSE/MNLP/Training Time (Column 1/2/3) as a function of the dimensionality of latent layers $D$. Row1/2/3: motion/flu/stock. The error bar is mean±standard deviation.

sharing a number of latent layers.

**Qualitative Results**: **(i)** For the *motion* dataset, we show that our particles are effective and efficient to estimate **y** in Fig. 5. The settings of our $\text{DGP}_{seq}$ in this figure are the number of latent layers $L = 2$, the dimensionality of latent layers $D = 10$, the number of particles $N_p = 100$, the size of active set $N_{AC} = 100$. **(ii)** For the *flu* dataset, we estimate its posterior using our $\text{HDGP}_{seq}$ with the number of latent layers $L = 2$, the dimensionality of latent layers $D = 5$, the number of particles $N_p = 100$, the size of active set $N_{AC} = 100$. As shown in Fig. 6, our $\text{HDGP}_{seq}$ successfully captures non-stationarity and heteroscedasticity of flu-activity-rate observations. **(iii)** For the *stock* dataset, we estimate its posterior using our $\text{HDGP}_{seq}$ with the number of latent layers $L = 2$, the dimensionality of latent layers $D = 10$, the number of particles $N_p = 100$, the size of active set $N_{AC} = 200$. As shown in Fig. 7, our $\text{HDGP}_{seq}$ successfully captures the strong heteroscedasticity of this data.

### 2.1.2 Dimensionality Reduction for Image Reconstruction

Low-dimensional latent representations provide a means to avoid the so-called "curse of dimensionality" and the low-dimensional latent layers of a DGP can be understood as a form of dimensionality reduction. Inspired by [5], we use DGP for dimensionality reduction to reconstruct images from noisy observations. Towards this goal we use the *Frey* face dataset[4] which is composed of 1900 images of size $20 \times 28 = 560$. We add Gaussian noise (std dev 0.1) to the images, and use these noisy images as both the input and output of the DGP model.

Figure 8 shows the RMSE between the noiseless image and the reconstruction as a function of the number of training episodes in our learning framework. Hyperparameter learning was performed after the first episode and the hyperparameters were fixed for the remaining

---

[4]http://www.cs.nyu.edu/~roweis/data.html

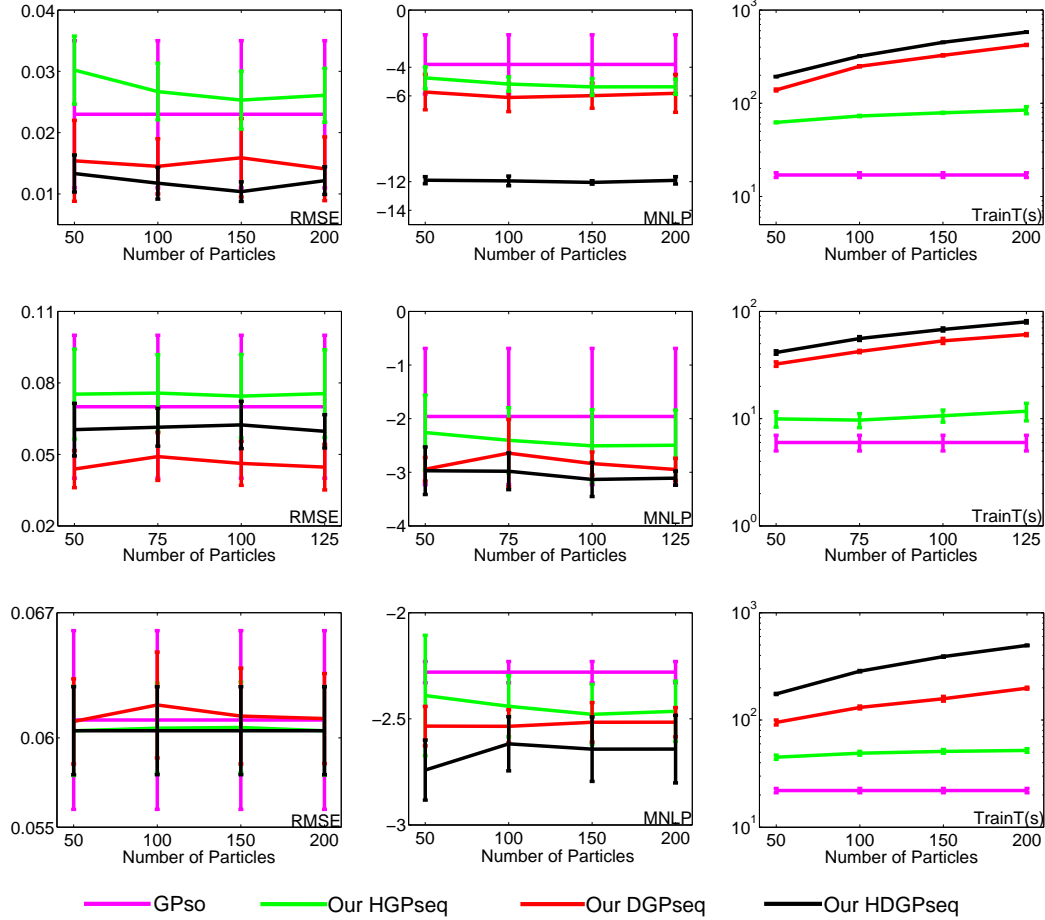**Yali Wang, Marcus Brubaker, Brahim Chaib-draa, Raquel Urtasun**

Figure 3: DGP as Deep Dynamical Prior (Properties of Our Inference Framework): RMSE/MNLP/Training Time (Column 1/2/3) as a function of the number of particles $N_p$. Row1/2/3: motion/flu/stock. The error bar is mean±standard deviation. Note that there is no $N_p$ in $GP_{so}$, hence $GP_{so}$ is a straight line for the plot of $N_p$.

episodes. We chose the number of particles $N_p$ and the size of active set $N_{AC}$ to be 100 in our $DGP_{seq}$. To reduce the randomness of sampling, we run our $DGP_{seq}$ five times for each episode and report the average reconstruction RMSE. As shown in Fig. 8, our $DGP_{seq}$ outperforms $DGP_{var}$. Furthermore, as expected, our $DGP_{seq}$ performs better when the number of training episodes increases, or the width and the depth of the deep structure also increases. Finally, Fig. 9 shows qualitatively that our $DGP_{seq}$ achieves better reconstruction than $DGP_{var}$.

## 2.2 DGP for Supervised Learning

In this section, we focus our evaluation in the supervised setting for large training sets and/or high-dimensional data sets.

### 2.2.1 Regression

We employ the *Parkinsons* telemonitoring dataset[5], which is a six-month biomedical voice recording from 42 people with early-stage Parkinson's disease for a total of 5875 input/output pairs. The input is a 16-dimensional biomedical voice feature vector and the output is a 2-dimensional score vector (motor UPDRS score and total UPDRS score). We randomly partition the data five times and choose each time training/test sets to be of size 5000/875. The basic setting of our $DGP_{seq}$ is $L = 2$, $D = 2$, $N_p = 500$, $N_{AC} = 100$. We run our $DGP_{seq}$ for one training episode with sequential inference and hyperparameter learning. We employ the same parameters of our $DGP_{seq}$ to $GP_{so}$ and $DGP_{var}$ whenever applicable. As shown in Table 2, our $DGP_{seq}$ outperforms $GP_{so}$ and $DGP_{var}$. Furthermore, compared to $DGP_{var}$, our $DGP_{seq}$ significantly reduces the computation.

---
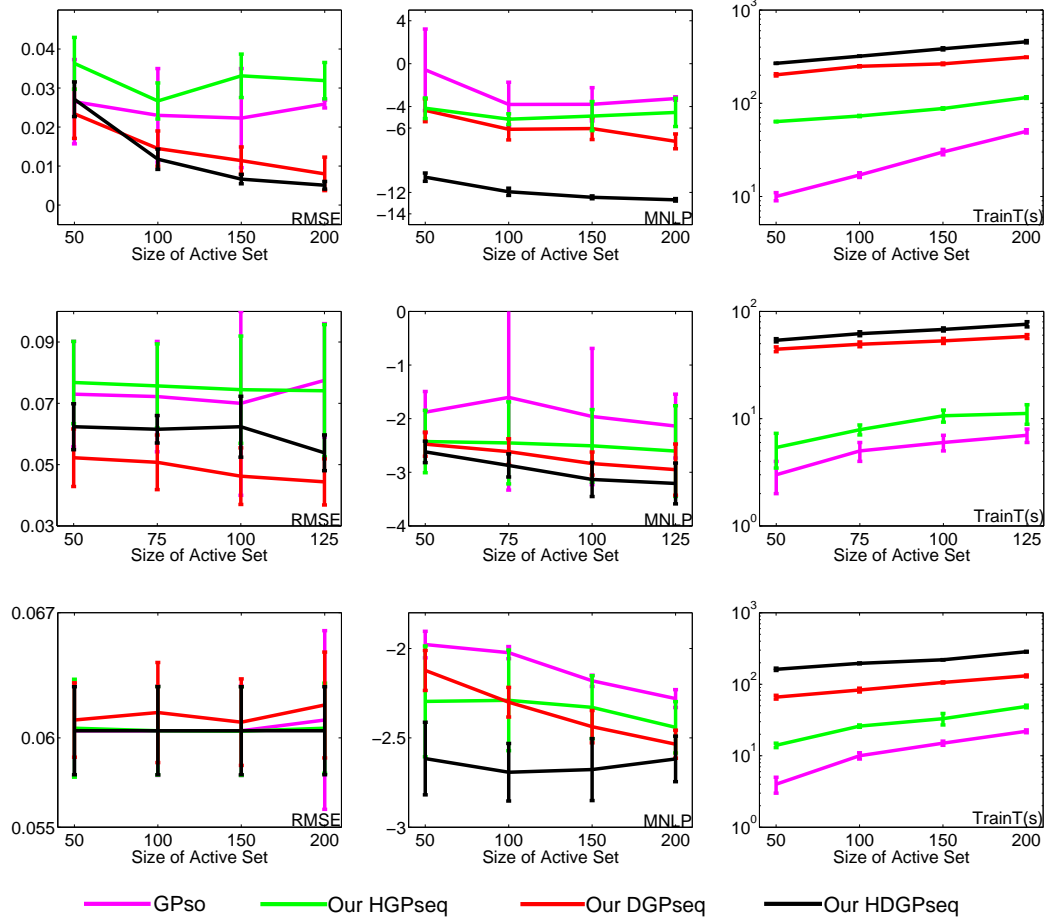
[5]http://archive.ics.uci.edu/ml/index.html

Figure 4: DGP as Deep Dynamical Prior (Properties of Our Inference Framework): RMSE/MNLP/Training Time (Column 1/2/3) as a function of the size of active set $N_{AC}$. Row1/2/3: motion/flu/stock. The error bar is mean±standard deviation.

### 2.2.2 Classification

We evaluate our approach for classification on the *Tumor* gene expression dataset[6] and the *MNIST* Digits dataset[7]. The tumor data consists of a predefined training and test sets of sizes 144 and 46 respectively. The input contains 16,063 tumor genes (16,063 dimensional input vector) and the output contains 14 cancer classes (14 dimensional binary output vector). For the tumor data, the basic settings of our $DGP_{seq}$ are $L = 3$, $D = 12$, $N_p = 200$. $N_{AC}$ is the size of the data as we did not employ sparsification. We use a linear kernel to reduce computation as the data is very high-dimensional. We run our $DGP_{seq}$ for five training episodes with sequential inference and hyperparameter learning. We employ the same parameters for our $DGP_{seq}$, $GP_{so}$ [1] and $DGP_{var}$ [3] whenever applicable. As shown in Table 3, our $DGP_{seq}$ achieves the best classification accuracy and also significantly out-

performs $DGP_{var}$ [3] in terms of computation.

For the MNIST data, there are 60,000 training images and 10,000 test images. The input is a $28 \times 28$ image (784 dimensional input vector) and the output contains 0-9 digits (10 dimensional binary output vector). Classification accuracy on the test set (i.e., the percentage of correctly classified examples) is evaluation metric. The basic settings of our $DGP_{seq}$ are $L = 1$, $D = 400$, $N_p = 100$ and $N_{AC} = 1000$. The covariance function of GP is chosen as the 1st-order arc-cosine kernel [6]. Additionally, instead of using the latent state $h$ as the output of each layer, we use the logistic activation function of $h$, i.e., $(1 + e^{-h})^{-1}$ as the output of each layer, which mimics neural networks for classification. Note that this modification is simple to accommodate in our inference algorithm by simply propagating the particles through the activation function when sampling. We run our $DGP_{seq}$ for one training episode. We employ the same parameters for $GP_{so}$ [1] whenever applicable. As this data

---

[6]http://www.genome.wi.mit.edu/

[7]http://yann.lecun.com/exdb/mnist/

| Data | Methods | RMSE(%) | MNLP | Train_T(s) |
|---|---|---|---|---|
| Motion | $\text{GP}_{so}$ | 1.45 | -4.81 | 20 |
| | our $\text{DGP}_{seq}$ | 1.0±0.1 | -4.95±0.75 | 202±2 |
| | our $\text{HDGP}_{seq}$ | 1.1±0.3 | -5.15±0.20 | 236±8 |
| Flu | $\text{GP}_{so}$ | 5.3 | -2.42 | 8 |
| | our $\text{DGP}_{seq}$ | 4.0±0.6 | -3.09±0.32 | 50±2 |
| | our $\text{HDGP}_{seq}$ | 3.0±0.7 | -3.16±0.11 | 108±6 |
| Stock | $\text{GP}_{so}$ | 6.4 | -1.95 | 28 |
| | our $\text{DGP}_{seq}$ | 6.4±0.1 | -2.18±0.07 | 305±2 |
| | our $\text{HDGP}_{seq}$ | 6.6±0.3 | -2.41±0.04 | 413±4 |

Table 1: Missing Data in Multi-Task Learning (Motion/Flu/Stock).

| Data | Methods | Acc | Train_T(s) |
|---|---|---|---|
| *Tumor* | $\text{GP}_{so}$ | 0.65 | 1 |
| | $\text{DGP}_{var}$ | 0.50 | 2256 |
| | our $\text{DGP}_{seq}$ | **0.73** | 53 |
| *MNIST* | $\text{GP}_{so}$ | 0.9500 | 60mins |
| | $\text{DGP}_{var}$ | - | - |
| | DVI-GPLVM | 0.9405 | 20mins* |
| | our $\text{DGP}_{seq}$ | 0.9424 | 180mins |

Table 3: Classification. For MNIST, the results of DVI-GPLVM are from [7]. (*) Distributed implementation using an unreported number of cores.

| Methods | RMSE | MNLP | Train_T(s) |
|---|---|---|---|
| GP | 22.81±3.96 | 9.22±0.34 | 0.24±0.07 |
| NGP | 22.80±3.95 | 9.22±0.33 | 0.45±0.09 |
| WGP | 22.88±4.18 | 8.73±0.90 | 1.79±0.19 |
| MLHGP | 22.89±3.60 | 8.44±1.07 | 1.03±0.22 |
| VHGP | 22.88±3.76 | 8.45±0.41 | 1.95±0.28 |
| our $\text{HGP}_{seq}$ | 22.79±3.98 | **8.25±0.52** | 0.93±0.05 |

Table 4: Comparison with State-of-the-art Heteroscedastic GP Approaches (Motorcycle Data).

set is large-size, $\text{DGP}_{var}$ [3] is infeasible. Instead, we compare our $\text{DGP}_{seq}$ to the state-of-the-art scalable GPLVM with distributed variational Inference (DVI-GPLVM) [7]. As shown in Table 3, our $\text{DGP}_{seq}$ outperforms DVI-GPLVM in terms of accuracy with a comparable amount of training time (considering that their reported timing was parallelized while our implementation is currently serial). $\text{GP}_{so}$ outperforms both DVI-GPLVM and our $\text{DGP}_{seq}$ on this dataset, however as noted in [7], the primary purpose here is to demonstrate the scalability of the approach. These results demonstrate the ability of our approach to scale to both large and high dimensional datasets.

## 2.3 Further Comparison for Extensions

The heteroscedastic GPs and multi-task GPs have been well studied in the GP community [8, 9, 10, 11, 12]. Hence, we use the benchmark heteroscedastic *motorcycle* dataset [13] and the multi-output *Jura* dataset [12] in order to compare with state-of-the-art approaches in these two domains.

### 2.3.1 Comparison with State-of-the-art Heteroscedastic GP Approaches

We use a benchmark motorcycle dataset [13] to evaluate heteroscedasticity. This data consists of 94 time/acceleration pairs, where there are three noise regions including a flat low noise region, a curved region

and a flat high noise region [14].

We compare our $\text{HGP}_{seq}$ with a number of state-of-the-art heteroscedastic GP approaches. Namely, the standard GP with ARD kernel (GP); the standard GP with the non-stationary kernel (NGP) compounded by ARD kernel, neural network kernel and linear kernel [15]; warped GP (WGP) [16]; most likely heteroscedastic GP (MLHGP) [8]; variational heteroscedastic GP (VHGP) [9].

For each of five date partitions, 60/34 data pairs are used for training/test. For our $\text{HGP}_{seq}$, $N_p = 100$ and $N_{AC}$ is the size of training set (as we did not use sparsification). The number of training episodes is two. The average results in Table 4 show that our $\text{HGP}_{seq}$ outperforms other GP approaches with a better MNLP and a competitive NMSE & training time. Additionally, we evaluate the posterior of the motorcycle data by using our $\text{HGP}_{seq}$ with the same setting before. As shown in Fig. 10, our $\text{HGP}_{seq}$ outperforms the standard GP due to the heteroscedastic noise modelling.

### 2.3.2 Comparison with State-of-the-art Multi-Task GP Approaches

We choose a benchmark Jura data[8] to show our comparison with state-of-the-art multi-task GP approaches. In the *Jura* dataset, the input is a 2D location and the output is the measurement of cadmium,

---

[8]http://home.comcast.net/ pgoovaerts/book.html

| Data | Methods | RMSE | MNLP | Train_T(s) |
|------|---------|------|------|------------|
| | $GP_{so}$ | 9.33 ±0.15 | 7.45 ±0.02 | 37±2 |
| Parkinsons | $DGP_{var}$ | 9.26±0.41 | 8.51±1.35 | 33486±1370 |
| | our $DGP_{seq}$ | **8.86±0.23** | **7.23±0.08** | 546±18 |

Table 2: DGP as Regressor (Parkinsons Data): Prediction Error & Training Efficiency.

| Methods | MAE | Train_T(s) |
|---------|-----|------------|
| GP | 0.5739±0.0003 | 74 |
| CMOGP | 0.4552±0.0013 | 784 |
| SLFM | 0.4578±0.0025 | 792 |
| GPRN | 0.4040±0.0006 | 1040 |
| GPRN-NPV1 | 0.4147±0.0001 | 130 |
| our $DGP_{seq}$ | 0.4150±0.0061 | 21 |

Table 5: Comparison with State-of-the-art Multi-task GPs (the Jura dataset). The results of GP, CMOGP, SLFM, GPRN are from [12].

nickel and zinc concentrations. The total number of data pairs is 359, where 100 measurements of cadmium are missing. We follow the experimental settings of [12] and use mean absolute error (MAE) for evaluation, choose the number of latent layers $L = 1$, the dimensionality of latent layers $D = 2$. Both the number of particles and the size of active set in our $DGP_{seq}$ are 200. We run our $DGP_{seq}$ five times for one training episode and compare the results with the standard GP, convolved multiple outputs GP (CMOGP)[11], semiparametric latent factor model (SLFM)[10], GP regression networks (GPRN)[12], and GPRN with nonparametric variational inference (mode 1, GPRN-NPV1) [17]. As shown in Table 5, our $DGP_{seq}$ achieves a competitive MAE but with much less training time indicating that our sequential inference is efficient and able to capture the correlations between multiple outputs.

# References

[1] L. Csató, M. Opper, Sparse Online Gaussian Processes, Neural Computation 14 (3) (2002) 641–668.

[2] S. V. Vaerenbergh, M. Lázaro-Gredilla, I. Santamaría, Kernel Recursive Least-Squares Tracker for Time-Varying Regression, IEEE Transactions on Neural Networks and Learning Systems 23 (8) (2012) 1313 – 1326.

[3] A. C. Damianou, N. D. Lawrence, Deep Gaussian Processes, in: AISTATS, 2013.

[4] N. D. Lawrence, A. J. Moore, Hierarchical Gaussian Process Latent Variable Models, in: ICML, 2007.

[5] J. Hensman, N. Lawrence, Nested Varitaional Compression in Deep Gaussian Processes, in: http://arxiv.org/pdf/1412.1370.pdf, 2014.

[6] Y. Cho, L. K. Saul, Kernel Methods for Deep Learning, in: NIPS, 2009.

[7] Y. Gal, M. van der Wilk, C. E. Rasmussen, Distributed Variational Inference in Sparse Gaussian Process Regression and Latent Variable Models, in: NIPS, 2014.

[8] K. Kersting, C. Plagemann, P. Pfaff, W. Burgard, Most Likely Heteroscedastic Gaussian Process Regression, in: ICML, 2007.

[9] M. Lázaro-Gredilla, M. K. Titsias, Variational Heteroscedastic Gaussian Process Regression, in: ICML, 2011.

[10] Y. W. Teh, M. Seeger, M. I. Jordan, Semiparametric Latent Factor Models, in: AISTATS, 2005.

[11] M. A. Alvarez, N. Lawrence, Computationally efficient convolved multiple output gaussian processes, JMLR 12 (2011) 1459–1500.

[12] A. G. Wilson, D. A. Knowles, Z. Ghahramani, Gaussian Process Regression Networks, in: ICML, 2012.

[13] B. W. Silverman, Some Aspects of the Spline Smoothing Approach to Non-Parametric Regression Curve Fitting, Journal of the Royal Statistical Society. Series B (Methodological) 47 (1) (1985) 1–52.

[14] C. E. Rasmussen, Z. Ghahramani, Infinite Mixture of Gaussian Process Experts, in: NIPS, 2002.

[15] C. E. Rasmussen, C. K. I. Williams, Gaussian Process for Machine learning, MIT Press, 2006.

[16] E. Snelson, C. E. Rasmussen, Z. Ghahramani, Warped Gaussian Processes, in: NIPS, 2004.

[17] T. V. Nguyen, E. V. Bonilla, Efficient Variational Inference for Gaussian Process Regression Networks, in: AISTATS, 2013.
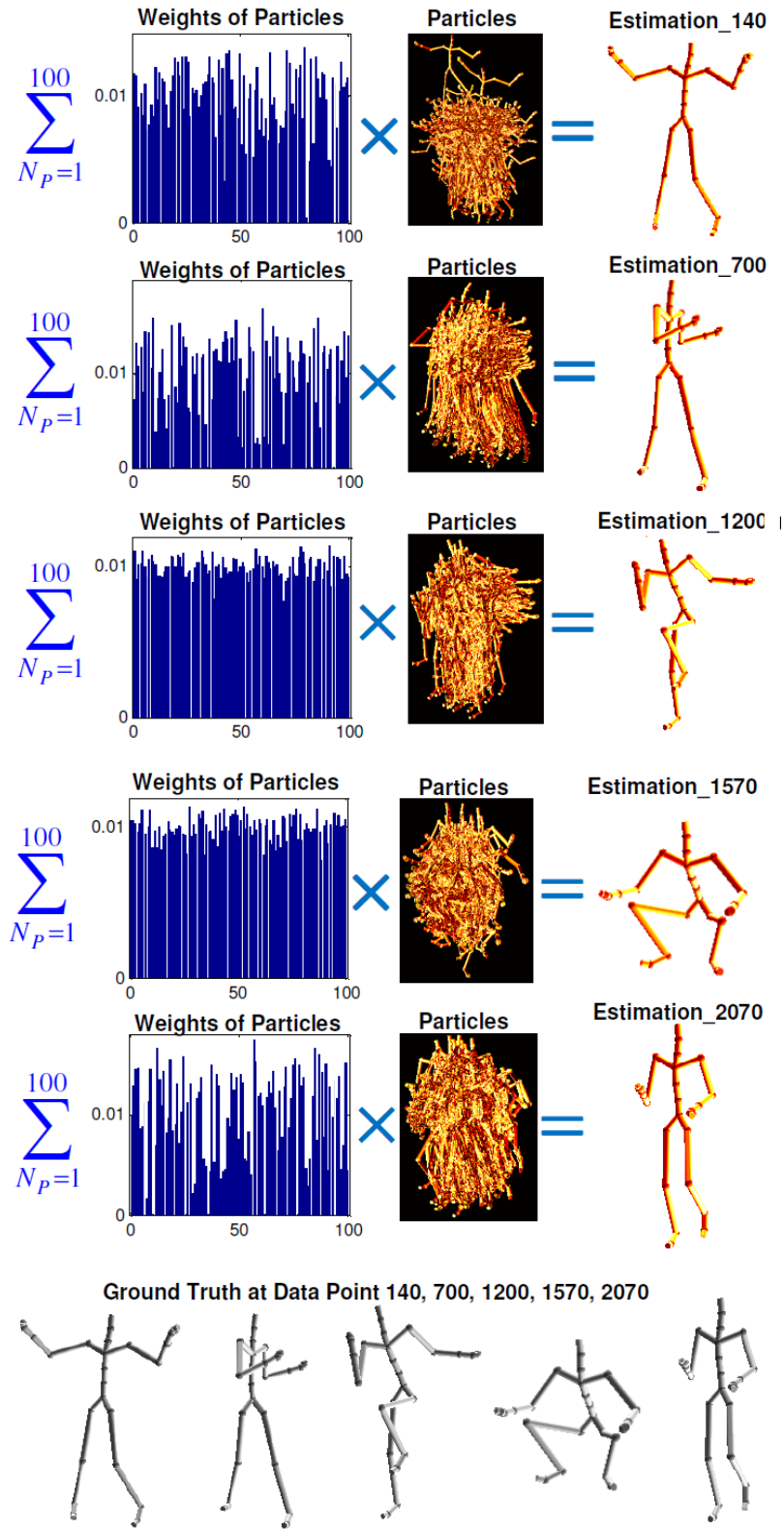
Figure 5: DGP as Deep Dynamical Prior (the motion dataset). Particle Estimation of **y** at the time steps 140, 700, 1200, 1570, 2070 where the number of particles $N_p$ is 100.
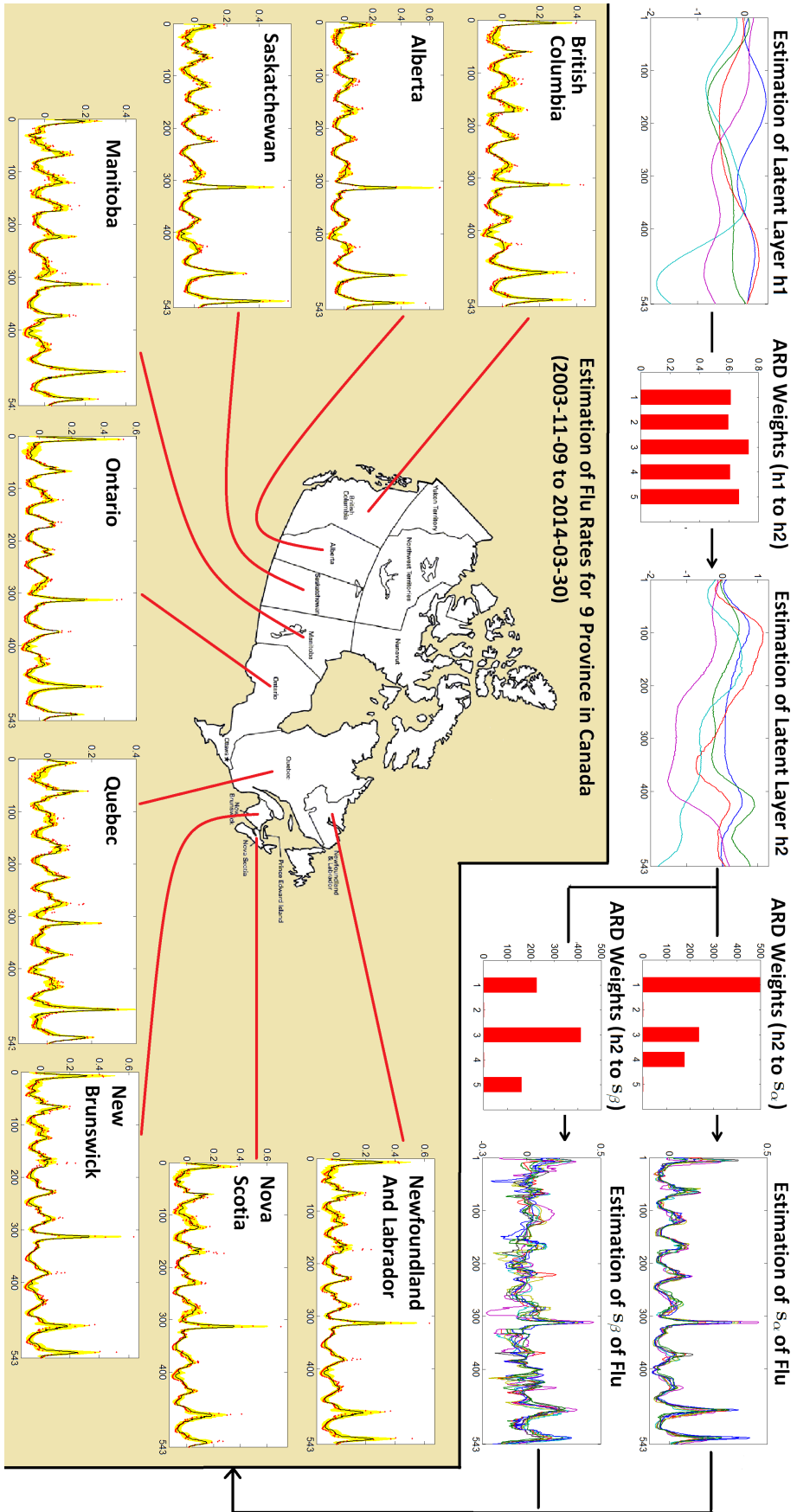
Figure 6: DGP as Deep Dynamical Prior (the flu dataset). Our HDGP $_{seq}$ has 2 latent layers, the dimensionality of each layer is 5. In the plots for the estimation of flu rates (9 provinces in Canada), the red dots are the observations. The black lines are the estimated means. The yellow regions are the 95% confidence interval.
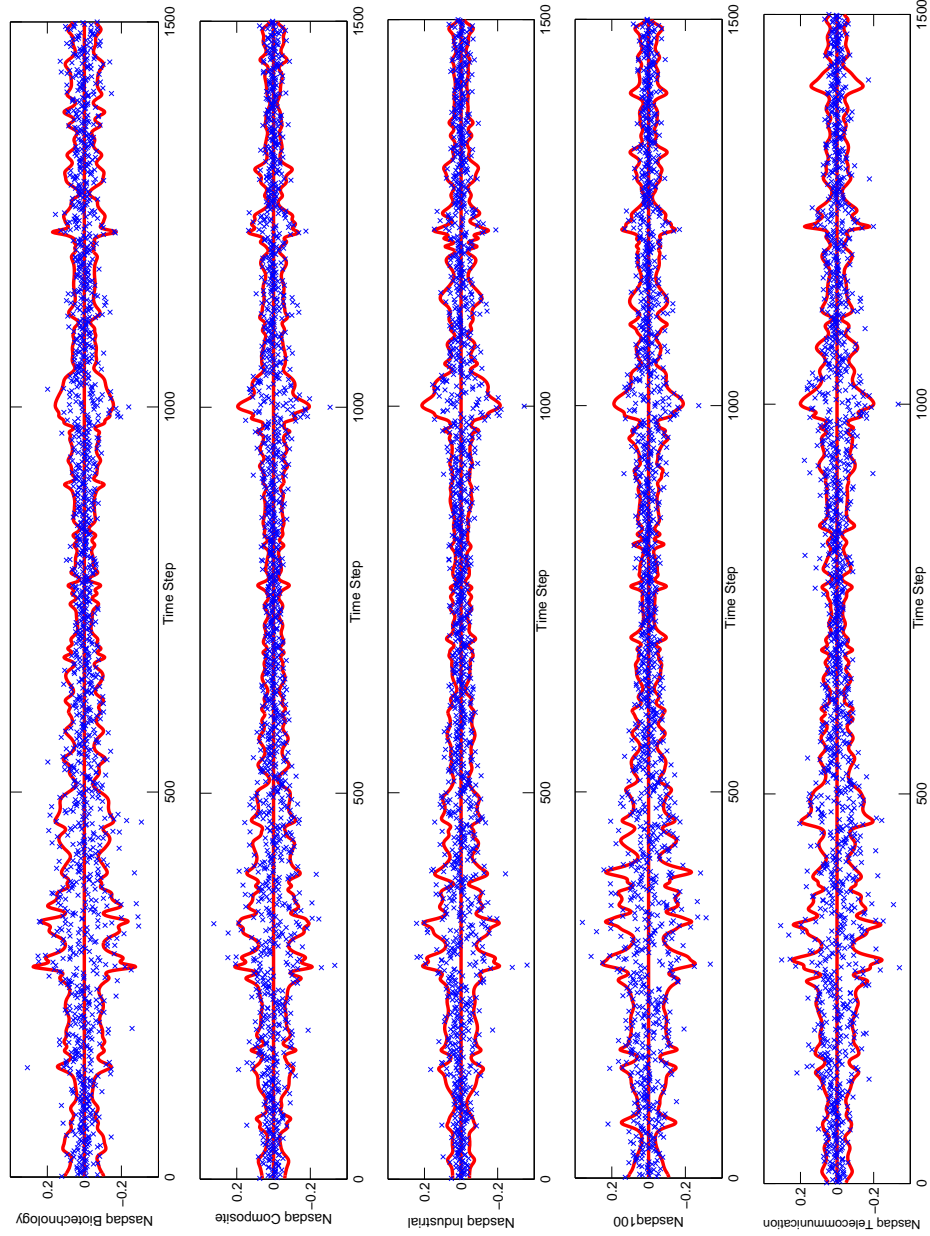
Figure 7: DGP as Deep Dynamical Prior (the stock dataset). The observations are blue crosses, the posterior mean with 95% confidence interval of our HDGP$_{seq}$ are red lines.
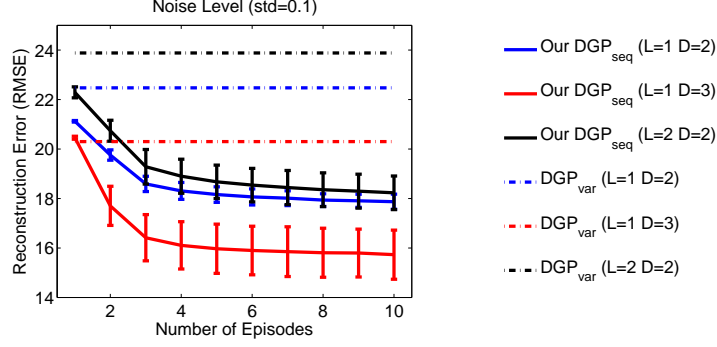
Figure 8: Dimensionality Reduction for Image Reconstruction (Face Data): the reconstruction error as a function of the number of episodes. Note that there is no sequential inference & hyperparameter learning episode in $DGP_{var}$. The results of $DGP_{var}$ are lines. The error bar is mean$\pm$standard deviation.
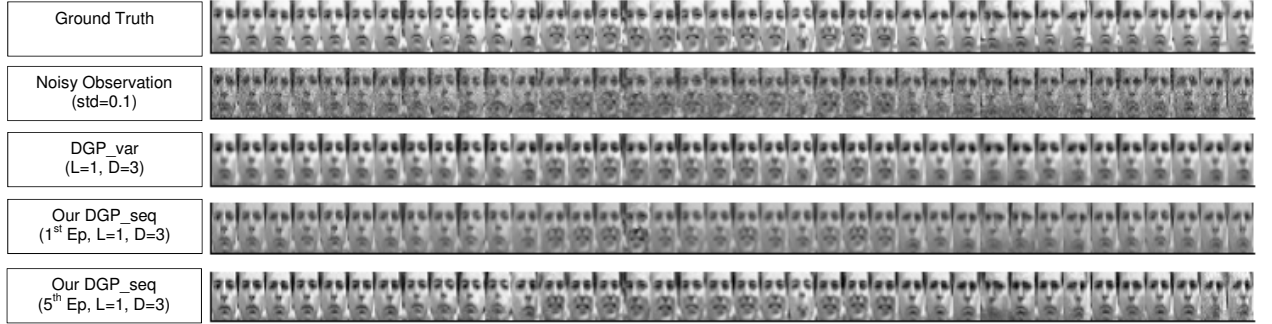


Figure 9: Dimensionality Reduction for Image Reconstruction (Face Data): the reconstruction image comparison at the time step 50:50:1900.
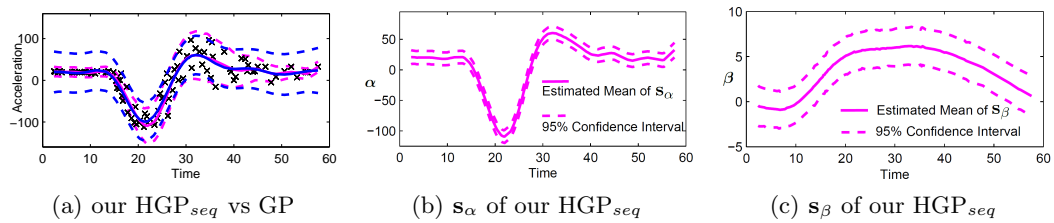


(a) our $HGP_{seq}$ vs GP          (b) $\mathbf{s}_\alpha$ of our $HGP_{seq}$          (c) $\mathbf{s}_\beta$ of our $HGP_{seq}$

Figure 10: Heteroscedastic GP Modeling (Posterior of Motorcycle Data). In Plot10(a), the observations are black crosses, the mean of our $HGP_{seq}$ / GP are pink / blue lines, 95% confidence interval of our $HGP_{seq}$ / GP are pink / blue dashed lines.