









- (#) Perform the following procedure:
  - (Prediction step) After drawing  $\{\mathbf{v}_t^{(i)}\}_{i=1}^N$ , generate  $N$  prediction samples of state variables  $\{\mathbf{x}_{t|t-1}^{(i)}, \boldsymbol{\theta}_{t-1}^{(i)}\}_{i=1}^N$  according to  $\mathbf{x}_{t|t-1}^{(i)} = \mathbf{f}(\mathbf{x}_{t-1|t-1}^{(i)}, \boldsymbol{\theta}_{t-1}^{(i)}) + \mathbf{v}_t^{(i)}$  for  $i = 1, \dots, N$ .
  - (★) (Filtering step) Execute the following procedure:
    - (i) Compute likelihood function of each particle  $p(\mathbf{y}_t | \mathbf{x}_{t|t-1}^{(i)})$ ,  $i = 1, \dots, N$ .
    - (ii) Calculate the normalizing weights  $q_t^{(i)} \propto p(\mathbf{y}_t | \mathbf{x}_{t|t-1}^{(i)})$  for  $i = 1, \dots, N$  such that  $\sum_{i=1}^N q_t^{(i)} = 1$ .
    - (iii) Generate  $\{\mathbf{x}_{i|t}, \boldsymbol{\theta}_t^{(i)}\}_{i=1}^N$  by the resampling of  $\{\mathbf{x}_{i|t-1}, \boldsymbol{\theta}_{t-1}^{(i)}\}_{i=1}^N$  with the probabilities  $\{q_t^{(i)}\}_{i=1}^N$ .
- If  $t$  is equal to  $T$ , stop the algorithm. Otherwise, set  $t \leftarrow t + 1$  and go to (#).

Here, the likelihood function  $p(\mathbf{y}_t | \mathbf{x}_t)$  is defined by the probability density function, e.g. normal distribution, corresponding to the observation model (3). The final draws of  $\{\boldsymbol{\theta}_T^{(i)}\}_{i=1}^N$  generated from this procedure are approximately distributed according to the augmented posterior distribution  $p(\mathbf{x}_{0:T}, \boldsymbol{\theta} | \mathbf{y}_{1:T})$ . Indeed, it can be proved that the empirical distributions given by  $N$  realization  $\{\boldsymbol{\theta}_T^{(i)}\}_{i=1}^N$  approach to the true distribution  $p(\boldsymbol{\theta} | \mathbf{y}_{1:T})$  with probability one as  $N$  goes to infinity.<sup>13</sup>

As the number of unknown parameters becomes large, the method requires to generate an exceedingly large number of Monte Carlo samples in order to cover the overall parameter space. For example, the transcription regulation model of circadian clock that will be shown in Section 3 contains 45 unknown parameters in total. Particularly, the problem of “*particle degeneracy*” is caused by a small number of Monte Carlo samples (Chapter 14 of Ref. 13). This is due to the fact that during the repeated resampling of initially generated  $\{\boldsymbol{\theta}_0^{(i)}\}_{i=1}^N$  drawn from  $p(\boldsymbol{\theta})$ , the number of distinct values in the resampled particles can diminish drastically as the time step increases. Hence, in order to estimate a large number of parameters, it is essential to use a fairly large amount of Monte Carlo samples.

Both the time complexity and space complexity of the prediction step are linear order of the number of particles. On the other hand, the time complexity of the filtering step is  $O(N \log N)$  because the resampling algorithm requires  $O(\log N)$  evaluations of weight table to yield one filtered particle (see Appendix A for details). We employ this type of resampling procedure in the numerical experiment. Rather than performing the ran-

dom resampling algorithms, the filtering step can be replaced by the resampling algorithm with the proportional-allocation that yields the filtered samples by increasing and decreasing the number of prediction samples proportional to their weights. In that case, the time complexity can be reduced to linear order.

Here, we briefly discuss the low coupling property of the particle filtering. In the prediction step, each predictive sample is updated independently. In the likelihood calculation (i) at the filtering step, each likelihood can be also computed independently. As a result, the tasks in each step can be divided into parallel processors. On the other hand, the processes of normalization (ii) and resampling (iii) require additional parallel computing technique in order to manage the communication among  $N$  Monte Carlo samples. In the step (ii), it is necessary that the likelihood values of  $N$  samples are collected into a root CPU. To implement a parallel computing for such interacted tasks, we can use the divide and conquer algorithm.<sup>17</sup> By the aid of it, time complexity of normalization and resampling can be reduced up to  $O((N/M) \log M)$  where  $M$  CPU processors are available, which indicates scalability of the particle filtering in parallel computing.

### 3. Numerical Experiments

#### 3.1. Implementation

Circadian clocks of living cells are controlled by oscillating feedback loop in the transcription circuits of clock-control genes.<sup>2,19</sup> The model that we consider was constructed by incorporating the regulatory mechanisms among the five mRNAs (per, cry, rev/erb, clock, bmal), the translated proteins (PER, CRY, REV/ERB, CLOCK, BMAL) and the two protein complexes (PER/CRY, CLOCK/BMAL), which were reported by Ref. 18. It consists of the 12 first-order stochastic difference equations as,

$$\begin{aligned}
 x_{1,t+1} &= x_{1,t} - k_{d_1}x_{1,t} + k_{tc_1}I(x_{5,t} < s_{1,5})I(x_{12,t} > s_{1,12}) + k_{c_1} + v_{1,t+1}, \\
 x_{2,t+1} &= x_{2,t} - k_{d_2}x_{2,t} + k_{tl_2}x_{1,t} - k_{b_{2,4}}x_{2,t}x_{4,t} + v_{2,t+1}, \\
 x_{3,t+1} &= x_{3,t} - k_{d_3}x_{3,t} + k_{tc_3}I(x_{5,t} < s_{3,5})I(x_{12,t} > s_{3,12}) + k_{c_3} + v_{3,t+1}, \\
 x_{4,t+1} &= x_{4,t} - k_{d_4}x_{4,t} + k_{tl_4}x_{3,t} - k_{b_{2,4}}x_{2,t}x_{4,t} + v_{4,t+1}, \\
 x_{5,t+1} &= x_{5,t} - k_{d_5}x_{5,t} + k_{b_{2,4}}x_{2,t}x_{4,t} + v_{5,t+1}, \\
 x_{6,t+1} &= x_{6,t} - k_{d_6}x_{6,t} + k_{tc_6}I(x_{5,t} < s_{6,5}) + k_{c_6} + v_{6,t+1}, \\
 x_{7,t+1} &= x_{7,t} - k_{d_7}x_{7,t} + k_{tl_7}x_{6,t} + v_{7,t+1}, \\
 x_{8,t+1} &= x_{8,t} - k_{d_8}x_{8,t} + k_{c_8} + v_{8,t+1},
 \end{aligned}$$

$$\begin{aligned}
x_{9,t+1} &= x_{9,t} - k_{d_9}x_{9,t} + k_{tl_9}x_{8,t} - k_{b_{9,11}}x_{9,t}x_{11,t} + v_{9,t+1}, \\
x_{10,t+1} &= x_{10,t} - k_{d_{10}}x_{10,t+1} + k_{tc_{10}}I(x_{5,t} > s_{10,5})I(x_{7,t} < s_{10,7}) \\
&\quad + k_{c_{10}} + v_{10,t+1}, \\
x_{11,t+1} &= x_{11,t} - k_{d_{11}}x_{11,t} + k_{tl_{11}}x_{10,t} - k_{b_{9,11}}x_{9,t}x_{11,t} + v_{11,t+1}, \\
x_{12,t+1} &= x_{12,t} - k_{d_{12}}x_{12,t} + k_{b_{9,11}}x_{9,t}x_{11,t} + v_{12,t+1}, \tag{5}
\end{aligned}$$

in order corresponding to per, PER, cry, CRY, PER/CRY, rev/erb, REV/ERB, clock, CLOCK, bmal, BMAL, and CLOCK/BMAL, respectively. Here,  $v_{i,t}$ ,  $i = 1, \dots, 12$ , stand for the Gaussian noises with mean 0 and variance  $2.0 \times 10^{-7}$ . Of interest is then to estimate the 26 rate constant parameters  $k$ ., the 7 threshold parameters  $s$ ., and the initial concentrations of 12 reactants  $x_{i,0}$ ,  $i = 1, \dots, 12$ . Note that we did not estimate values of  $k_{d_s}$  and  $k_{c_s}$  for the reason shown below.

Ueda *et al.*<sup>19</sup> studied the transcription regulations of mammalian circadian clocks based on the time course gene expression data which were generated with Affymetrix mouse high-density oligonucleotide probe array. Time course data were measured at equally-spaced 12 time points during 44 hours (about two days). We extracted gene expression measurements of the five mRNAs (per, cry, rev/erb, clock, bmal) from Figure 1a and the supplementary data ([http://sirius.cdb.riken.jp/MouseLiver/MouseLiverCCG\(020707\).html](http://sirius.cdb.riken.jp/MouseLiver/MouseLiverCCG(020707).html)<sup>19</sup>), where the set of 12 time points were allocated to the 700 simulation time steps by the conduction of time alignment. Note that in the current model (5), there are no regulatory mechanisms which yield oscillating concentrations of clock mRNA. As a possible solution, we decided to treat the measurements of clock mRNA concentrations as an input sequence. To this end, we interpolated the observed 12 data points of clock by using spline functions, and then, the interpolated values were substituted into the corresponding state variables which remained to be fixed during the parameter estimation process.

Figure 1 shows the observed time course of gene expression profiles and the result of simulation which was generated according to (5) with the pre-determined parameters without inclusion of system noises  $v_{i,t}$ s. This figure shows inconsistency in periods, phases and amplitude among the simulation paths and observed profiles. The objective of this numerical experiment is to capture the observed oscillating patterns of circadian rhythm by the application of the particle filtering where the currently obtained parameters have the bias in the cyclic behavior.

Before proceeding to the particle filtering, the observation noise  $w_t$  was set as the normal distribution  $N(0, 0.64)$  for per mRNA and  $N(0, 0.25)$

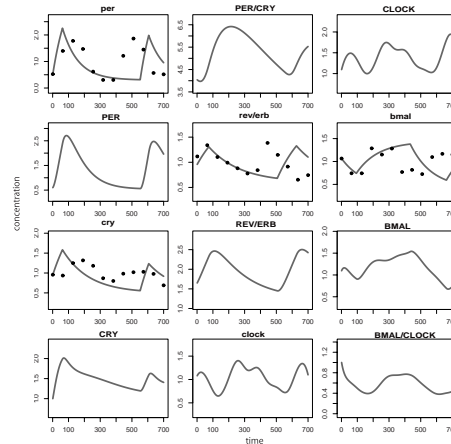


Figure 1. Result of initial simulation which was generated according to (5) with the pre-determined parameters without inclusion of system noises  $v_{i,t}$ s. The observed 12 data points are indicated by dots.

for the other reactants. We adopted the values that are employed in the simulation shown in Figure 1 as the prior mean of the parameter  $\theta$  and the initial state  $x_0$ . Because the parameters and the states should be non-negative, the prior distributions were set to truncated normal distributions. The prior variances of initial concentration  $x_{\cdot,0}$ , coefficients  $k$  and threshold values  $s$  were set to 0.01, 0.001 and 0.05, respectively.

In order to realize the implementation of the particle filtering with  $10^8$  samples, all particles were divided into 20 sets of  $5 \times 10^6$  particles, and allocated into 20 computational units. After repeating the particle filtering with  $5 \times 10^6$  samples for each computational unit, which yields  $\{\theta_T^{(i,k)}\}_{i=1}^{10^6}$  for  $k = 1, \dots, 20$ , we can compute the estimator by  $\hat{\theta} = \frac{1}{10^8} \sum_{k=1}^{20} \sum_{i=1}^{5 \times 10^6} \theta_T^{(i,k)}$ . Due to the independence of the computational units, the integrated estimator  $\hat{\theta}$  still preserves the statistical consistency.

### 3.2. Results

Figure 2 shows the histograms of Monte Carlo samples drawn from the posterior distribution of the parameters  $k_{d_i}$  ( $i = 1, \dots, 12$ ), where the sizes of Monte Carlo samples were set to  $10^8$  and  $10^5$ , respectively. It is observed from the right panel of Figure 2 that the number of distinct values declines drastically for the estimation with  $10^5$  particles. Obviously, we would not obtain any efficient estimates by using such degenerated particles. On the



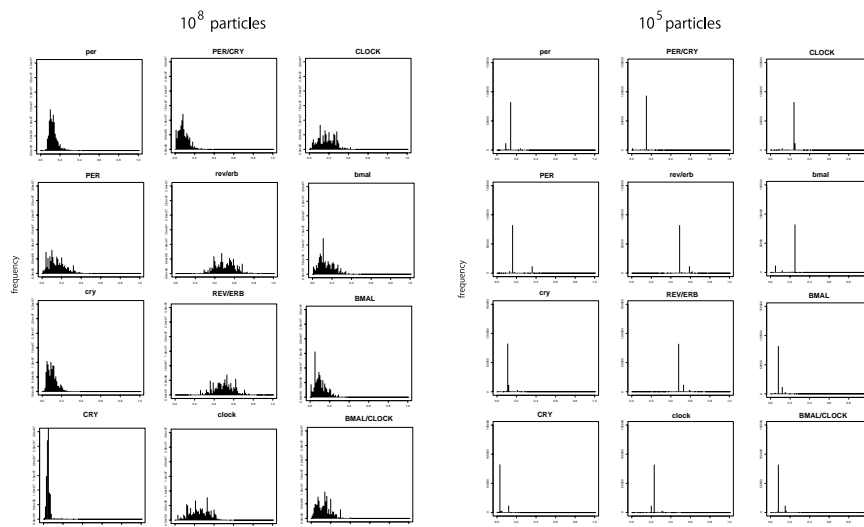


Figure 2. Approximated posterior distributions of the 12 “degradation rate”  $k_d$ , which were generated by using  $10^8$  (left) and  $10^5$  (right) particles.

other hand, the ability of the particle filtering was improved greatly by the use of  $10^8$  particles as the histograms of the Monte Carlo samples reveal the shapes of posterior distributions clearly as shown in the left panel of Figure 2.

Figure 3 shows the results of simulations based on the estimated parameters and initial state variables, which were computed as the averages of  $10^8$  and  $10^5$  particles, respectively. In the case of  $10^8$  particles, simulation paths were generated from the estimated parameters in each of the 20 computational units, and then, the generated 20 paths were averaged in order to compute the final estimates shown in Figure 3. The observed data points are also shown in the Figure 3. The simulated concentrations corresponding to  $10^5$  particles preceded the observed periodicity of circadian rhythm whereas the difference in periodicity has been reduced greatly by using  $10^8$  particles compared with the initially-chosen parameters (prior means).

### 3.3. Estimation of computational time

We evaluate computational time for execution of large-scale particle filtering on upcoming peta-scale supercomputers.

We measured and estimated CPU times in the following three cases:

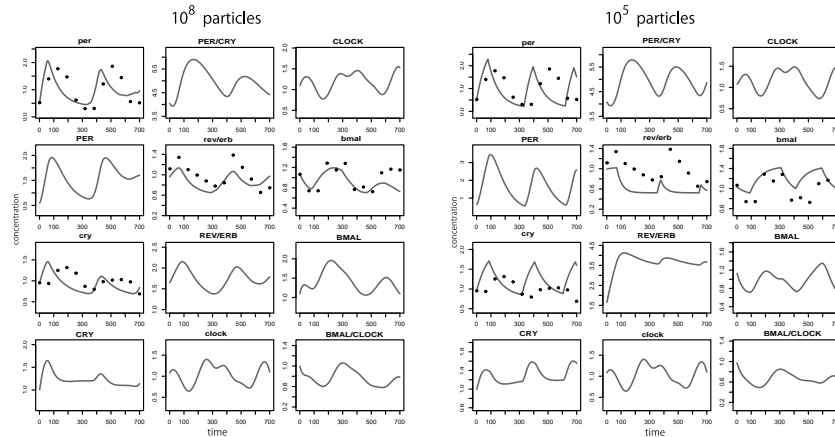


Figure 3. State transition estimated by using  $10^8$  (left) and  $10^5$  (right) particles.

- $10^8$  particles are divided into 20 sets of 5 million particles, and then, repeat the particle filtering 20 times on Opteron 2200 (about 5 gigaflops).
- Particle filtering with  $10^8$  (100 million) samples is performed on one petaflops machine.
- Particle filtering with  $10^9$  (one billion) samples is performed on one petaflops machine.

For the evaluations of (b) and (c), the CPU times were estimated based on the observed CPU time of the experiment (a). The results of the tests are summarized in Table 1. The current implementation that divides  $10^8$  particles into 20 computational subunits requires about eight days on average, which can be reduced to 90 seconds by the aid of one petaflops computer. Even with one billion particles, the estimated CPU time can be within 20 minutes.

Table 1. Observed and estimated CPU times (second). The left column shows the measured CPU time on Opteron 2200 and 20 sets of 5 million particles. The middle and right columns show the estimated CPU times for 100 million and one billion particles.

	Opteron, 5mil. $\times$ 20	1 petaflops, 100 mil.	1 petaflops, 1 billion
Ave.	$6.8 \times 10^5$	90	$1.0 \times 10^3$
Min.	$3.7 \times 10^5$	50	$5.6 \times 10^2$
Max.	$2.0 \times 10^5$	$2.7 \times 10^2$	$3.0 \times 10^3$

#### 4. Conclusion

This paper has shown the promising approach based on large scale particle filtering towards exhaustive explorations for kinetic parameters of *in silico* biological pathways. The approach could offer an overall view of posterior distributions defined over the 45 unknown parameters within a reasonable computing time. To our knowledge, the use of one hundred million particles is the first attempt not only in systems biology but also in the other scientific fields. Because the idea is too obvious, such a direct and straightforward approach may have been overlooked so far.

A limitation of our approach is the sensitivity of the estimation results to the choice of the prior distributions, that is, if the prior distributions are fairly inappropriate, we would not obtain reasonable posterior estimations. In order to overcome such a difficulty, we have to develop a methodology for incorporating the existing knowledge on biochemistry and systems biology.

According to the evaluation of computational time, it has been found that we can execute the particle filtering even with one billion particles within 20 minutes by the aid of one petaflops computer that will be accessible in Japan. The results of the experiments indicate that the large scale particle filtering can be a promising approach to parameter estimation of *in silico* biological pathways in the near future.

#### Appendix A. Resampling Algorithm

The filtering step ( $\star$ ) of the algorithm that has  $O(N \log N)$  time complexity is achieved as follows:

- Compute likelihood function of each particle  $p(\mathbf{y}_t | \mathbf{x}_{t|t-1}^{(i)})$ ,  $i = 1, \dots, N$ . Time complexity of this step is given by  $O(N)$ .
- Calculate the normalizing weights  $q_t^{(i)} = p(\mathbf{y}_t | \mathbf{x}_{t|t-1}^{(i)})$  for  $i = 1, \dots, N$  such that  $\sum_{i=1}^N q_t^{(i)} = 1$ . Time complexity of this step is given by  $O(N)$ .
- Compute cumulative weight table  $\hat{q}_t^{(i)}$  for  $i = 0, \dots, N$  by  $\sum_{j=1}^i \hat{q}_t^{(j)}$  ( $\hat{q}_t^{(0)} = 0$ ). Time complexity is given by  $O(N)$ .
- For each  $j = 1, \dots, N$ 
  - Generate random number  $u$  from uniform distribution  $U(0, 1]$ . Time complexity is  $O(N)$ .
  - Search the number  $k$  where  $\hat{q}_t^{(k-1)} < u \leq \hat{q}_t^{(k)}$  by binary search. Time complexity is  $O(N \log N)$ .
  - Set  $(\mathbf{x}_{t|t}^{(j)}, \boldsymbol{\theta}_t^{(j)}) \leftarrow (\mathbf{x}_{t|t-1}^{(k)}, \boldsymbol{\theta}_{t-1}^{(k)})$ , which requires  $O(N)$ .

Only the binary searching requires  $O(N \log N)$ . In order to deal with much larger number of particles, we can replace the binary search algorithm by the resampling with the proportional-allocation that yields the filtered samples by increasing and decreasing the number of prediction samples proportional to their weights.

## References

1. K. C. Chen, L. Calzone, A. Csikasz-Nagy, F. R. Cross, B. Novak and J. J. Tyson, *Mol Biol Cell.*, **15**, 3841-3862 (2004).
2. H. Matsuno, R. Murakami, R. Yamane, N. Yamasaki, S. Fujita, H. Yoshimori and S. Miyano, *Pac Symp Biocomput.*, **8**, 152-163 (2003).
3. A. Doi, S. Fujita, H. Matsuno, M. Nagasaki and S. Miyano, *In Silico Biol.*, **4**, 271-291 (2004).
4. A. Doi, M. Nagasaki, H. Matsuno and S. Miyano, *In Silico Biol.*, **6**, 1-13 (2006).
5. K. Fajarewicz, M. Kimmel, T. Lipniacki and A. Swierniak, *IEEE/ACM Transaction on Computational Biology and Bioinformatics*, **4**, 322-335 (2007).
6. M. Quach, N. Brunel and F. d'Alché-Buc, *Bioinformatics*, **23**, 3209-3216 (2007).
7. I. Nachman, A. Regev and N. Friedman, *Bioinformatics*, **20**, i248-i256 (2004).
8. G. Koh, H. F. C. Teong, M-V ment, D. Hsu and P.S. Thiagarajan, *Bioinformatics*, **22**, e271-e280 (2006).
9. B.-E. Perrin, L. Ralaivola, A. Mazurie, S. Bottani, J. Mallet, F. d'Alché-Buc, *Bioinformatics*, **19**, ii386-ii349 (2003).
10. G. Kitagawa, *J. Computational and Graphical Stat.* **5**, 1-25 (1996).
11. N. J. Gordon, D. J. Salmond and A. F. M. Smith *IEE Proceedings-F*, **140**, 107-113 (1993).
12. T. Higuchi and G. Kitagawa, *IEICE Tras. Inf. Syst.*, **E83-D**, 36-43 (2000).
13. A. Doucet and N. De Freitas, Springer-Verlag (2001).
14. M. Nagasaki, R. Yamaguchi, R. Yoshida, S. Imoto, A. Doi, Y. Tamada, H. Matsuno, S. Miyano and T. Higuchi, *Genome Inform.* **17(1)**, 46-61 (2006).
15. S. Tasaki, M. Nagasaki, M. Oyama, H. Hata, K. Ueno, R. Yoshida, T. Higuchi, S. Sugano, S. Miyano, *Genome Inform.* **17(2)**, 226-238 (2006).
16. G. Kitagawa, *J. Amer. Stat. Assoc.*, **93**, 1203-1215 (1998).
17. Peter Pacheco, *Parallel Programming with MPI*, Morgan Kaufmann (1996).
18. Y. Fujii, Y. Okitsu, H. Matsuno, S. Miyano and S. T. Inoue, A new regulatory interactions suggested by simulations for circadian genetic control mechanism in mammals, *Genome Inform.*, available at <http://www.jsbi.org/journal/GIWO>, (2004).
19. H. R. Ueda, W. Chen, A. Adachi, H. Wakamatsu, S. Hayashi, T. Takasugi, M. Nagano, K. Nakahama, Y. Suzuki, S. Sugano, M. Iino, Y. Shige-yoshi and S. Hashimoto, *Nature*, **418**, 534-539 (2002).