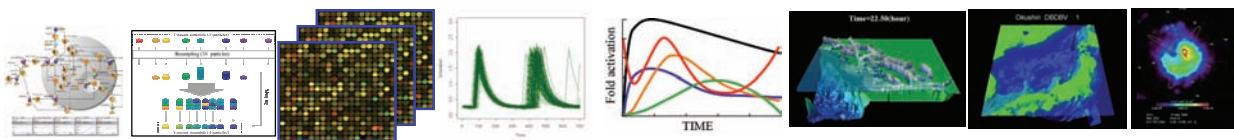


2010/01/25 10:10—11:10 The 3rd Workshop on Integration of EFD and CFD

Introduction to Sequential Data assimilation methods: Their mathematical basis and recent development

Tomoyuki Higuchi

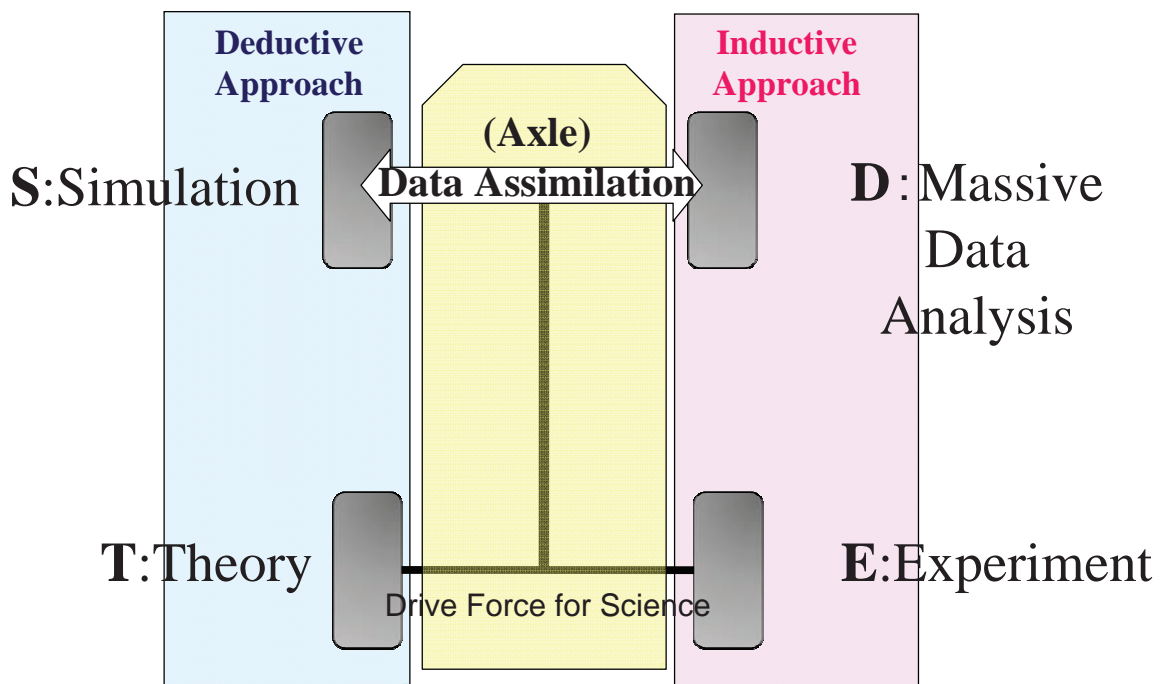
Research Organization of Information and Systems
The Institute of Statistical Mathematics/JST CREST



1/42



TESD: Four Kinds of Methodology of Science



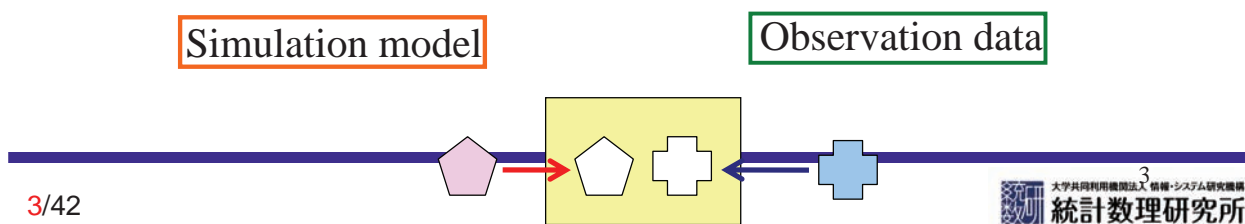
2/42

Red color indicates a slide used in the last year's presentation
(2nd. Workshop on Integration of EFD and CFD)



What is Data Assimilation?

- Emerging subject in meteorology and oceanography.
 - Methodology to synthesize numerical simulation model and observed data
 - **Simulation model** can not reflect real physics accurately.
 - (e.g.) Accurate weather forecast needs good initial conditions.
 - Uncertainty in the model (boundary condition, initial condition, unknown parameters, unknown dynamics...) exists.
 - **Observation data** have some physical/budgetary restrictions.
- ➔ Correct variables in numerical simulation model using observation data. = Data Assimilation

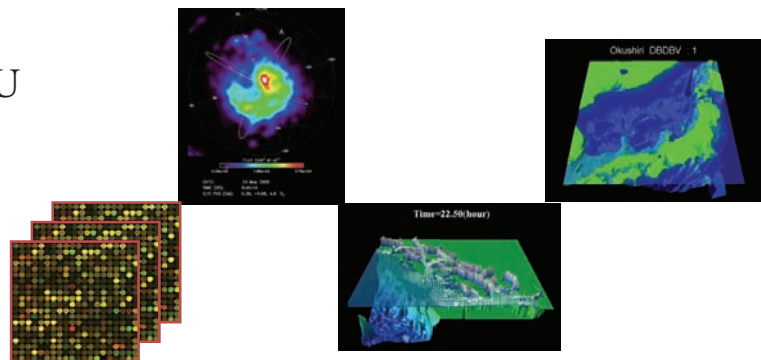


Objects of Data Assimilation from a viewpoint of Meteorology and Oceanography

- [1] To produce the best (better) **initial condition** for forecasting. It is actually realized in the real weather forecast (ex., Japan Meteorological Agency).
- [2] To find the best (better) **boundary condition** in constructing a simulation model. This procedure includes a setting of appropriate boundary conditions necessary for dealing with a coupled phenomena.
- [3] To attain an optimal **parameter** vector that appears in an empirical law (scheme) employed for describing complicated phenomena with the different time and spatial scales. A **validation** of the empirically given values is regarded as this problem.
- [4] To inter/extrapolate (estimate) a physical quantity at times and locations without observations based on a numerical simulation model. This procedure is called “a **generation of re-analysis dataset** (product)”. This dataset is used to discover a new scientific findings by general geophysical researchers.
- [5] To conduct an experiment with a virtual observation network and perform a **sensitivity analysis** in an attempt to construct an effective observation network system with less budgetary cost and less consuming time.

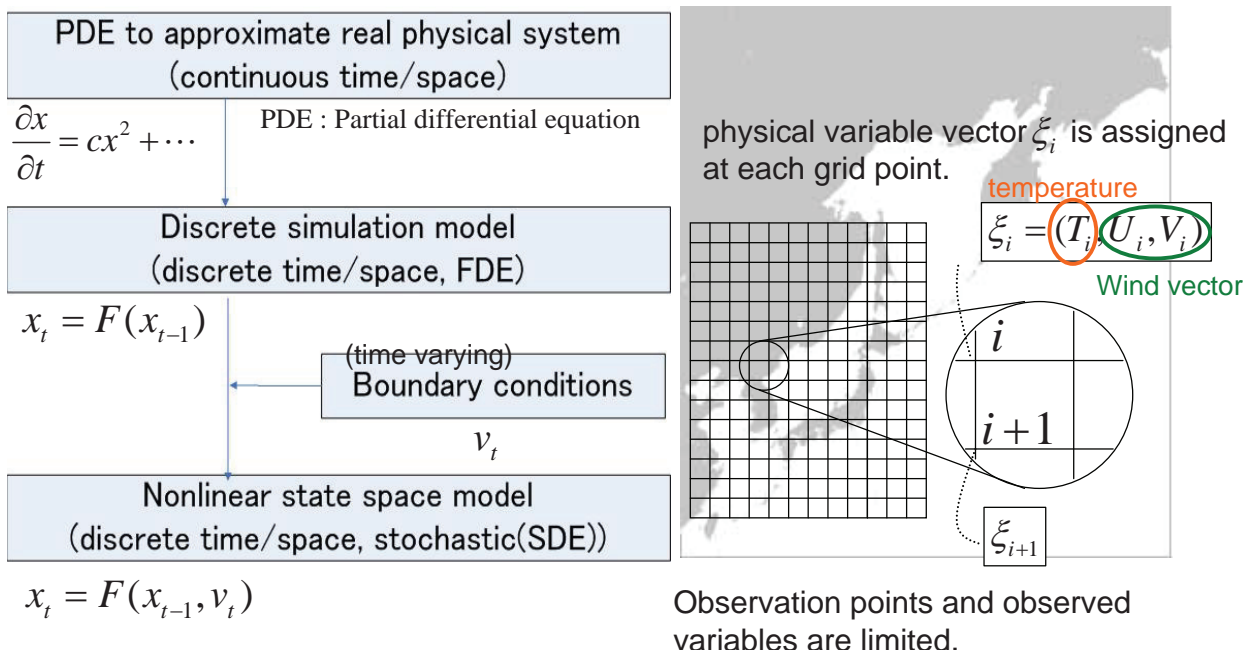
Outline

- Mathematical basis and Bayesian computation
- Sequential data assimilation
- Ensemble-based nonlinear filtering method
 - Particle filter (PF)
- Advanced methods for PF
 - Merging PF
 - Meta PF
 - PF with GPGPU
- Conclusions

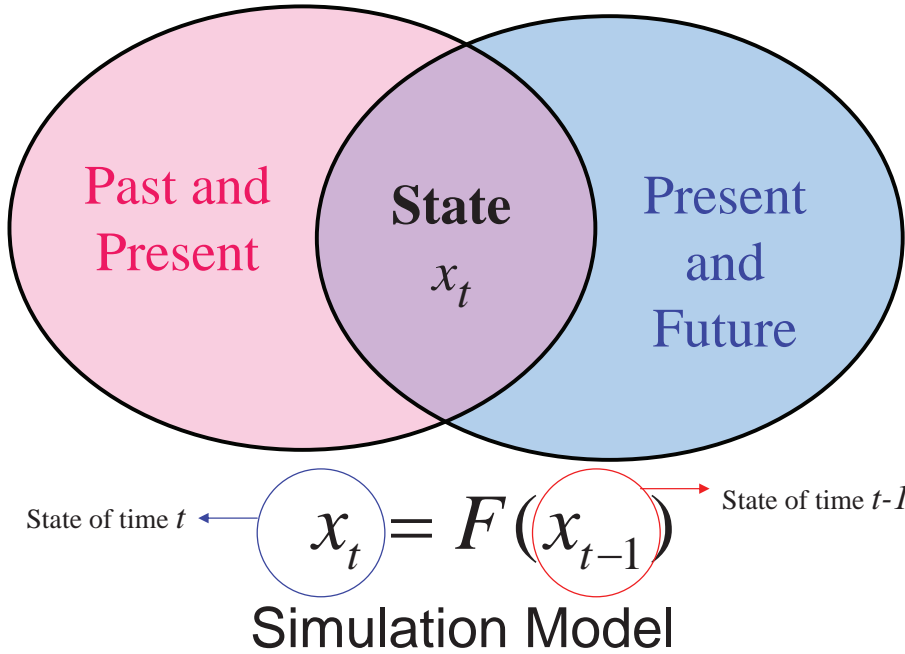


Construction of Simulation Model

(simplified meteorological model around Japan)



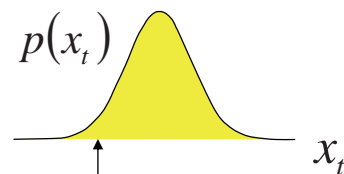
State Vector : Contact point between past and future



From one path to PDF (=Probability Distribution Function)

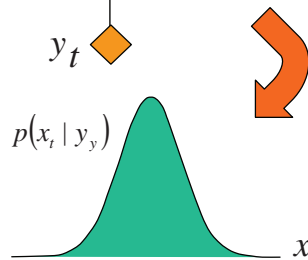


Simulation model with uncertainty

$$x_t \approx F_t(x_{t-1})$$


Relation to observed data

$$y_t \approx h_t(x_t)$$



Next slide

Conditional distribution $p(\mathbf{x}_t | \mathbf{y}_t)$

DA = Estimate $p(x_t | y_t)$

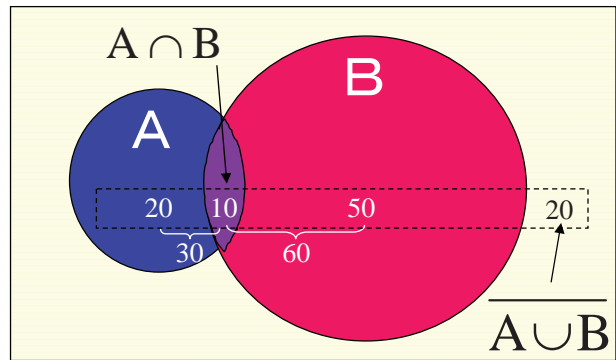
Conditional and Joint Probabilities

$p(A) \equiv$ Probability of **A**
 $p(A | B) \equiv$ Probability of **B** given **A** ← Conditional Probability
 $p(A, B) \equiv$ Probability of **A** and **B** ← Joint Probability

$p(A=1) = \frac{30}{100}$, $p(A=1, B=1) = \frac{10}{100}$, $p(B=1|A=1) = \frac{10}{30} = \frac{10}{20+10}$

Total: 100
 # of consumers to buy a bag of coffee grounds : 30
 # of consumers to buy milk: 60
 # of consumers to buy a coffee bag and milk: 10

A=1: Buy a bag of coffee grounds
 =0: Not buy
 B=1: Buy milk
 =0: Not buy



9/42

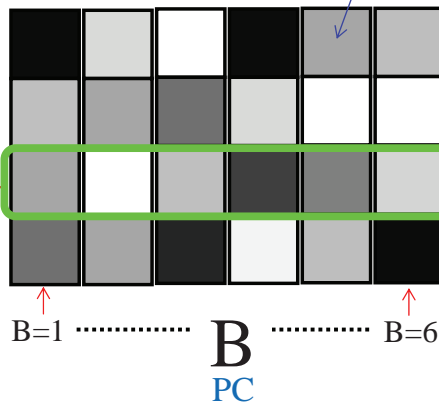
Marginalization

$$p(A = 3) = \sum_{j=1}^6 p(A = 3, B = j)$$

Joint probability $p(A = i, B = j)$

A=1: Yahoo
 =2: Google
 =3: Microsoft
 =4: Others

A A=3 →
 Web Search Engine



$$\sum_{i=1}^4 \sum_{j=1}^6 p(A = i, B = j) = 1$$

B=1: NEC, =2: Fujitsu, =3: Dell
 =4: Toshiba, =5: Apple, =6: Others

$$p(A = A_i) = \sum_{j \in \text{possible } B_j} p(A_i, B_j)$$



$$p(A) = \int p(A, B) dB$$

10/42

Bayes' Theorem

$$p(B | A) = \frac{p(A, B)}{p(A)} \Rightarrow p(A, B) = p(B | A)p(A)$$

It is **easy** to calculate.

$$p(\underline{A} | \underline{B}) = \frac{p(A, B)}{p(B)} = \frac{p(A, B)}{\sum_{A \in \mathcal{A}} p(A, B)} = \frac{p(\underline{B} | \underline{A})p(A)}{\sum p(\underline{B} | \underline{A})p(A)}$$

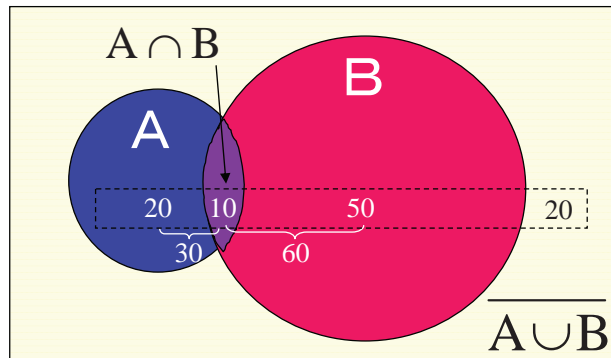
A ∈ A のとる可能性

A=1: Buy a bag of coffee grounds (Search Engine type)

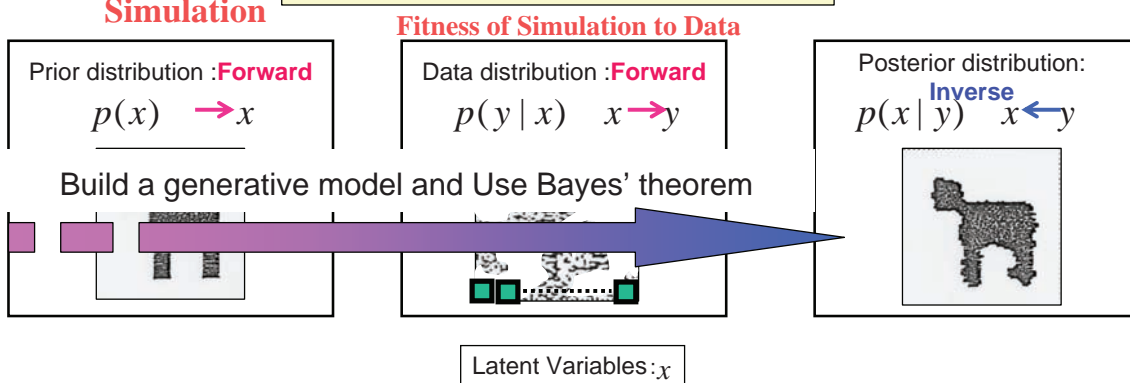
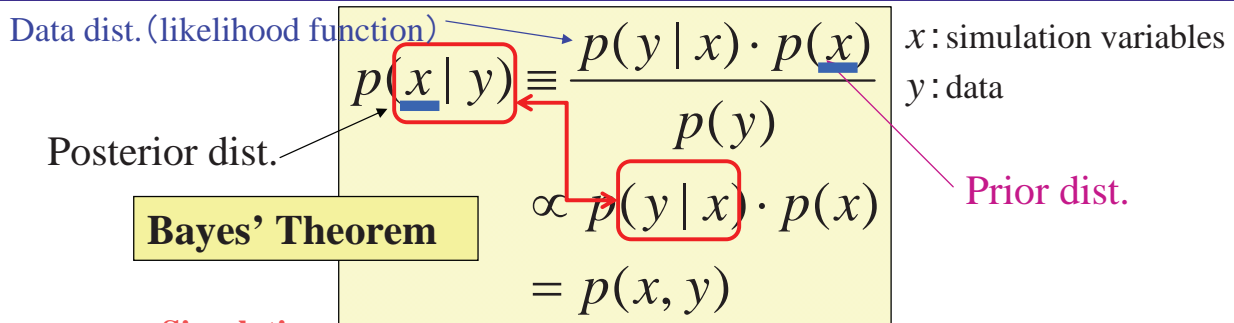
B=1: Buy milk (PC type)

$$p(A=1|B=1) = \frac{p(B=1|A=1)p(A=1)}{p(B=1|A=1)p(A=1)+p(B=1|A=0)p(A=0)}$$

$$= \frac{10 \cdot 30}{30 \cdot 100} = \frac{10}{10+50}$$

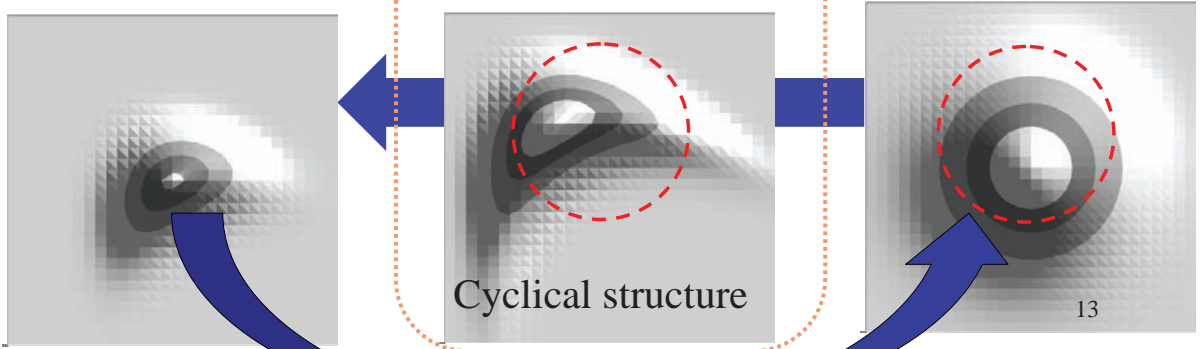
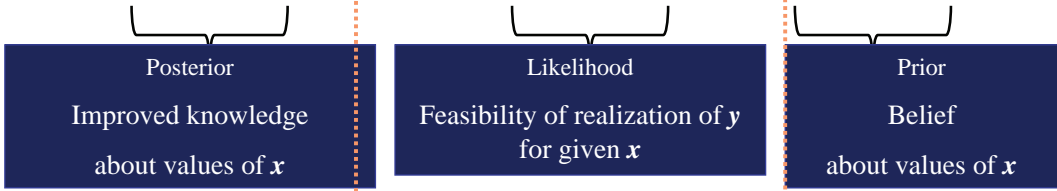


Generative Model, Inversion with Bayes' theorem, and Data Assimilation



Bayesian estimation

$$p(\mathbf{x} | \mathbf{y}) \propto p(\mathbf{y} | \mathbf{x}) \cdot p(\mathbf{x})$$



13/42

Data Assimilation in Generalized State Space Model

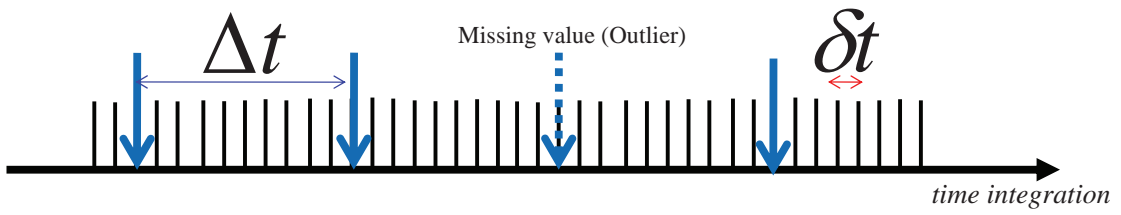
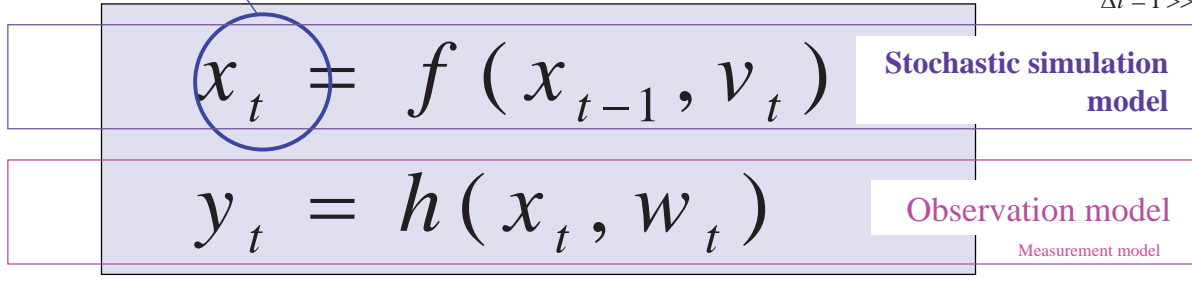
State Vector (Simulation variables)

$L \Rightarrow L$: nonlinear map

Δt : sampling time of observations

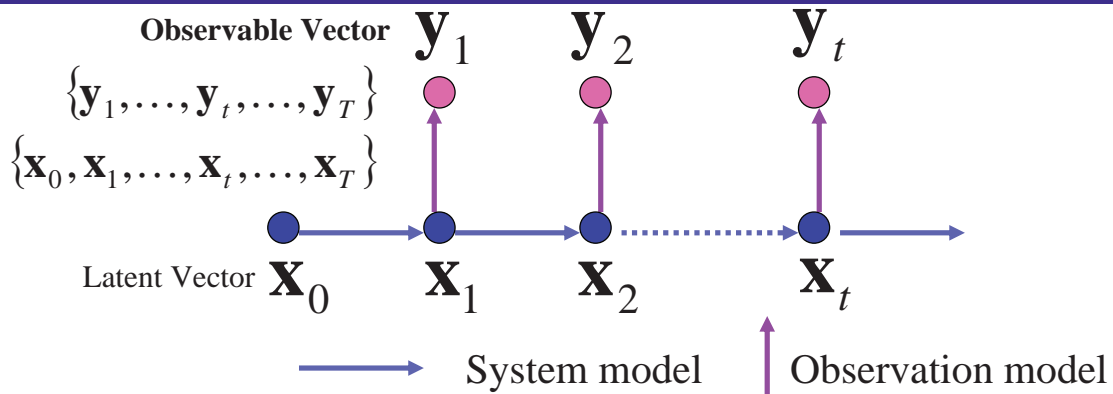
δt : simulation time step

$\Delta t = 1 \gg \delta t$



14/42

Chain Structure Graphical Model



Suppose a statistical inference problem on a daily economic status given daily stock market data

Today's economic situation given yesterday's stock market data $p(\mathbf{x}_t | \mathbf{y}_{1:t-1} \equiv [y_1, y_2, \dots, y_{t-1}])$

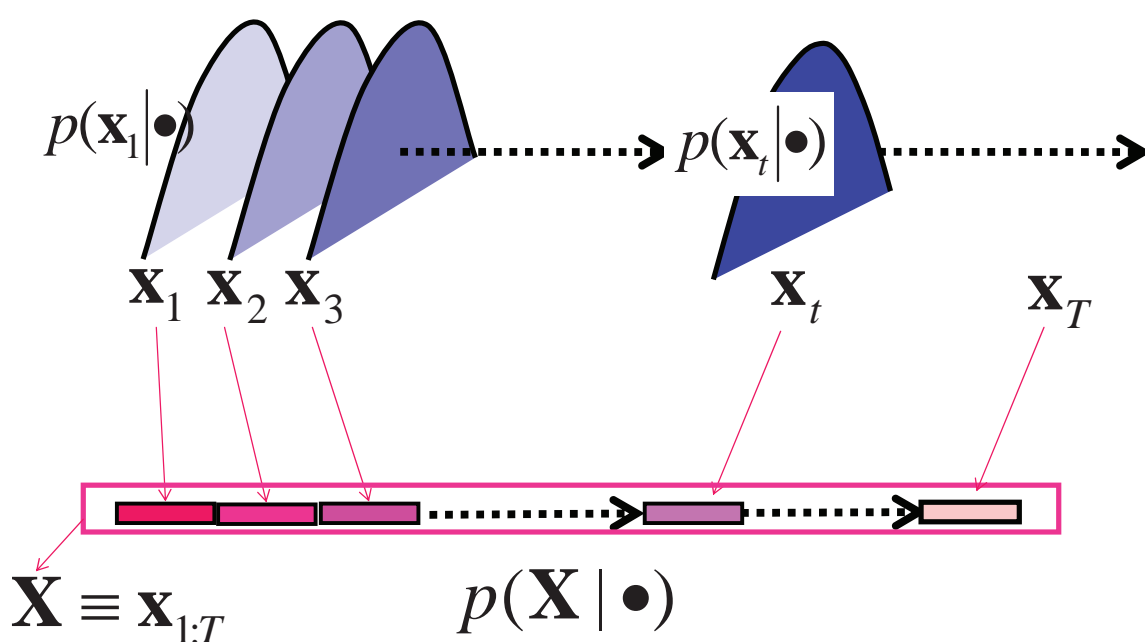
Today's economic situation estimated by the stock market data up to today $p(\mathbf{x}_t | \mathbf{y}_{1:t} \equiv [y_1, y_2, \dots, y_t])$

Today's economic situation analyzed by using all available data when we look back on the today in future $p(\mathbf{x}_t | \mathbf{y}_{1:T} \equiv [y_1, y_2, \dots, y_T])$

$p(\mathbf{x}_{1:T} | \mathbf{y}_{1:T})$

15/42

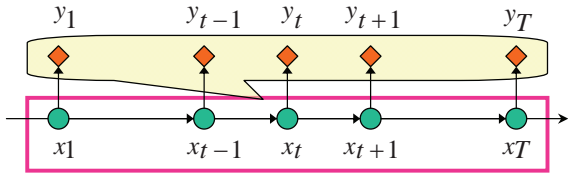
State Vector and Concatenated State Vectors



16/42

Two ways of DA method

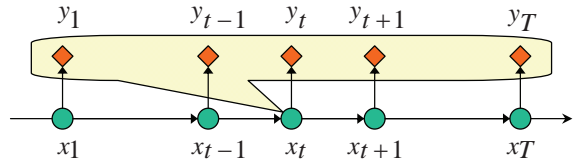
Find $\mathbf{X}=[x_1, \dots, x_T]$ that maximize
 $p(\mathbf{X} | y_{1:T})$



Variational DA method

- Adjoint method (4DVar)
- Representer method

For $t = 1, \dots, T$:
 estimate $p(x_t | y_{1:T})$
 Smoothing dist.



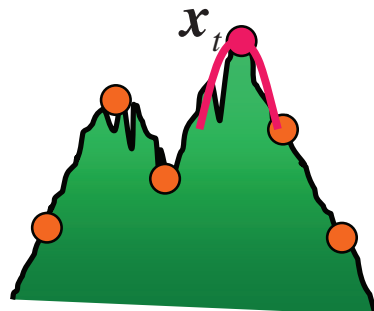
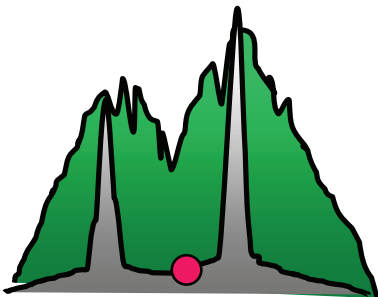
Sequential DA method

- Kalman filter (KF), smoother
- Extended KF (EKF)
- Ensemble KF (EnKF)
- Particle filter (PF)

Optimization and Statistical Inference

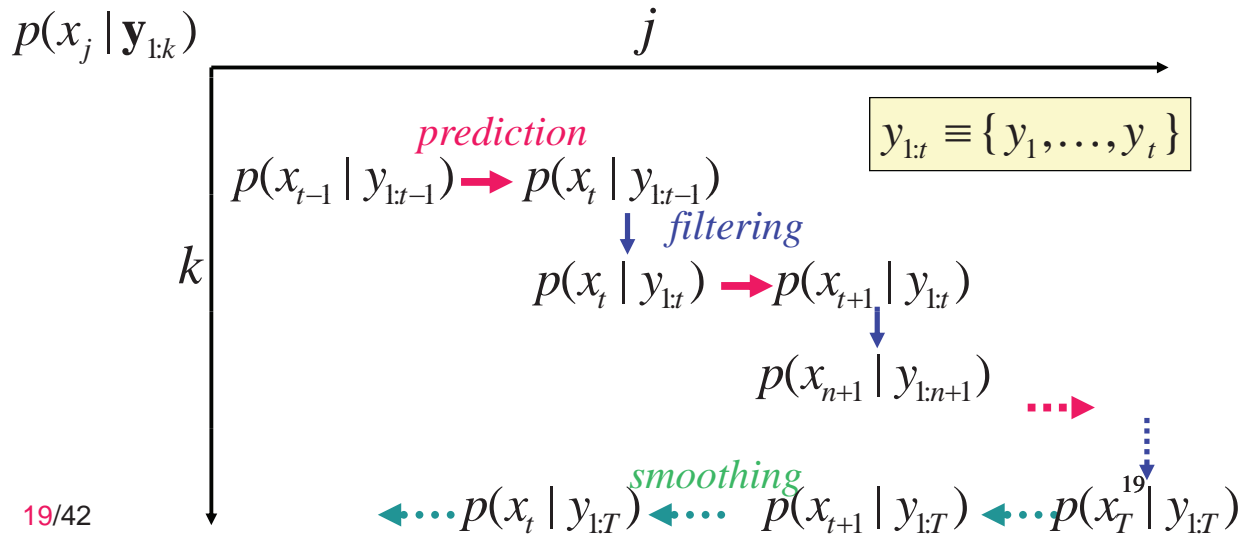
$f(\mathbf{X})$ Dimension of \mathbf{X} is huge
 $\hat{\mathbf{X}} = \max \{-f(\mathbf{X})\}$

$-f(x_t) \Rightarrow p(x_t)$
 $\hat{x}_t = \int p(x_t) \cdot x_t dx_t$



Conditional Distribution **Recursive formula**

<p>predictive density: $p(x_t \mathbf{y}_{1:t-1})$</p>	<p>Today's economic situation given yesterday's stock market data</p>
<p>filter density: $p(x_t \mathbf{y}_{1:t})$</p>	<p>Today's economic situation estimated by the stock market data up to today</p>
<p>smoother density: $p(x_t \mathbf{y}_{1:T})$</p>	<p>Today's economic situation analyzed by using all available data when we look back on the today in future</p>



19/42

Prediction

$$\begin{aligned}
 & p(x_t | y_{1:t-1}) \\
 &= \int p(x_t, x_{t-1} | y_{1:t-1}) dx_{t-1} \\
 &= \int p(x_t | x_{t-1}, y_{1:t-1}) p(x_{t-1} | y_{1:t-1}) dx_{t-1} \\
 & \quad \boxed{p(x_t | x_{t-1}, y_{1:t-1}) = p(x_t | x_{t-1})} \quad \text{Markov property (1)} \\
 &= \int p(x_t | x_{t-1}) \boxed{p(x_{t-1} | y_{1:t-1})} dx_{t-1}
 \end{aligned}$$

Filter pdf at time $t-1$

19-1/42

20


filtering

$$\begin{aligned}
 p(x_t | y_{1:t}) &= p(x_t | y_t, y_{1:t-1}) \\
 &= \frac{p(x_t, y_t | y_{1:t-1})}{p(y_t | y_{1:t-1})} \\
 &= \frac{p(y_t | x_t, y_{1:t-1}) \cdot p(x_t | y_{1:t-1})}{p(y_t | y_{1:t-1})} \\
 &= \frac{p(y_t | x_t) \cdot p(x_t | y_{1:t-1})}{p(y_t | y_{1:t-1})} \\
 &= \frac{p(y_t | x_t) \cdot p(x_t | y_{1:t-1})}{\int p(y_t | x_t) \cdot p(x_t | y_{1:t-1}) dx_t}
 \end{aligned}$$

Posterior, Belief

Markov Property (2)

$p(y_t | x_t, y_{1:t-1}) = p(y_t | x_t)$



19-2/42
21

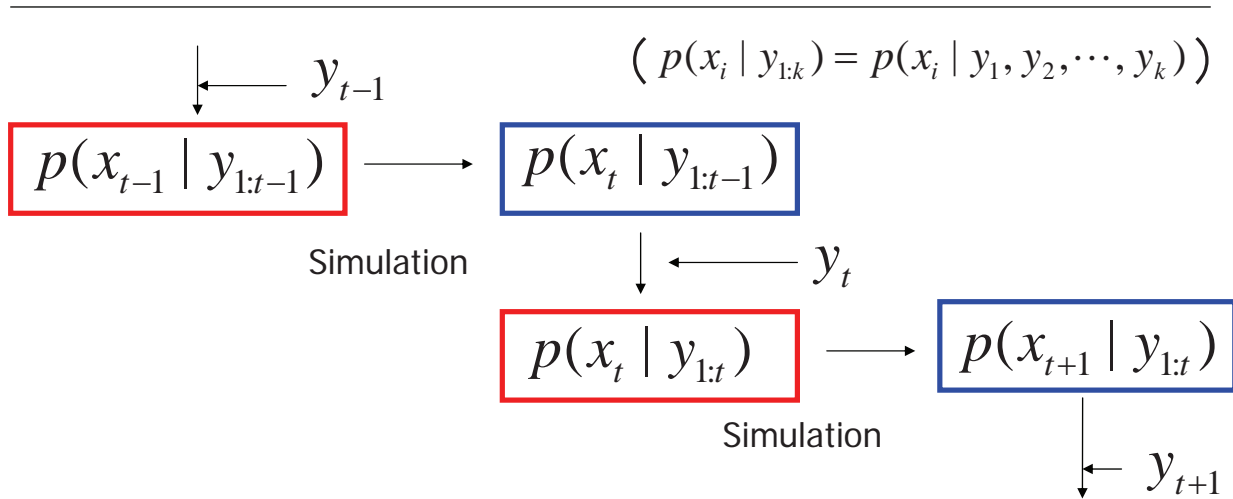
Smoothing

$$\begin{aligned}
 p(x_t | y_{1:T}) &= \int p(x_t, x_{t+1} | y_{1:T}) dx_{t+1} \\
 &= \int p(x_t | x_{t+1}, y_{1:T}) \cdot p(x_{t+1} | y_{1:T}) dx_{t+1} \\
 &= \int p(x_t | x_{t+1}, y_{1:t}) \cdot p(x_{t+1} | y_{1:T}) dx_{t+1} \\
 &= \int \frac{p(x_t, x_{t+1} | y_{1:t})}{p(x_{t+1} | y_{1:t})} \cdot p(x_{t+1} | y_{1:T}) dx_{t+1} \\
 &= \int \frac{p(x_t | y_{1:t}) \cdot p(x_{t+1} | x_t, y_{1:t})}{p(x_{t+1} | y_{1:t})} \cdot p(x_{t+1} | y_{1:T}) dx_{t+1} \\
 &= \underbrace{p(x_t | y_{1:t})}_{\text{Filter Dist.}} \cdot \int \underbrace{\frac{p(x_{t+1} | x_t)}{p(x_{t+1} | y_{1:t})}}_{\text{Prediction Dist.}} \cdot \underbrace{p(x_{t+1} | y_{1:T})}_{\text{Smoothing Dist.}} dx_{t+1}
 \end{aligned}$$

19-3/42
22

Sequential Data Assimilation

Estimate PDF of state vector x_t or its moments (mean, variance, ...) sequentially on each observation



Challenging problem: Huge dimension and inversion

- Data Assimilation = Estimation problem of state vector x_t :

(system model) $x_t = F_t(x_{t-1}, v_t | x_0)$

(observation model) $y_t = H_t x_t + w_t$ or $y_t = h_t(x_t) + w_t$

- x_t : All variables in simulation model
- y_t : All observed variables
- v_t : Stochastic part to represent uncertainty of model (boundary condition, ...)
- w_t : Observation error
- v_t, w_t : Normally Gaussian x_0 : Initial condition

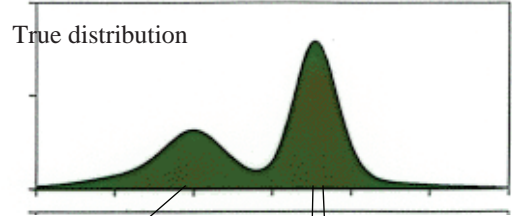
dimension $x_t : 10^4 \sim 10^6$ $y_t : 10^2 \sim 10^5$ $\dim(x_t) \gg \dim(y_t)$

Numerical representation of distribution

$$p(x_t | y_{1:t-1}), p(x_t | y_{1:t}), p(x_t | y_{1:T})$$

Monte Carlo approximation

Represent pdf by the actual realizations.



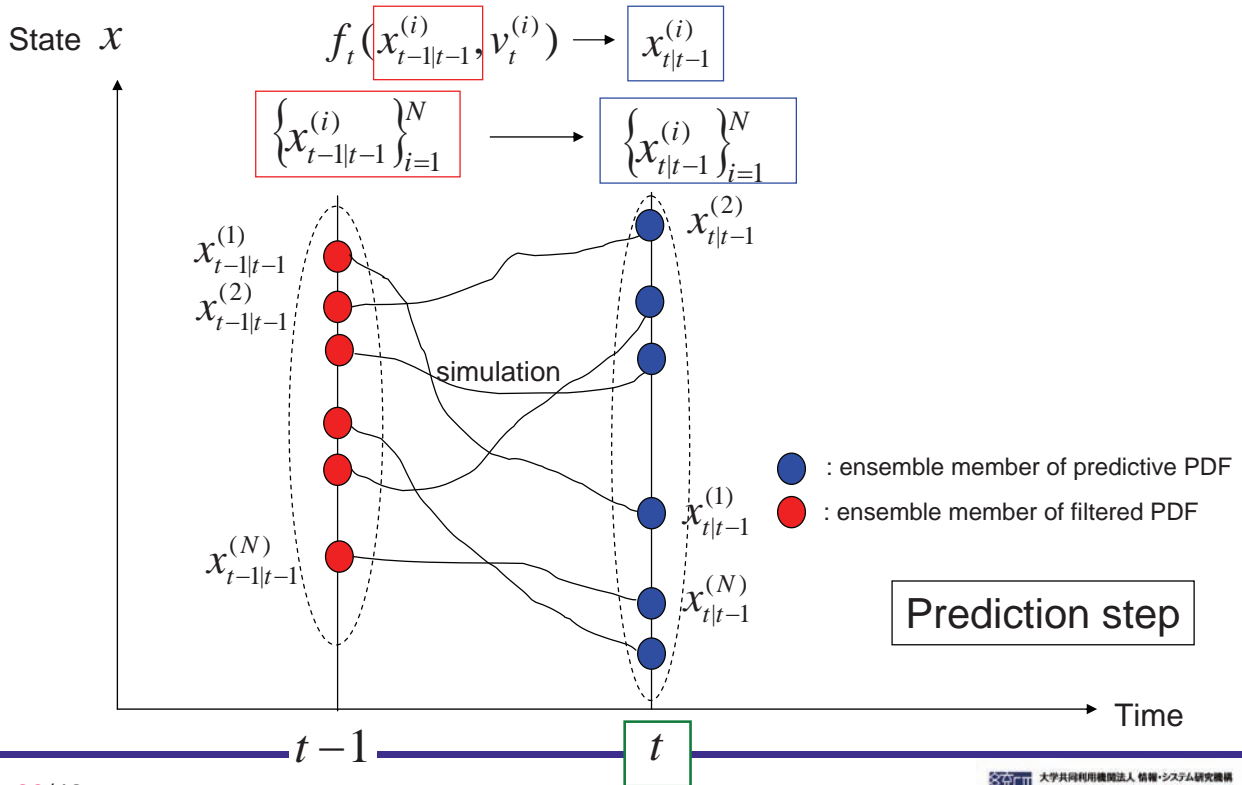
N : # of particles

$$p(x_t | y_{1:t-1}) \cong X_{t|t-1} \equiv [x_{t|t-1}^{(1)}, x_{t|t-1}^{(2)}, \dots, x_{t|t-1}^{(N)}]$$

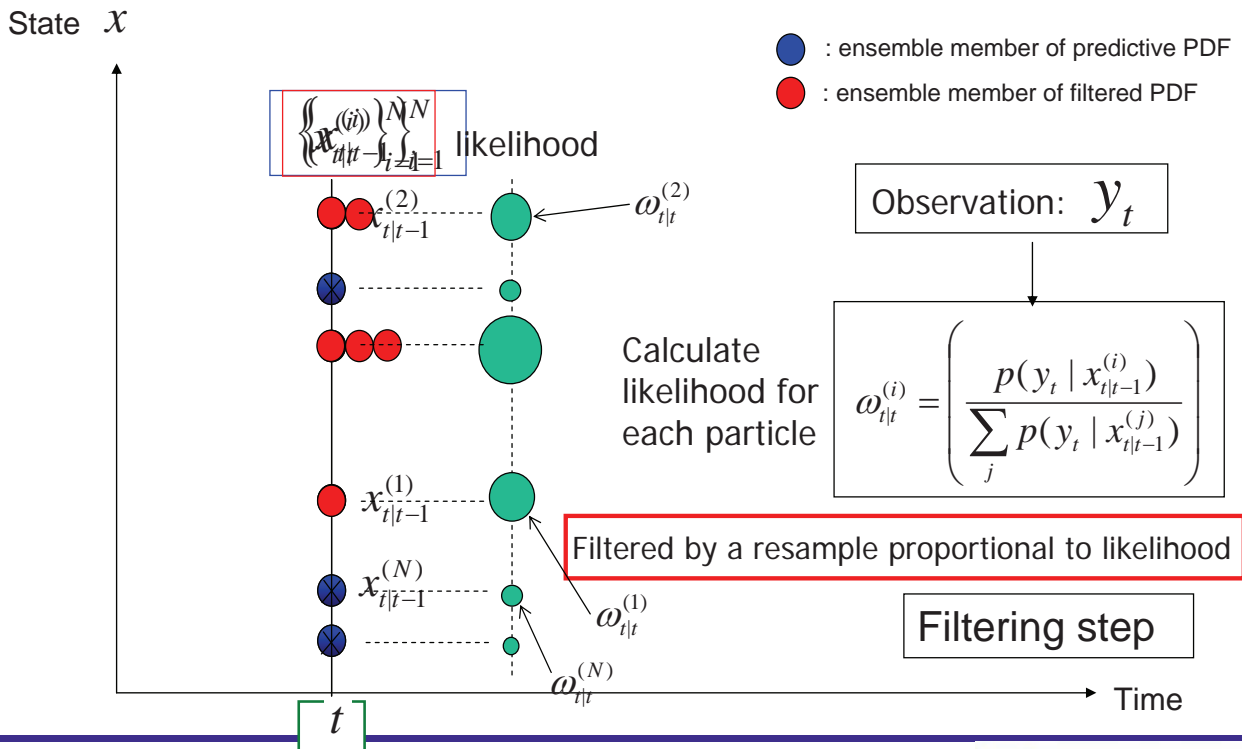
$$p(x_t | y_{1:t}) \cong X_{t|t} \equiv [x_{t|t}^{(1)}, x_{t|t}^{(2)}, \dots, x_{t|t}^{(N)}]$$



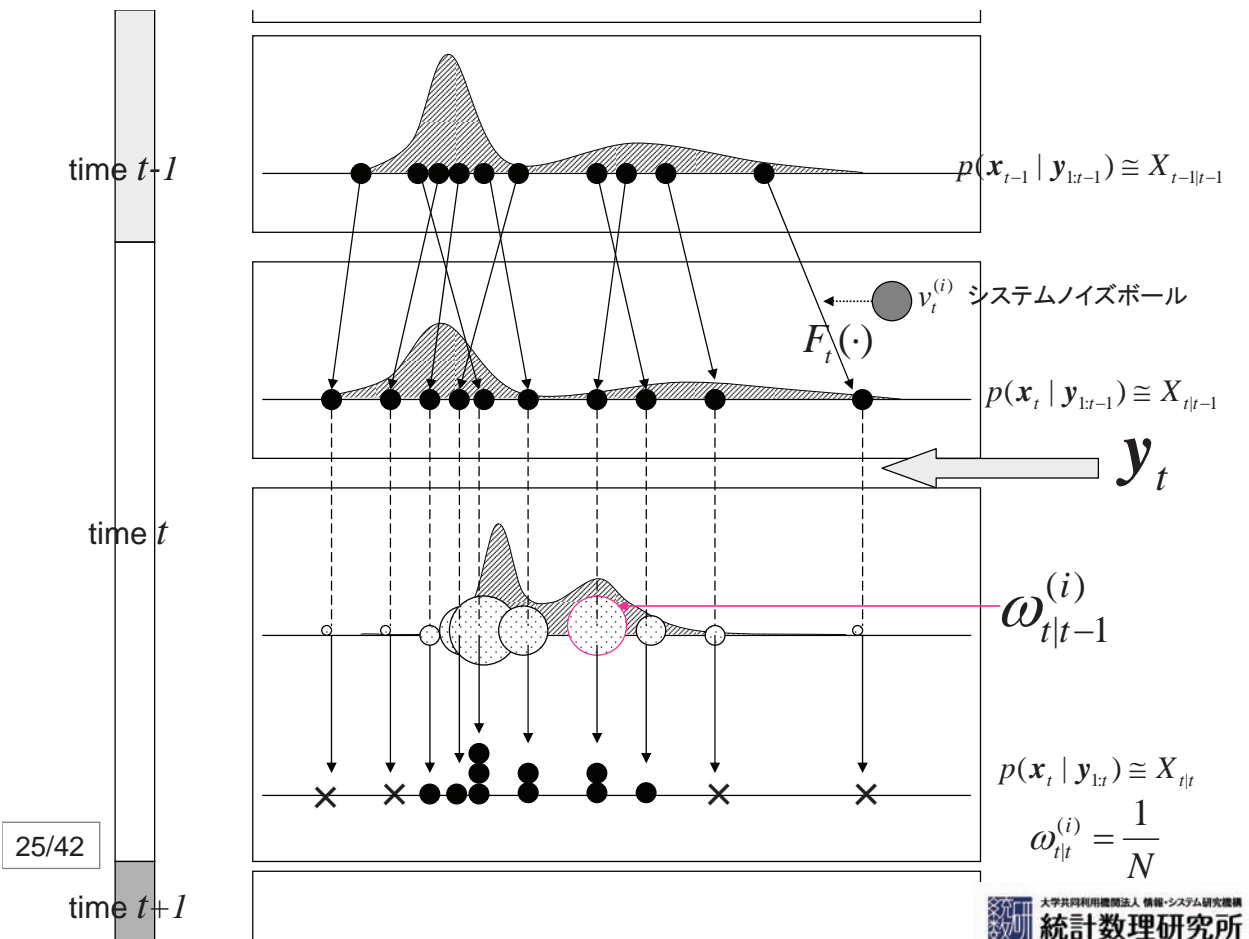
Prediction Step (Common in EnKF and PF)



Filtering Step of PF



24/42



Family of particle filters

- Kalman filter (EKF: Extended Kalman filter)
- EnKF: Ensemble Kalman filter (Evensen 1994)
- Particle filter
 - SIR filter (e.g., Gordon et al. 1993, Kitagawa 1993)★
 - Gaussian particle filter (e.g., Kotecha and Djuric 2003)
 - Kernel filter (e.g., Hurzeler and Kunsh 1998)
 - Merging particle filter (MPF)

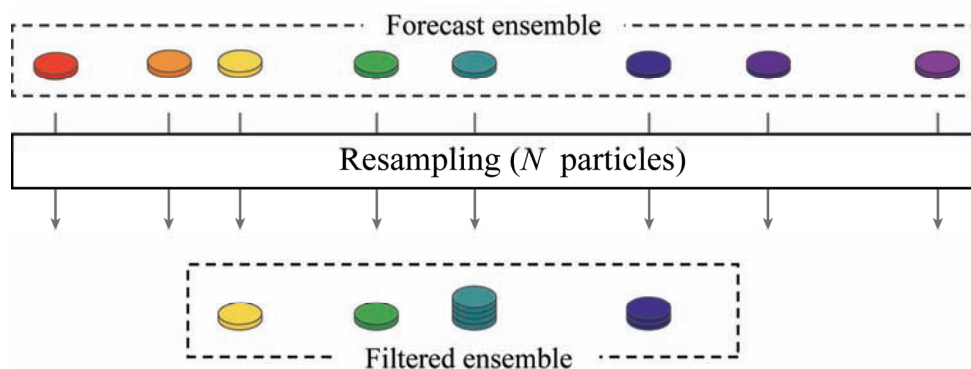
★ It encounters a problem called ‘degeneration’ in applying to high-dimensional models. (i.e., the diversity of an ensemble is lost after repeating resampling procedures.)

※Gaussian PF: 1) Each particle in filtered ensemble is drawn from a Gaussian function with the mean and covariance of the forecast ensemble. 2) It requires high computational cost due to a factorization of a high-dimensional covariance matrix in generating Gaussian samples.

26/42

29
 大学共同利用機関法人 情報・システム研究機構
 統計数理研究所

Particle filter (SIR: Sequential Importance Resampling)

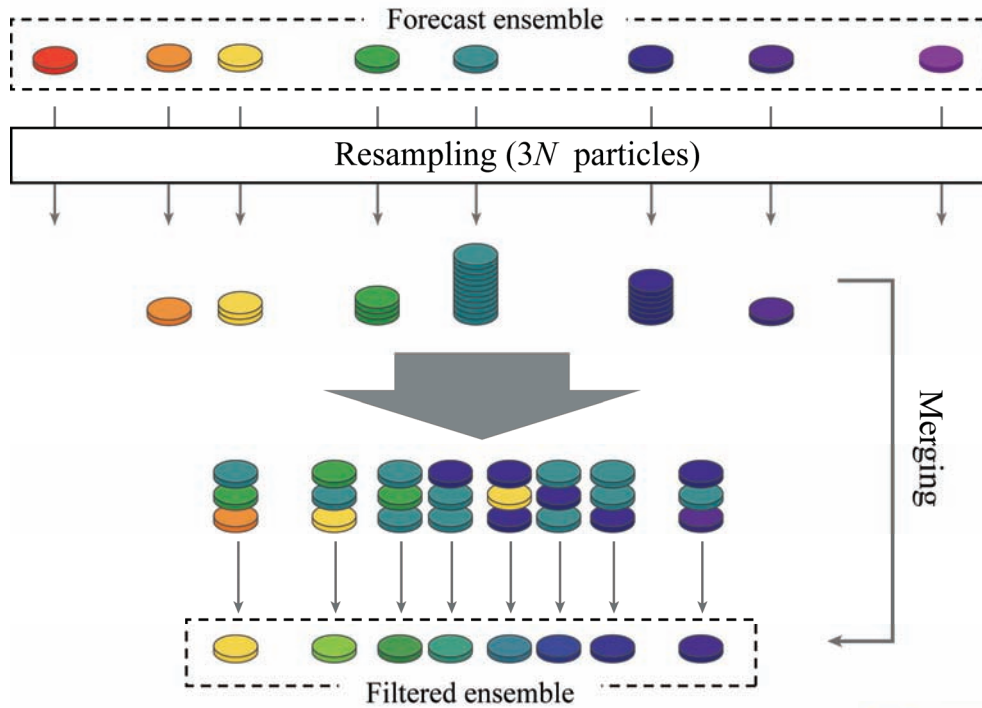


A posterior (filtered) ensemble is obtained by resampling the forecast ensemble with weights of likelihood. Thus, an ensemble member is duplicated in the filtered ensemble according to its likelihood.

27/42

大学共同利用機関法人 情報・システム研究機構
 統計数理研究所

Merging particle filter (MPF)



28/42

MPF algorithm (1)

We draw $n \times N$ samples from the forecast ensemble with weights of w_i , and obtain an ensemble: $\{\hat{\mathbf{x}}_{t|t}^{(1,1)}, \dots, \hat{\mathbf{x}}_{t|t}^{(n,1)}, \dots, \hat{\mathbf{x}}_{t|t}^{(1,N)}, \dots, \hat{\mathbf{x}}_{t|t}^{(n,N)}\}$.

A subset $\{\hat{\mathbf{x}}_{t|t}^{(j,1)}, \dots, \hat{\mathbf{x}}_{t|t}^{(j,N)}\}$ from this $n \times N$ samples satisfies

$$p(\mathbf{x}_t | \mathbf{y}_{1:t}) \approx \frac{1}{N} \sum_{i=1}^N \delta(\mathbf{x}_t - \hat{\mathbf{x}}_{t|t}^{(j,i)})$$

because it is a filtered ensemble obtained using normal PF.

Each number of a filtered ensemble is generated as a weighted sum of n samples from the $n \times N$ sample set as:

$$\mathbf{x}_{t|t}^{(i)} = \sum_{j=1}^n \alpha_j \hat{\mathbf{x}}_{t|t}^{(j,i)}$$

29/42

MPF algorithm (2)

In order to ensure that the newly generated ensemble approximately preserves the mean and covariances of the filtered PDF, the merging weights α_j are set to satisfy

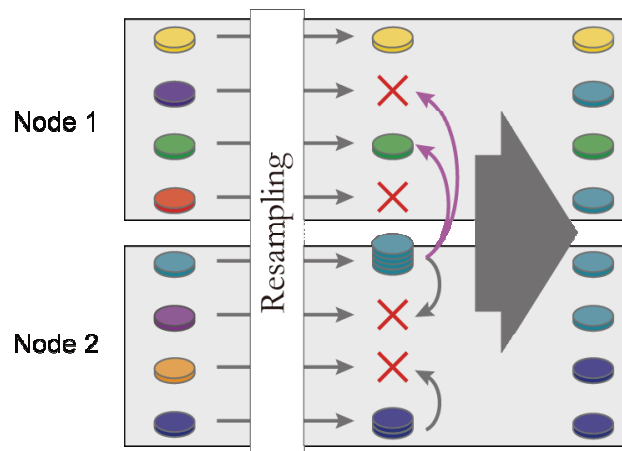
$$\sum_{j=1}^n \alpha_j = 1, \quad \sum_{j=1}^n \alpha_j^2 = 1 \quad (n \leq 3 \text{ such that } \alpha_j \neq 0 \text{ for all } j)$$

where each α_j is a real number.

Then, a new ensemble approximation of the filtered PDF $p(\mathbf{x}_t | \mathbf{y}_{1:t})$ is obtained as

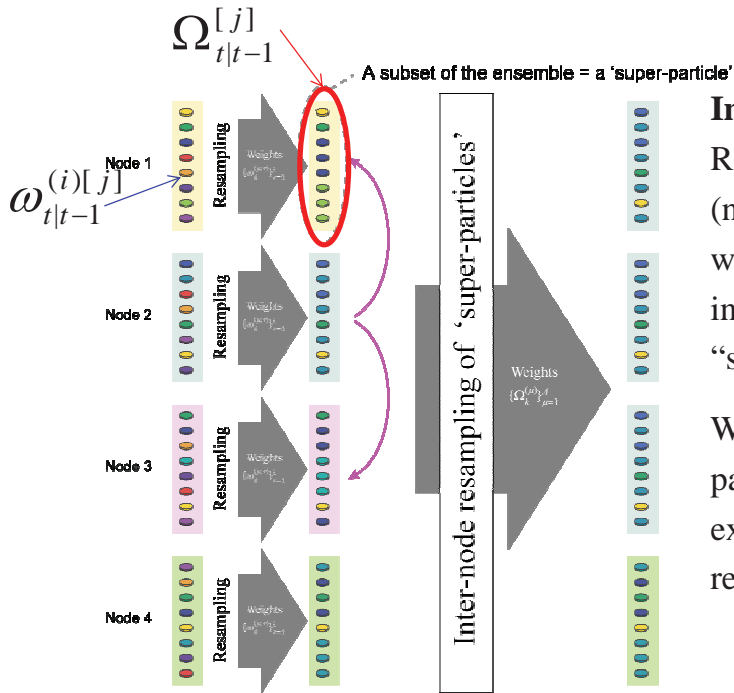
$$p(\mathbf{x}_t | \mathbf{y}_{1:t}) \approx \frac{1}{N} \sum_{i=1}^N \delta(\mathbf{x}_t - \hat{\mathbf{x}}_{t|t}^{(j,i)})$$

Flowchart of PF



A concept of the “Islands” in GA is similar, but different.

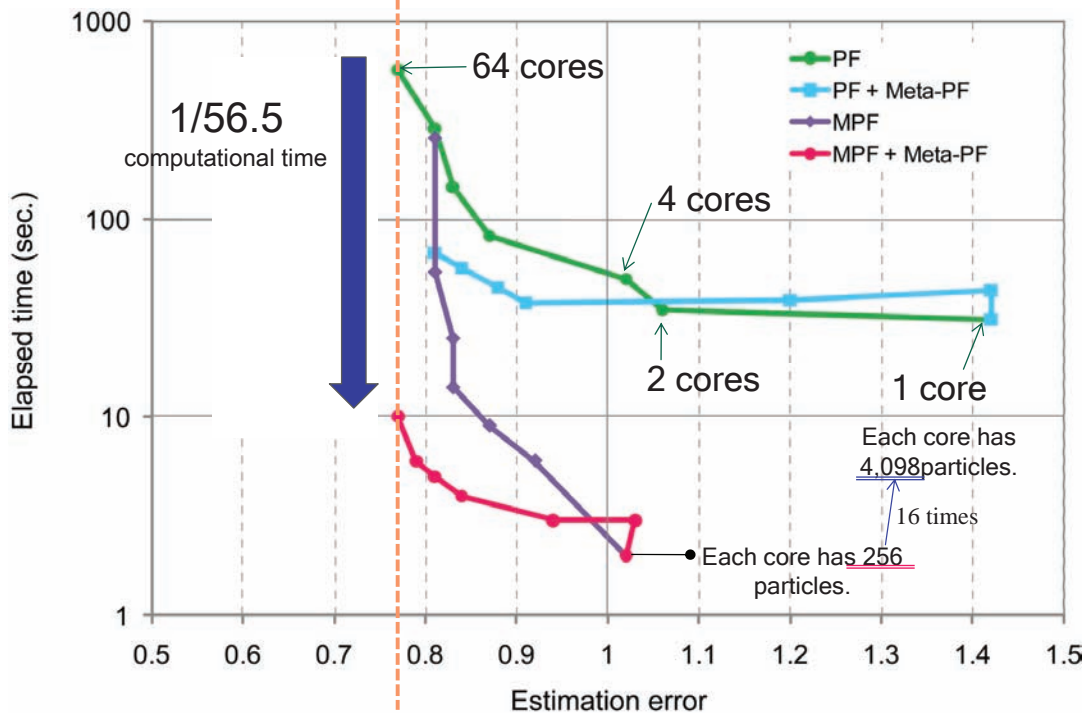
Meta particle filter (DPF: Distributed particle filter)



Inter-node resampling:
Resampling between ensembles (not ensemble member) each of which consists of many particles in a node. This ensemble is called "super-particle".

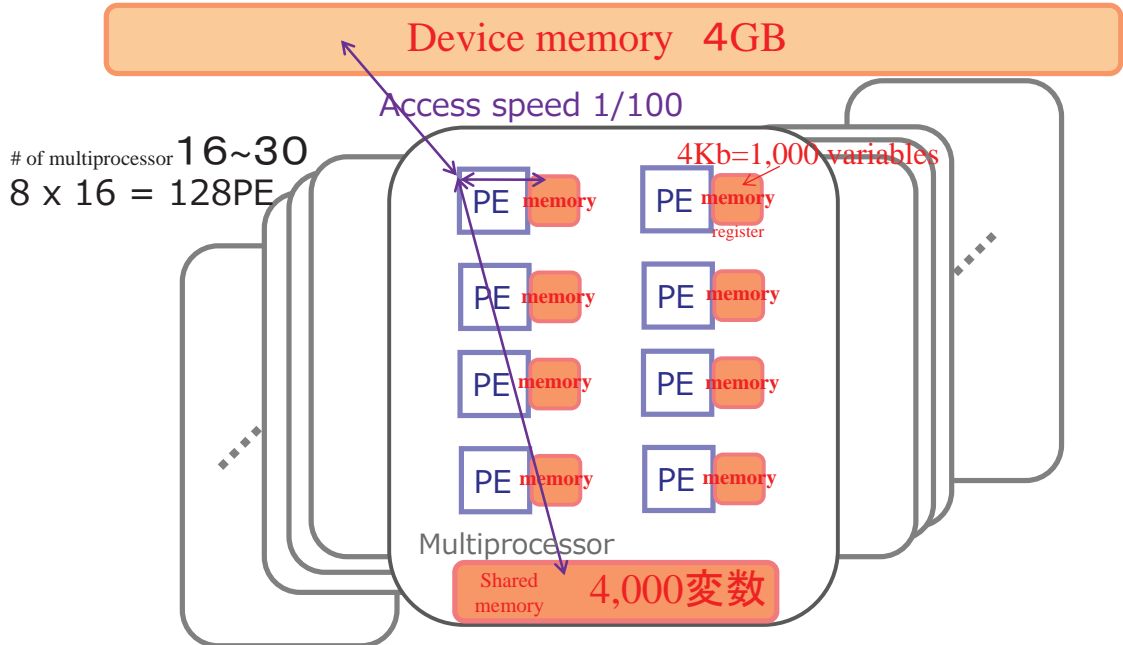
When a weight for any super-particles (i.e., nodes) $\Omega_{t|t-1}^{[j]}$ exceeds 0.3, the inter-node resampling procedure is applied.

Result



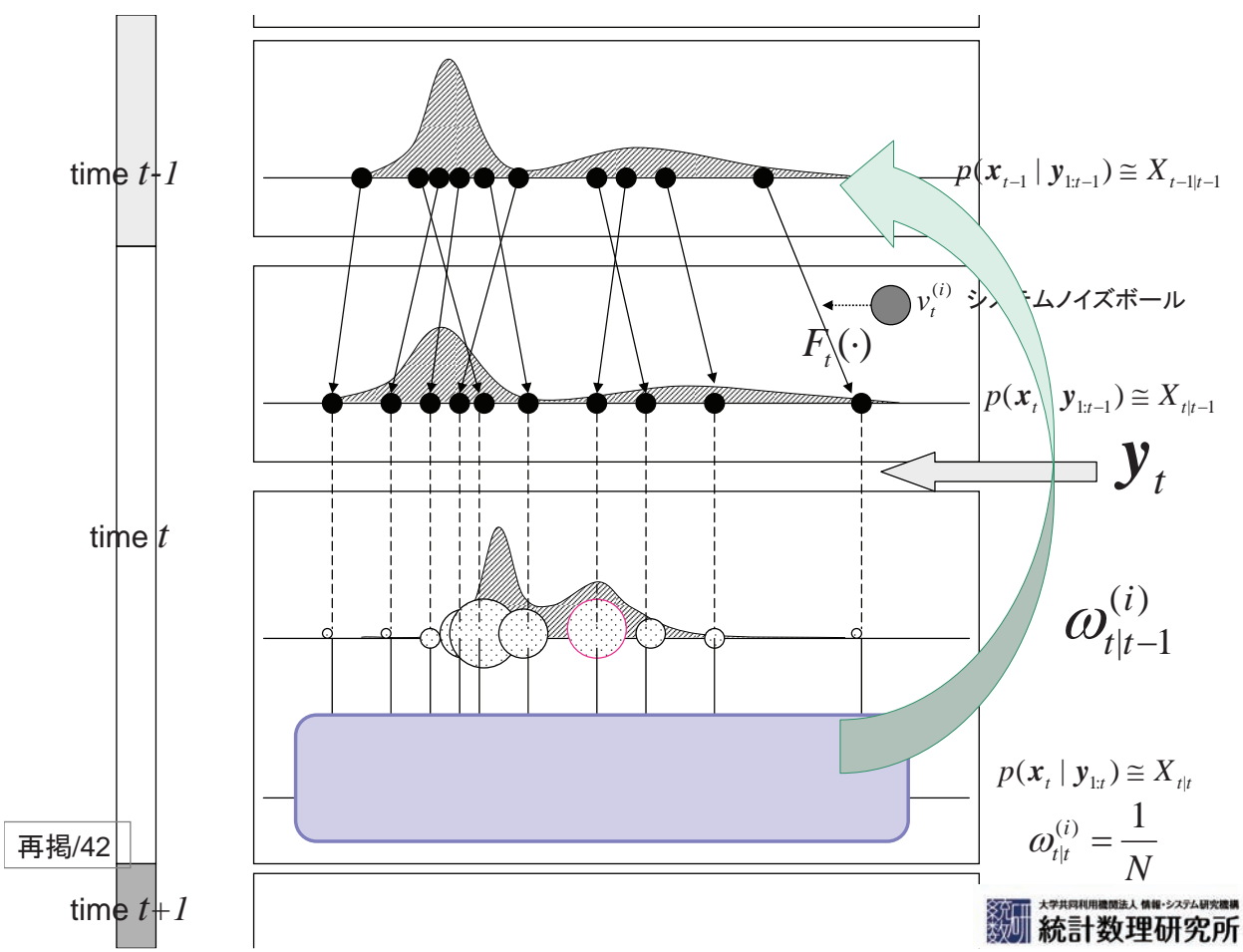
Inside GPU: 128~240 parallel computing

GPGPU's power is equivalent to a computer with 128 1.2GHz processors

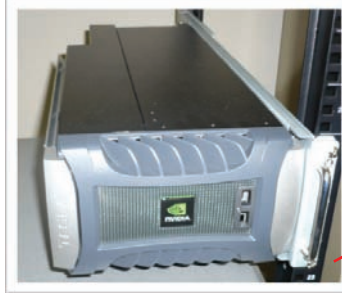


PE: Processing Element

34/42



Opteron + Tesla + ClearSpeed@ Higuchi Lab.



Tesla (GPU/Cuda)

nVIDIA S1070 1U
GPU computing server
(1.6GiB/1.6GHz/410GB/s)



Host machine (CPU)

HP xw9400 Workstation
Dual-Core AMD Opteron(tm) Processor 2220



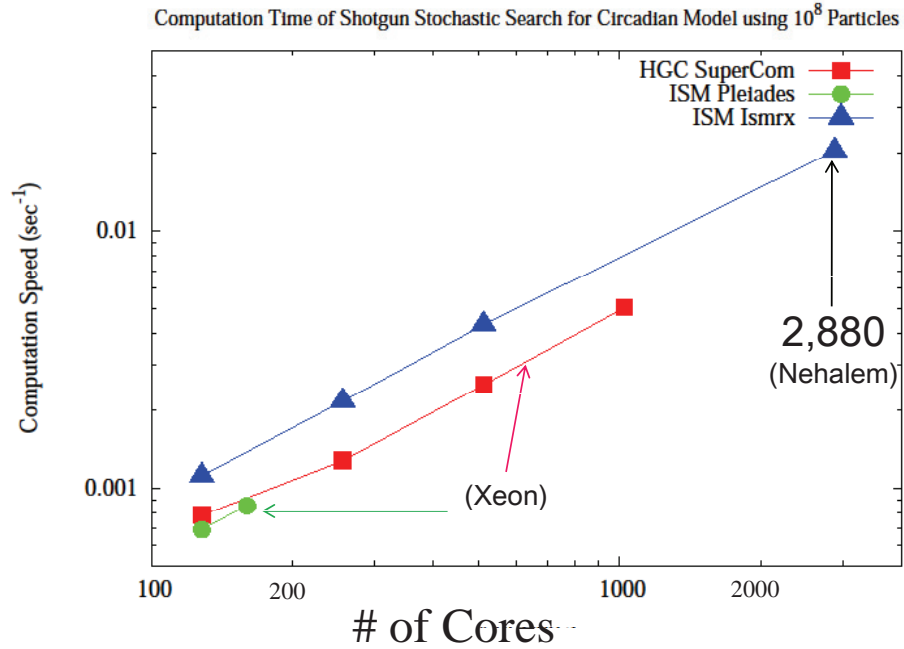
ClearSpeed e710

SIS meets GPGPU.

- SIS on GPGPU designed for parameter estimation.
 - Simulation is carried out on GPGPU.
 - Parameter estimation is carried out on CPU.

particles	PF Opteron 2220 1core (Nakamura et al., 2009)	PF Opteron2220 1core + GPGPU(Tesla C870)	SIS Opteron2220 1core + GPGPU(Tesla C870)
100,000,000 (1億)	8Days (6.8×10^5 sec)	12Hours (4.5×10^4 sec)	3Hours (1.0×10^4 sec)
	1	× 15	× 67
1,000,000,000 (10億)	2.6 Months?	5 Days?	28Hours (1.0×10^5 sec)

HPC of our group




37/42

Next-Generation of Supercomputer in Japan at Kobe



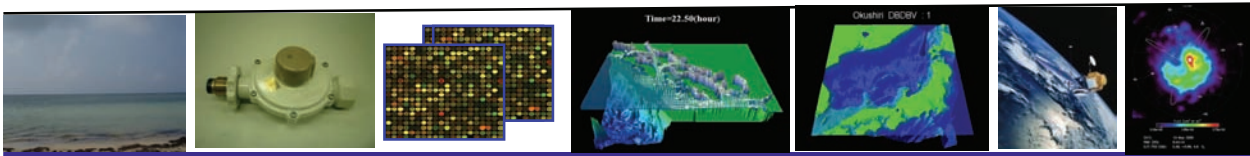
Japanese Government will spend more than 1 billion US\$ for this national project. It has more than **600,000** cores.

- Grand Challenge:  ISLAM
Integrated Simulation of Living Matter
- Nanotech (Institute for Molecular Science)
- Life Science (RIKEN)

38/42

Research projects in progress by our group

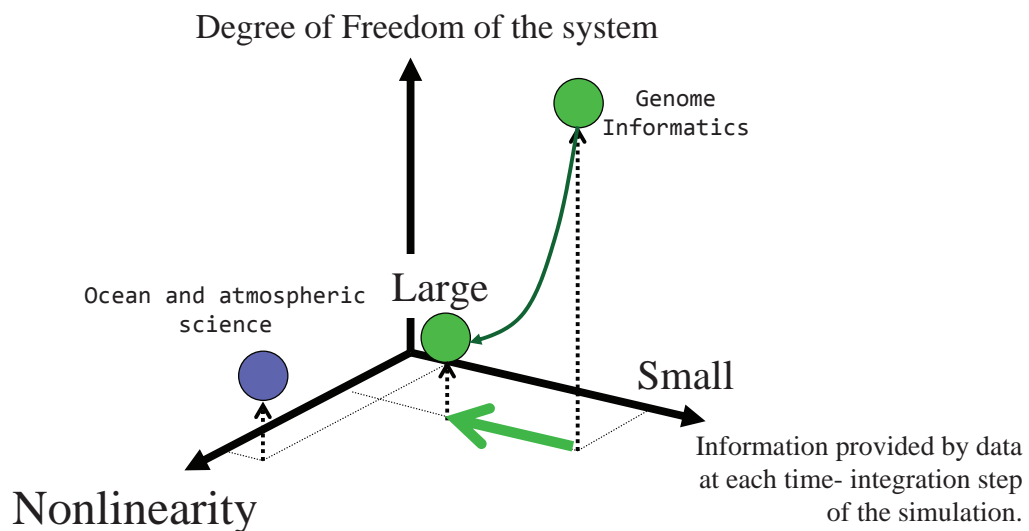
- Coupled Ocean-Atmosphere model
- Tsunami model
- Ocean tide
- 3D structure of ring current
- Genome informatics
- Marketing (with agent simulations)



39/42

43

TIPS: A choice of the data assimilation methods Reduction of degree of freedom



40/42

Contact

Email: higuchi@ism.ac.jp

Homepage:
<http://daweb.ism.ac.jp/>