

HPCとともに進化する大規模データ同化

統計数理研究所

樋口知之

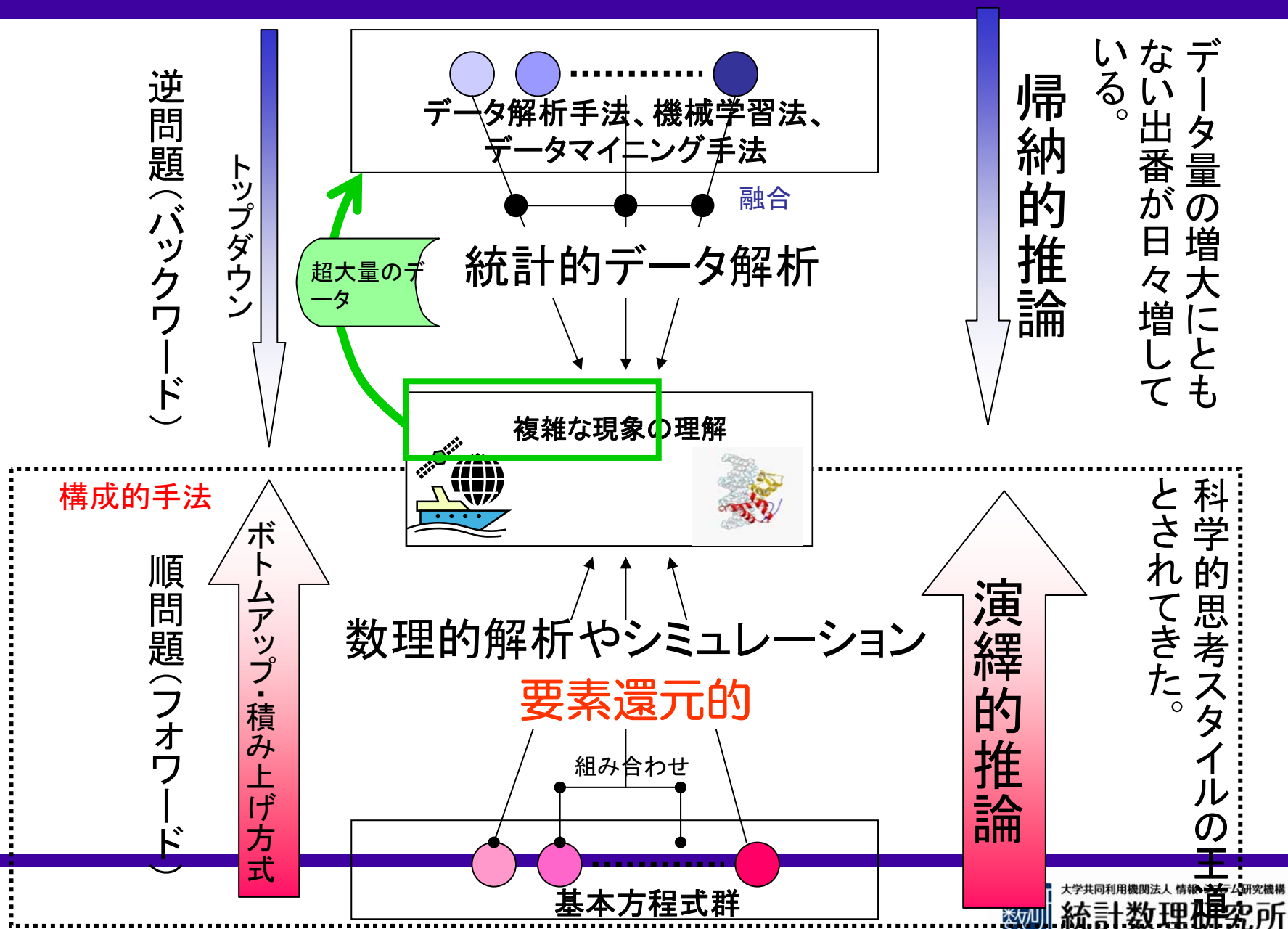
予測発見戦略研究センター データ同化グループ

樋口、上野准教授、吉田助教、中野助教
ポスドク研究員： 斉藤、林、井元、長尾、才田

次世代シミュレーションNOEグループ

田村教授（主幹）、中野教授（センター長）、樋口
佐藤准教授、上野准教授、吉田助教、中野助教

帰納的アプローチ



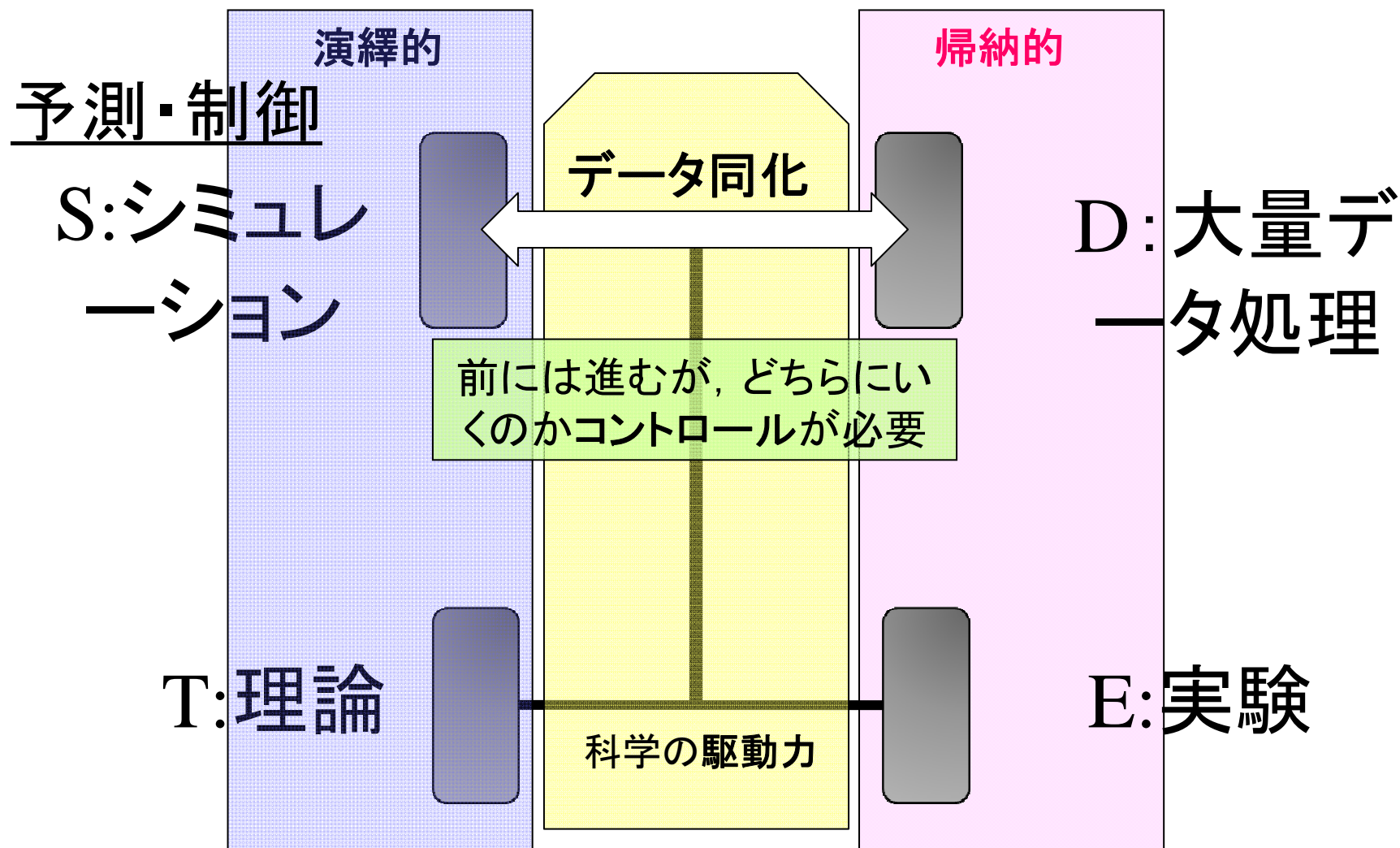
社会的要請: Personalization

- 背景:
 - 無駄を省く(低価格化, 低コスト)
 - 資源の有効利用のために選択と集中
 - 価値観の多様化
 - “コ”(個人, 個性, 固有, 個別)に特化

大量生産・大量消費をめざした20世紀→
個人に焦点をあわせる科学へ

オーダーメイド医療, 副作用の研究, マイクロマーケティング,
One-to-One *, Situation *, 環境に優しい

両者をつなぐのは禁手とされてきた



話の流れ

- ・ データ同化 (Data Assimilation (DA)) のコンセプト
- ・ シミュレーションと一般状態空間モデル
- ・ 逐次ベイズ計算
- ・ アンサンブルベース逐次データ同化
 - アンサンブルカルマンフィルタ (EnKF: Ensemble Kalman filter)
 - 粒子フィルタ (PF: Particle Filter)
- ・ 次世代スーパーコンピュータ実装にむけたアルゴリズム開発のTips
- ・ 未来デザインの道具へ

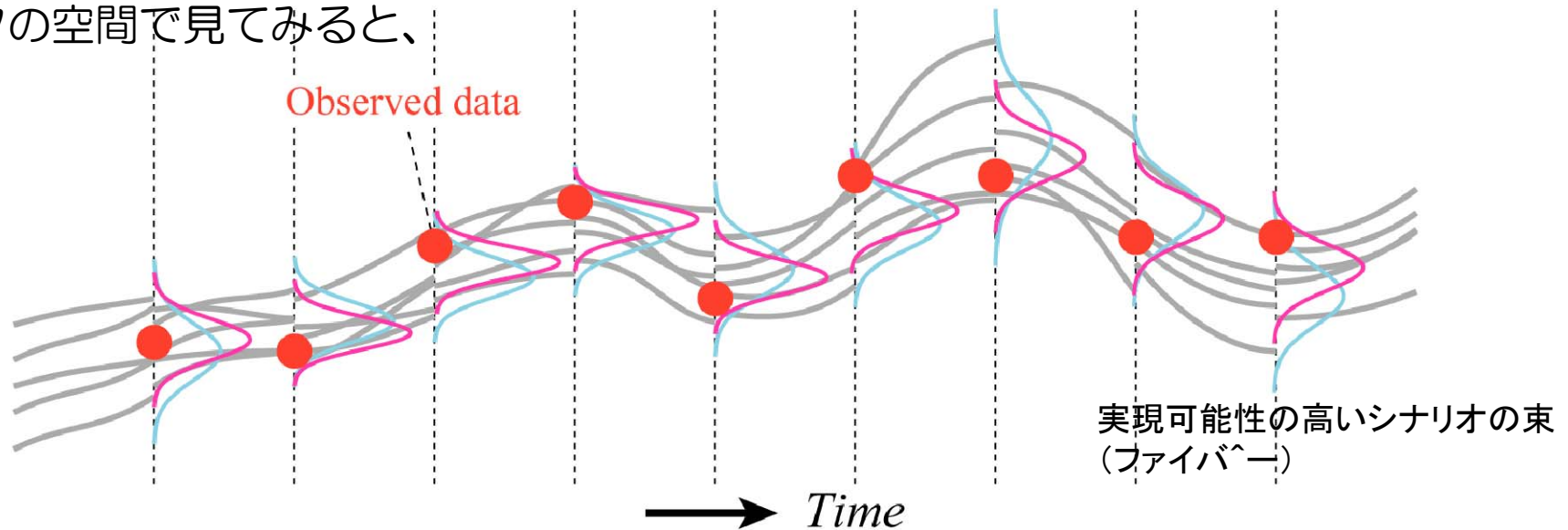
データ同化の目的: 気象・海洋学の観点から

- [1] 予報を行うための最適な初期条件を求める。これは既に、現業の天気予報で実用化されていることである。
- [2] シミュレーションモデルを構成する際の最適な境界条件を求める。連成現象を取り扱う際の適応的な境界条件設定もこの作業に含まれる。
- [3] スケールが異なるシミュレーションモデル間の橋渡しを行うスキーム内に含まれる諸パラメータの最適な値を求める。経験的に与えられるモデル内のパラメータ値の検証も一つの具体例である。
- [4] シミュレーション(物理)モデルにもとづいた、観測されていない時間・空間点における観測値の補間を行う。この作業は再解析データセットの生成とも呼ばれる。このデータセットから新しい科学的発見をもくろむ。
- [5] 時間・経費を節約できる効率的な観測システムを構築するための仮想観測ネットワークシミュレーション実験や感度解析を行う。

データ同化のイメージ

データが解空間を制約

データの空間で見ると、



A physical simulation model can provide various **scenario** about the temporal evolution of the system of interest if we change initial conditions, boundary conditions, parameter settings, and so forth.

Data assimilation aims at producing the **most likely scenario** by incorporating observations into the physical model.

シミュレーションモデルと状態ベクトル

page1

(日本周辺の簡易化した気象モデルの例を用いて説明)

実システムに対応した偏微分方程式
(連続時間・空間)

$$\frac{\partial x}{\partial t} = cx^2 + \dots$$

数値シミュレーションモデル
(離散時間・空間, 有限差分方程式)

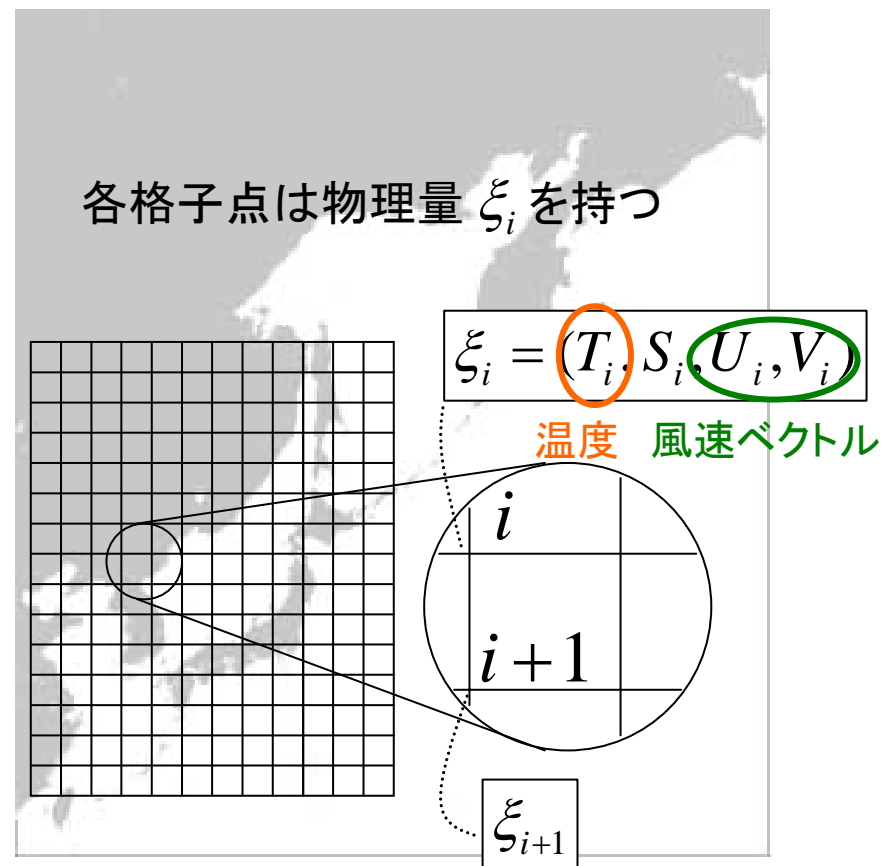
$$x_t = f_t(x_{t-1})$$

境界条件・モデル化誤差
由来の不確かさ

v_t

非線形状態空間表現のシステムモデル
(離散時間・空間, 確率差分方程式)

$$x_t = f_t(x_{t-1}, v_t)$$



観測点・観測される物理量の種類は限定

シミュレーションモデルと状態ベクトル

page2

(日本周辺の簡易化した気象モデルの例を用いて説明)

実システムに対応した偏微分方程式
(連続時間・空間)

$$\frac{\partial x}{\partial t} = cx^2 + \dots$$

数値シミュレーションモデル
(離散時間・空間, 有限差分方程式)

$$x_t = f_t(x_{t-1})$$

境界条件・モデル化誤差
由来の不確かさ

v_t

非線形状態空間表現のシステムモデル
(離散時間・空間, 確率差分方程式)

$$x_t = f_t(x_{t-1}, v_t)$$

$$x_t = \begin{bmatrix} \xi'_1 \\ \vdots \\ \xi'_i \\ \xi'_{i+1} \\ \vdots \\ \xi'_N \\ \theta \end{bmatrix}$$

$S_i(U_i, V_i)$

風速ベクトル

類は限定

シミュレーションとデータ同化

State Vector (Simulation variables)

$L \Rightarrow L$: nonlinear map

Δt : sampling time of observations

δt : simulation time step

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

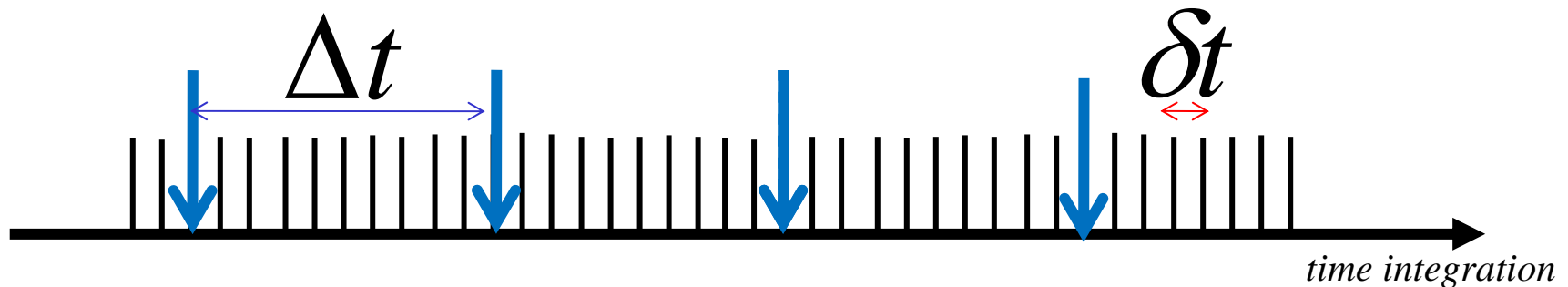
$$x_t = f(x_{t-1}, v_t)$$

Stochastic simulation model

$$y_t = h(x_t, w_t)$$

Observation model

Measurement model



時系列モデルの拡張

統計数理研究所の伝統の強み (Akaike (1980), Kitagawa (1987, 1996))

Linear-Gaussian

State Space Model [SSM]

$$x_t = Fx_{t-1} + Gv_t$$

$$y_t = Hx_t + e_t$$



Non-linear Non-Gaussian Model

$$x_t = f(x_{t-1}, v_t)$$

$$y_t = h(x_t, e_t)$$



Generalized State Space Model [GSSM]

$$x_t = f(x_{t-1}, v_t)$$

$$y_t \sim r(\cdot | x_t)$$

条件付確率と同時確率

$p(A) \equiv A$ が起きる確率

$p(A | B) \equiv B$ が起きたもとで A が起きる確率 ← 条件付確率

$p(A, B) \equiv A$ と B が同時に起きる確率 ← 同時確率

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

消費者(お客の)数全体: 100

コーヒー豆を買った人の数: 30

ミルクを買った人の数 60

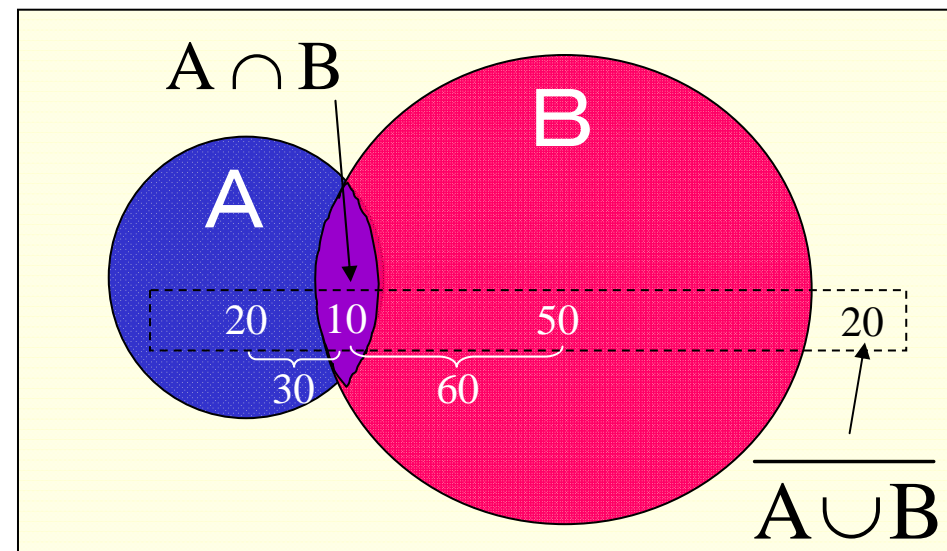
コーヒー豆もミルクも買った人の数: 10

A=1: コーヒー豆も買った

=0: 買わなかった

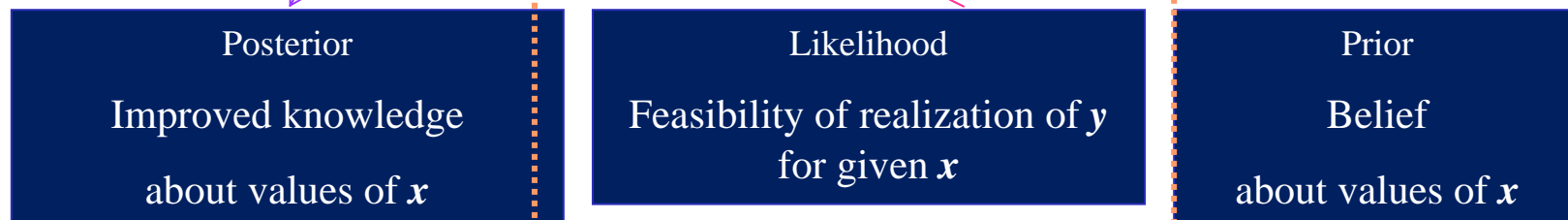
B=1: ミルクを買った

=0: 買わなかった

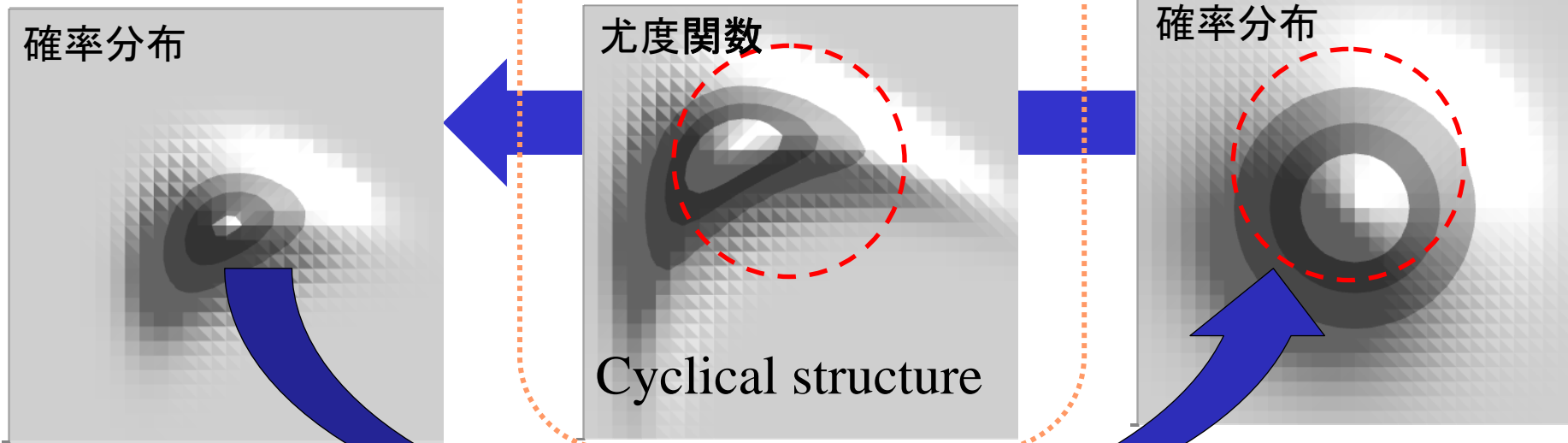


ベイズの定理と情報循環

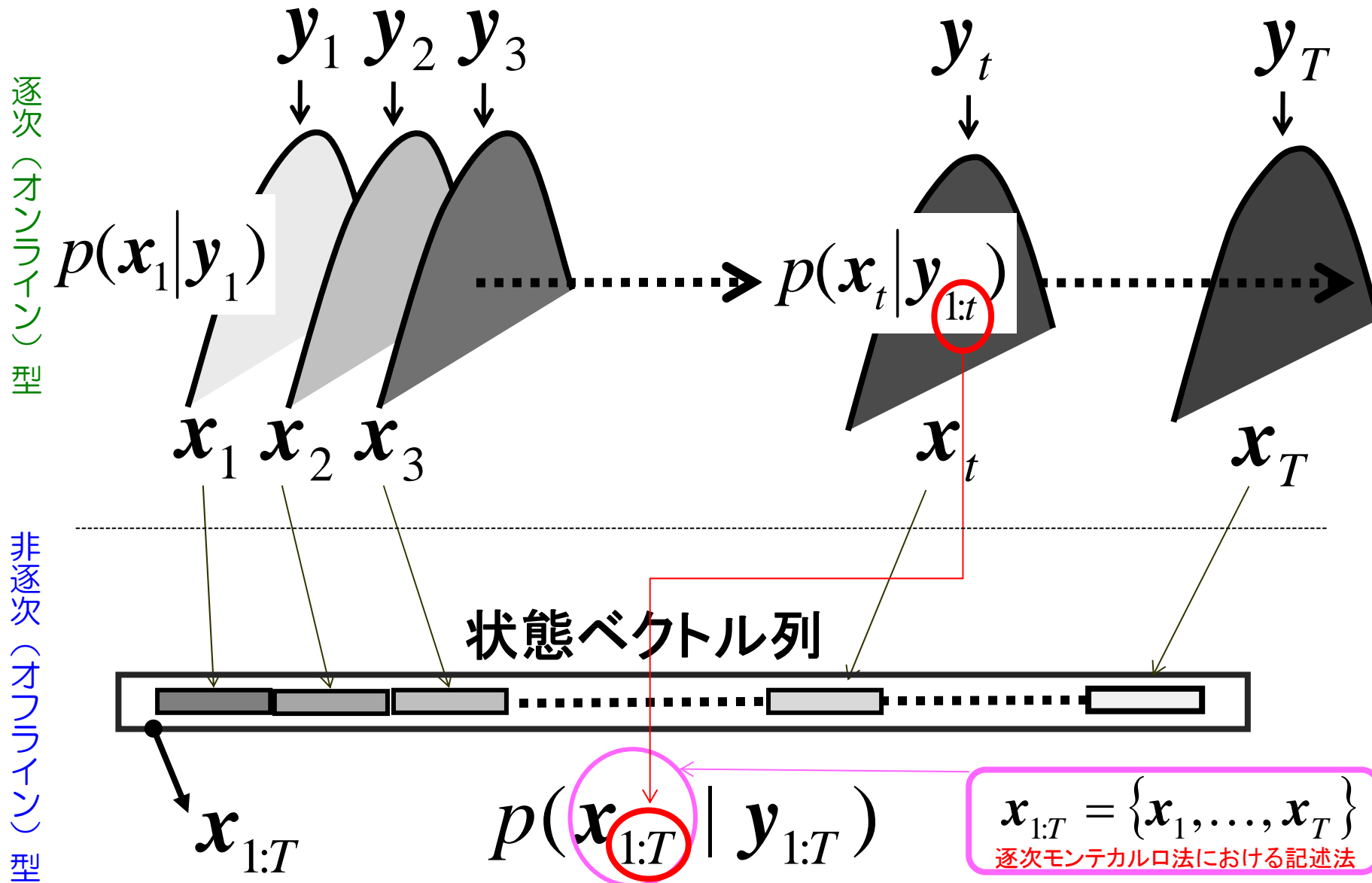
$$p(\mathbf{x} | \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{x}) \cdot p(\mathbf{x})}{p(\mathbf{y})} \propto p(\mathbf{y} | \mathbf{x}) \cdot p(\mathbf{x})$$



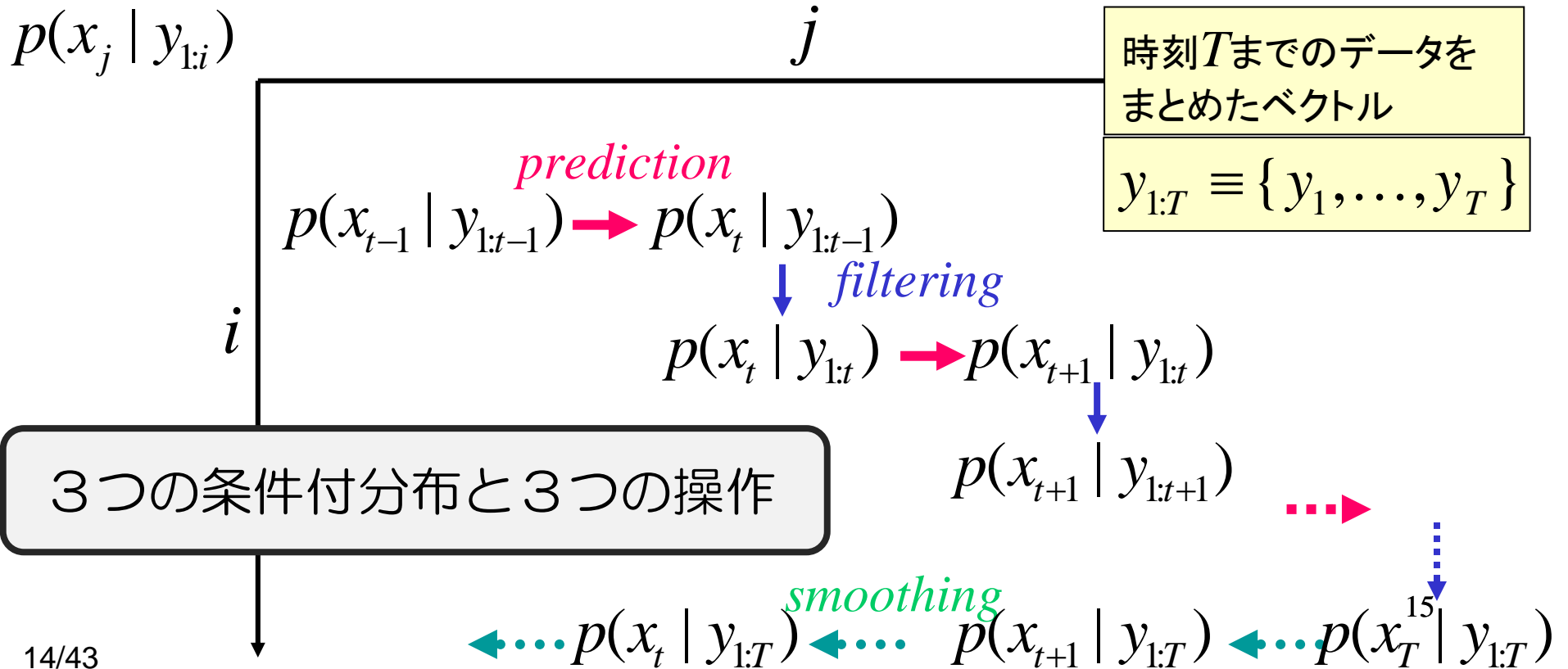
\mathbf{x} の空間



周辺分布と同時分布

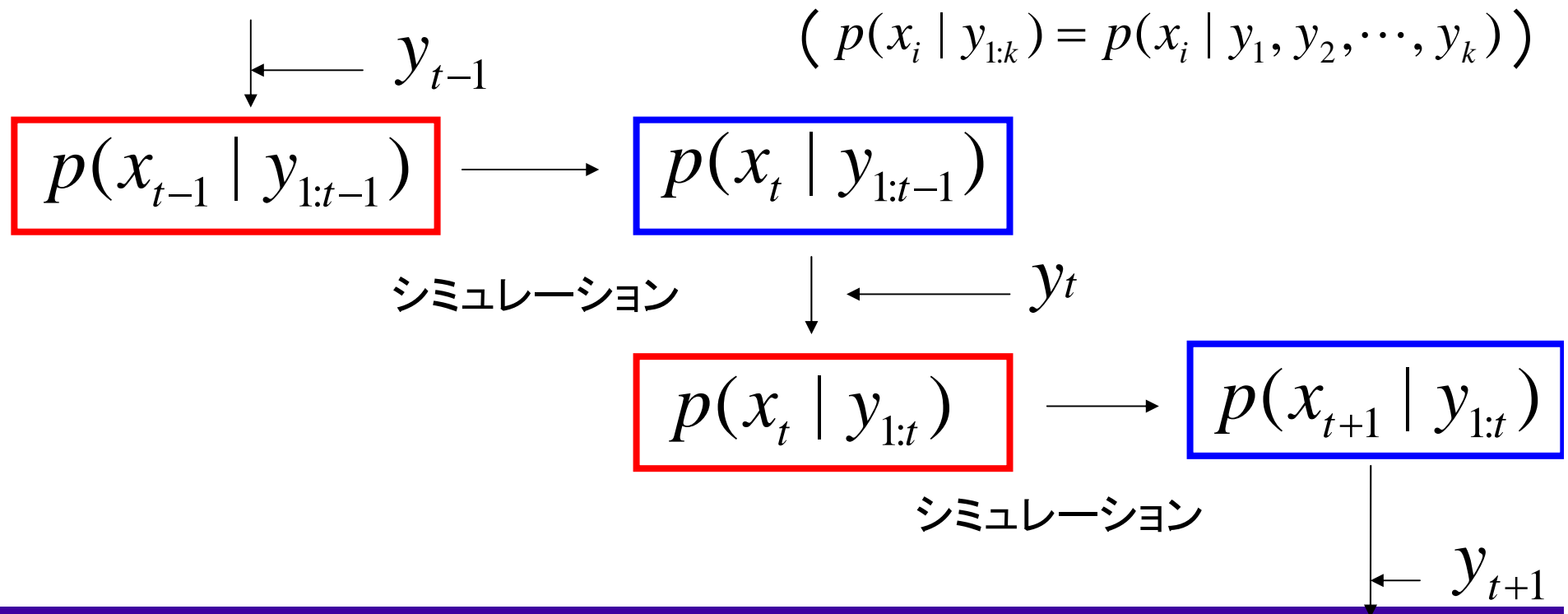


| | |
|---|-----------------------------------|
| <p>predictive density: $p(x_t y_{1:t-1})$</p> | <p>きのうまでのデータに基づく今日の状態</p> |
| <p>filter density: $p(x_t y_{1:t})$</p> <p>フィルタ分布</p> | <p>今日までのデータに基づく今日の状態</p> |
| <p>smoother density: $p(x_t y_{1:T})$</p> <p>平滑化分布</p> | <p>数年後、データをすべて得たもとで振り返った今日の状態</p> |



逐次データ同化のアルゴリズム

逐次データ同化では観測を得るたびに確率変数 x_t の分布または値の推定を行う



アンサンブルベース逐次データ同化の手法

page1

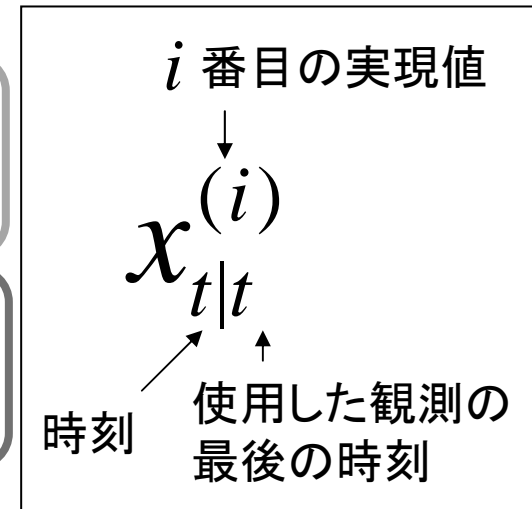
条件付分布(cPDF)の近似を実現値の集合(アンサンブル)として保持

$$p(x_t | y_{1:t-1}) \cong \frac{1}{N} \sum_{i=1}^N \delta(x_t - x_{t|t-1}^{(i)})$$

$$\left\{ x_{t|t-1}^{(i)} \right\}_{i=1}^N$$

$$p(x_t | y_{1:t}) \cong \frac{1}{N} \sum_{i=1}^N \delta(x_t - x_{t|t}^{(i)})$$

$$\left\{ x_{t|t}^{(i)} \right\}_{i=1}^N$$



○ 粒子フィルタ (PF: Particle Filter)

- ・理論的にはアンサンブルはアンサンブルメンバー数無限大で条件付分布の理論値に収束
- ・有限粒子数では退化の問題をどう回避するかがカギ

○ アンサンブルカルマンフィルタ (EnKF: Ensemble Kalman Filter)

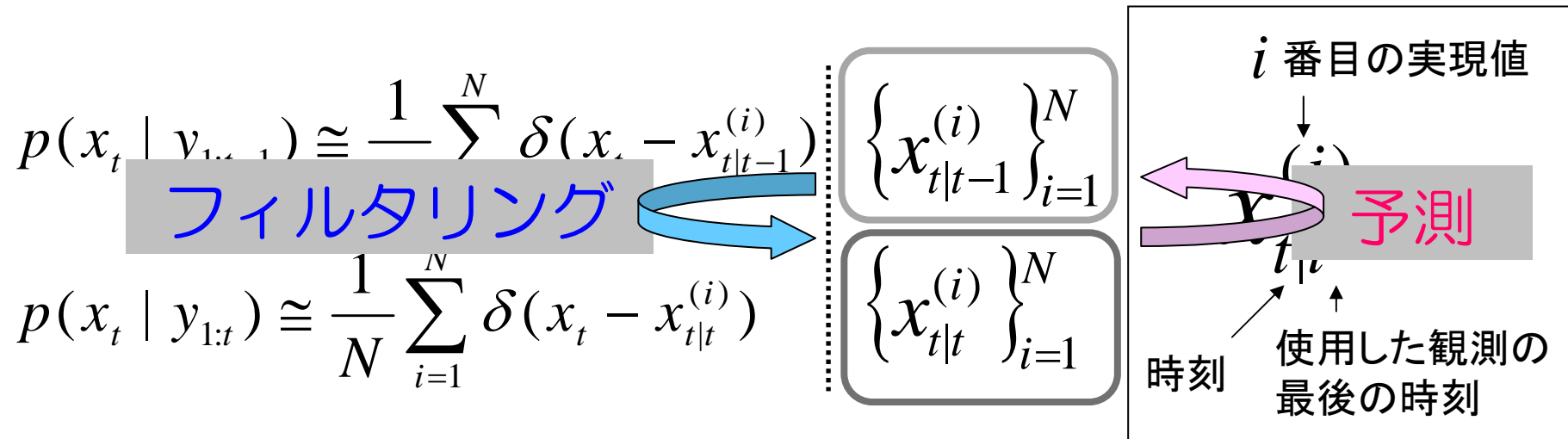
- ・cPDFの理論値には収束しない。アンサンブルはcPDFとは似ていて別物。
ただし、2次モーメントまでを近似。
- ・退化の問題は起こらない。

Swarm Filter と呼ぶべき

アンサンブルベース逐次データ同化の手法

page2

条件付分布(cPDF)の近似を実現値の集合(アンサンブル)として保持



○ 粒子フィルタ (PF: Particle Filter)

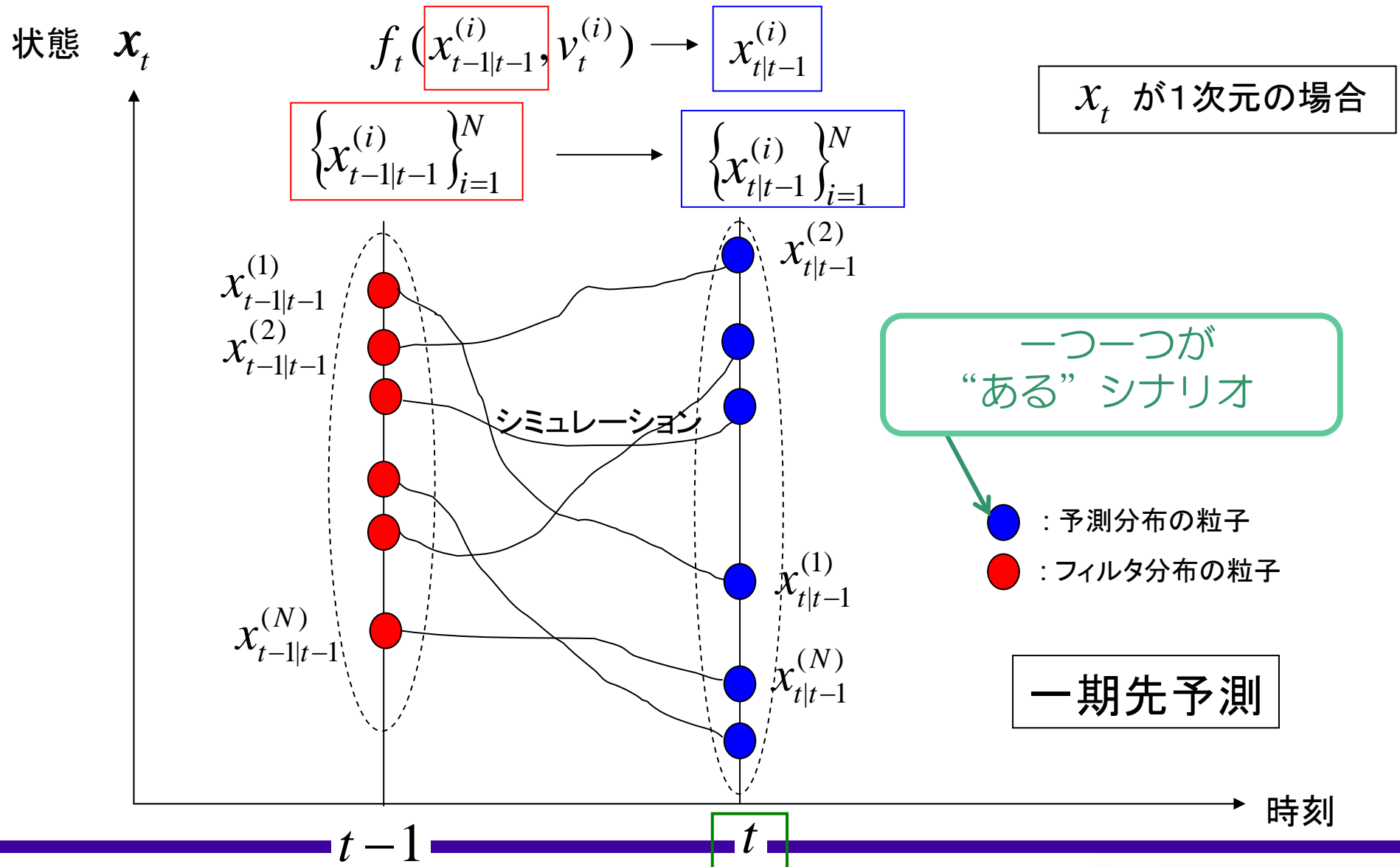
- ・理論的にはアンサンブルはアンサンブルメンバー数無限大で条件付分布の理論値に収束
- ・有限粒子数では退化の問題をどう回避するかがカギ

○ アンサンブルカルマンフィルタ (EnKF: Ensemble Kalman Filter)

- ・cPDFの理論値には収束しない。アンサンブルはcPDFとは似ていて別物。
ただし、2次モーメントまでを近似。
- ・退化の問題は起こらない。

Swarm Filter と呼ぶべき

EnKFとPFにおける一期先予測



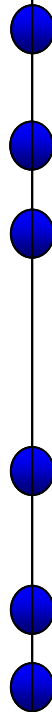
PFにおけるフィルタリング

page1

状態 x_t



$$\left\{ x_{t|t-1}^{(i)} \right\}_{i=1}^N$$



t

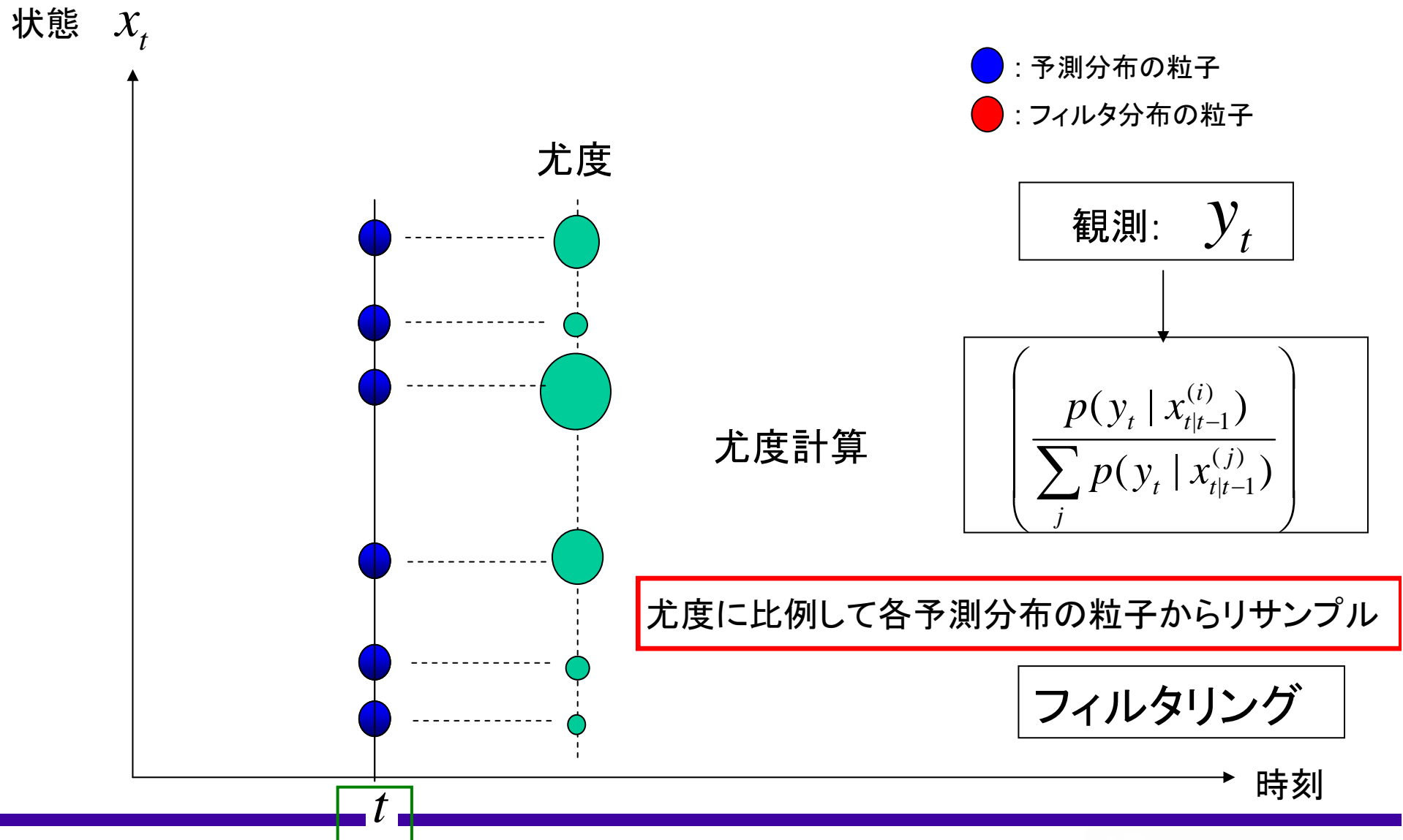
- : 予測分布の粒子
- : フィルタ分布の粒子

観測: y_t

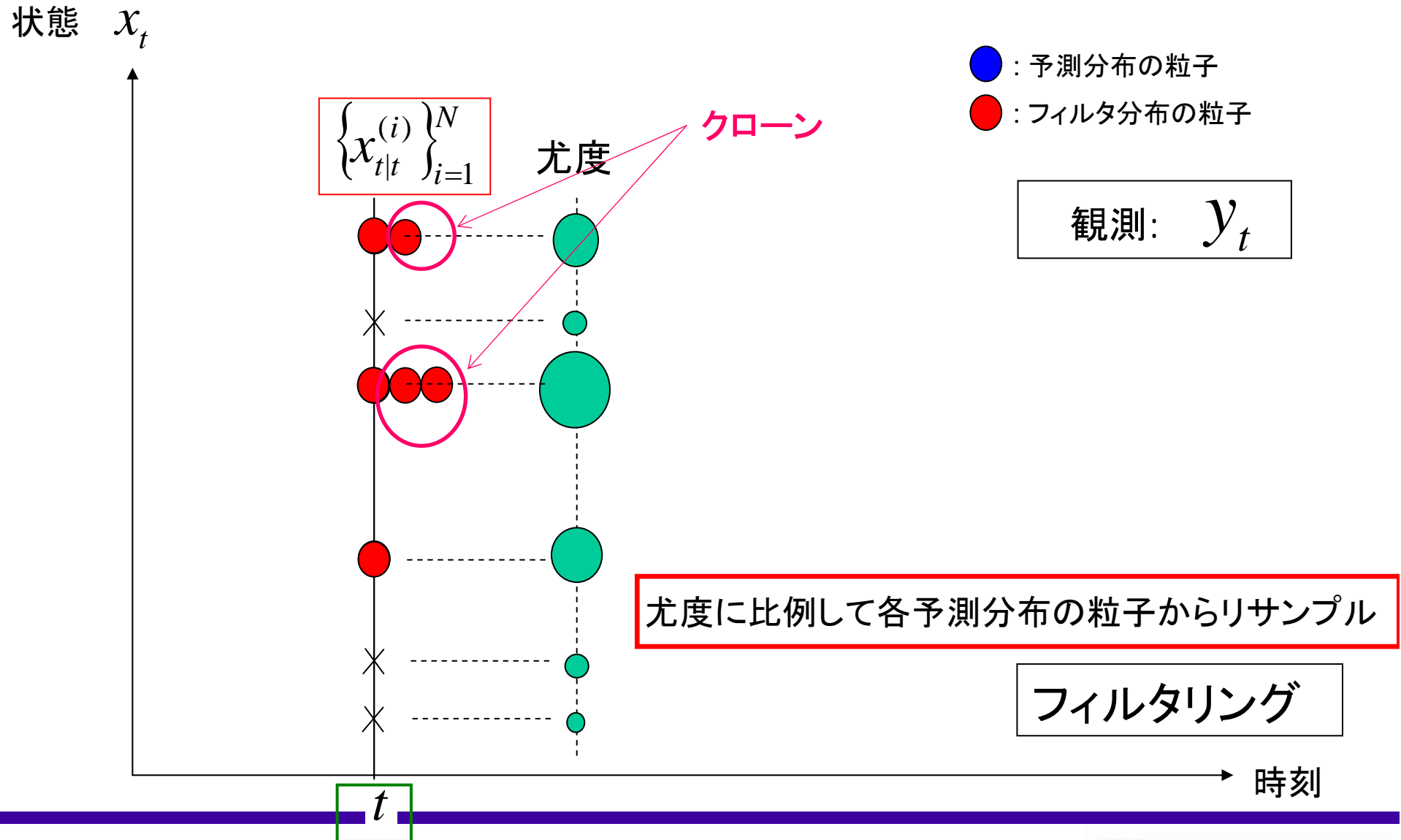
フィルタリング

時刻

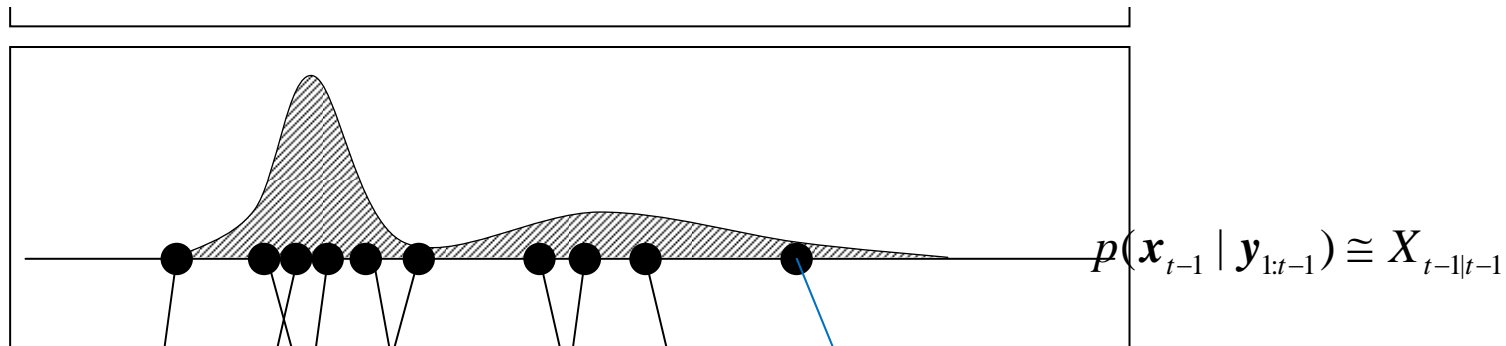
PFにおけるフィルタリング



PFにおけるフィルタリング



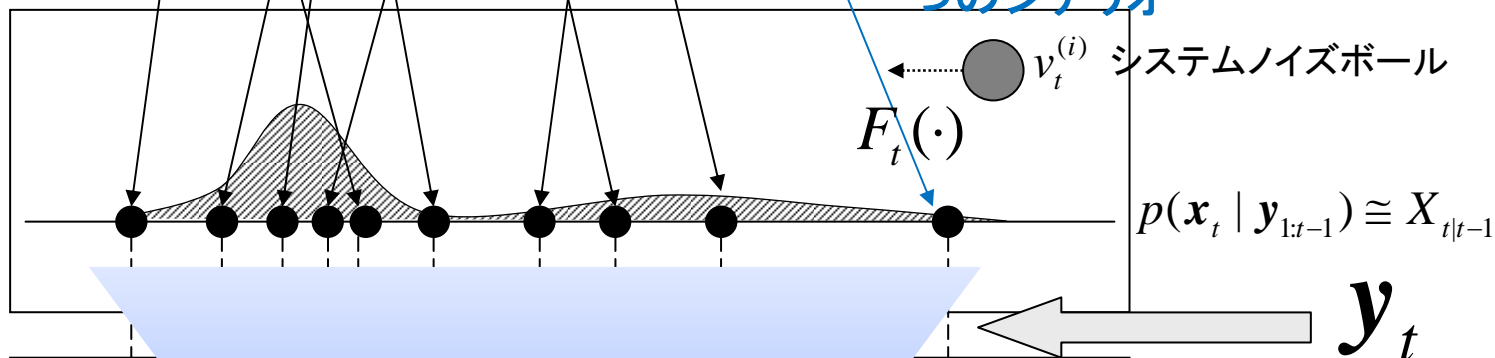
時刻 $t-1$



一つのシナリオ

$v_t^{(i)}$ システムノイズボール
 $F_t(\cdot)$

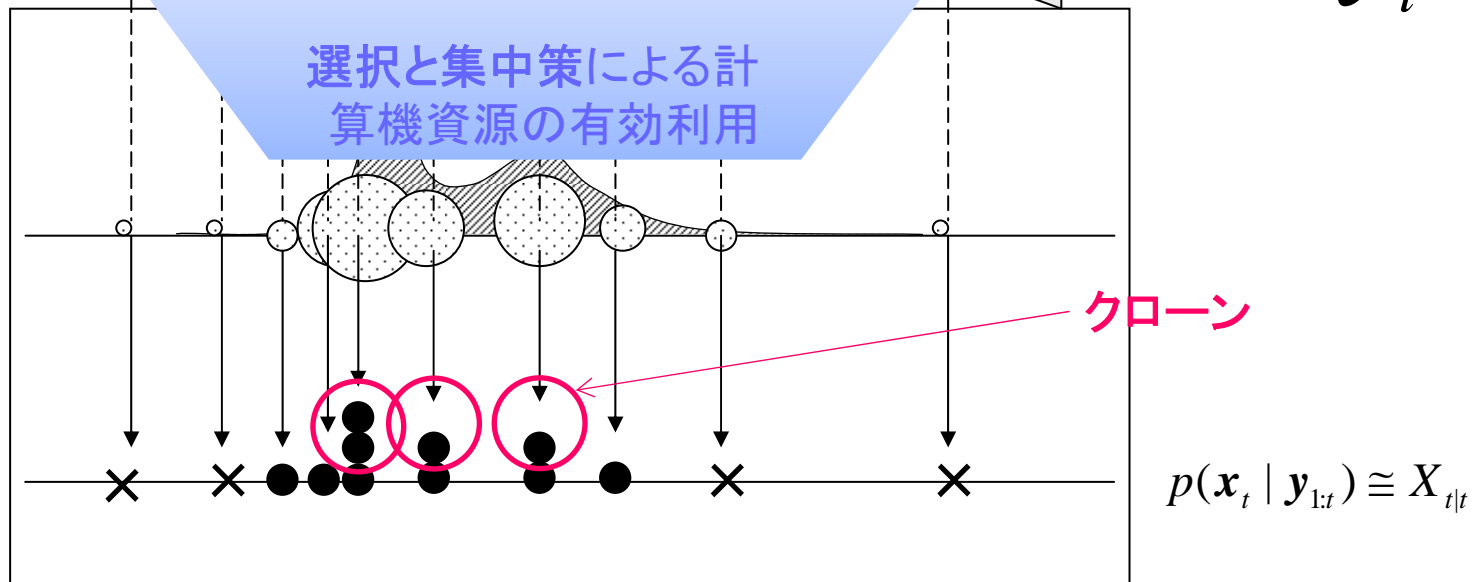
時刻 t



選択と集中策による計算機資源の有効利用

クローン

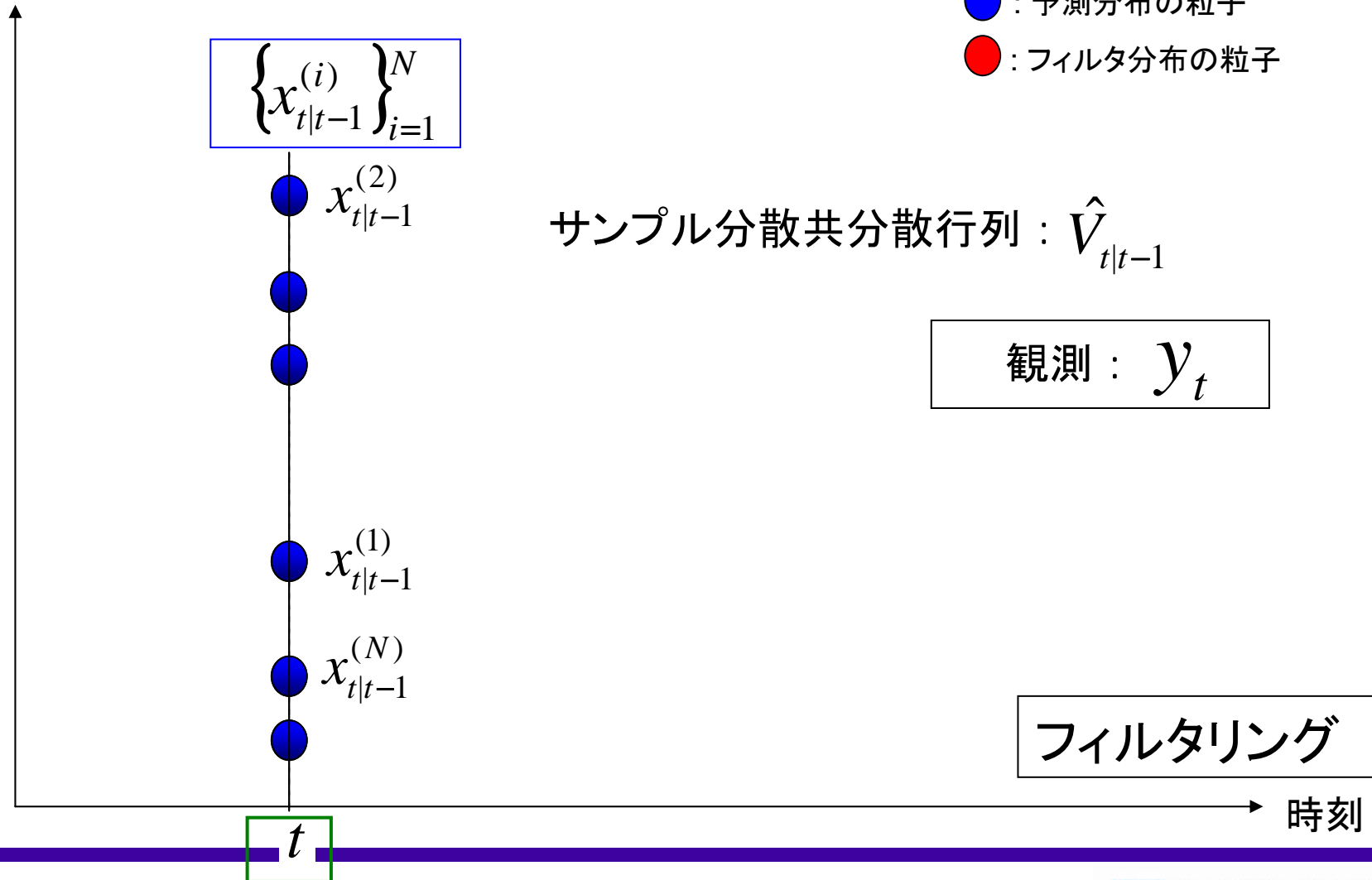
時刻 $t+1$



EnKFにおけるフィルタリング

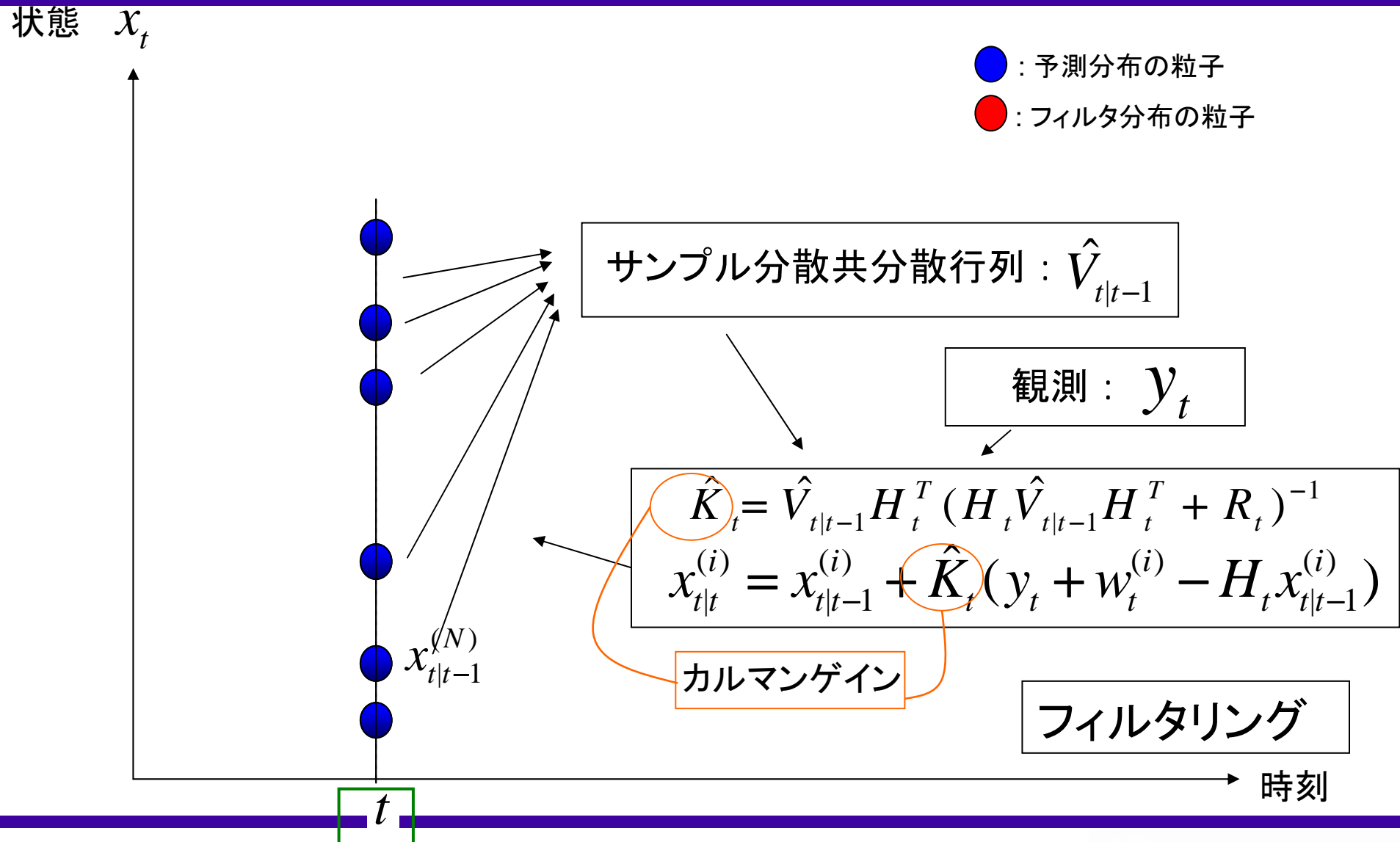
page1

状態 x_t



EnKFにおけるフィルタリング

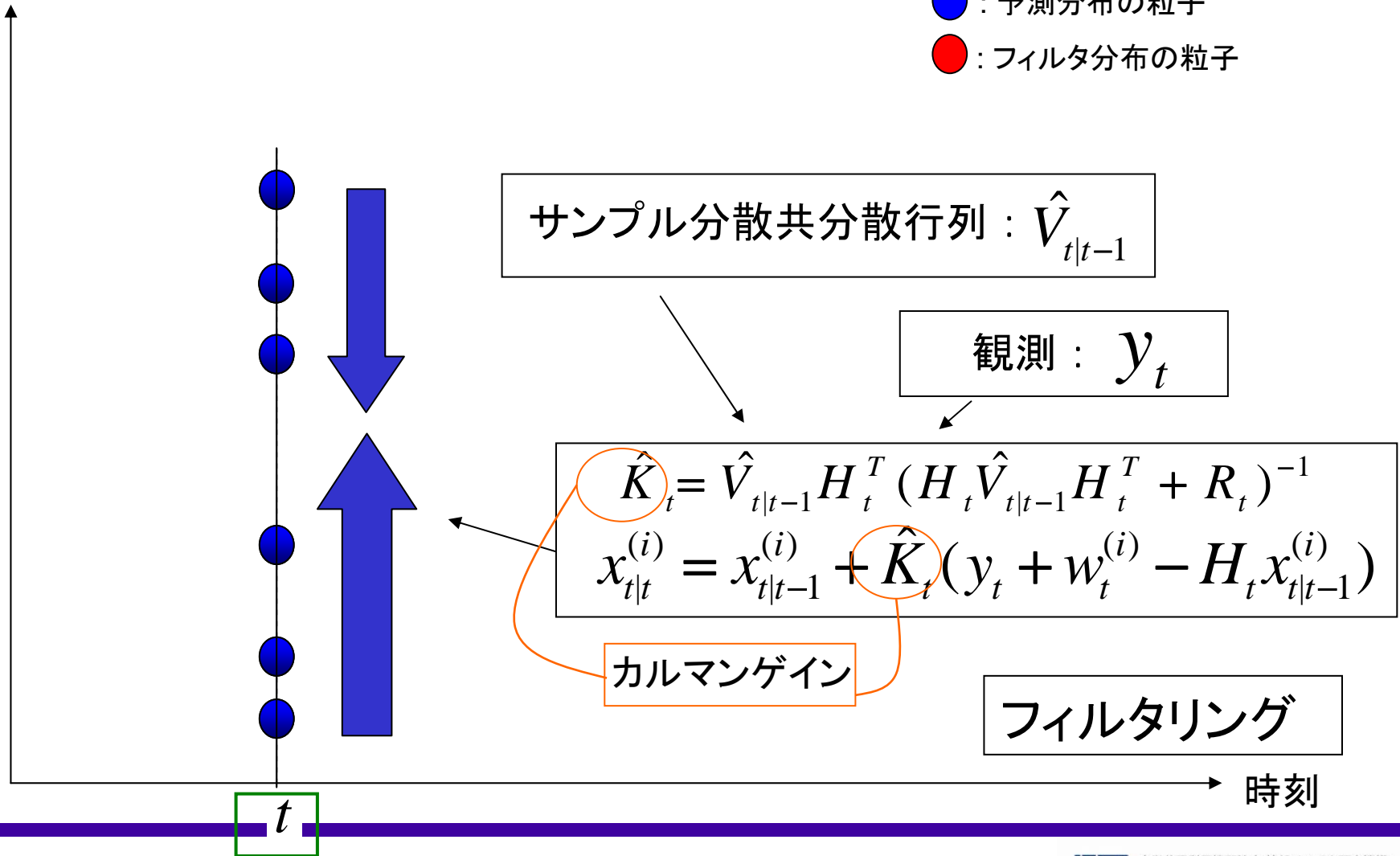
page2



EnKFにおけるフィルタリング

状態 x_t

- : 予測分布の粒子
- : フィルタ分布の粒子



EnKFにおけるフィルタリング

状態 x_t

- : 予測分布の粒子
- : フィルタ分布の粒子

$$\{x_{t|t}^{(i)}\}_{i=1}^N$$

サンプル分散共分散行列 : $\hat{V}_{t|t-1}$

観測 : y_t

$$\hat{K}_t = \hat{V}_{t|t-1} H_t^T (H_t \hat{V}_{t|t-1} H_t^T + R_t)^{-1}$$
$$x_{t|t}^{(i)} = x_{t|t-1}^{(i)} + \hat{K}_t (y_t + w_t^{(i)} - H_t x_{t|t-1}^{(i)})$$

カルマンゲイン

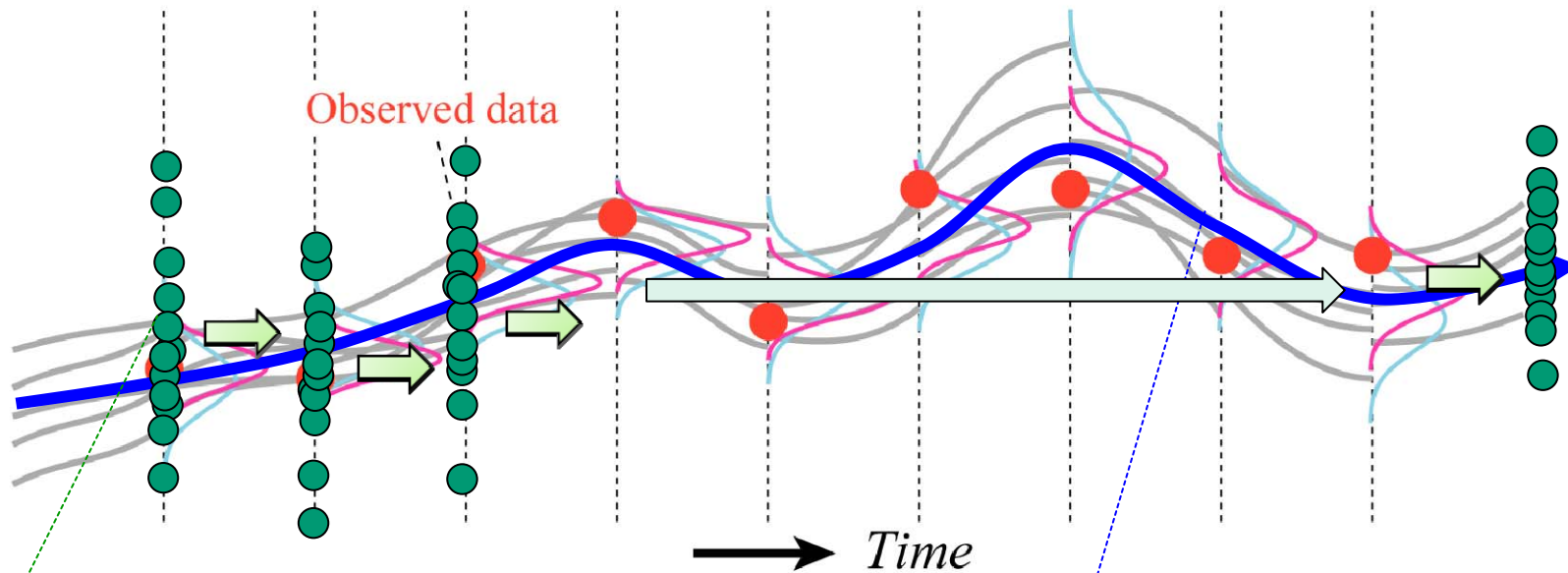
フィルタリング

時刻

t

逐次 vs. 非逐次

再掲: データ同化のイメージ



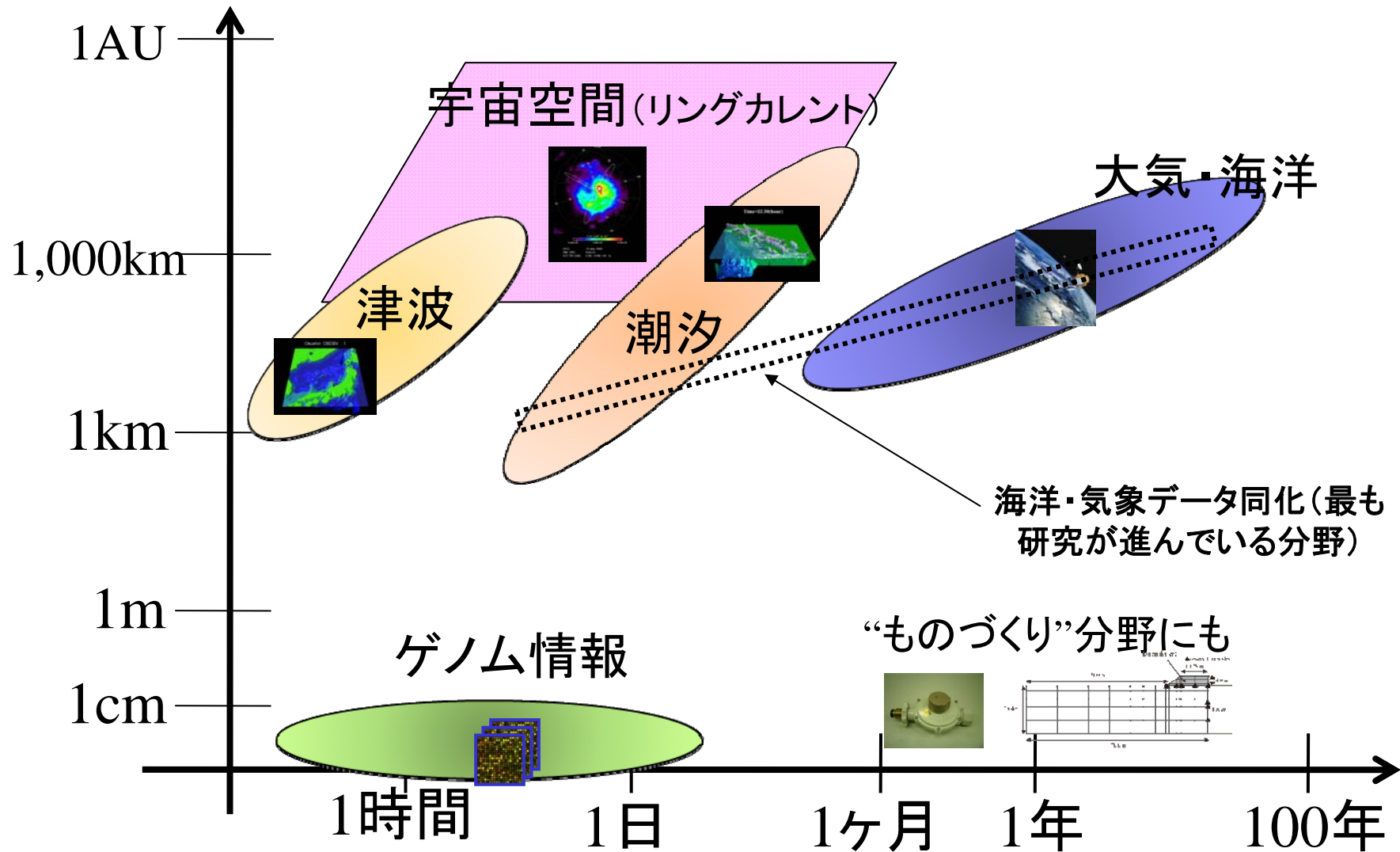
逐次(オンライン)型: 集団の時間発展を追う。つまり、Swarm Filter

代表例: EnKF

非逐次(オフライン)型: ベストなパスを求める

代表例: 4次元変分法(Adjoint法)

時間スケール・空間スケール



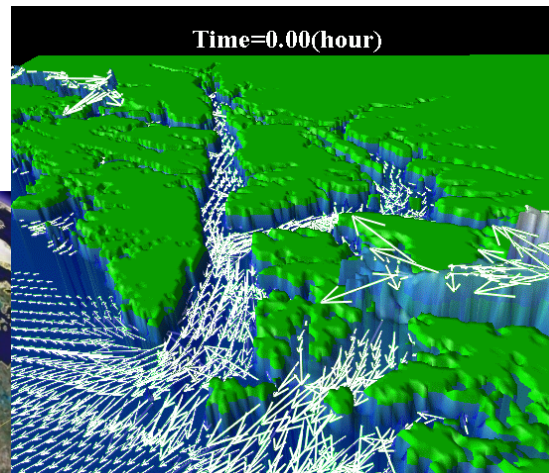
“個”にマッチしたシミュレーション: 境界条件の設定機能をパーソナライズする

“個”によって異なる形状, 形態情報をシミュレーション
モデルに取り込む 『メタシミュレーションモデル』

運動方程式: $\mathbf{v} \frac{\partial \mathbf{v}}{\partial t} + (\mathbf{v} \cdot \nabla) \mathbf{v} + \mathbf{f} \times \mathbf{v} = -g \nabla \eta - \underbrace{\gamma_b}_{\text{海底摩擦係数 (地域依存性)}} \frac{\mathbf{v}|\mathbf{v}|}{\underbrace{H}_{\text{水深}}} + A_H \nabla^2 \mathbf{v}$

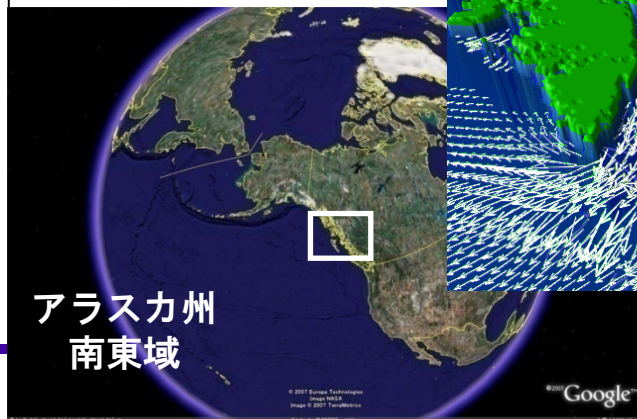
連続式: $\frac{\partial \eta}{\partial t} + \nabla \cdot (\mathbf{v}H) = 0$

我々研究チームによる,
潮汐シミュレーションの例



\mathbf{v} : 水平(2次元)流速ベクトル
 η : 海面水位
 H : 水深, \mathbf{f} : コリオリパラメータ

多品種少量生産を基本とする製品開発現場でのステップの簡略化や、患者一人一人に合った治療サービスの提供と期間の短縮化

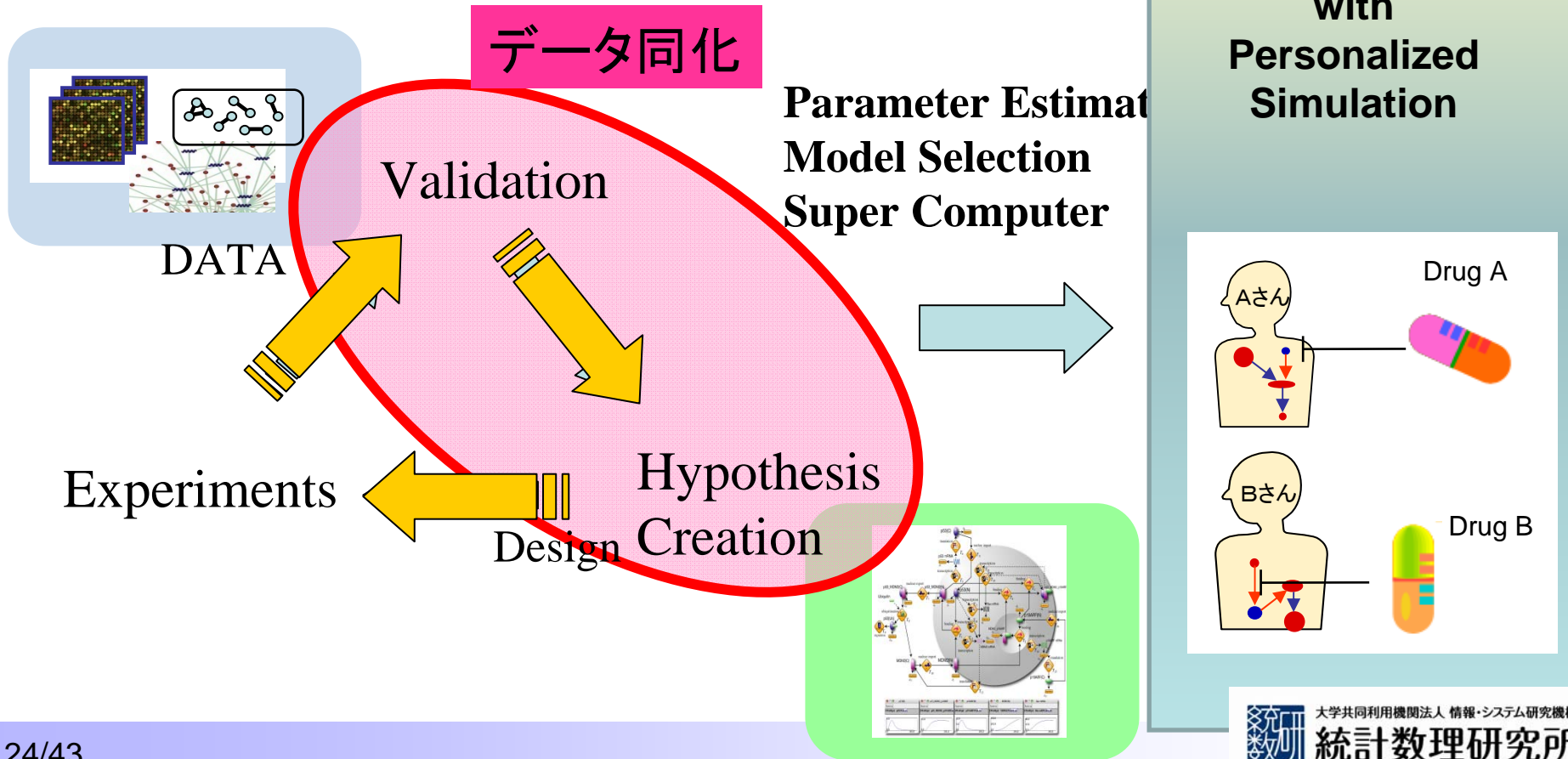


アラスカ州
南東域

オーダーメイド医療・創薬に向けたPersonalized Simulation

データ同化

- 生体システムを統合的に理解・予測するための有望な方法論
- 実験科学と統計・シミュレーション科学の融合





Life Science Data Assimilation System

Computational Systems Biology in New Dimensions

Molecular biology has now become a data-driven science with technological innovations of experimental biology and rapid increase in available genomic information. Over the past decades, a range of data science technologies has been developed in bioinformatics and computational systems biology, especially in order to unveil a complex world of biochemical reaction networks in living organisms. Life Science Data Assimilation System (LiSDAS) is an integrated computing platform which embodies ‘In Silico Statistical Analysis’ of biochemical networks. Our technology brings a range of data science in molecular biology and medical science, involving biosimulations coupled with Data Assimilation technology, dynamic system analyses of biochemical reaction networks with genomic data on diverse scales from the level of molecules to -omics.

計算機内で生化学反応系の仮説モデルを複数生成し、パラメータ・ネットワーク構造のショットガンサーチを行い、高品質の「イン・シリコモデルを創る」ためのデータ同化システム。生化学反応系シミュレータと実験生物学を融合するための新しい「データ駆動型イン・シリコ解析」アプリケーションを開発。

Computer-Driven Production of BioSimulators

代謝反応経路，転写因子ネットワーク，シグナル伝達系の生化学反応シミュレータの設計・開発支援ツール。ユーザーによって設計された基本モデルを起点に，タンパク質相互作用や転写因子データベースに登録されている(仮説)反応経路を取り込み，生化学反応系を大量に生成する。SBMLやCSMLによる反応表現をサポート。

Model Checking and Analysis of Network Dynamics

生化学反応モデルの動的特性を解析するためのプログラム群。反応物質の擬似ノックダウン実験やモデルの確率的摂動によるロバスト性解析，反応パラメータと生化学反応ダイナミクスのマッピングなど，イン・シリコシステム解析を行うための計算環境を提供。

Shotgun Stochastic Search for High-Performance Models via Data Assimilation

SBMLやCSMLで表現された大量の生化学反応系から，実験データや既存の生化学的知見に整合的な仮説モデルをショットガンサーチで列挙するためのデータ同化プログラム：

- (1) シミュレータと実験データからモデルの最適パラメータを計算するための大規模ベイズ逆問題の数値解法（粒子フィルタ，アニーリングのハイブリッド大規模計算システム）
- (2) ベイズ型モデル評価規準にもとづく，候補反応モデルのスクリーニング

統計数理研究所データ同化Gが利用している並列計算機

| System | CPU | Clock frequency | # of nodes | # of cores | Memory |
|-----------------------------|---------------------|-----------------|-------------------|--------------------|-----------|
| ismrx (ISM) | Intel Xeon X5570 | 2.93 GHz | 360 nodes | 2,880 cores | 32GB/node |
| pleiades (Higuchi Lab) | Intel Xeon E5440 | 2.83 GHz | 24 nodes | 192 cores | 32GB/node |
| sheep (Higuchi Lab) | Intel Xeon E5550 | 2.66 GHz | 13 nodes | 104 cores | 24GB/node |
| Type-B (HGC/Univ. Tokyo) | Intel Xeon E5450 | 3.00 GHz | 768 nodes | 6,144 cores | 32GB/node |
| RICC (RIKEN) | Intel Xeon X5570 | 2.93 GHz | 1,024 nodes | 8,192 cores | 12GB/node |
| Kei* (RIKEN) | SPARC64 VIIIfx | 2.00 GHz | > 80,000 nodes | > 640,000 cores | 16GB/node |

JST/GREST

Fujitsu

* the Next-Generation Supercomputer

ClockTime: Cloud Computing Kernel for Time-series Modeling Engine



multivariate analysis データ同化のアプリではない

http://sheep.ism.ac.jp/CloCK-TIME/service/index.html - Windows Internet Explorer

http://sheep.ism.ac.jp/CloCK-TIME/service/index.html

CloCK-TIME

Cloud Computing Kernel for Time-series Modelling Engine

Home TimeSeries

Table of Contents

- Overview
- Parameter Settings
 - 2.1 Data
 - 2.2 Analysis Mode
 - 2.3 Observation Model
 - 2.4 Trend Component
 - 2.5 Seasonal Component
 - 2.6 AR Component
 - 2.7 Number of Hyper Parameters
- Analysis Results

1. Overview

"Cloud Computing Kernel for Time-series Modelling Engine" (CloCK-TIME) is an online analyzer for multivariate time-series data, which runs on a PC cluster in a cloud computing system. CloCK-TIME decomposes given multivariate time-series data into trend, seasonal, autoregressive (AR), and observation noise components by a hybrid method consisting of the particle filter and the Markov Chain Monte Carlo (MCMC). All of functions such as data uploading and parameter settings can be controlled interactively through a user interface.

2. Parameter Settings

2.1 Data

2.1.1 Data Upload

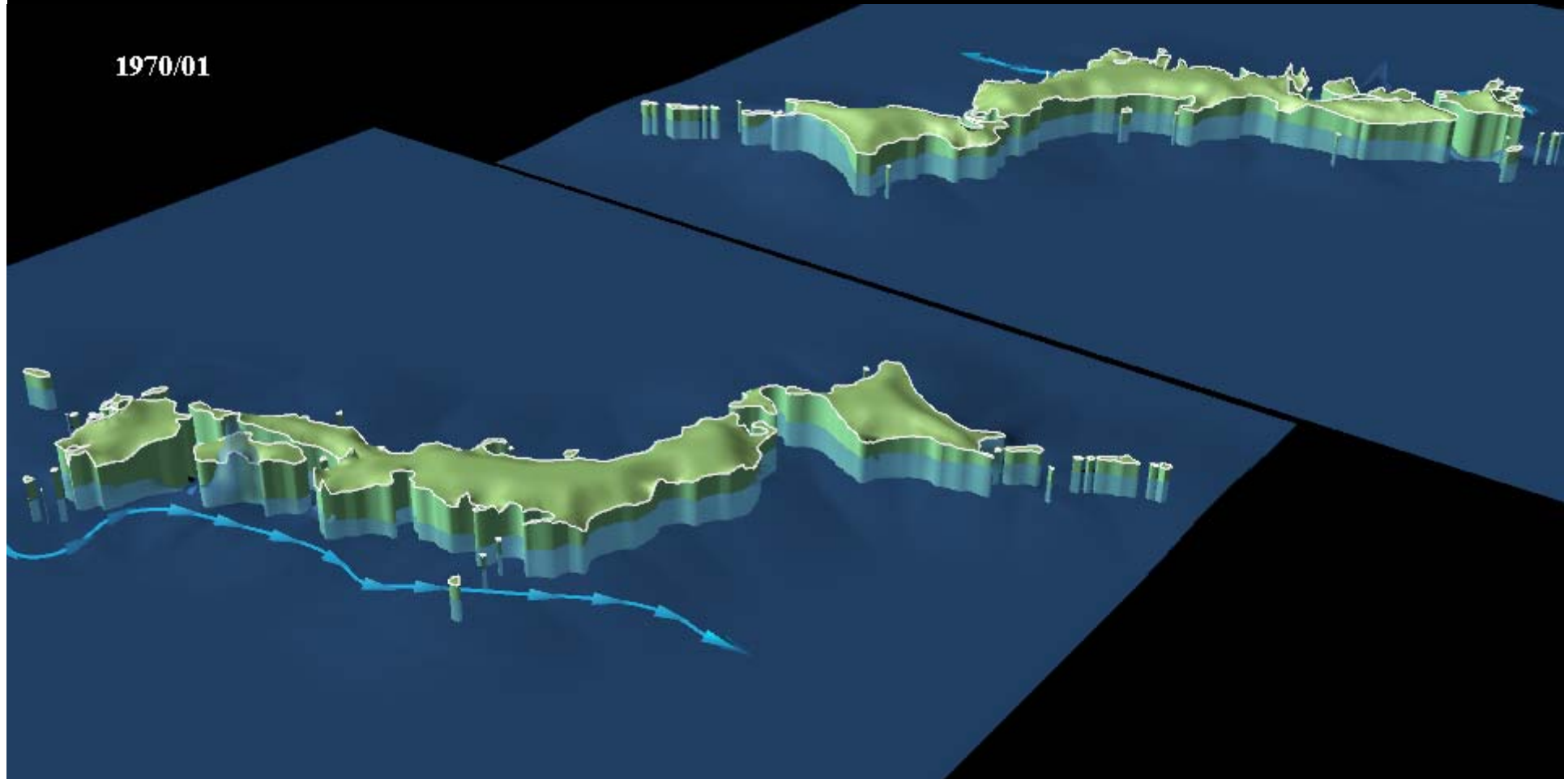
Upload your data file to the server. The data should be in a matrix form, in which each column corresponds to a time series of each channel. Each data in a row should be separated with space, tab, or comma. Some comments to the server can be described in the header.

| | A | B | C | D | E |
|----|---------------------------------------|------|------|---|---|
| 1 | # Sample input data (dimension = 3) | | | | |
| 2 | #@MissValue = -99999 | | | | |
| 3 | #@Nbin = 10 | | | | |
| 4 | 810 | 1851 | 1592 | | |
| 5 | 730 | 1783 | 1547 | | |
| 6 | 729 | 1767 | 1543 | | |
| 7 | 766 | 1780 | 1557 | | |
| 8 | 753 | 1778 | 1546 | | |
| 9 | 765 | 1805 | 1577 | | |
| 10 | 837 | 1854 | 1586 | | |
| 11 | 943 | 1969 | 1658 | | |
| 12 | 977 | 1997 | 1715 | | |
| 13 | 968 | 1990 | 1729 | | |
| 14 | 878 | 1944 | 1687 | | |
| 15 | 871 | 1927 | 1663 | | |
| 16 | 830 | 1888 | 1623 | | |
| 17 | 775 | 1797 | 1532 | | |
| 18 | 724 | 1743 | 1481 | | |
| 19 | 756 | 1756 | 1492 | | |
| 20 | 786 | 1792 | 1508 | | |
| 21 | 801 | 1851 | 1569 | | |
| 22 | 884 | 1915 | 1628 | | |
| 23 | 932 | 1969 | 1646 | | |
| 24 | 1003 | 2016 | 1760 | | |
| 25 | 951 | 2000 | 1741 | | |
| 26 | 892 | 1944 | 1692 | | |
| 27 | 853 | 1893 | 1634 | | |
| 28 | 822 | 1875 | 1599 | | |
| 29 | 760 | 1809 | 1565 | | |
| 30 | 705 | 1787 | 1494 | | |

data file

潮位月別データの分解

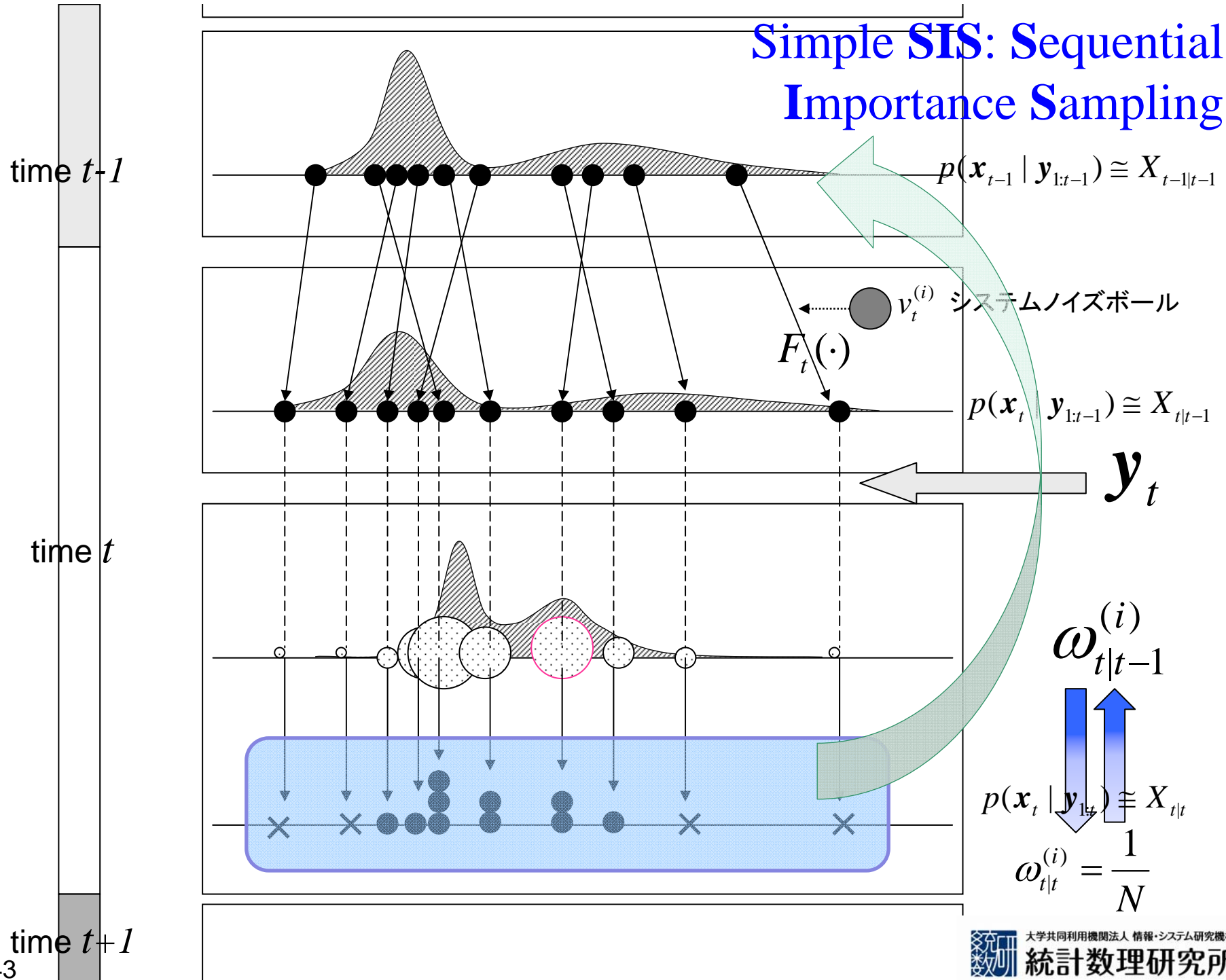
1970/01



赤い部分は、1970年からの海岸線の沈降量を示す。**青い→**は、黒潮の流れを示す。

(1)太平洋側、特に東海・東南海地域において最近は大震災が起こっておらず、プレート沈み込みによる応力エネルギーがどんどんと蓄積されている。(2)潮位の数年周期変動は、日本海側よりも太平洋側の方が大きく、かつ地域性がはっきりとしている。これは北太平洋における海洋大循環の影響と考えられている。

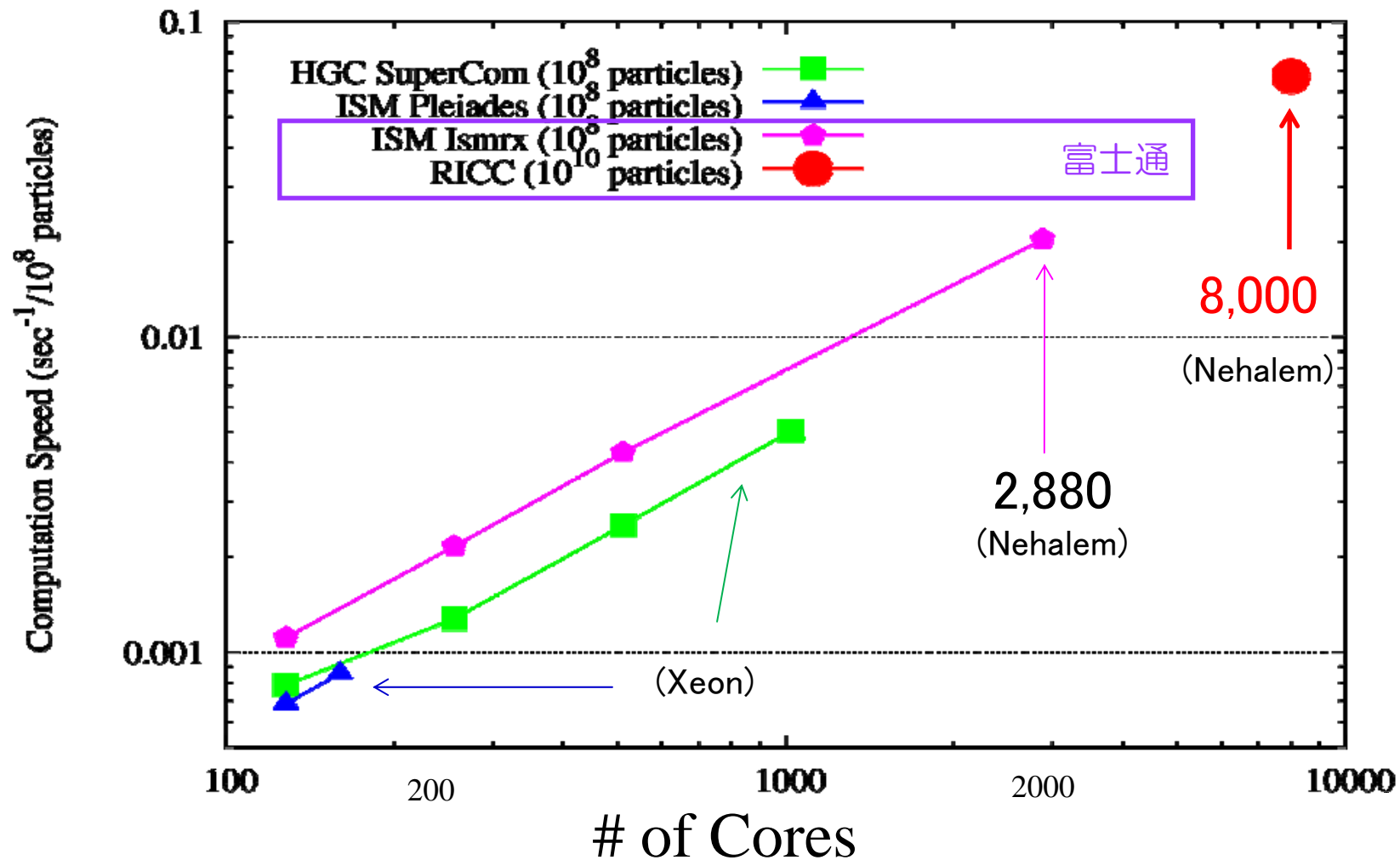
Simple SIS: Sequential Importance Sampling



Simple SISでスケーラビリティを確認



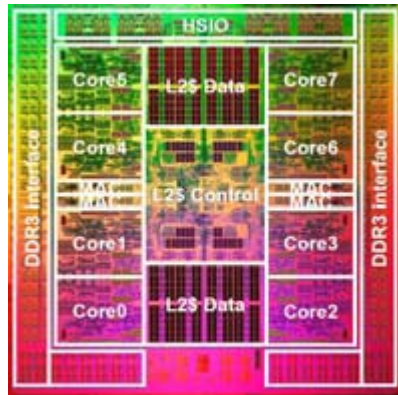
Computation Time for Biochemical Reaction Model of Mammalian Circadian Rhythm



並列計算機への粒子フィルタの実装

- リサンプリングの際に, node間通信が頻発.
- 任意の2 nodes間でほぼat randomに通信が発生し, 並列化しにくい.

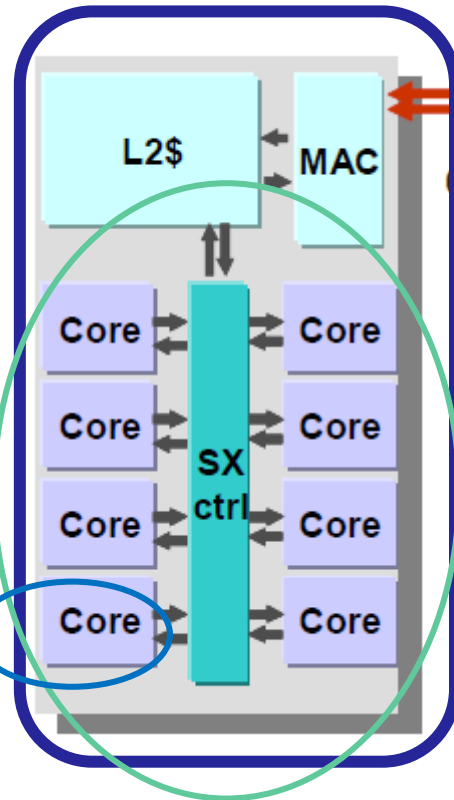
①階層構造をはじめから意識



SPARC64 VIIIfx

出典: 理研/富士通
のホームページから

[第一階層]
コア内: SIMD化(コンパイラが対応)

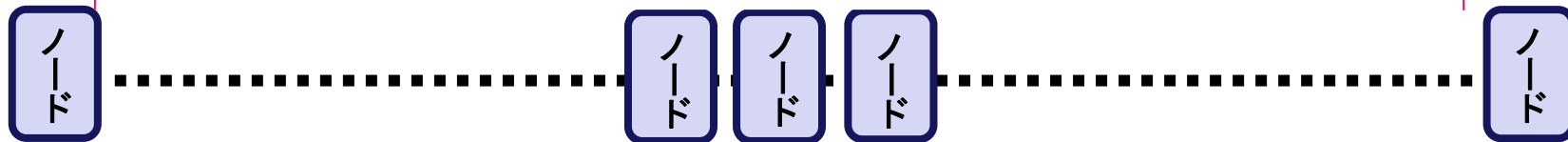


1ノード(CPU) = 8コア
16GFlops x 8 = 128 Gflops

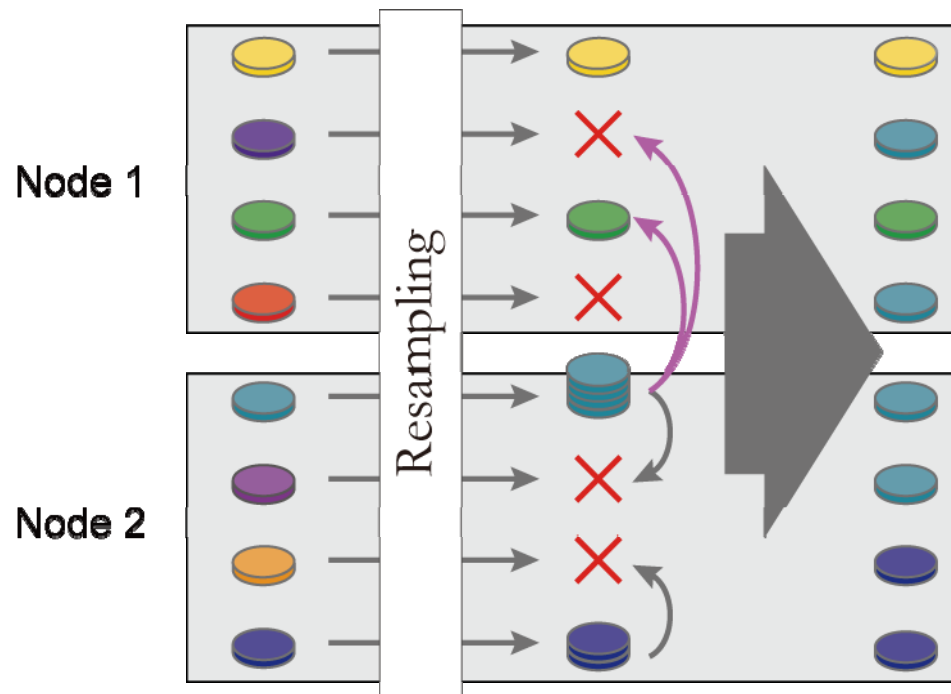
ノード数は8万以上

[第2階層]ノード内
OpenMP: スレッド並列

[第3階層] ノード間 MPI: プロセス並列

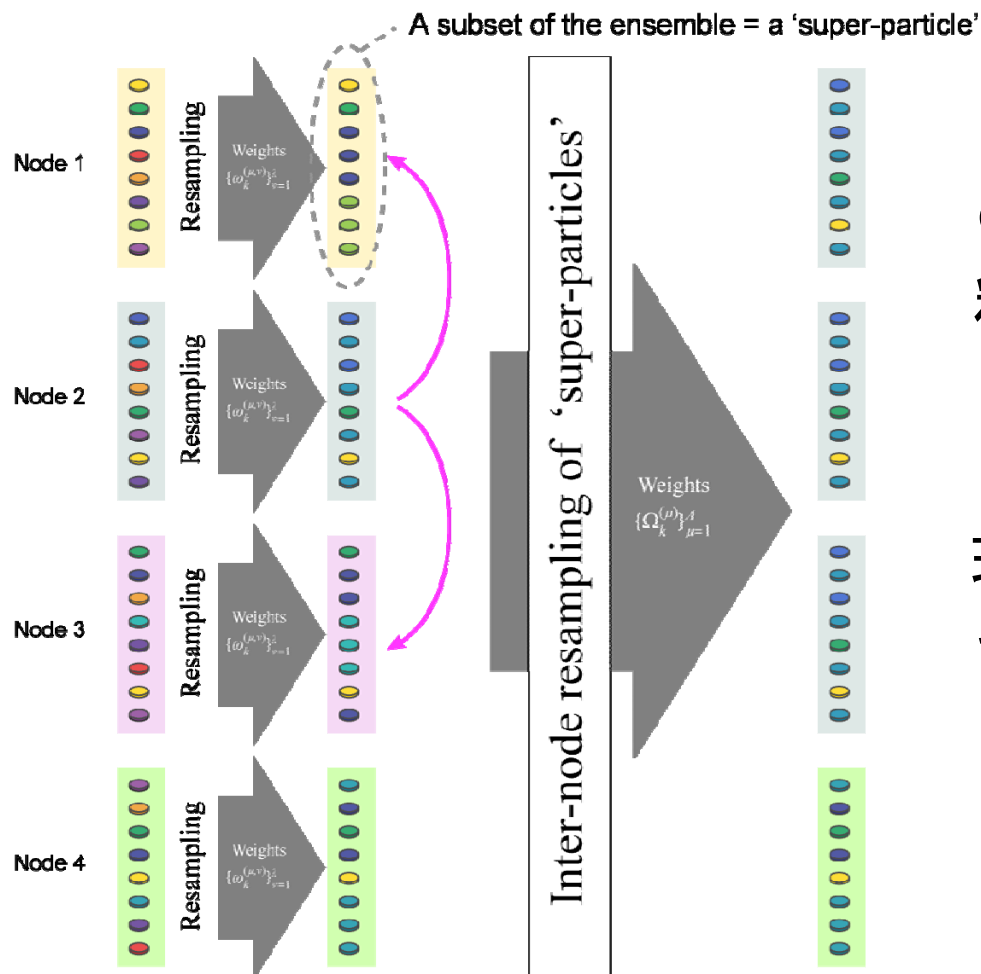


Flowchart of PF



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

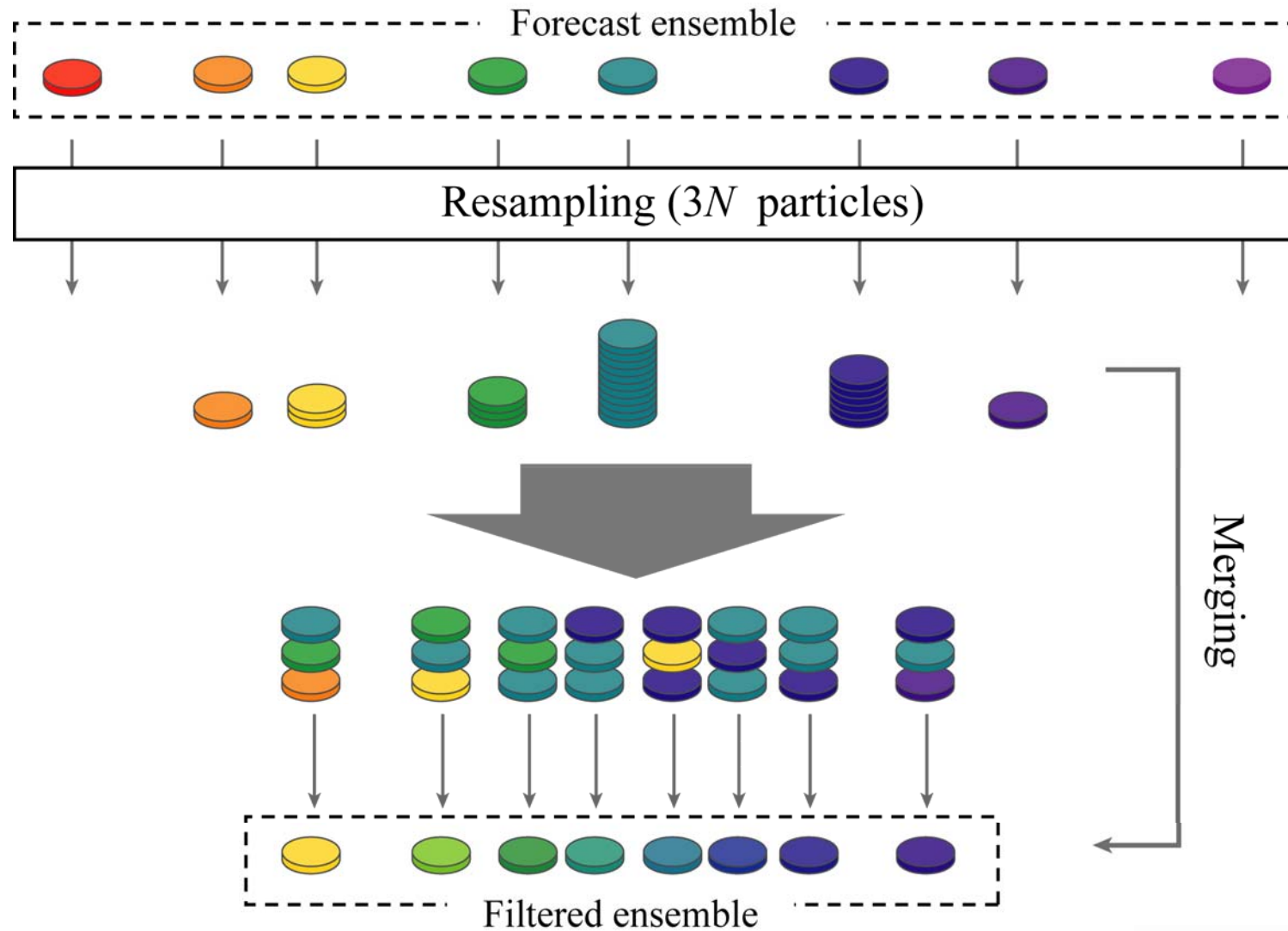
MetaPF: Meta-particle filter



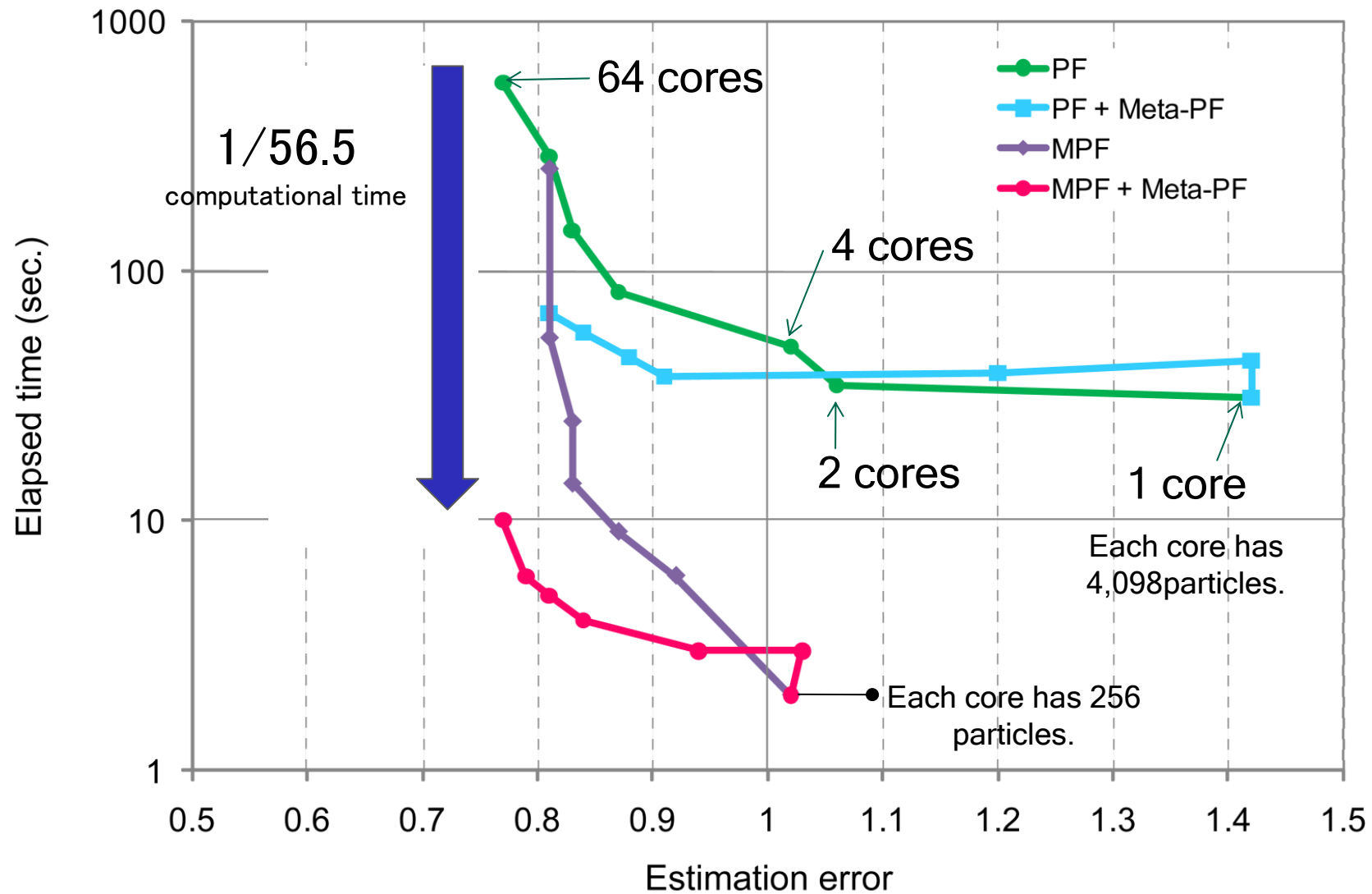
Nodeに割り当てられた ensemble の subset (=超粒子)をresamplingする.

現状では, nodeに割り当てられた重みがどれか1つでも0.3を超えたら resamplingすることになっている.

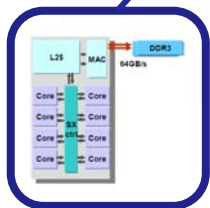
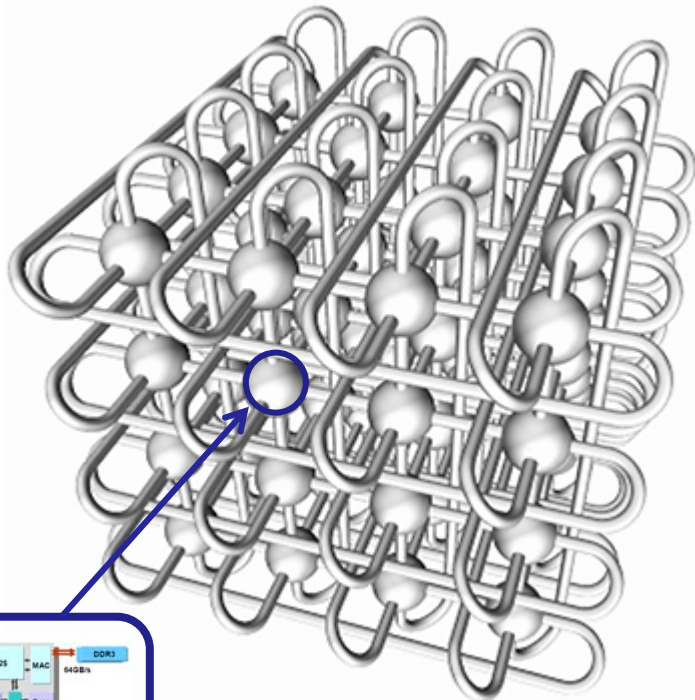
MPF: Merging particle filter



Result



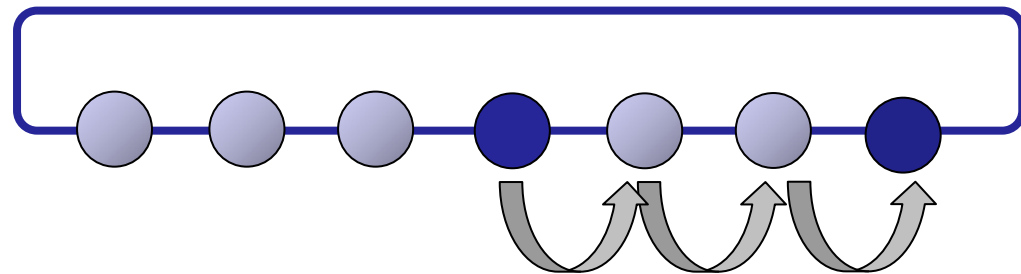
②ネットワーク構造をはじめから意識



出典: 理研のホームページから

IBM Sequoia (BlueGene/Q)
Cray XT5(Jaguar) も3次元トーラス

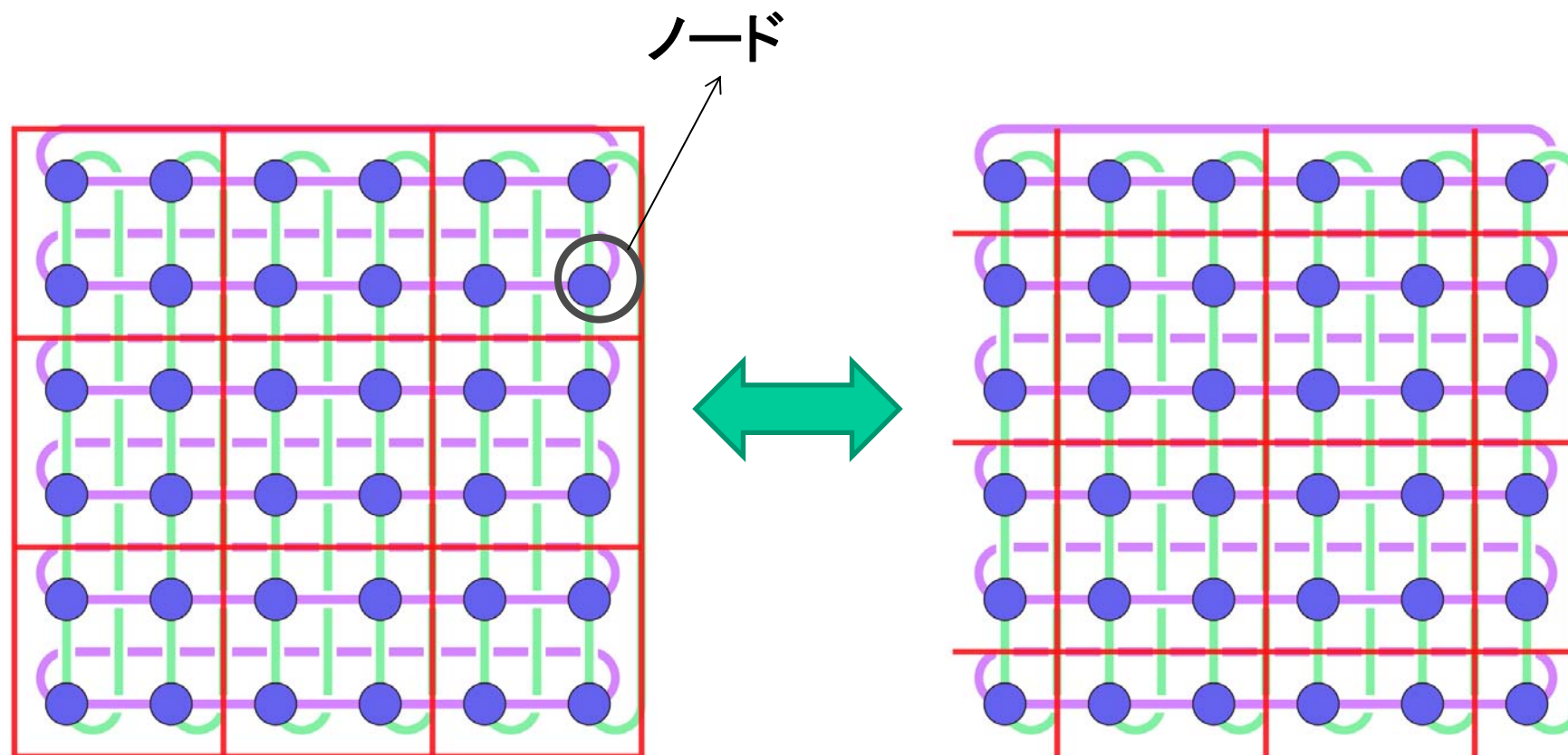
6次元メッシュトーラス: 3次元の直方体に配置したノードをそれぞれ6方向で結合し、各次元がそれぞれリング状に結合されるネットワーク構成



Q. ネットワークについて、三次元トーラスで隣り合うノードどうしの通信が速いとのことですが、隣り合わないノード間の通信はその通信回数分遅くなるのでしょうか。

A. 隣り合わないノード間の通信は、基本的にはバケツリレー方式で行われるので、経由するノードの数が多ければ多いほど、通信時間(レイテンシ)は長くなります。

Alternate lattice-pattern switching

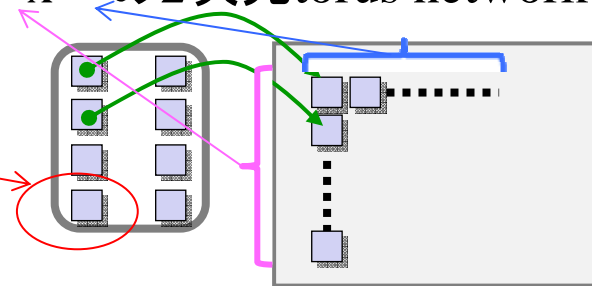


ここでは、上図のような2種類の格子パターンでノードをグループ化することを考え、2種類のパターンを毎回切り替えるものとする。

Alternate lattice-pattern switching

- 以下では，2種類のグループ分けパターンを交互に切り替え，各時点でのresamplingはその時の個々のグループの中で行う，という戦略を取る。
- 毎回グループの組み方を切り替えることで，GAのneighborhood modelのように，情報がうまく全体に行き渡り，退化が回避されるという効果も期待できる。
- 実際のnetwork topologyは気にせずとも，仮想的にノード間のつながりを定義すれば実装は可能なので，今回は，とりあえずプログラムを書いてみて，どの程度の精度が出るのかを調べる。
- また，今回は，3次元ではなく，* x * の2次元torus networkを設定して計算。

1コアが仮想的ノードに対応



実験結果 1 (ismrx:富士通)

Lorenz 96 model (Lorenz and Emanuel 1998)

$$\frac{dx_j}{dt} = (x_{j+1} - x_{j-2})x_{j-1} - x_j + 8 \quad \text{for } j = 1, \dots, 40$$

$(x_{-1} = x_{39}, x_0 = x_{40}, x_{41} = x_1.)$

| PEs | Particles | Global PF | | Alternative switching | |
|-------|-----------|-----------|---------------|-----------------------|---------------|
| | | RMSE | Elapsed times | RMSE | Elapsed times |
| 4x4 | 65536 | 0.85 | 01m12s | 0.84 | 00m40s |
| 6x6 | 147456 | 0.78 | 02m26s | 0.78 | 00m50s |
| 8x8 | 262144 | 0.77 | 03m55s | 0.77 | 00m54s |
| 10x10 | 409600 | 0.74 | 06m18s | 0.75 | 01m03s |
| 12x12 | 589824 | 0.74 | 08m48s | 0.73 | 01m11s |
| 16x16 | 1048576 | 0.73 | 17m18s | 0.73 | 01m45s |

- The experiments were performed on the PC cluster system (CPU: Xeon E5450 3.00GHz (Quad core) x2; Memory: 32GB).
- A virtual two-dimensional torus network is assumed. Each core in the real architecture is assigned to each node in the virtual torus network.

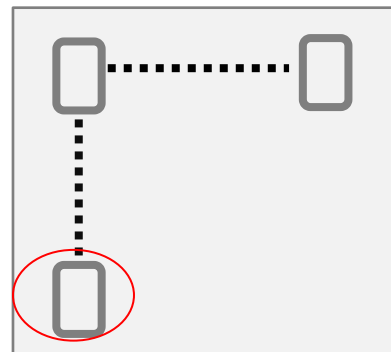
Alternate lattice-pattern switching

- 以下では、2種類のグループ分けパターンを交互に切り替え、各時点でのresamplingはその時の個々のグループの中で行う、という戦略を取る.
- 毎回グループの組み方を切り替えることで、GAのneighborhood modelのように、情報がうまく全体に行き渡り、退化が回避されるという効果も期待できる.

- Cray XT6m (1ノード=12コアx2ソケット、1~44ノード)で計算

1ノード当たりの粒子数を固定。従って、ノードを増やしても通信が発生しなければ、計算所要時間は一定。

ハード的に2次元トーラス



実験結果2 (統数研・Cray XT6m)

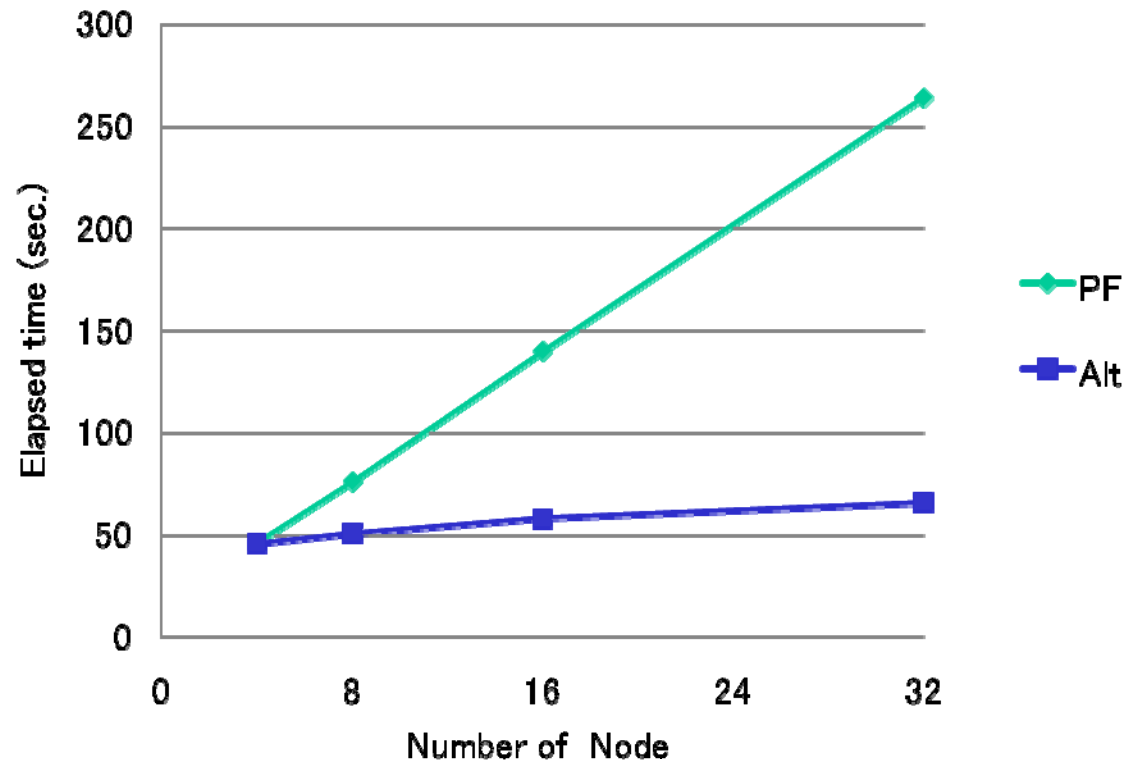
Lorenz 96 model (Lorenz and Emanuel 1998)

$$\frac{dx_j}{dt} = (x_{j+1} - x_{j-2})x_{j-1} - x_j + 8 \quad \text{for } j = 1, \dots, 40$$

$(x_{-1} = x_{39}, x_0 = x_{40}, x_{41} = x_1.)$

Cray XT6m の結果

ノード数が少ないので、ALの効果よりも、局所リサンプリングアルゴリズムの効果が大きい。



③I/Oまわり+ストリーミング計算が肝

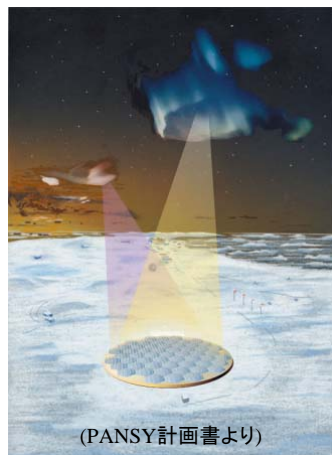
page1

あらゆる分野で発生する問題

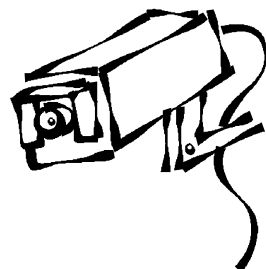
計測・観測現場



(イルミナ社HP)



(PANSY計画書より)



計算の現場

データマイニング、
機械学習、統計数理

+

スーパーコンピュータ



出典: 理研/富士通
のホームページから

データの次元が計算できる
サイズを超えている。

- ・次世代シーケンサーのデータはもはやインターネットでは送れない。
- ・プライバシーをどう守るのか？
- ・I/Oの遅さをどう克服するのか？

③I/Oまわり+ストリーミング計算が肝

page2

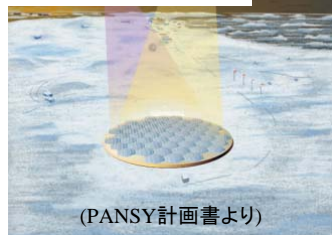
あらゆる分野で発生する問題

計測・観測現場

1テラ級の データ！

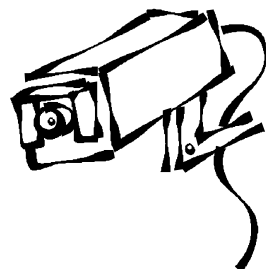


(イルミナ社HP)



(PANSY計画書より)

現場で取得されるデータ量の爆発スピードが、ネットワークや計算機インフラの技術発達スピードを圧倒的に凌駕してしまった。



新たな枠組み
(方法)が必要

計算の現場

データマイニング、
機械学習、統計数理

+

スーパーコンピュータ



出典：理研/富士通
のホームページから

データの次元が計算できるサイズを超えている。

- ・次世代シーケンサーのデータはもはやインターネットでは送れない。
- ・プライバシーをどう守るのか？
- ・I/Oの遅さをどう克服するのか？

データ同化は「モデルを創る」

地球規模の複雑な現象の高精度予測のために、時空間観測・計測データとシミュレーションモデルを統合し、初期値・境界値やパラメータを効率的に求めるための技術。

● シミュレーション科学のプロトコルは…

- 基本モデルの作成
- モデルパラメータのチューニング
- シミュレーション実験
- これまでの知見や観測データとの整合性を検証
- リモデリング

データ同化はシミュレーション科学に
統計科学のエッセンスを注入する。

● データ同化では…

- 基本モデルの作成
- **観測データにもとづくモデルパラメータの学習**
- シミュレーション実験
- **統計的規準にもとづくモデルの性能評価・ロバスト解析**
- **リモデリング**

ご静聴ありがとうございました。



Email: higuchi@ism.ac.jp

Homepage:

<http://www.ism.ac.jp/~higuchi/>

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