

アンサンブルベース逐次データ同化と HPC

樋口知之 (統計数理研究所)

強風予報列車の運行スムーズに

3分〜数十分後に強風が吹きそくた。風による列車運行への支障を少なくするため、JR東日本が「強風予報システム」を開発した。過去数年間の風速の傾向をもとにはじき出す仕組み。来年度から運用方法の検討を始める。実用化されるかどうかは未定だが採用されれば列車の安全性が高まるうえ、運転を休止する時間も短縮できるとい。

JR東日本が開発

る。

同社は1年前から強風に見舞われる路線を持つ輸送指令室にモニターを置き、京葉線、大湊線、仙山線など16カ所でデータ収集やシステムの信頼性の確認をしてい

システムは、3分間隔で記録した風速データを統計学の手法で作った数式にあてはめ、その後の風速をコンピュータに計算させる。数式は「風が安定した時」「荒れ模様の時」など6パターンあり、過去数年の実績で最も

用する。

30分後までの予測は、ほぼ実用段階に達したという。現在の同社の運転基準(在来線の多くで採用)では、風速30km以上が実際に吹いたら

30分間列車を止め、この間に再び30km以上の風が吹かなければ運転再開ができる。新システムが実用化すると、3分

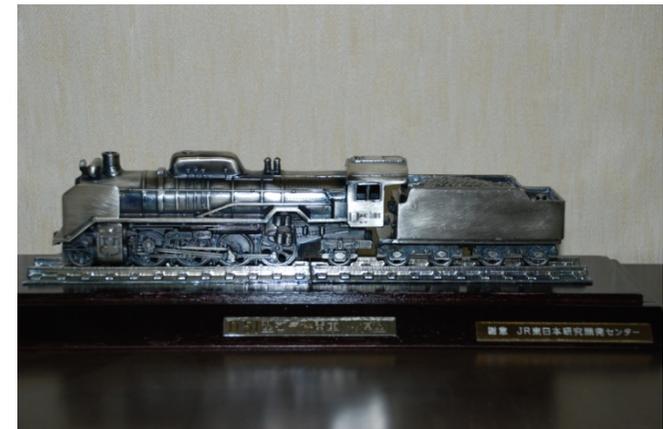
後に30kmを超えると算出された時点で運転を休止するが、風速30km以上にならないと予測できた時点で運転再開できるようにするという。

実用化にこぎつけた場合、運転休止時間は、2〜3割減らせると見込んでいる。新年度から、開発者と輸送指令の担当者が議論を重ねて、どのように利用すれば実用化できるか検討していくという。

実用化なら運休時間2、3割減

新年度から、開発者と輸送指令の担当者が議論を重ねて、どのように利用すれば実用化できるか検討していくという。

JR東 安全研との共同研究



帰納の論理

●論理学者 vs. 数学者 『数学者の厳密性の甘さかげんに我慢がならない』

●物理学者:『60という数はすべての数で割りきれぬ』

1,2,3,4,5,6 OK. “勝手にとったところの” 10,12,15,20,30,...

●数学者 『すべての素数は素数』

例を使って考えるなどといっているのではない。
むしろそうしろと言っているのだ。

よって、これは測定の誤差

個々の特殊な事例を調べると、そこに共通する性質や関係から、一般的な命題や法則を推論できることがある。これを帰納という。帰納の結果が誤りに導くこともいかに多いかをわれわれは知っている。しかし、注意深く、感性を研ぎ澄ましてかかれば、帰納がすばらしい結果に導くことがあるということもわれわれは数々の例から知っているのである。

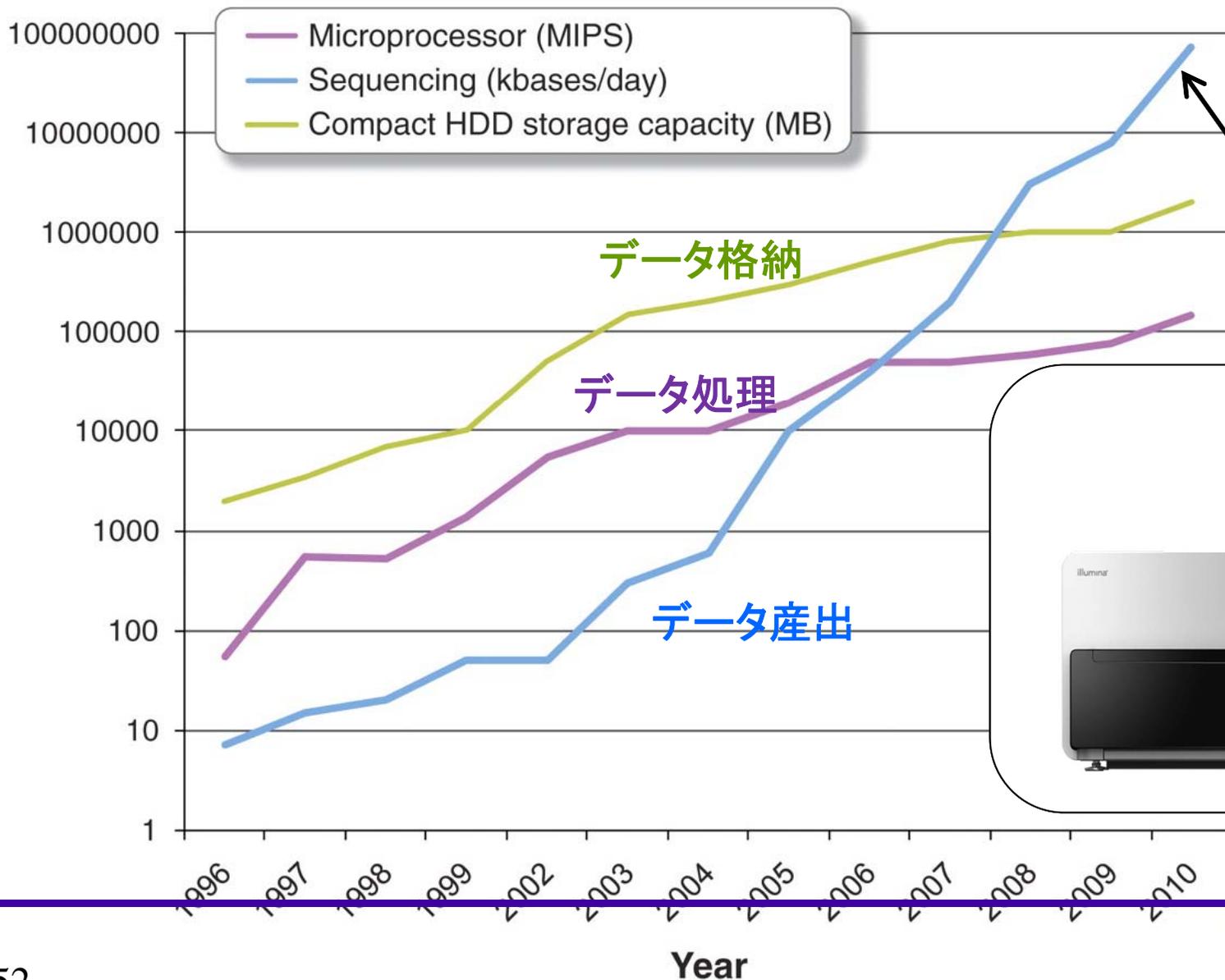


G. ポリア著、「帰納と類比」から
(実際は、金出武雄著「素人のように考え、玄人として実行する」から)

データ中心主義の時代へ

Sequencing Progress vs Compute and Storage

Moore's and Kryder's Laws fall far behind

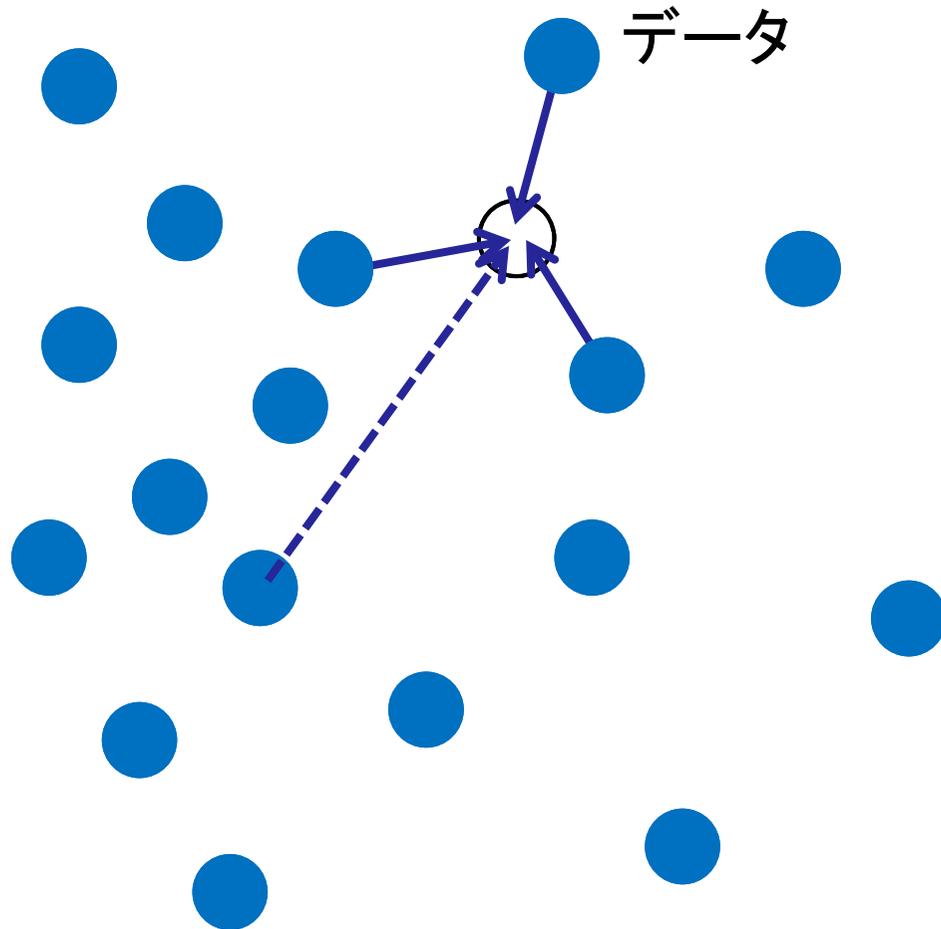


100Gb/day /sequencer



内挿と外挿問題

—得意と不得意—



データ無しの領域



CPSとは

- Center for Planetary Science



予測能力の向上に集中！

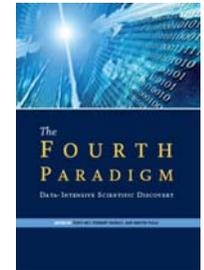
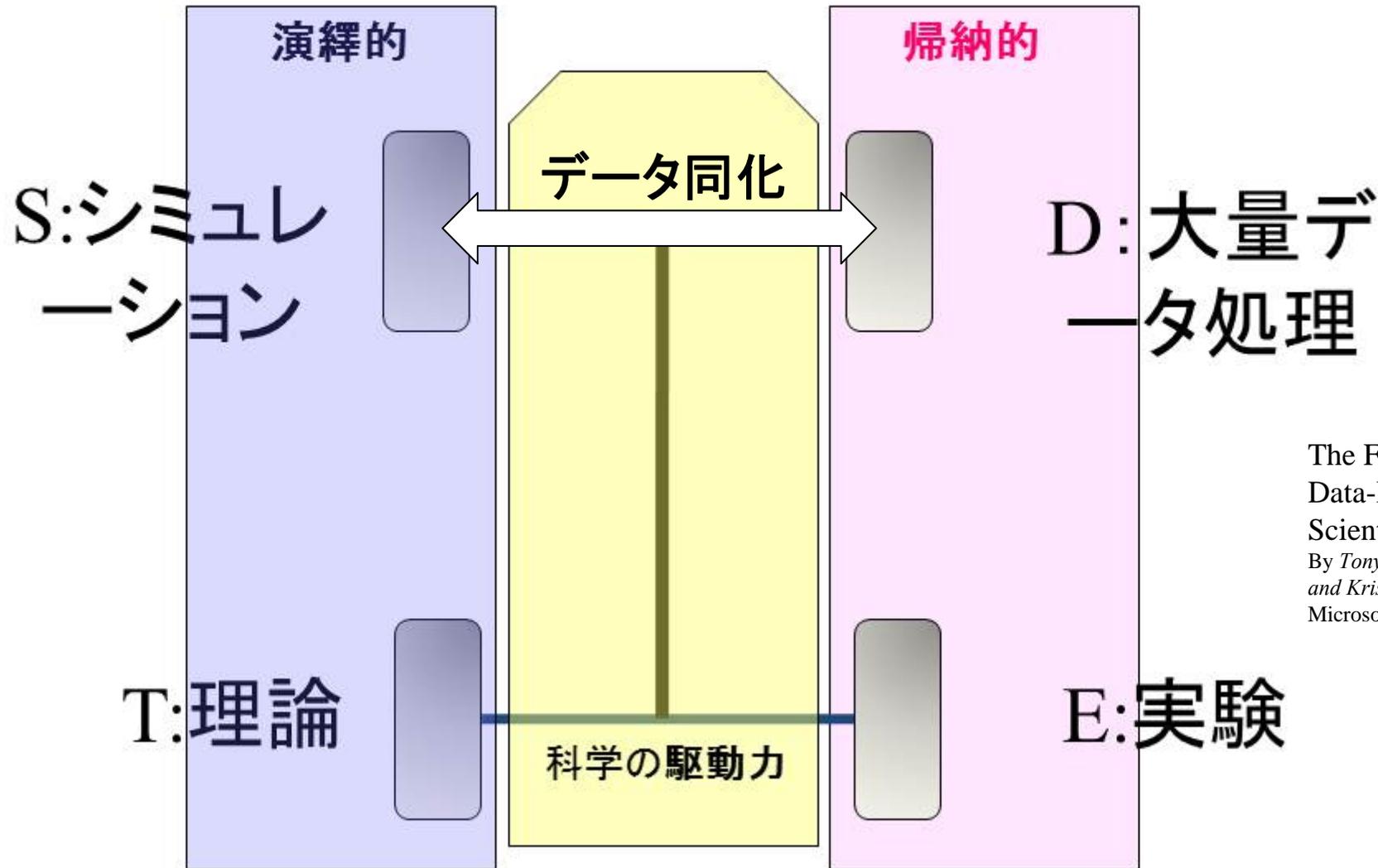
二つの機能の性能の合成

フォワード(前向き)計算モデルの記述力
+
対象の現状態(現況)を捉える認識力
= 予測能力

計測手法のイノベーションと直結

地球や生命体のような複雑なシステムの理解と制御においては、対象に関する知識は常に不完全であることを前提に、現象の予測能力でもって研究の進め方を評価し、修正する方策が有効

つなぐ：データ同化



The Fourth Paradigm:
Data-Intensive
Scientific Discovery
By Tony Hey, Stewart Tansley,
and Kristin Tolle, Eds.
Microsoft Research, 2009.

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。

(参考文献：蒲地 他、「統計数理」、54(2), 223-245, 2006.)

シミュレーションモデルの構成 (1)

(Tsunami simulation model in Japan Sea)

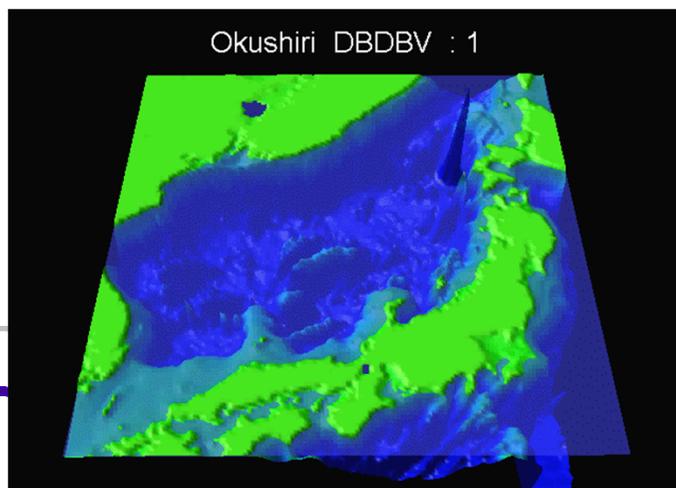
PDE to approximate real physical system
(continuous time/space)

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

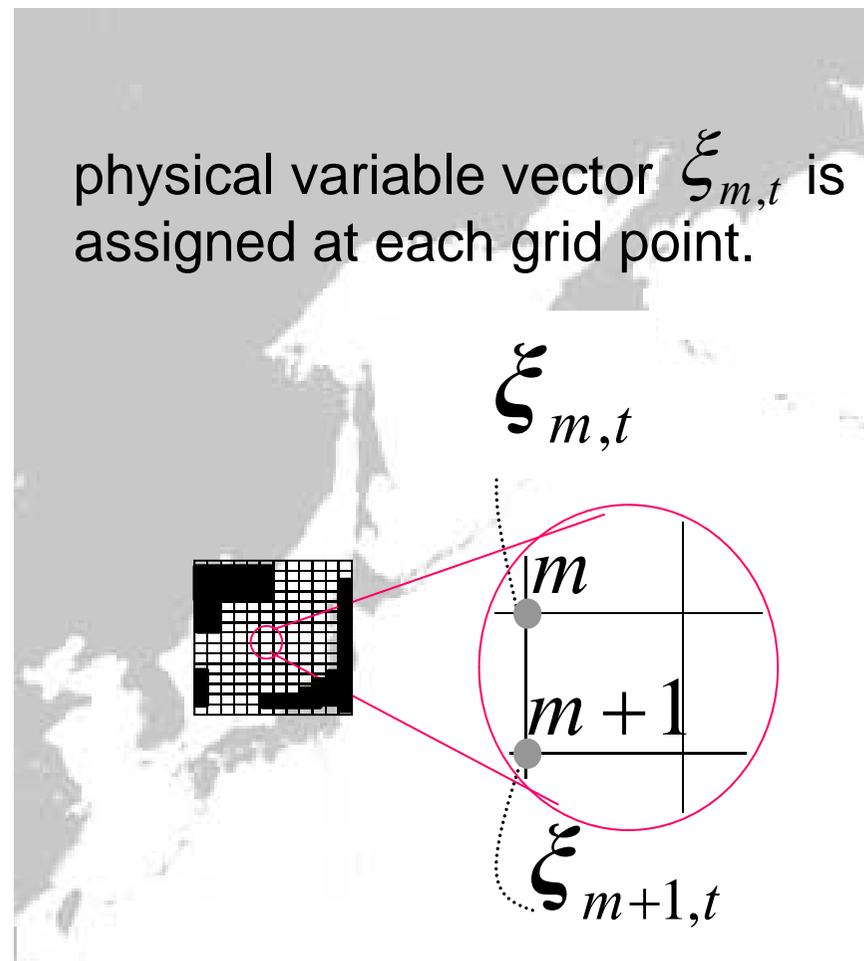
PDE : Partial differential equation

Discrete simulation model
(discrete time/space, FDE)

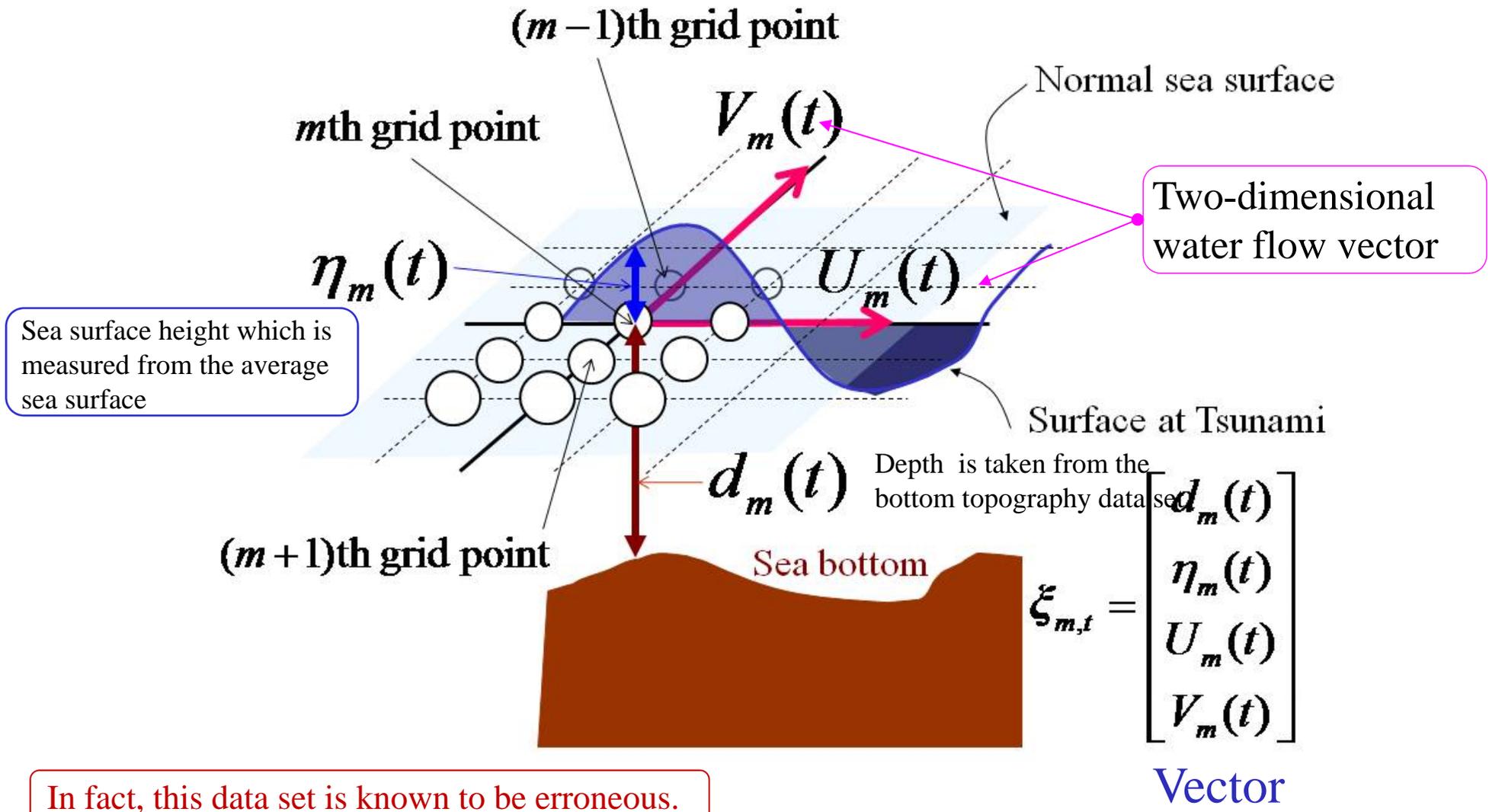
Suppose a case where we conduct a two-dimensional simulation experiment for understanding the flow of shallow water such as the tsunami.



physical variable vector $\xi_{m,t}$ is assigned at each grid point.

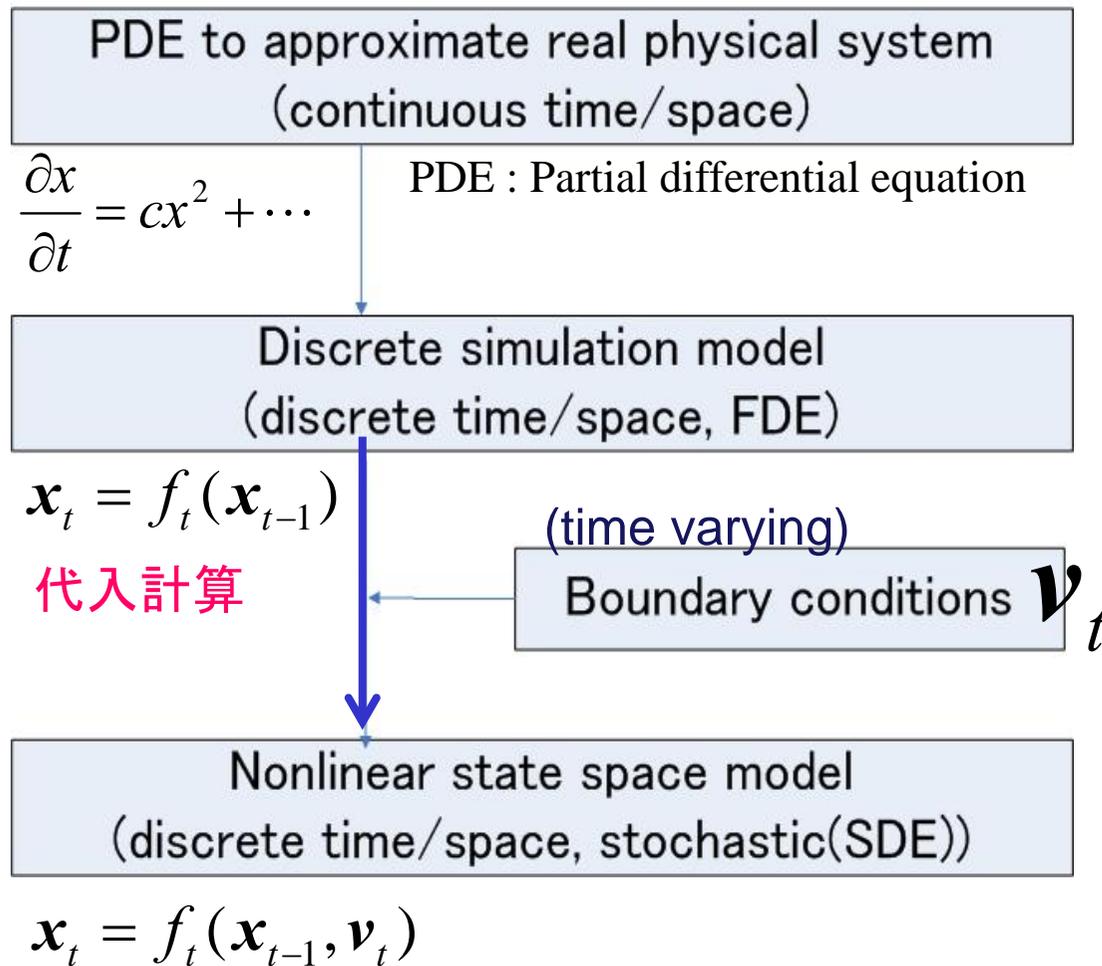


シミュレーションモデルの構成 (2)



システムモデルとしてのシミュレーションモデル

(simplified meteorological model around Japan)



State Vector

$$\mathbf{x}_t = \begin{bmatrix} \xi_{1,t} \\ \vdots \\ \xi_{m,t} \\ \xi_{m+1,t} \\ \vdots \\ \xi_{M,t} \\ \theta \end{bmatrix}$$

データ同化と一般状態空間モデル

State Vector (Simulation variables)

システムモデル

Stochastic simulation model

$$\left\{ \begin{array}{l} \mathbf{x}_t = f_t(\mathbf{x}_{t-1}, \mathbf{v}_t), \quad \mathbf{v}_t \sim p(\mathbf{v} | \boldsymbol{\theta}_{\text{sys}}) \\ \mathbf{y}_t = h_t(\mathbf{x}_t, \mathbf{w}_t), \quad \mathbf{w}_t \sim p(\mathbf{w} | \boldsymbol{\theta}_{\text{obs}}) \end{array} \right.$$

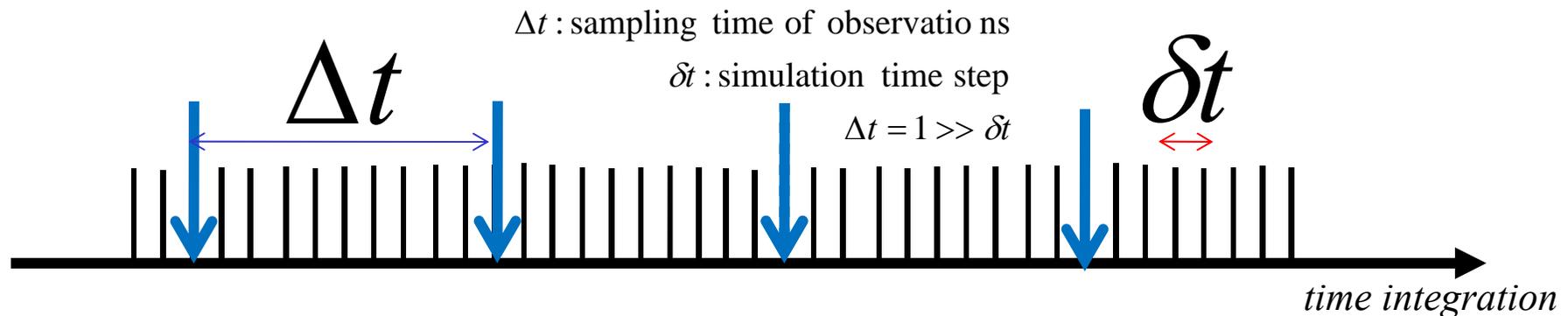
気象・海洋のデータ
同化の枠組み

$$\mathbf{y}_t = H_t \mathbf{x}_t + \mathbf{w}_t, \quad \mathbf{w}_t \sim N(0, R_{\text{obs}})$$

Observation model

Measurement model

観測モデル



逐次ベイズ計算

--- 日次株価データを考えると ---

条件付き分布 **predictive density:** $p(x_t | y_{1:t-1})$
 予測分布

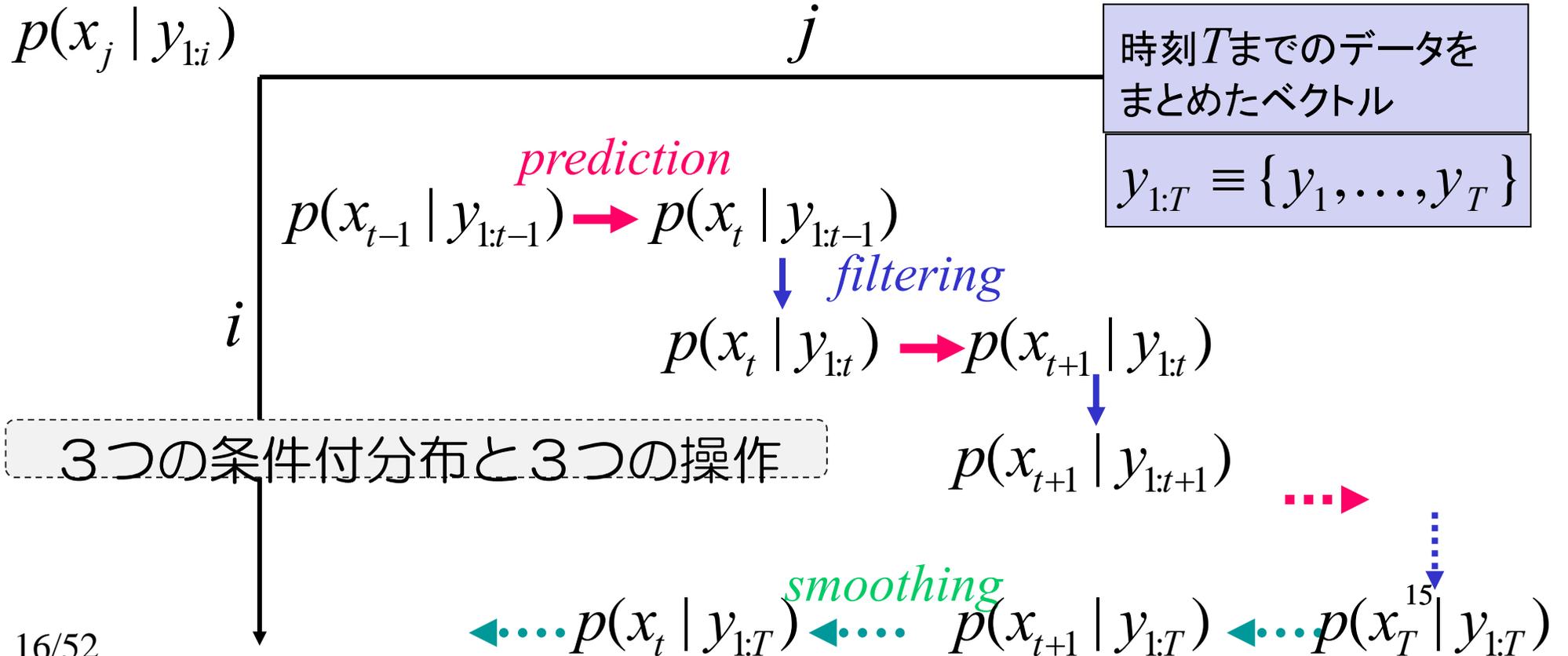
きのうまでのデータに
 基づく今日の状態

filter density: $p(x_t | y_{1:t})$
 フィルタ分布

今日までのデータに基
 づく今日の状態

smoother density: $p(x_t | y_{1:T})$
 平滑化分布

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



フィルタリングとベイズの定理

\mathbf{x}_t : 時刻 t の興味のある対象 $\mathbf{x}_{1:T} \equiv \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$

\mathbf{y}_t : データ $\mathbf{y}_{1:T} \equiv \{\mathbf{y}_1, \dots, \mathbf{y}_T\}$

フィルタリング
フィルタ分布

$$p(\mathbf{x}_t | \mathbf{y}_{1:T}) = \frac{\overset{\text{観測モデル}}{p(\mathbf{y}_t | \mathbf{x}_t)} \overset{\text{予測分布}}{p(\mathbf{x}_t | \mathbf{y}_{1:t-1})}}{\sum_{\mathbf{x}_t} p(\mathbf{y}_t | \mathbf{x}_t) p(\mathbf{x}_t | \mathbf{y}_{1:t-1})}$$

データの尤度 $p(\mathbf{y}_{1:T} | \theta) = \prod_{t=1}^T p(\mathbf{y}_t | \mathbf{y}_{1:t-1}, \theta)$ 一期先尤度

パラメータの事後分布

$$p(\theta | \mathbf{y}_{1:T}) = \frac{\overset{\text{事後分布}}{p(\theta | \mathbf{y}_{1:T})} \overset{\text{事前分布}}{p(\theta)}}{\sum_{\theta} p(\mathbf{y}_{1:T} | \theta) p(\theta)}$$



ベイズの定理がなぜ今役立つのか？ 4つの理由

イギリスの牧師・数学者(1702 - 1761年)
1763年に発見

x : 興味のある対象

y : データ

2. 対象の特徴をとらえるセンサー性能の向上
高精度センサーのコモディティ(日用品)化

4. 高速(無線)インターネット網の整備

ベイズの反転公式

$$p(\underline{x} | \underline{y}) = \frac{p(\underline{y} | \underline{x}) p(\underline{x})}{\sum p(\underline{y} | \underline{x}) p(\underline{x})}$$

1. 膨大な数の積分(和)操作には高速な計算機が必要
コンピュータの性能向上

3. 対象の細かい情報を不確実性を含めて数値化。個人の情報を網羅的に収集
ストレージの廉価化

ベイズの定理と情報循環

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

Posterior

Improved knowledge
about values of \mathbf{x}

Likelihood

Feasibility of realization of \mathbf{y}
for given \mathbf{x}

Prior

Belief
about values of \mathbf{x}

\mathbf{x} の空間

確率分布

尤度関数

確率分布

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_n の分布または値の推定を行う

↓ ← y_{n-1}

$$(p(x_i | y_{1:k}) = p(x_i | y_1, y_2, \dots, y_k))$$

$$p(x_{n-1} | y_{1:n-1})$$

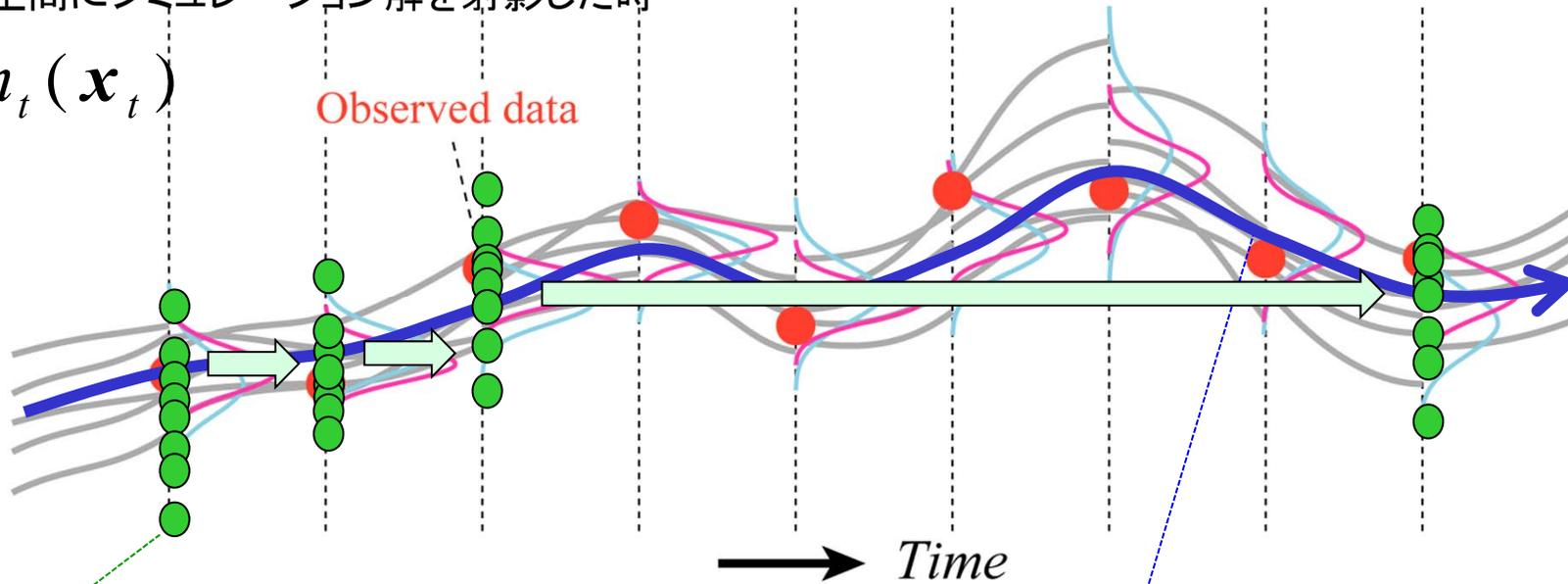
逐次（アンサンブル） vs. 非逐次（最適パス）

データ同化のイメージ

データ空間にシミュレーション解を射影した時

$$h_t(\mathbf{x}_t)$$

Observed data



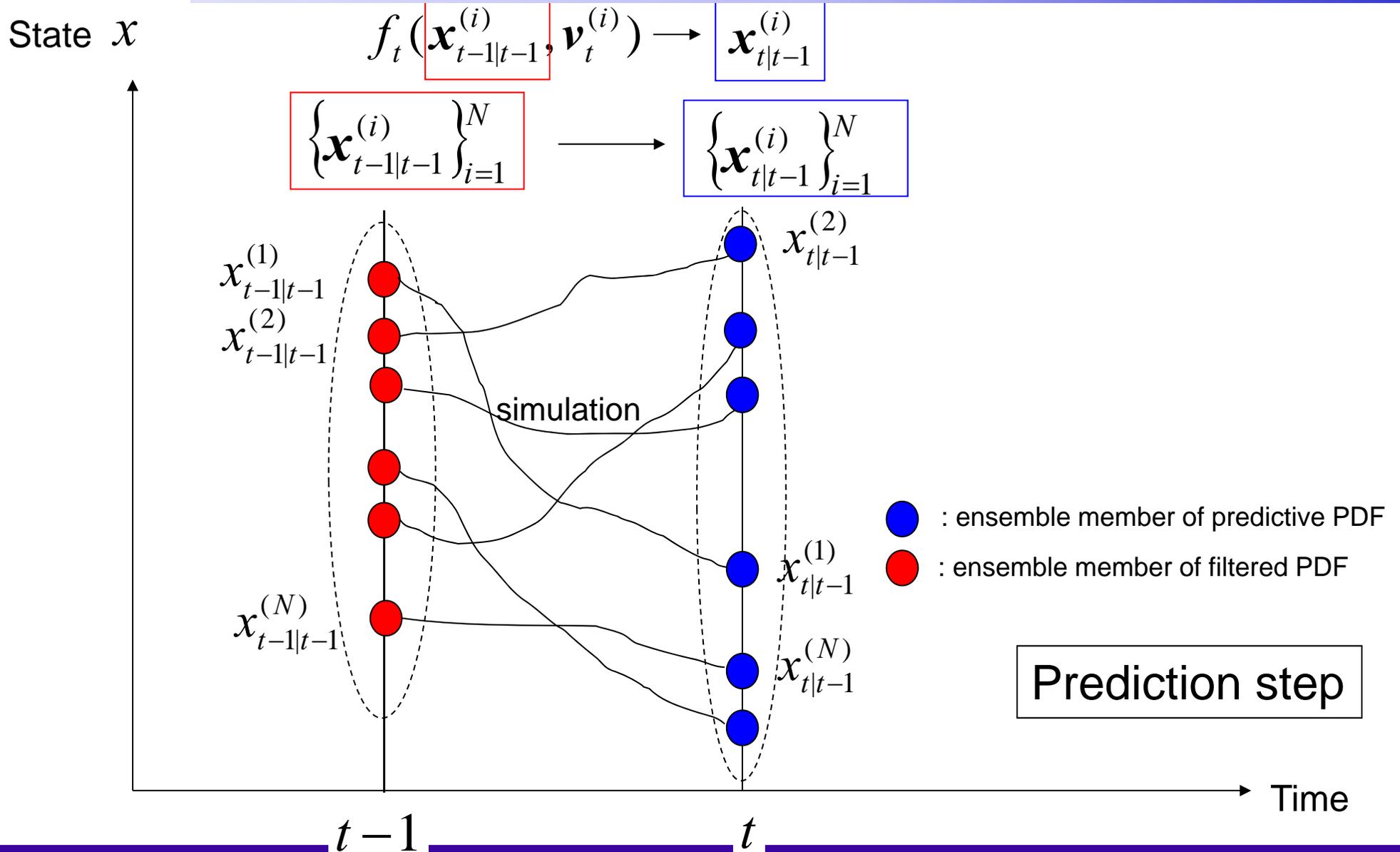
逐次(オンライン)型: 集団の時間発展を追う。つまり、**Swarm Filter**

代表例: EnKF (Ensemble Kalman Filter)

非逐次(オフライン)型: ベスト初期値をもつパスを求める

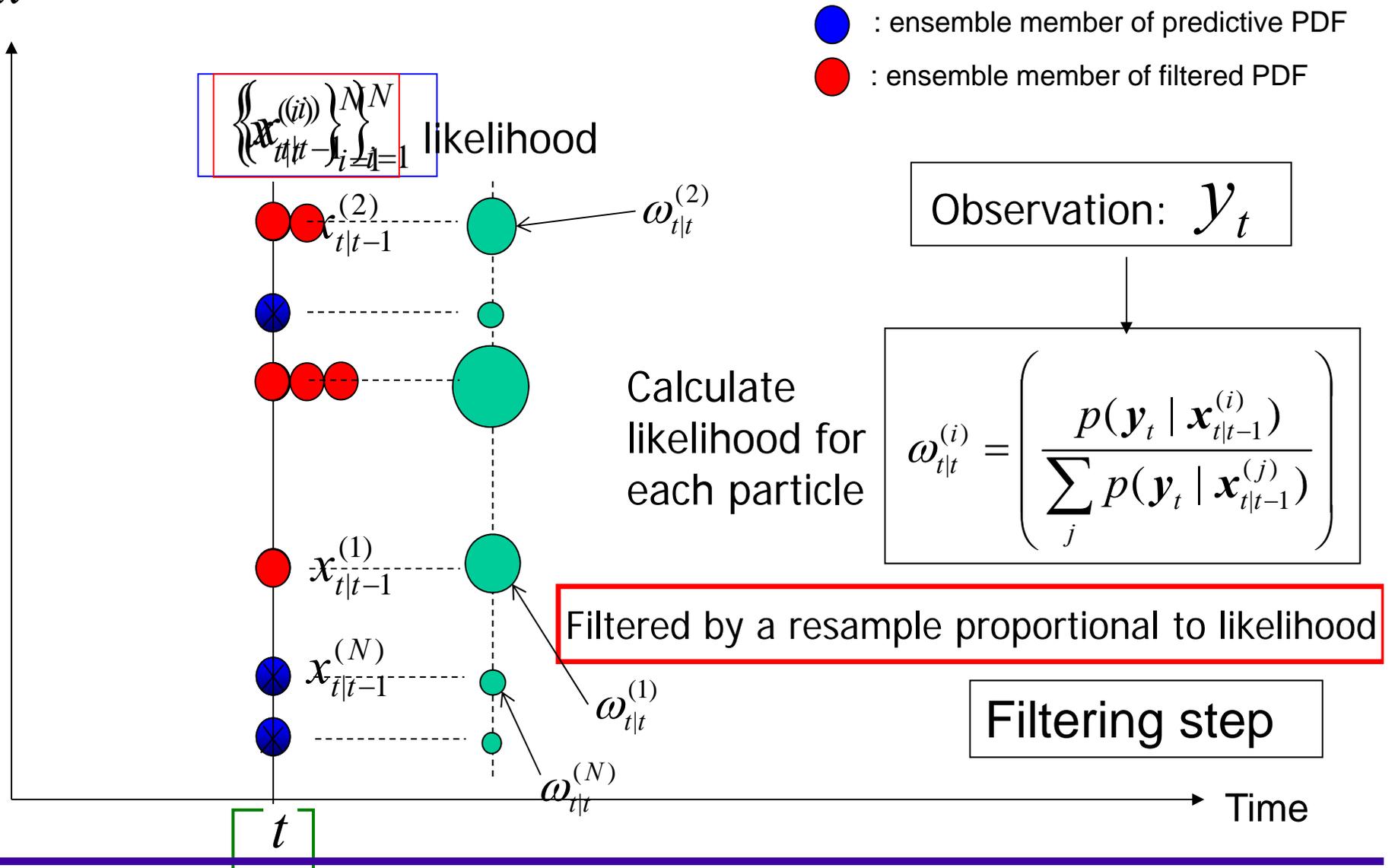
代表例: 4次元変分法 (Adjoint法)

予測のステップ (EnKFとPFで共通)

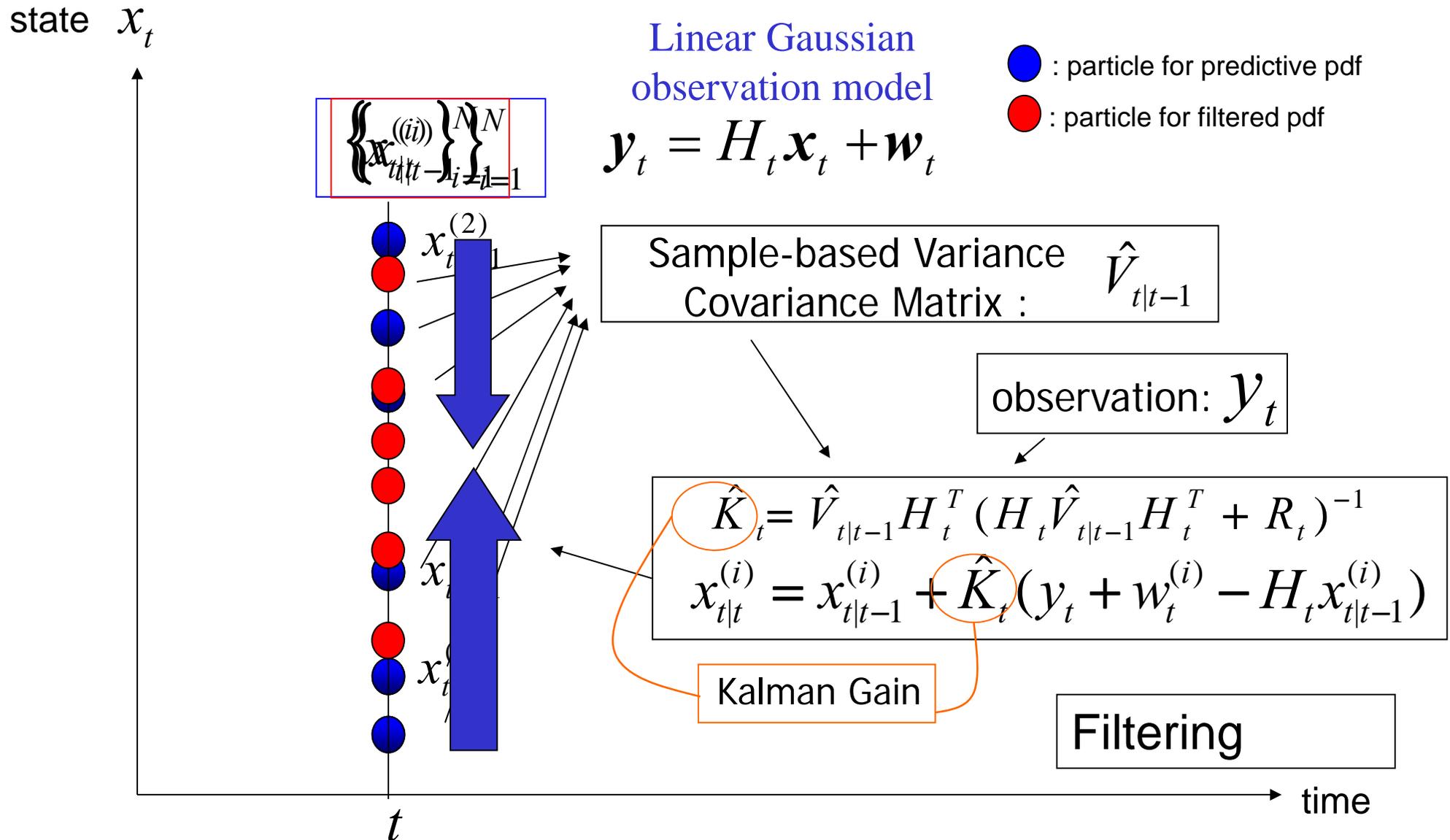


PFでのフィルタリングのステップ

State x

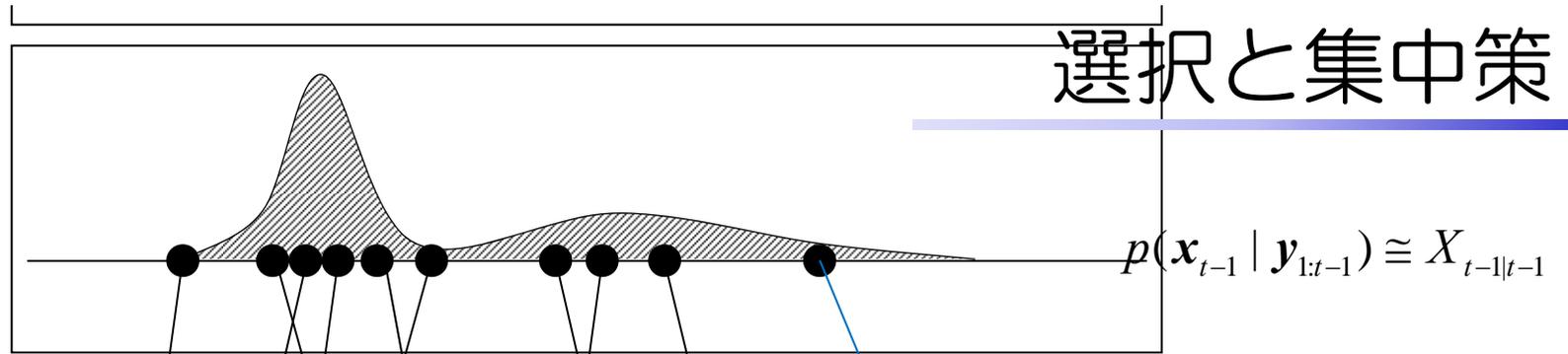


EnKFでのフィルタリングのステップ

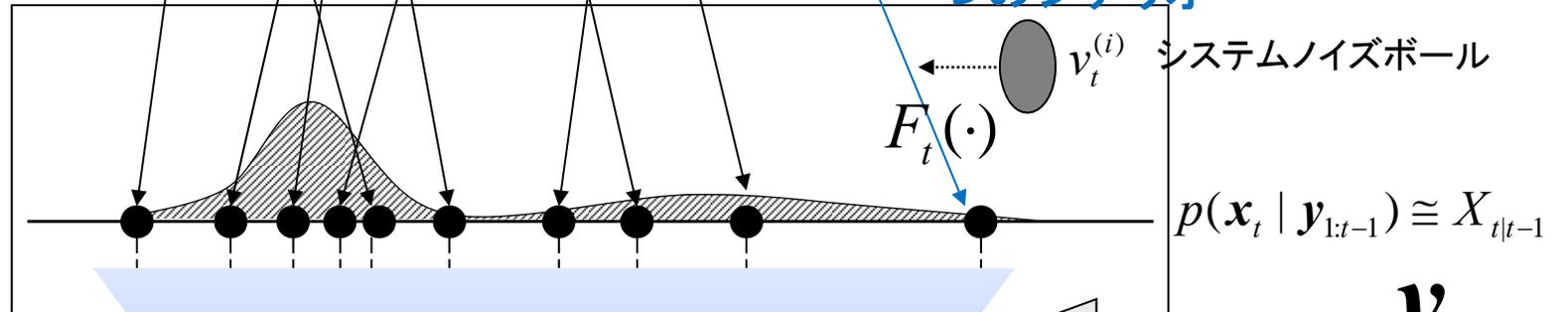


選択と集中策

時刻 $t-1$

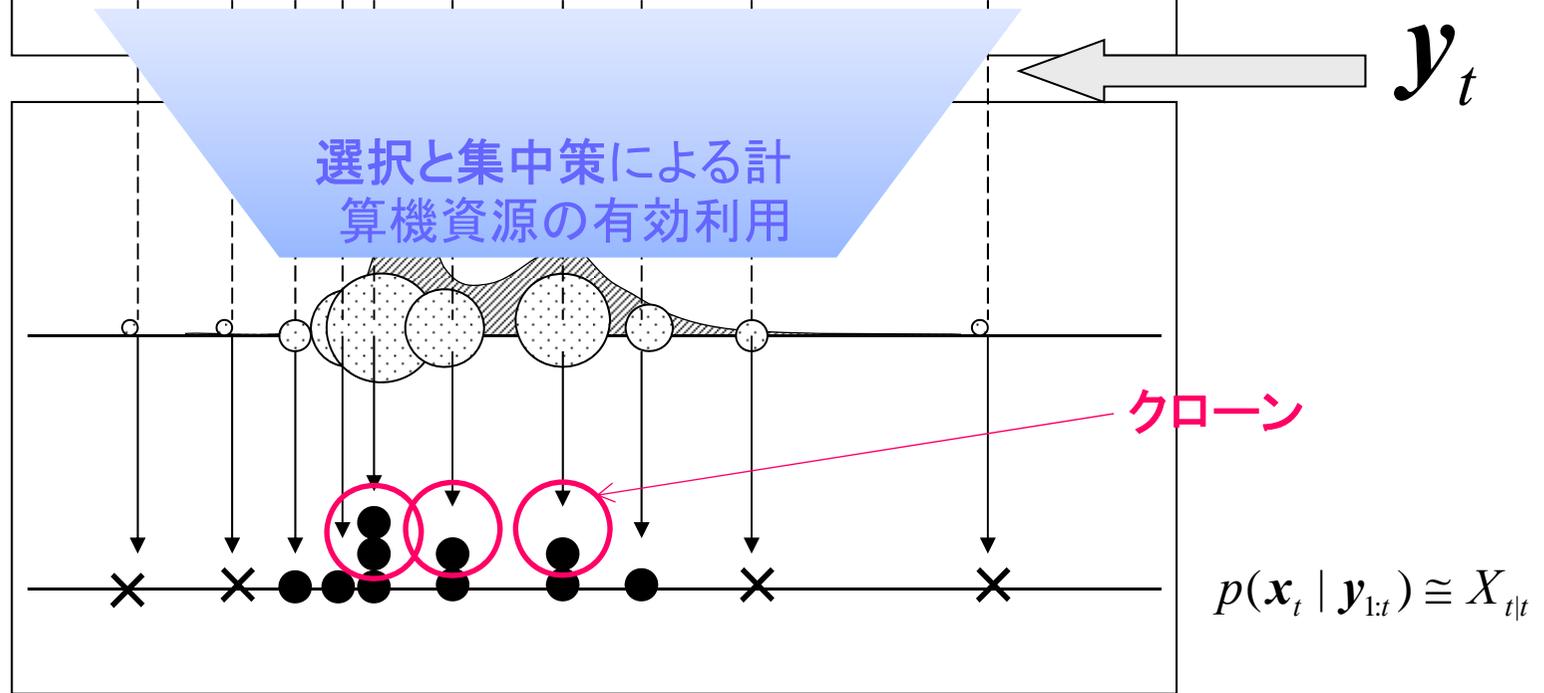


一つのシナリオ



\mathbf{y}_t

時刻 t



25/52 時刻 $t+1$

4次元変分法とベイズ統計

予測

システムモデル

$$\hat{\mathbf{x}}_t = f(\hat{\mathbf{x}}_{t-1})$$

システムノイズを入れたくない
(保存則が破れることを避けたい)

フィルタリング

$$p(\mathbf{y}_t | \mathbf{y}_{1:t-1}) = \int p(\mathbf{y}_t | \mathbf{x}_t) p(\mathbf{x}_t | \mathbf{y}_{1:t-1}) d\mathbf{x}_t = p(\mathbf{y}_t | \hat{\mathbf{x}}_t)$$

初期値のみ最適化

$$p(\hat{\mathbf{x}}_0^{(k)} | \mathbf{y}_{1:T}) \propto p(\mathbf{y}_{1:T} | \hat{\mathbf{x}}_0^{(k)}) p(\hat{\mathbf{x}}_0^{(k)} | \mathbf{y}_{*:0})$$

$$\hat{\mathbf{x}}_0^{(k+1)} \leftarrow \hat{\mathbf{x}}_0^{(k)}$$

パラメータの事後分布

$$p(\theta | \mathbf{y}_{1:T}) \propto p(\mathbf{y}_{1:T} | \theta^{(k)}) p(\theta^{(k)})$$

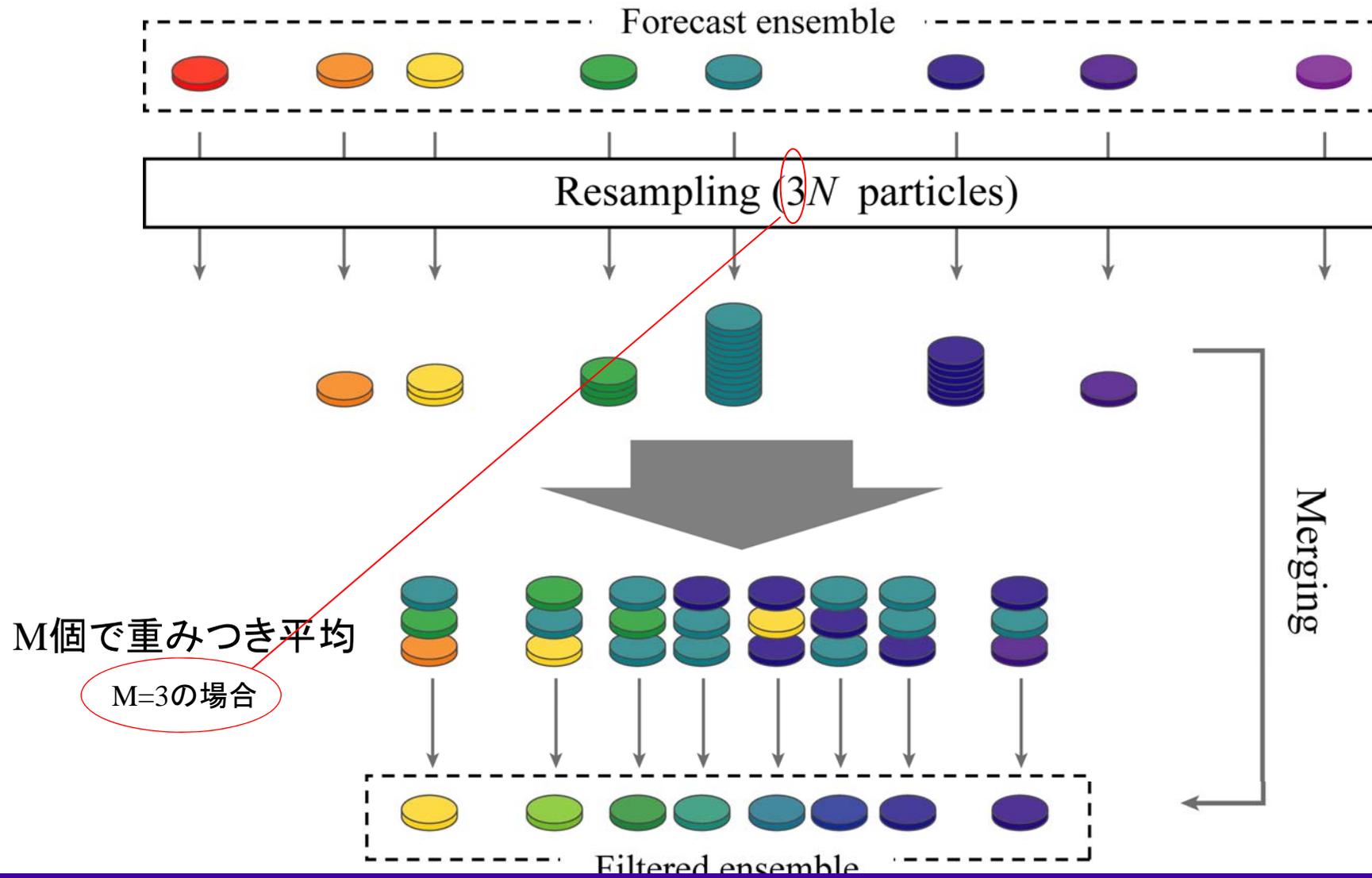
$$\theta^{(k+1)} \leftarrow \theta^{(k)}$$

Importance Sampling
MCMC

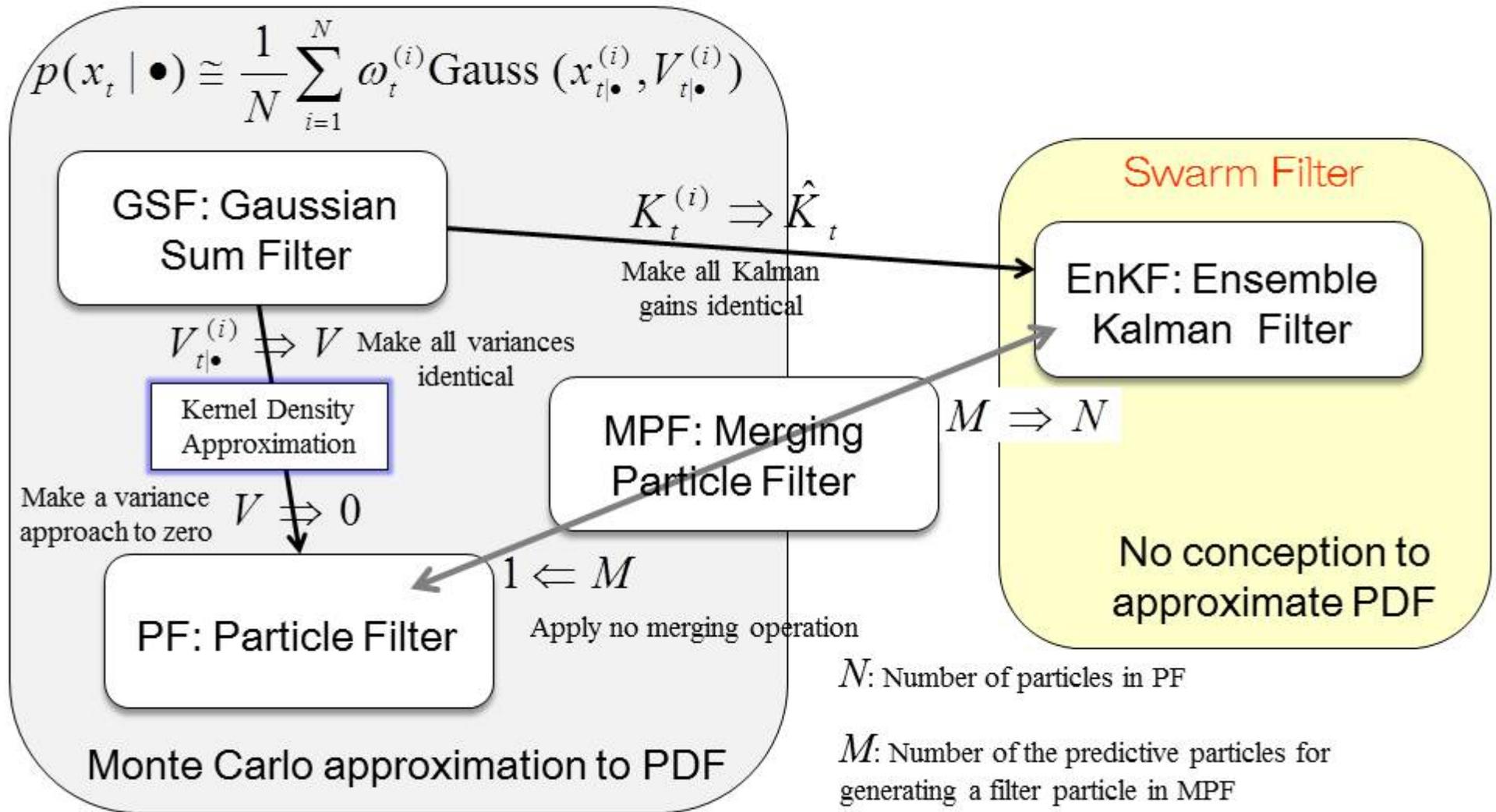
4DVar vs. アンサンブルベースデータ同化法

	非逐次型 4DVar <small>力学的バランスを重視(システムノイズ無し)</small>	逐次型	
		EnKF <small>簡便&安心</small>	原型PF <small>超簡便&原理的には万能</small>
Pros (○)	<ul style="list-style-type: none"> ・状態ベクトルの次元が非常に大きい、最大規模のシミュレーションモデルが取り扱える。 ・感度解析が可能 ・ベクトル計算向き 	<ul style="list-style-type: none"> ・実装が容易 ・退化現象(分布表現能力の減少。)が原理的におきない。 	<ul style="list-style-type: none"> ・実装が著しく容易 ・観測モデルが非線形の場合にも自然に対応可能 ・並列計算向き
Cons (×)	<ul style="list-style-type: none"> ・時間を遡るシミュレーションコードを書きおろす必要があるため、人的労力の負荷が高い。 ・統計モデルでないので、統一的な視点や基準でもってモデル解析ができない。 	<ul style="list-style-type: none"> ・共分散行列の更新ステップの計算コストが高い。 ・分布がガウスから大きく逸脱した時には誤った結果を導く。 ・超高次元状態ベクトルのシミュレーションモデルが取り扱えない。 	<ul style="list-style-type: none"> ・時不変パラメータ推定問題の場合は、退化現象がおきる。 ・厳密な平滑化アルゴリズムの実現が実質的に無理

Merging particle filter (MPF)



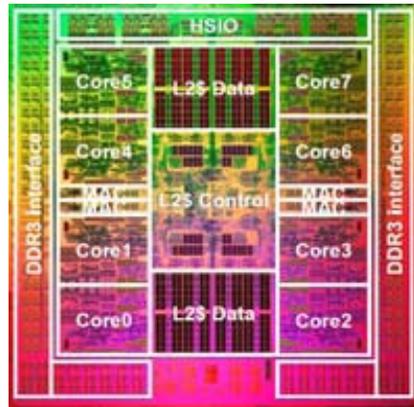
非線形フィルタリング間の関係



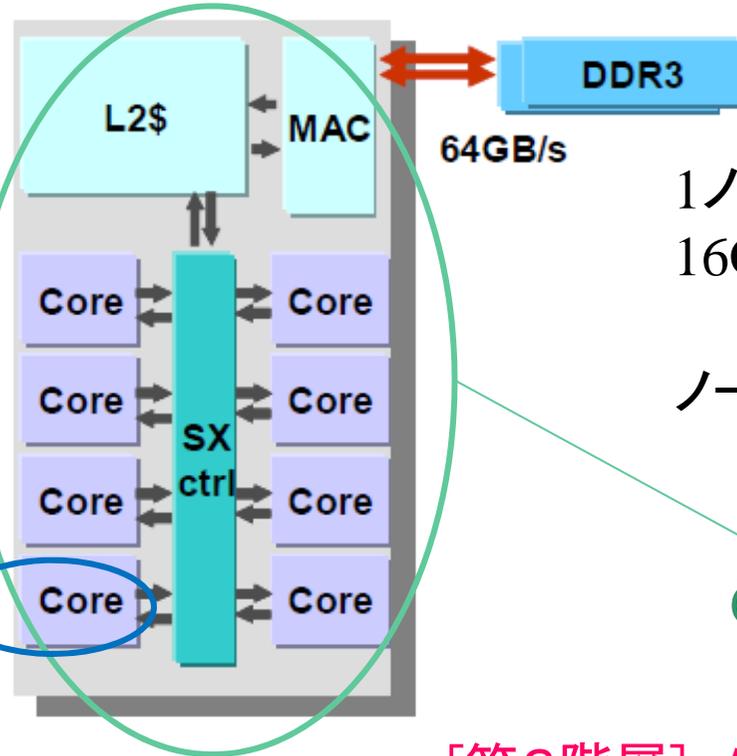
計算機への粒子フィルタの実装

- 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



64GB/s

DDR3

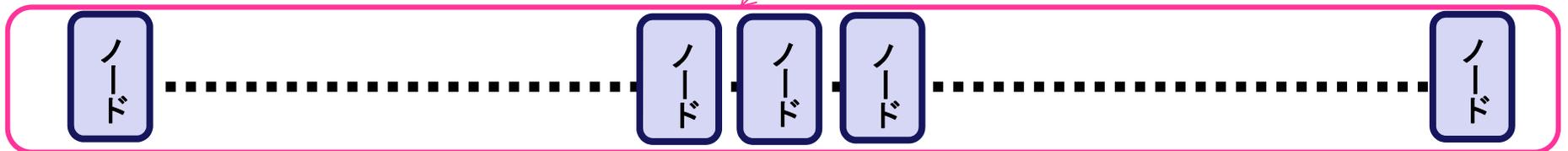
1ノード(CPU) = 8コア
16GFlops x 8 = 128 Gflops

ノード数は8万以上

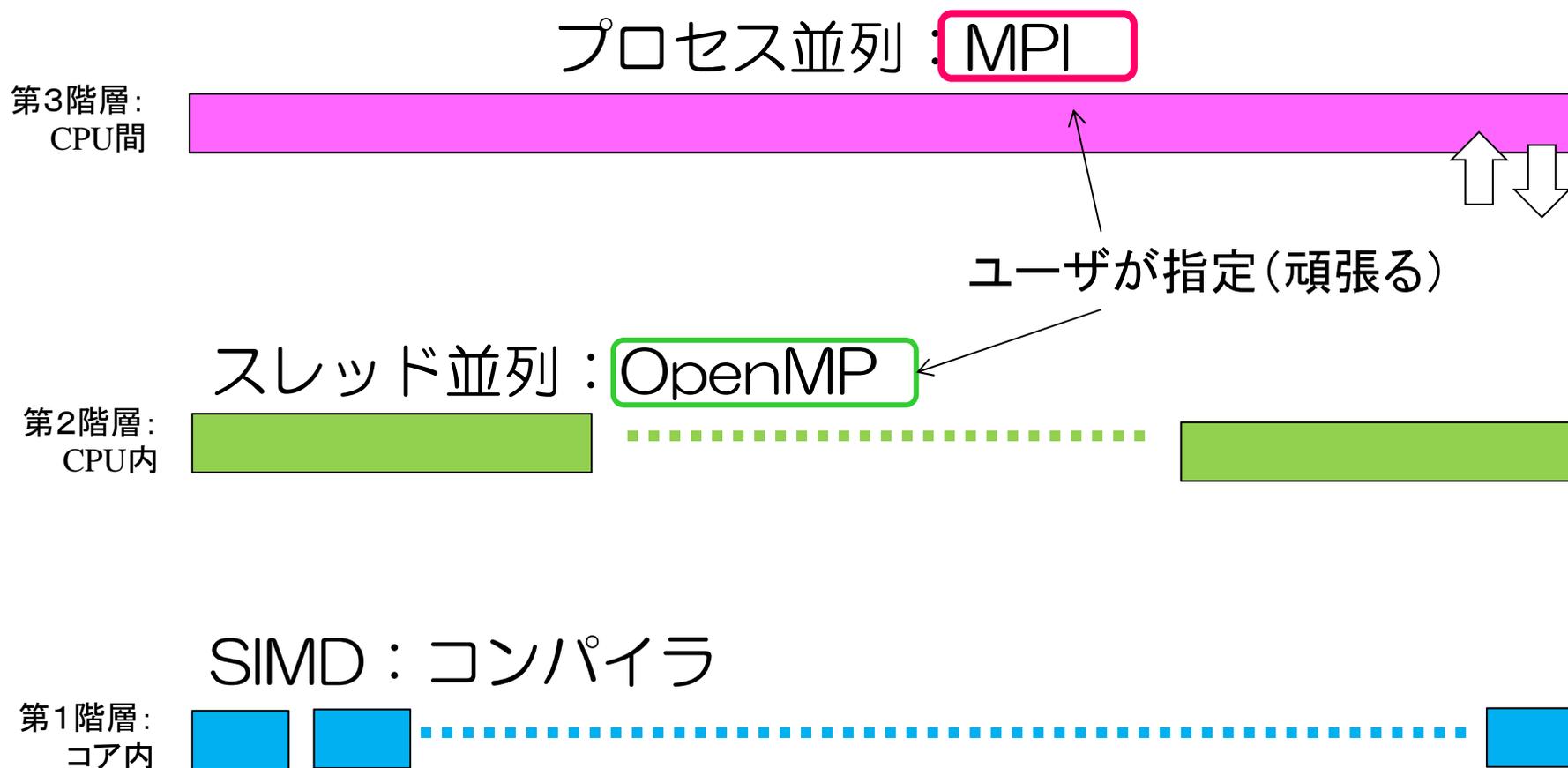
[第2階層]ノード内
OpenMP : スレッド並列

[第3階層] ノード間 MPI: プロセス並列

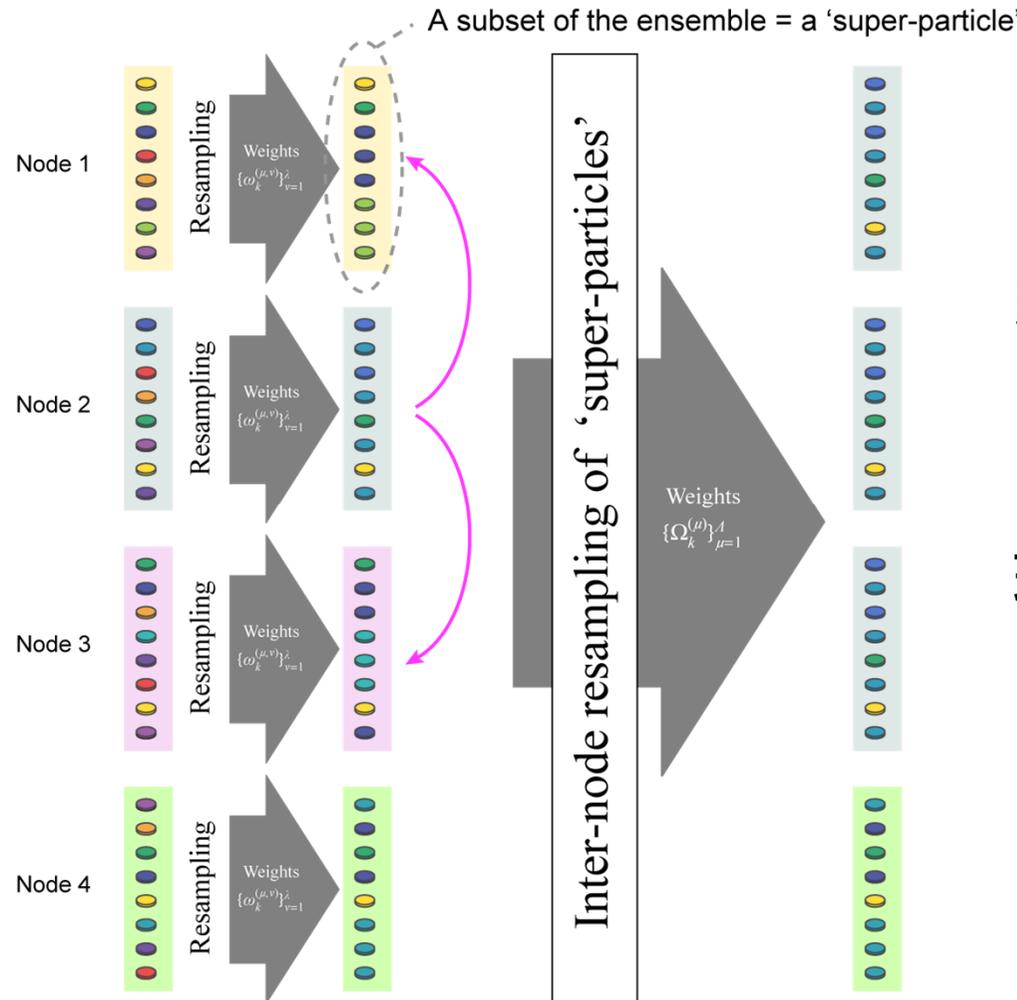
[第1階層]
コア内: SIMD化(コンパイラが対応)



①階層構造をはじめから意識



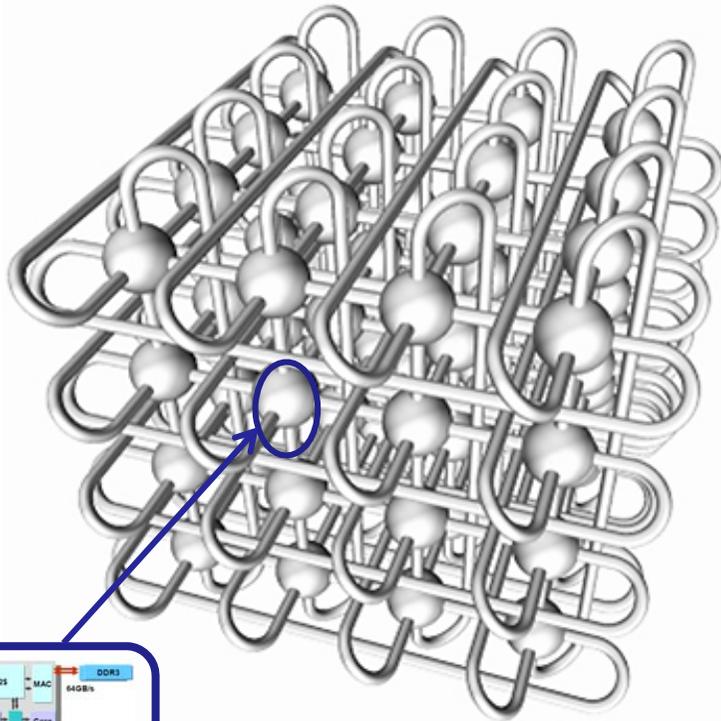
MetaPF: Meta-particle filter



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.

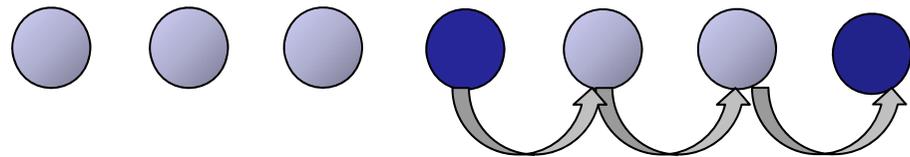
②ネットワーク構造をはじめから意識



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

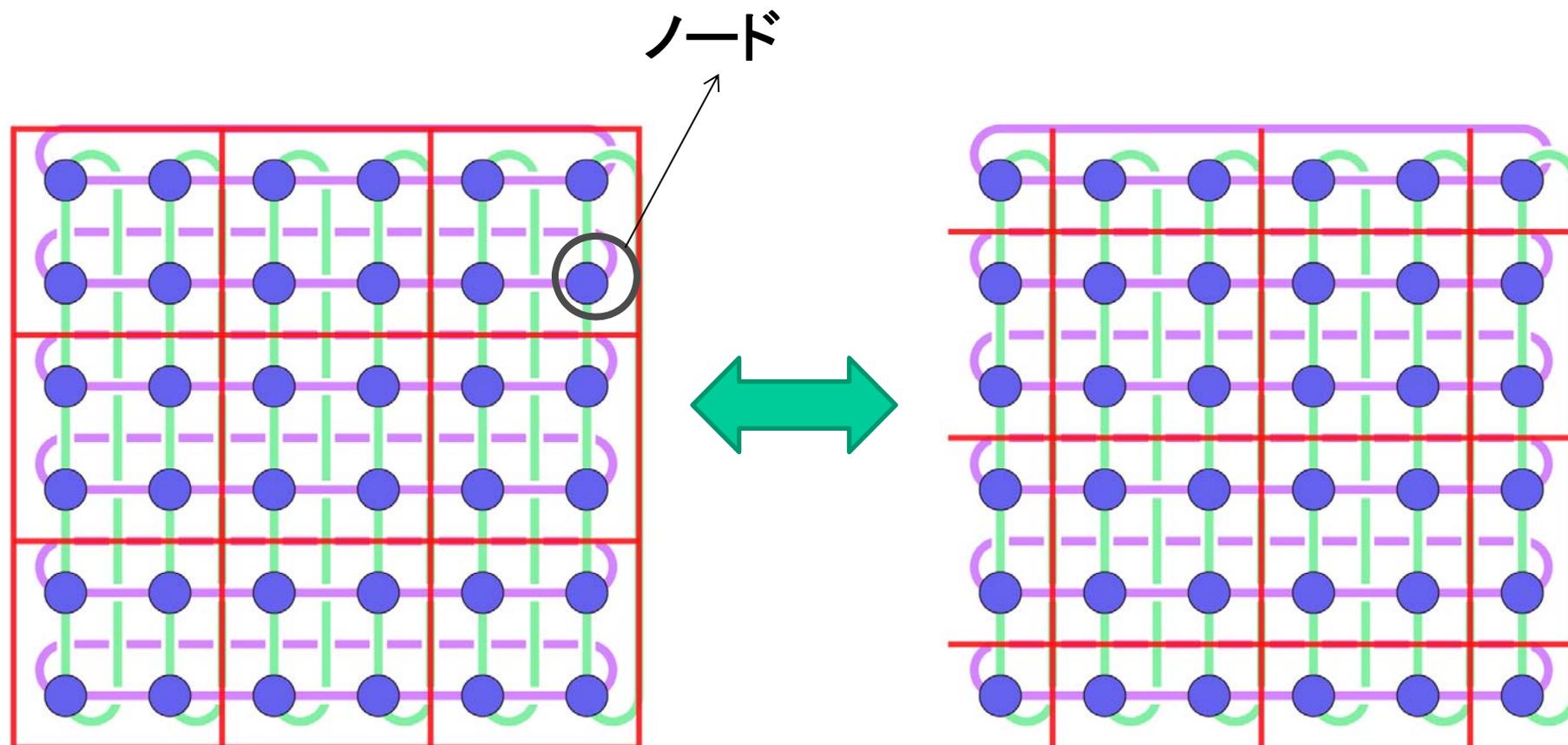
6次元メッシュトーラス: 3次元の直方体に配置したノードをそれぞれ6方向で結合し、各次元がそれぞれリング状に結合されるネットワーク構成



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。

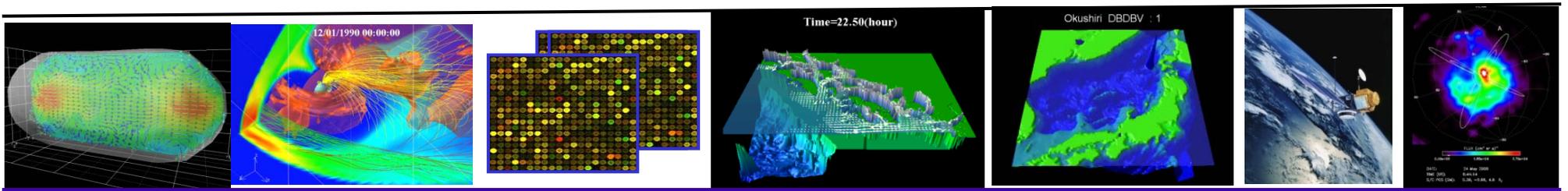
Alternate lattice-pattern switching



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。

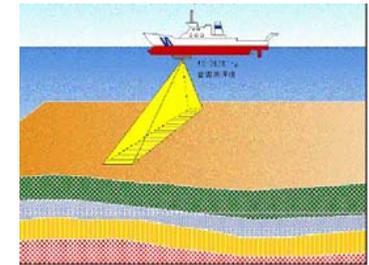
データ同化研究開発センターでの研究プロジェクト

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

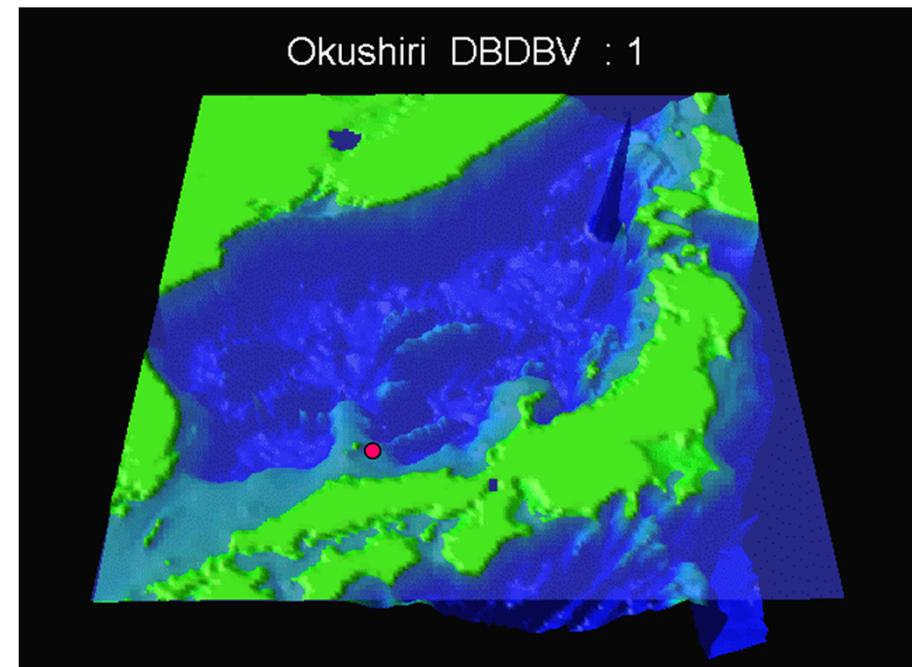
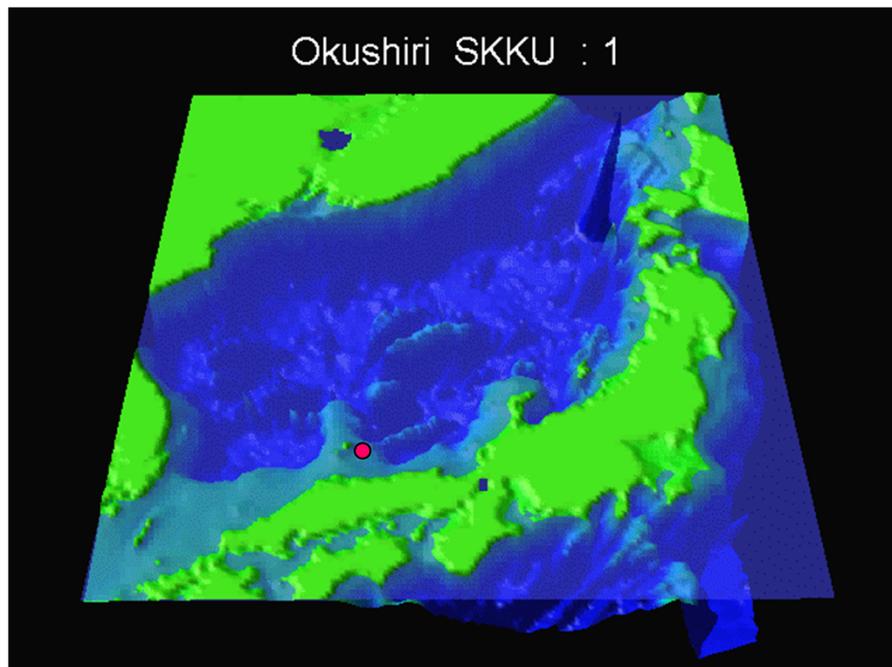


シミュレーション (同化結果でない)

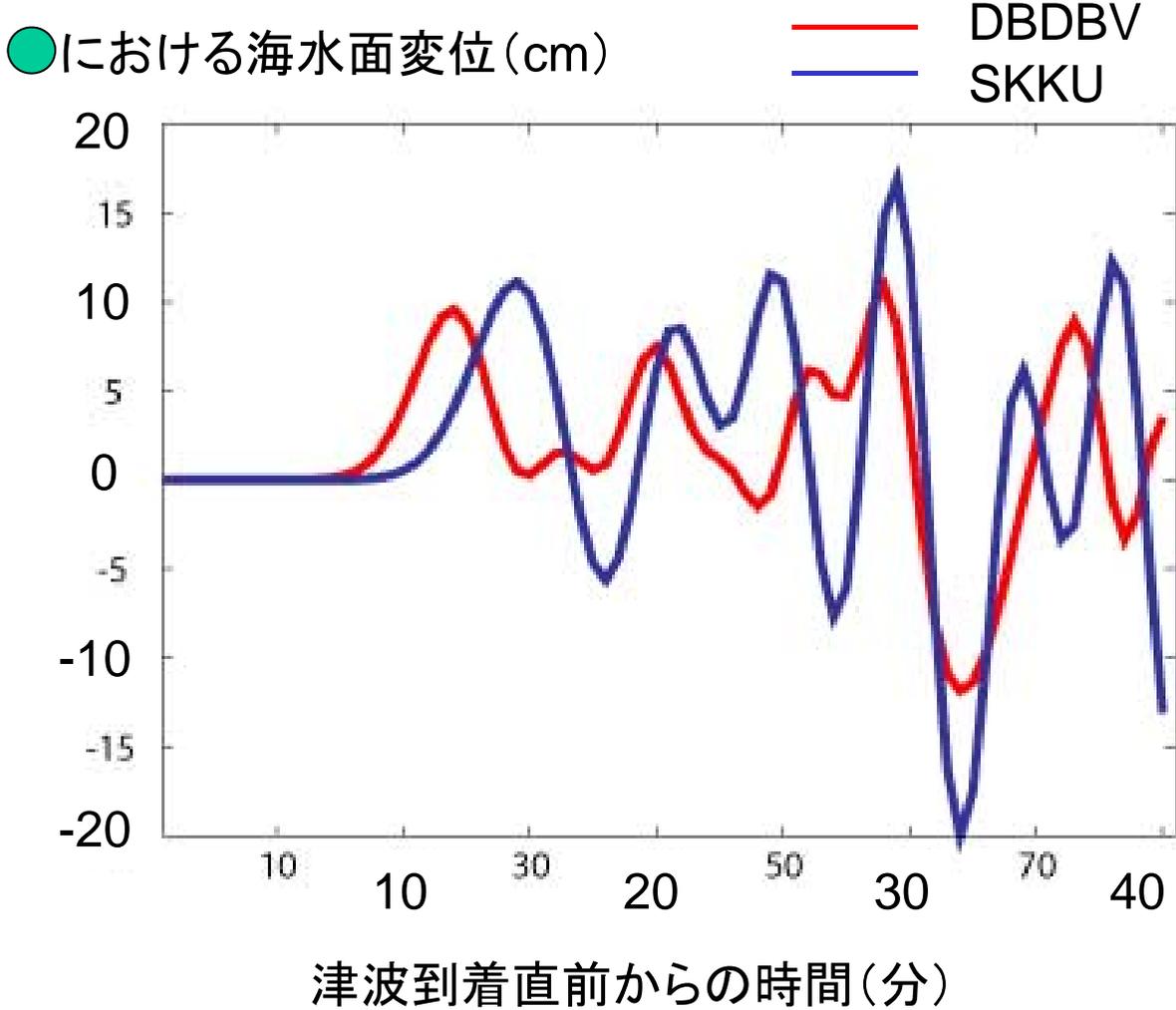
- Simulation of Okushiri Tsunami
 - Simulation based on topographies made by different organizations.
 - It looks similar, but time series of sea surface displacement at a point () is ...



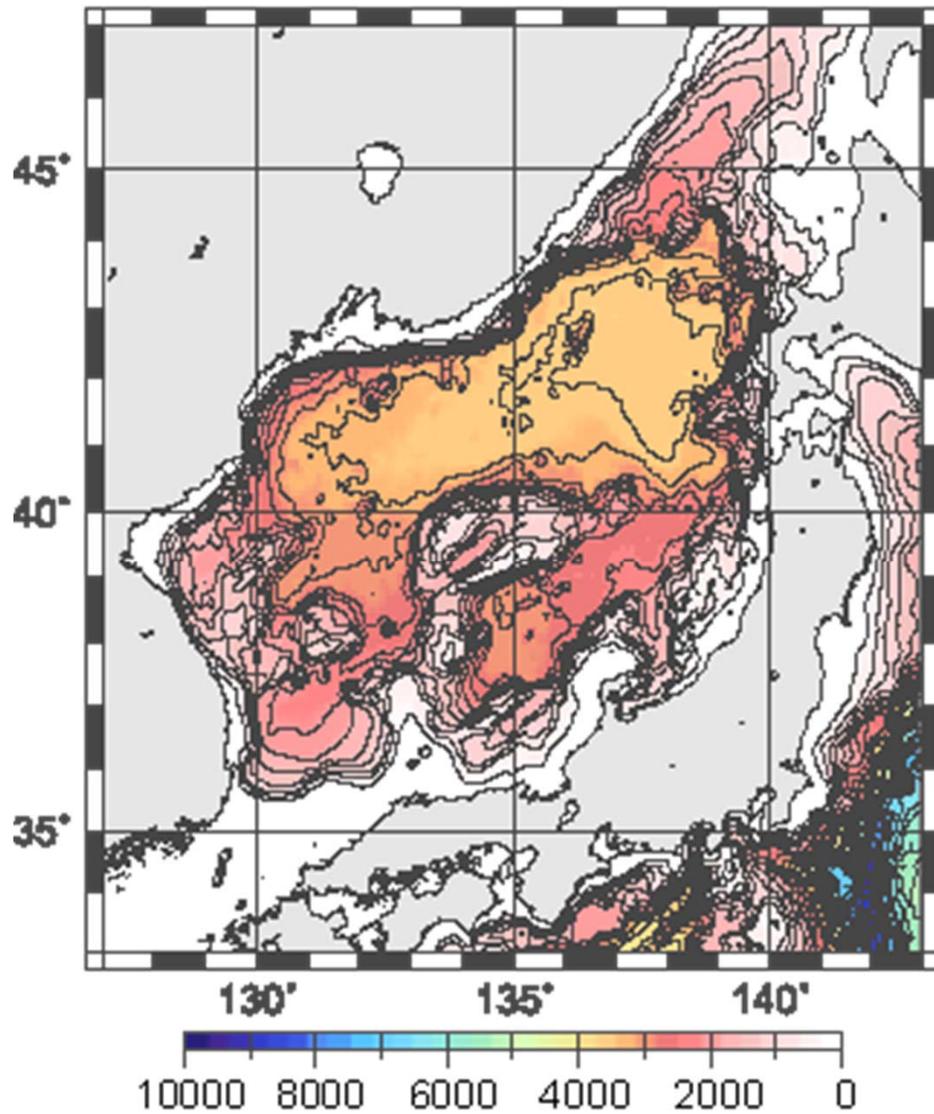
(From Japan Coast Guard)



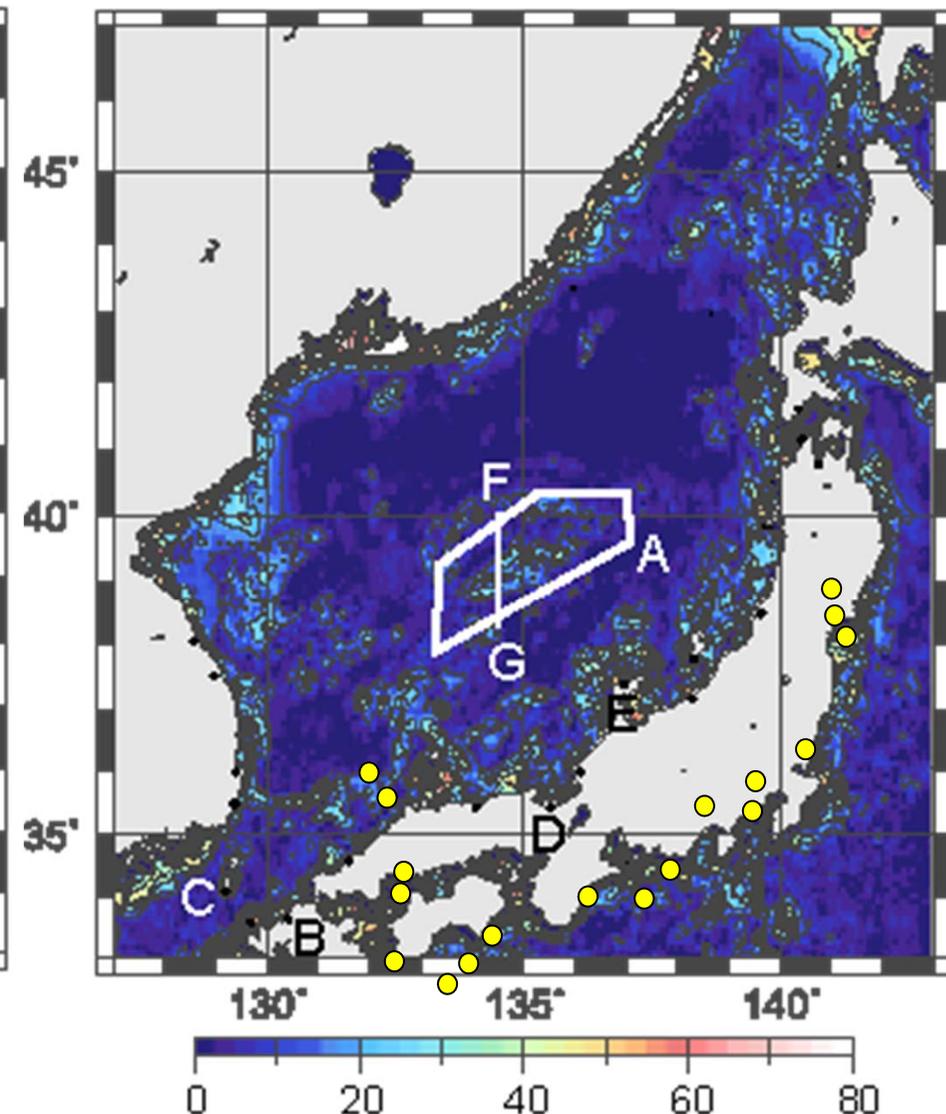
海水面変位



大きい相対的誤差： 深さに関する4つのデータベース間の比較



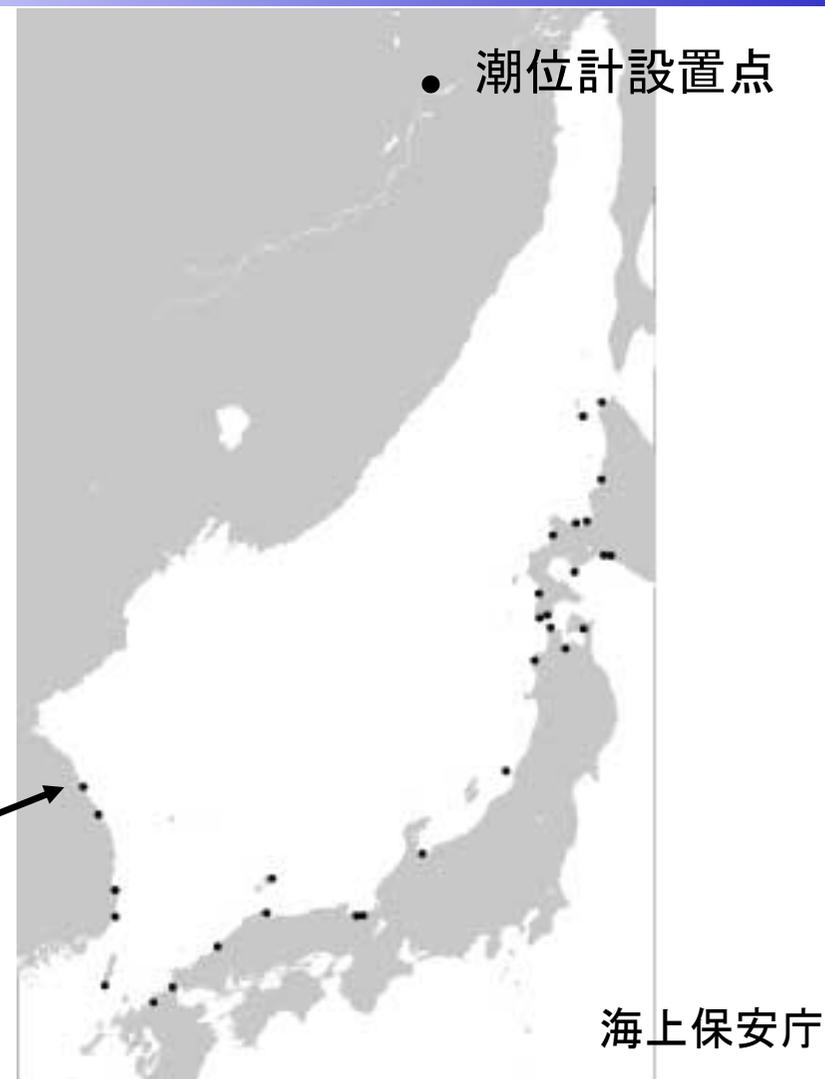
Depth (m)



SD/Depth(%)

● 潮位計設置点

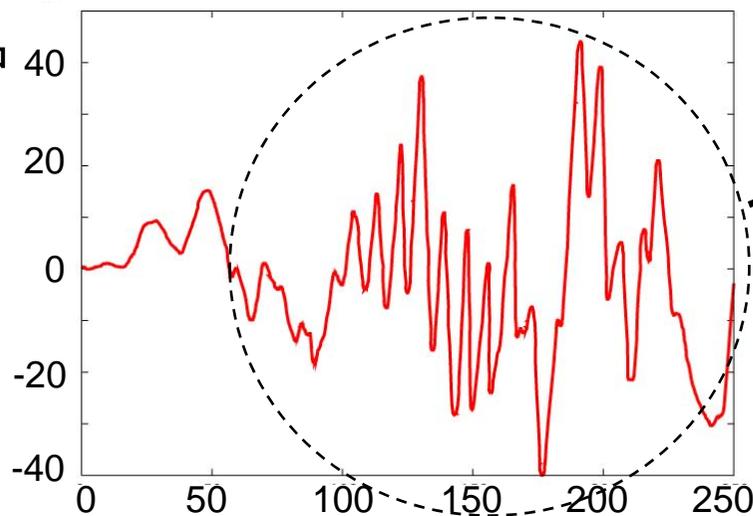
- 潮位計データ
 - 設置地点付近の海面変位を反映
 - 得られるデータは1次元時系列
- 設置地点は右図
 - 30点程度



例)

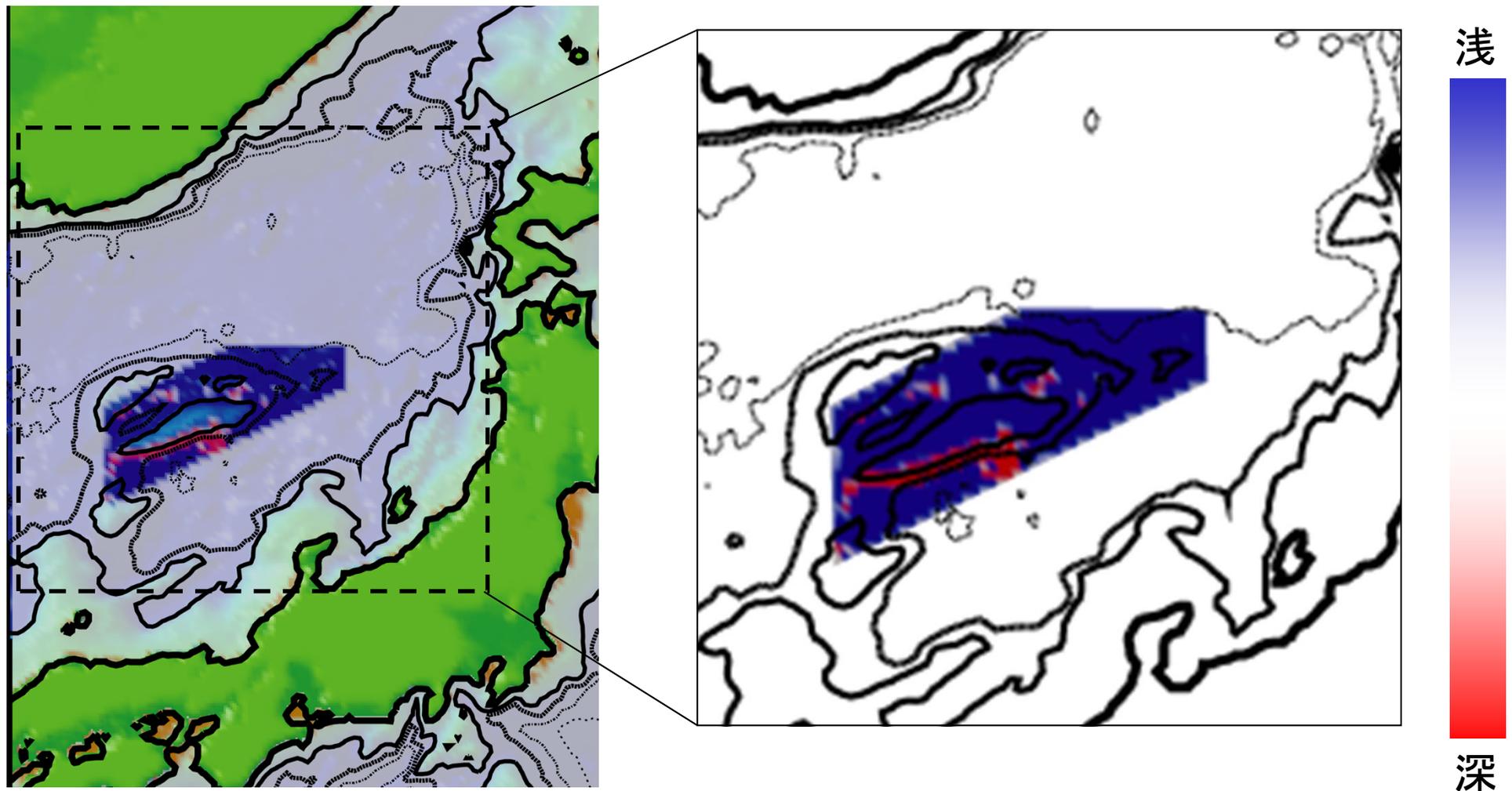
北海道南西沖
地震津波
(韓国・東草)

潮位 (cm)



津波到着直前からの時間(分)

同化結果



- 大和海嶺周辺は4種類の海底地形データの平均よりもやや浅いと推定される
- 南斜面に平均よりも深いと推定される部分がある

“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

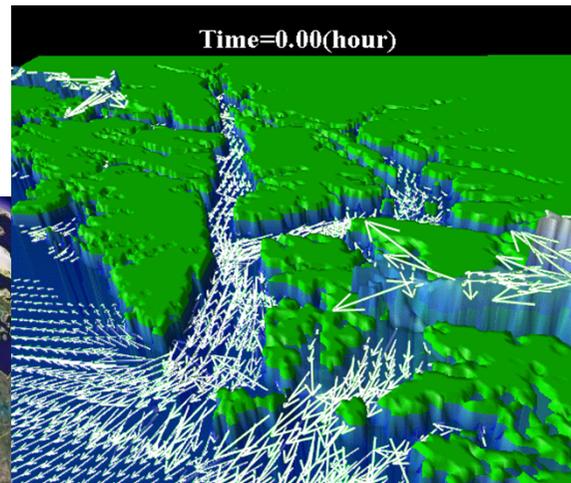
“個”によって異なる形状, 形態情報をシミュレーションモデルに取り込む 『メタシミュレーションモデル』

海底摩擦項

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

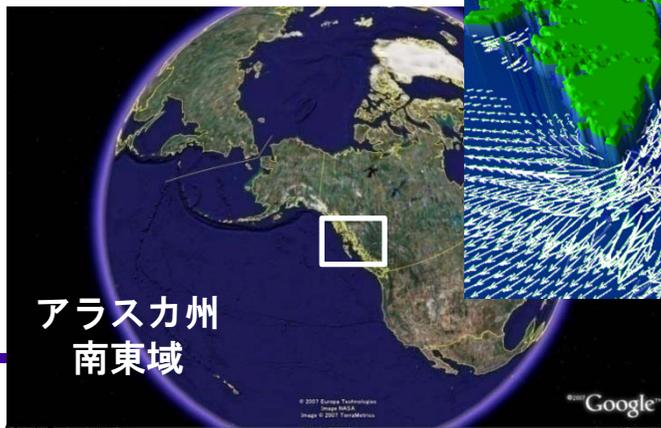
連続式:
$$\frac{\partial \eta}{\partial t} + \nabla \cdot (\mathbf{v} H) = 0$$

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



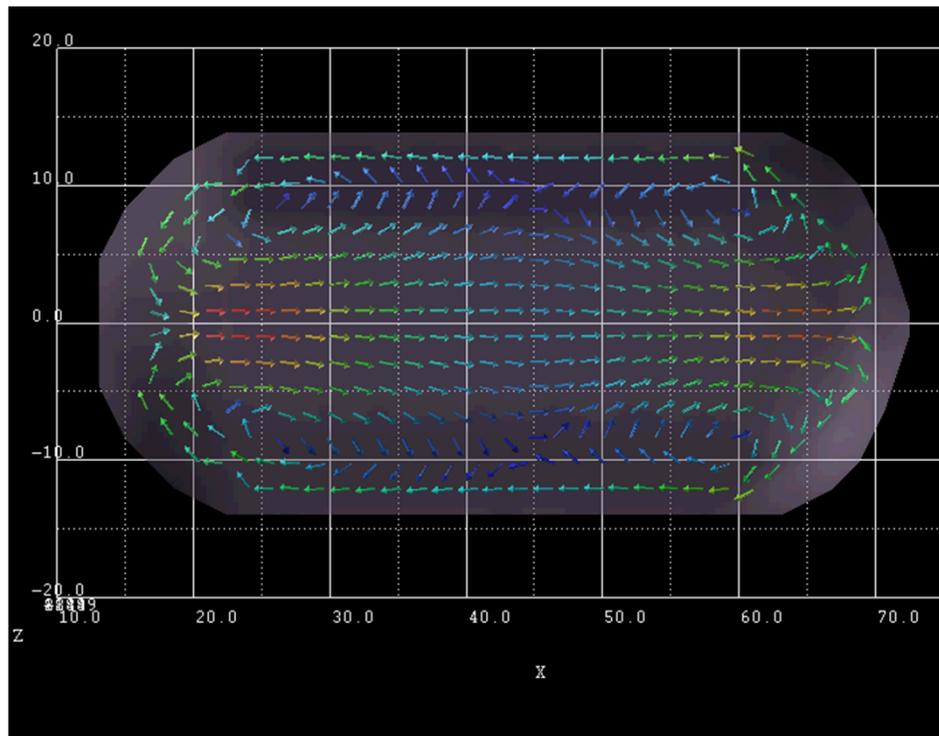
流動の原動力のパラメータ推定

国立遺伝学研究所・木村研究室と共同研究

細胞質流動(線虫)

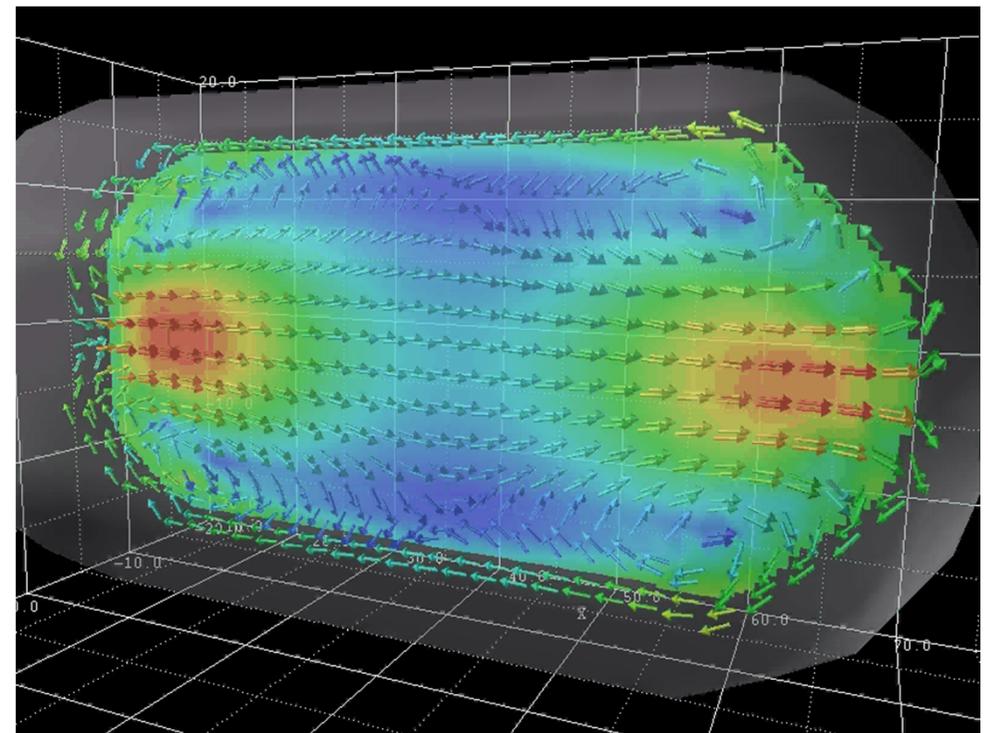
- ・境界条件の推定: ミオシンがつくる剪断応力値を推定

細胞質流動の定量データ



PIV(Particle Image Velocitmetry)計測

細胞質流動のシミュレーション



数値解法はMPS(Moving Particle semi-implicit)法を採用。連続体を多数の粒子で近似。計算格子を生成しなくていい。自由表面などの移動境界がある問題の扱いが簡単

東北地方太平洋沖地震（東日本大震災）に伴う 地震波・地震音波伝搬シミュレーション

物理モデル

- 地球中心から地表までの固体地球および地表から高度1000kmまでの大気を考慮した1次元地球構造モデル
- 気象庁が決定した震源解（モーメントテンソル解）を長さ150kmの断層に沿って配置した断層破壊モデル

各計算グリッドにおける地震波（表面波）および音波の応答波形をノーマルモード法（Kobayashi [2007]）によって計算

グリッド幅は緯度・経度方向0.1度および高度方向10km

断層を含む緯度・経度方向10度および高度方向120kmの範囲の可視化した結果を示す。

統計数理研究所が所有するスーパーコンピュータ（富士通Dシステム）を約4000ノード時間（512並列で約60時間）使用

津波よりも伝搬速度が速い音波を利用した将来型津波警報システムの構築

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

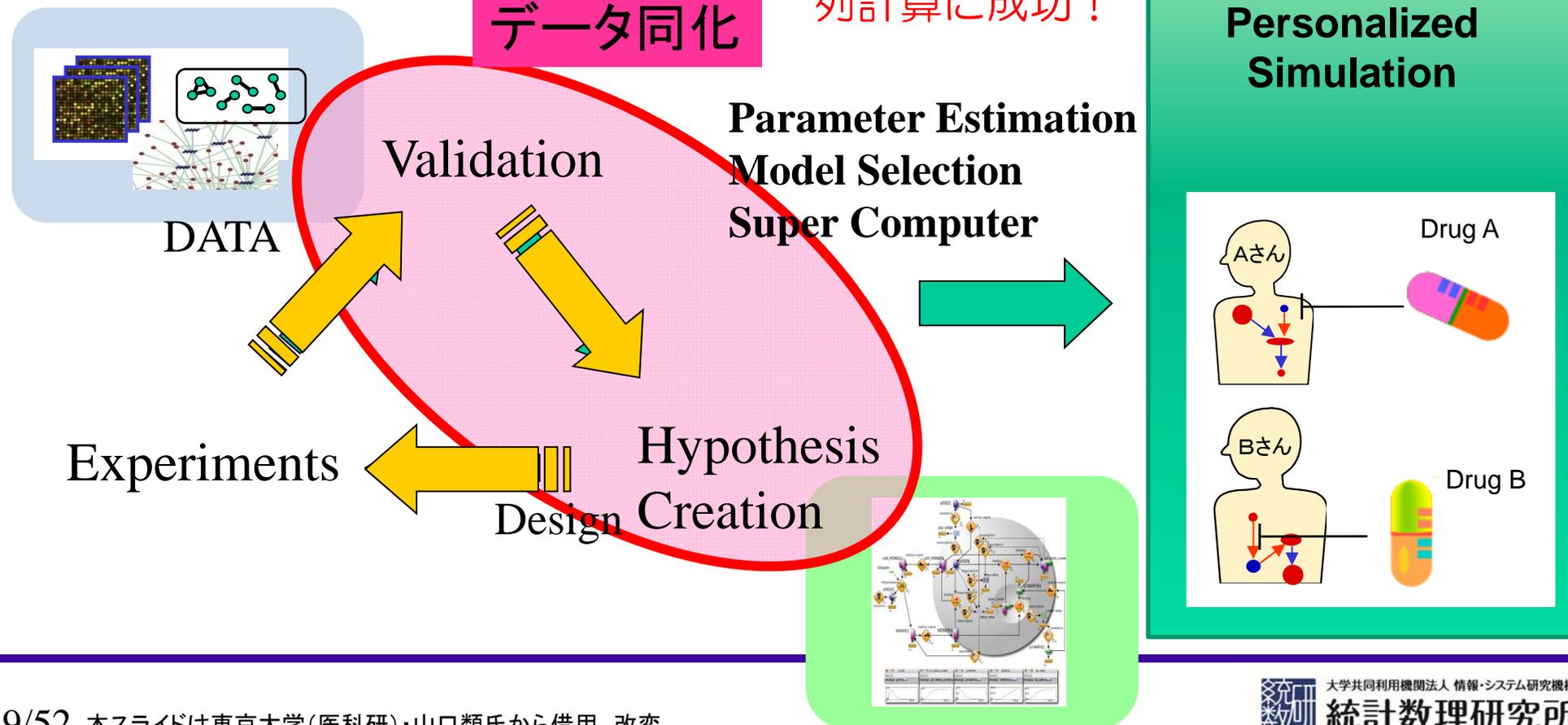
東京大学・医科学研究所・宮野研究室と共同研究

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



10万コアの超並列計算に成功!



地盤シミュレーションのデータ同化

京都大学(農学研究科)・村上教授グループとの共同研究

■実学的観点から: 地盤沈下の最終沈下量の予測

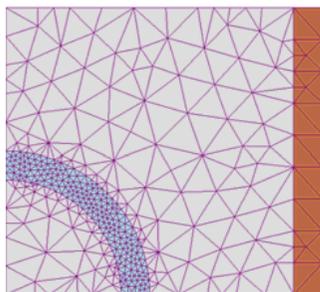
・関西国際空港の埋め立てにおける沈下量予測など、途中までの計測を用いて最終沈下量の予測をしたい。

・データ同化の観点からは、シミュレーション予報のための初期値生成作業。

■計算手法の観点から: 弾性体であれば、線形性・復元性によりデータ同化は簡単。(カルマンフィルタ(KF)が有効。実際の土は塑性をもつので、変形の経路依存性を考慮しなければならないので、KFは適用不可。

構造シミュレーション:

有限要素法、境界要素法



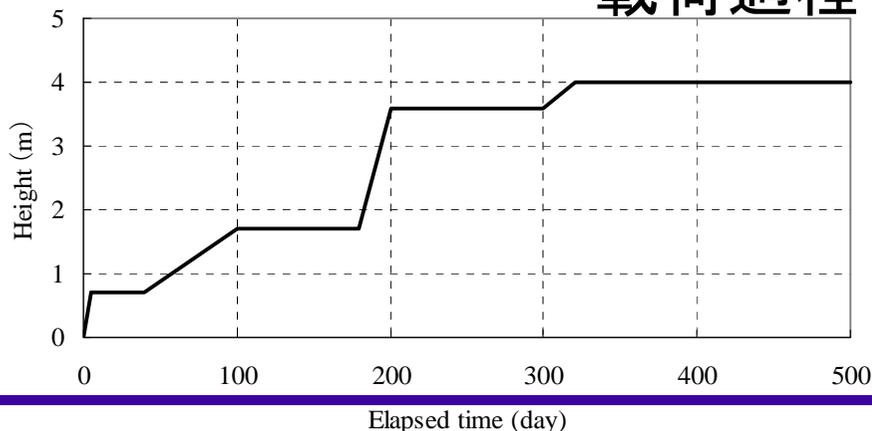
(Wikipediaから)



観測点

8点の観測(沈下量測定)を用いて、弾塑性パラメータを推定する

Embankment loading history 載荷過程



・塑性変形も加わると、非線形性に伴うシミュレーション計算の破綻(保存則を満たさない等)の問題が出てきてしまい、KFでは難しかったがSIS(逐次データ同化の一種)で解決できた。



(理研のホームページから)

前世紀： 物質（「もの」）を均質に大量に生産するシステム



21世紀： 個人化された情報サービスを提供するシステム

個人をターゲットにした商品・サービスの提供を効率的に行えるシステム

“コ” — 個人，個性，個別，固有一が大切！

ご静聴ありがとうございました。

バーチャルにリアルを埋め込むデータ同化： ベイズ、アンサンブル予測

ヒューマンモデリング、個人化情報サービス

(2011年4月刊行)



(2011年9月刊行)



(2011年11月刊行)



データ同化研究開発センターの紹介ビデオ(10分)
を作成しました。YouTubeにアップしています。