

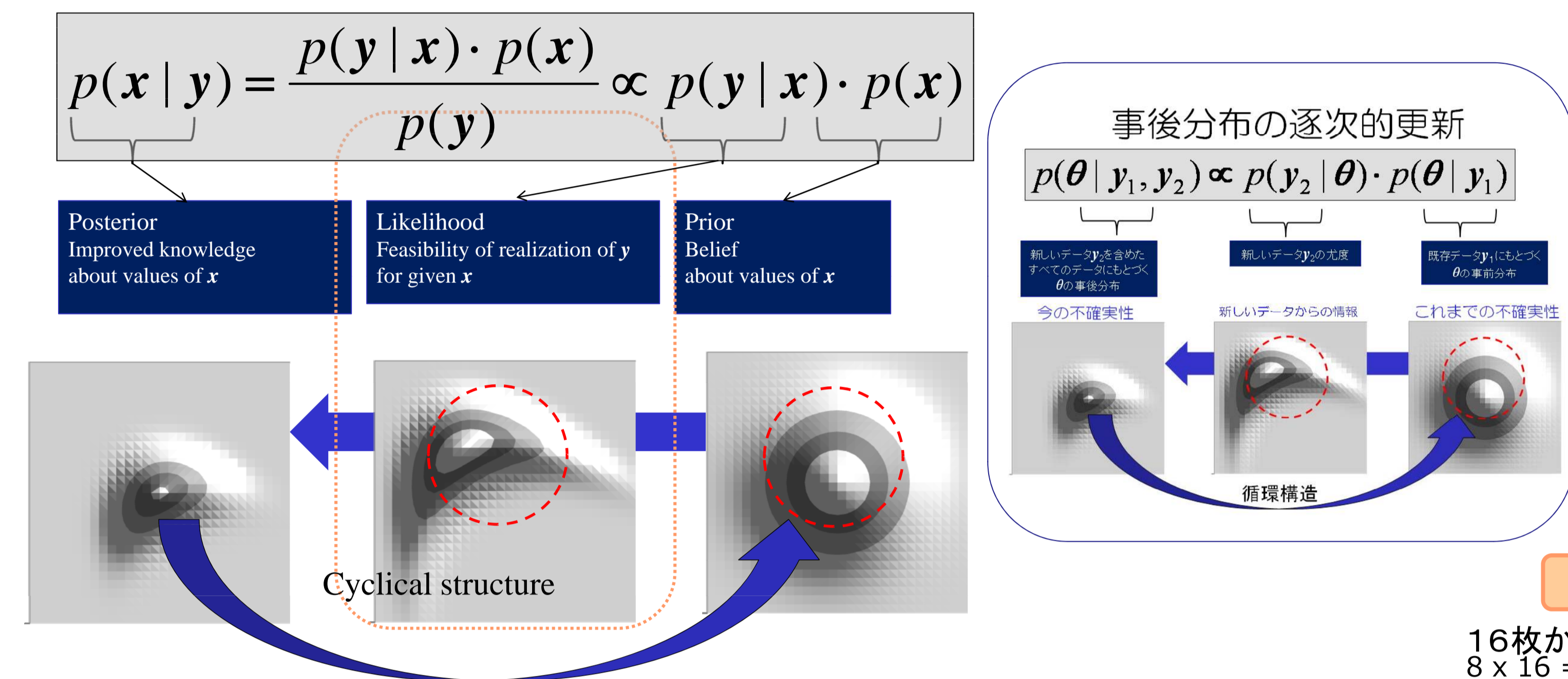
樋口 知之 モデリング研究系 教授

1. ベイズモデリングとは？

■動的かつ複雑な対象を理解する：対象から得られるデータは、さまざまな計測・観測条件に応じて時間的にも空間的にも多様な様相を呈する。この各状況の特性に即しつつデータもうまく説明できる表現方法をも手にしたのなら、予測や制御といった次のステップの作業が見通しよく完遂できる。通常我々はこの要求に、データyの生成メカニズムを、潜在変数ベクトルxを持つ統計モデルp(y|x)をもって近似することで解決を図る。

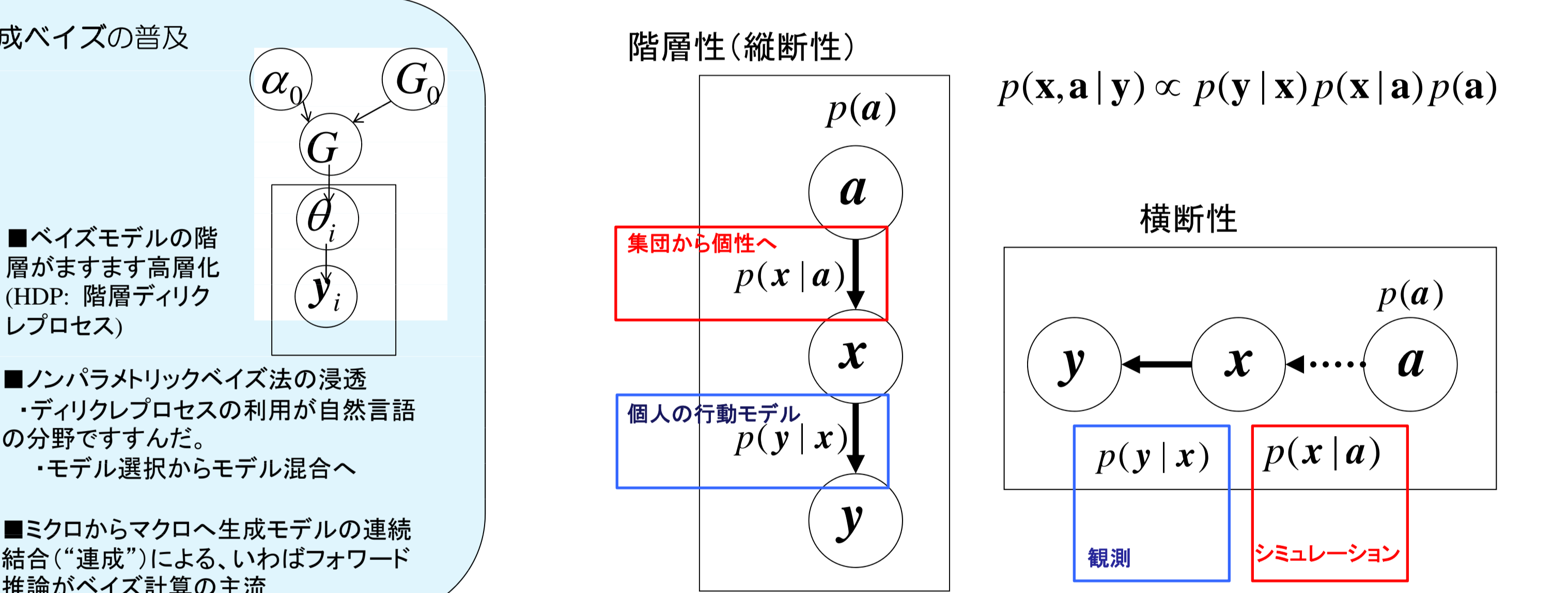
■事前分布：従来の統計的モデリングではxの次元をなるべく小さくすることが王道、さらにいうと美徳とされてきた。潜在変数を増やせば増やすほど統計モデルの記述能力は向上するが、汎化能力と呼ぶ将来のデータの予測能力が減少する。この問題への対策として、xについても統計モデルp(x)を想定するのがベイズモデルである。このp(x)をベイズ統計では事前分布と呼ぶ。

■ベイズの定理：事前分布の導入により以下のベイズの定理をもって、想定した事前分布p(x)がどのように修正されるのか、つまりxに関する不確実性がデータによりどの程度修正されたかを観察するのである。このp(x|y)を事後分布と呼ぶ。この仕組みにより多数のパラメータも安定して推論できるようになり、結果として高い予測能力とデータ記述能力を同時に持つ、総合的な統計モデルが構成できる。ここでは、この一連のモデル化行為をベイズモデリングと呼んでいる。



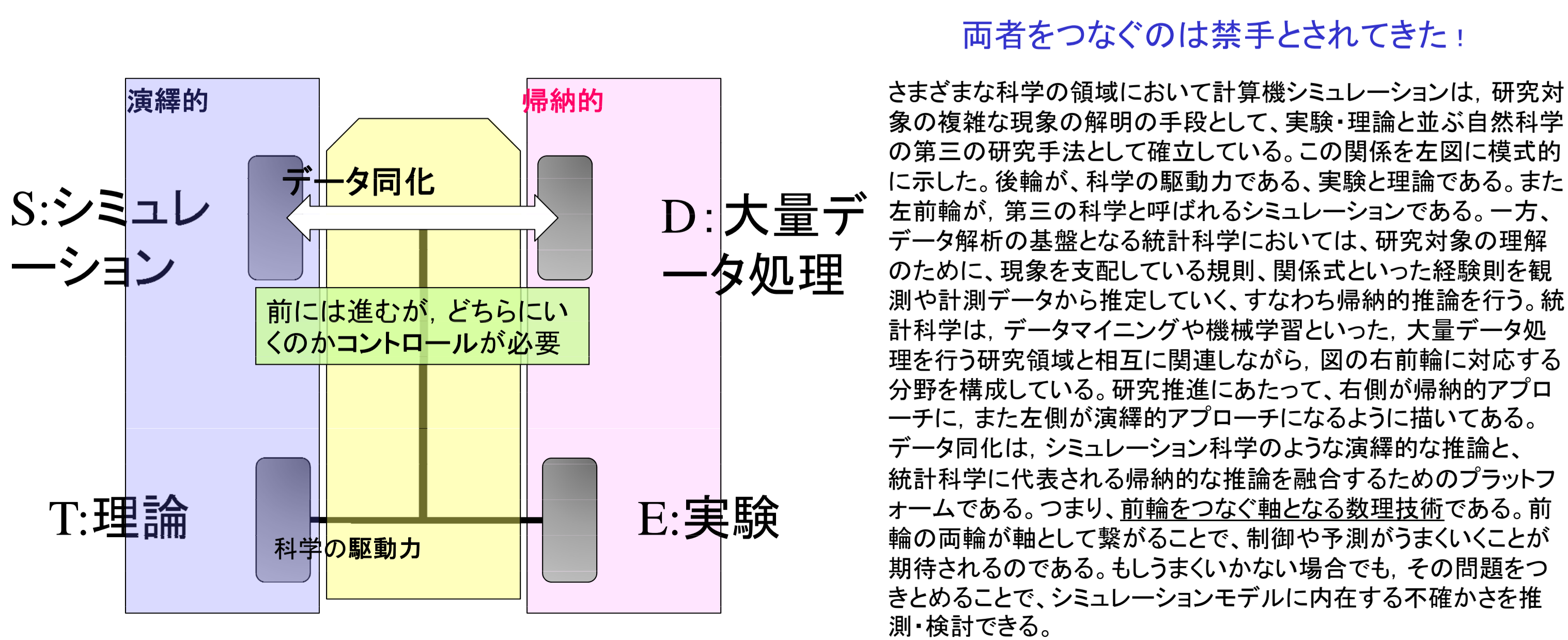
2. ベイズモデルの典型的な利用法

■階層ベイズモデル：事前分布の設定法、つまり事前分布への信念の置き具合はどのようにして決めるのかという疑問が沸く。これには、事前分布にもパラメータを導入することで自由度を残し、データ処理前の事前分布の決め打ちは避けることで対処する。つまり、事前分布をp(x|a)で与える。パラメータに対してさらに不確実性を許容し、p(a)を具体的に計算に導入し、さまざまな積分操作によって推論を行っていくのが階層ベイズ法である。



■対象の構造を自然にモデル化：上図の左側を使って、各個人の購買にまでモデル化を試みるマイクロマーケティングを例に階層ベイズモデルの利用法を考察する。最下層では個人の購買行動といったマイクロ単位の確率モデルを採用し、階層があがるにつれ、個人の特長、層別化された集団の特性、地域の特性、そして時代効果といったふうに、各demographicが階層ベイズモデルの階層に対応しているのが分かる。このように、マイクロマーケティングと階層ベイズモデルは親和性が高い。

■異種情報をつなぐ：階層構造を横にすることで、階層ベイズモデルは異種情報を統合するプラットフォームにもなり得る。このことを模式的に示したのが図の右側である。今、一番左端にあるように変数yとxの間の関係がp(y|x)で、またaの情報の不確実性をp(a)でモデル化されている状況を想定する。このとき、xとaの間の関係をもし何らかの形で確率的に表現し、それを統計モデルp(x|a)で与えることができたならば、データyに基づいたaの不確実性の評価が可能になる。図では、点線で示されているp(x|a)の関係を明示的に与えたならば(実線にしたならば)、いままで分離していたyに関する情報と、aに関する情報を統合することができる。この仕組みを利用した例がデータ同化(下図)やゲノムデータ解析である。



3. データ同化

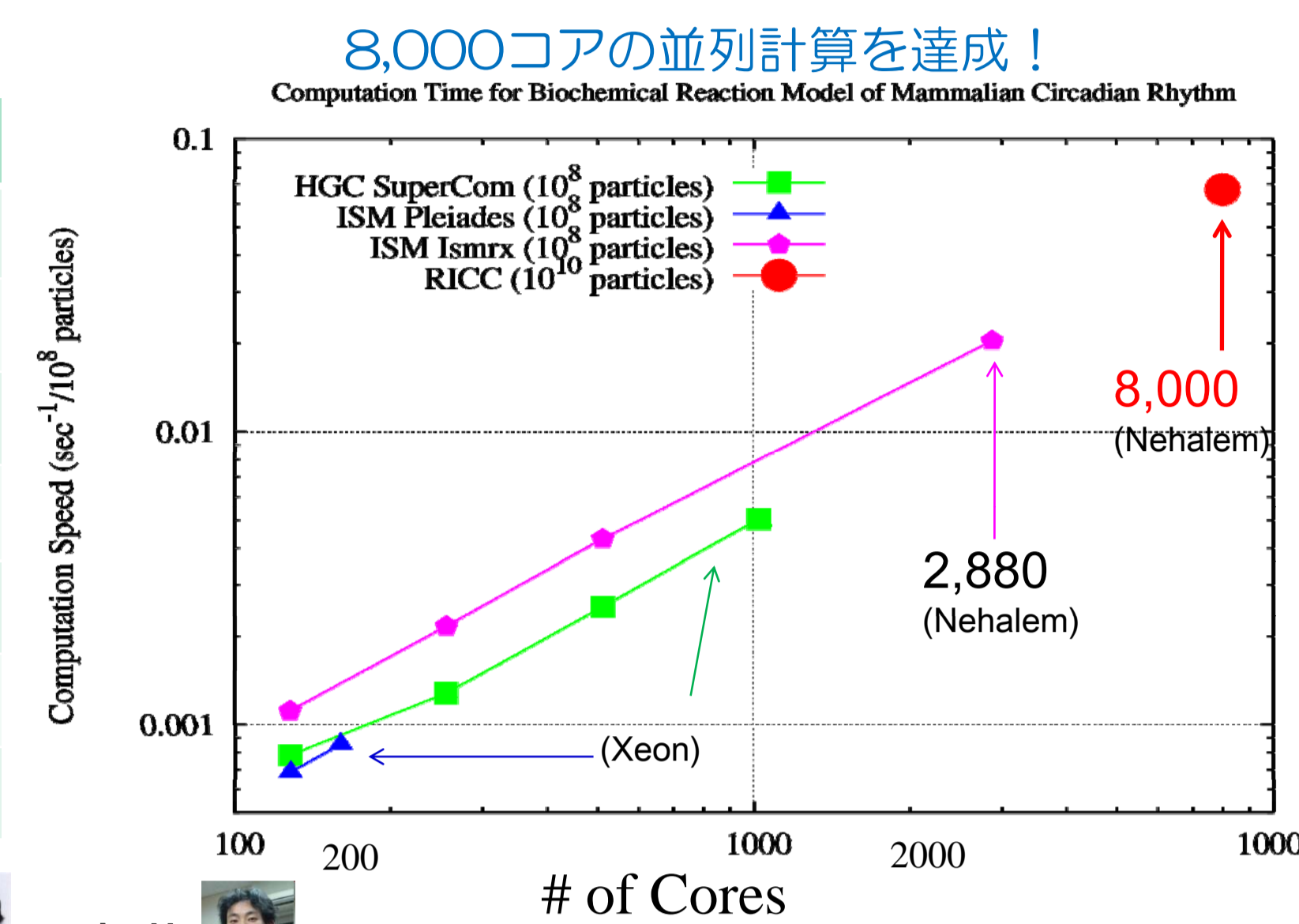
■手法：我々はデータ同化手法を統計科学の枠組みで正確に定義し、統計科学や情報科学の分野において蓄積されてきたアルゴリズムやモデリングに関する研究成果を利用しながら、これからの計算機インフラを視野に入れつつ、実装が平易かつ広い分野に適用できる逐次データ同化手法の開発を行っている。

※1：次世代スーパーコンピュータプロジェクトへ参加 (吉田、中野(慎)、長尾、齊藤らと。情報・システム研究機構では我々のみ国家プロジェクトに参加！)

汎用CPUを超並列につないだ分散メモリ型にHPCのプラットフォームは確実に移っていく。このHPC環境の中で、計算機の集約的利用に基づく統計計算、つまり計算集約型統計計算の計算効率を高めるためには、HPCのハード特性を最初から意識した手法の開発が肝要となるであろう。考慮すべきハード特性としては、階層性とネットワーク構造の二つが特に重要である。

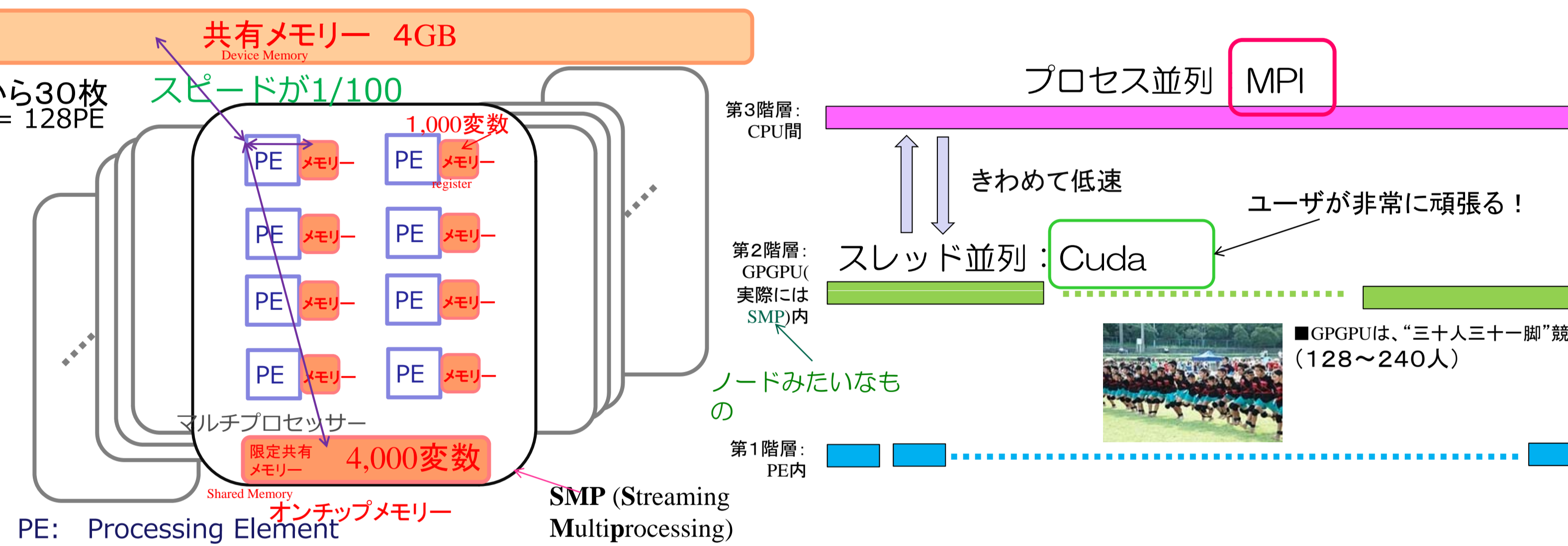
データ同化Gが利用している(予定の)スパコン

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
PetaCom* (RIKEN)	SPARC64 VIIIfx	2.00 GHz	> 80,000 nodes	> 640,000 cores	16GB/node
M System (JAXA)	SPARC64 VII	2.50 GHz	3,008 nodes	12,032 cores	32GB/node

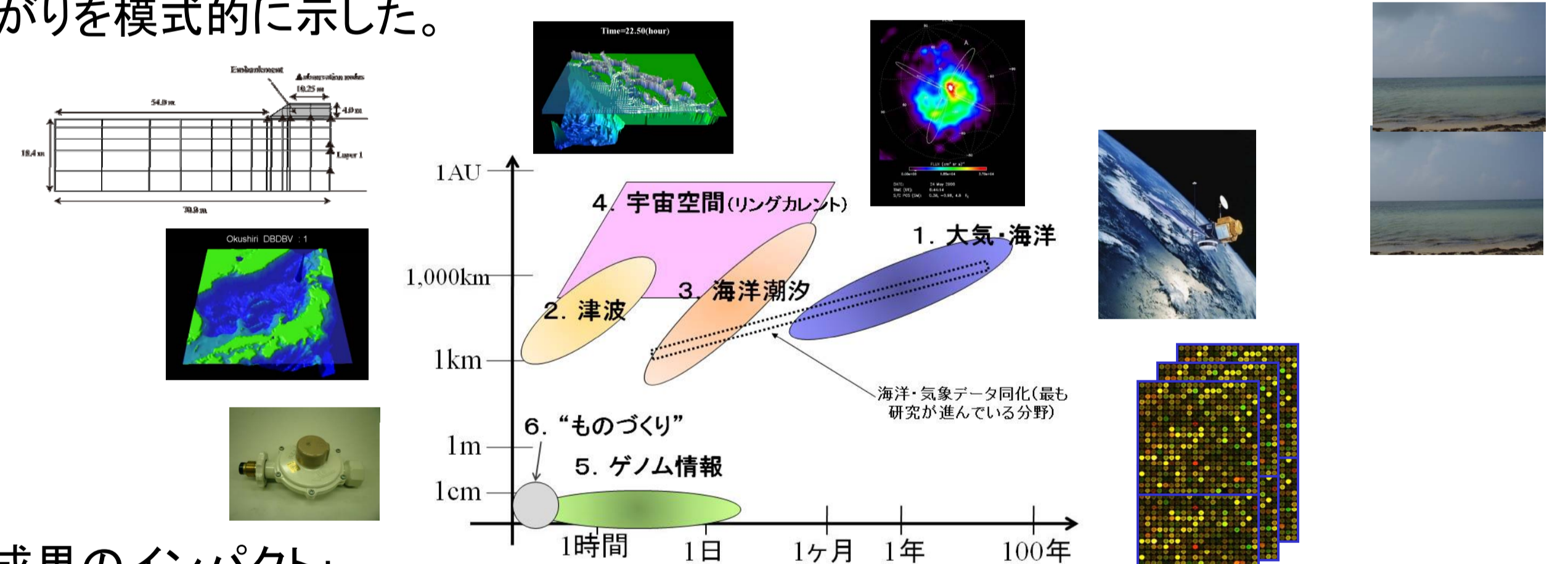


※2：GPGPUを用いた革新的高速化 (林、齊藤らと)

GPGPUとはGeneral-purpose computing on graphics processing unitsの略で、画像処理を専門とする補助演算装置であるGPUを画像処理以外の一般的な計算に用いる技術の総称である。GPUのクロック数自体はCPUと比較すると劣る(1.3GHz程度。CPUの約1/2~1/3)が、一つのインストラクションを多数の異なるデータ(統計でのデータの意味でなく、計算科学でのデータの意味)に同時に適用する操作、つまりSIMD型タイプとしては著しく性能が良い。この特性をいかした統計計算の高速化に取り組んでいる。



■適用分野：データ同化手法の研究を行うには、具体的なシミュレーションモデルとデータセットの二つ、つまり具体的テーマの選定が必要である。我々は主に、大気・海洋、津波、海洋潮汐、宇宙空間(リングカレント)、ゲノム情報、ものづくり、環境経済の7つの領域における新しいデータ同化実験に取り組んでいる。下図に、各応用領域における典型的な現象の時間スケールと空間スケールの広がりを模式的に示した。



■研究成果のインパクト：

- 1) 物理化学等の法則でなく、むしろデータを参照することを“第一原理”とした、初期条件・境界条件・パラメータの最適化と計算資源の集中化の研究を行うことで、逆問題解析の革新的アルゴリズムの開発に貢献。
- 2) 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等の、シミュレーション技術の個人化技術への発展的進化に貢献。

4. マイクロマーケティング

■個人化技術：個人化技術の実現に必要な基本的要素技術は、ケース数Nよりも属性変数の数Pが圧倒的に大きい状況、つまり『新NP問題』の克服法である。この困難を減じるためには、ある特徴量(属性)で似た値をとるものは、他の特徴量でも似た値をとることが期待できるといったようなさまざまな先見的知識を活用し、スパースな情報空間を確率的に埋めていくことが肝である。

■応用例：スーパーマーケットのPOSデータの解析と消費者行動のモデル化を通して、個人に特化したサービスが提供できるシステム構築の研究を行っている。

