

分野連携による新しい科学の創出

ビッグデータの有効利用 ビッグデータの有効利用例②: ゲノム解析

目標・目的、克服すべき学術的課題

- 計算機科学との連携による、飛躍的に増加した大量実験データ(ハイスループットデータ)から生物学的な発見を行おうとする研究手法(バイオインフォマティクス)の高度化・発展
- 遺伝子発現データから複雑な計算によって遺伝子間の関係を予測・推定する遺伝子ネットワーク解析の発展

従来の研究

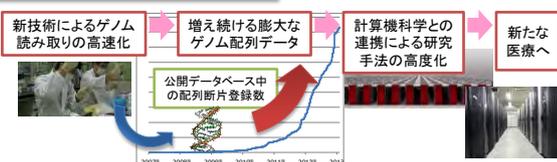
- 新技術により急速に蓄積、短時間で大量に出力されるデータを用いたバイオインフォマティクス
- 特定の条件下で観測された比較的小規模なデータから遺伝子ネットワークを推定し解析

分野連携の方策

- 多種多様で膨大なデータを組み合わせた解析

大規模計算で実現されること

- ゲノム情報の検索と効率的な解析法の開発
- 個人のゲノム情報に基づく最適な医療(個人ゲノム医療)のコストダウンによる一般的な治療化
- 解析対象をそれほど絞らずに数百程度のデータセット(数万サンプル)を抽出し、それに対して網羅的に推定ソフトウェアを適用するというアプローチ



ゲノム解析において今後、必要となる計算機性能を下表に示す。

課題	要求性能 (PFLOPS)	要求メモリ/バンド幅 (PB/s)	要求ファイル/O性能 (TB/s)	メモリ量/ケース (PB)	ストレージ量/ケース (PB)	計算時間/ケース (hour)	ケース数	総演算量 (EFLOP)	概要・計算手法	問題規模	備考
個人ゲノム解析	100	50	0.5	9	100	700	200	50,000,000	がんゲノム解析 200,000人分のマッピングおよび変異同定		変異検出に最低1000人の解析が必要なのでそれを1ケースとした。メモリ量は140Knode X 64GBで計算
疾患遺伝子発見のための統計的解析	500	500	5	200	2	140	5	1,300,000	ゲノムワイド連鎖解析(GWAS)	ヒトゲノム3Gbp x 200,000人分・1ケース4万人	メモリ量は800GB/node・ノード数25万を仮定

※ 本見積もりは、9月末日での見積もりである。未だ精査の余地があり、最終版では、より精度の高い数値を記載する予定である。

3.3 大規模実験施設との連携

(1) X線自由電子レーザー施設 SACLA 等の大型研究施設との連携

XFEL (X線自由電子レーザー) は、波の位相がきれいにそろったレーザーの性質を持つ超高輝度の X 線を発生させることのできる光源であり、これまで構造を解くのが難しかった非結晶粒子や微細結晶の構造解析に威力を発揮すると期待されている。2011年3月に完成した XFEL 施設 SACLA (Spring-8 Angstrom Compact free electron LAsEr) では1日に最大で約 500 万枚の回折パターンが得られるため大量のデータ処理を行う必要があり、計算機科学技術との連携の重要性が謳われている。特に、生体粒子の動態をも含めた 4 次元イメージングに期待が高まっている。

また、SACLA による非結晶試料のイメージング実験では、10 ナノメートルからマイク