

## 生命科学と情報科学の融合：一塩基多型（SNP）データベースの開発とその解析

理化学研究所 播磨研究所 メンブレンダイナミクス研究グループ

石野 洋子

ここ数年来のヒトや微生物も含めた動植物のゲノムプロジェクトの成果として、既に多くのゲノム塩基配列が公開されている。これに伴い、塩基配列の解析手法自体に力点が置かれていたこれまでの生命科学の研究も、得られた塩基配列情報をどう利用するかという方向へと移行しつつある。そこで最近注目を浴びているのが、生命科学と情報科学の融合した学問である、“バイオインフォマティクス”や“計算生物学”という分野である。

バイオインフォマティクスはもともと生物学のデータ管理に情報技術を応用させたもので、主にデータベースに関連した技術を指していた。しかし、最近では、生物システムの特性を数学的あるいは物理的モデルへ抽象化することや、データ解析のための新しいアルゴリズムを導入することなど、より広い範囲を指すようになった。疾病遺伝子の探索や遺伝子ネットワークの推定などは、最近の主なトピックスである。ゲノム塩基配列の大量なデータを解析し、有用な情報や知識を抽出するというバイオインフォマティクスの一連の研究の流れには、近年の計算機技術の進歩やインターネットの普及・発展が寄与するところが大きいのは論を待たない。さらにこのような研究を遂行するためには、情報技術の理解のみならず生物学的そして化学的知識も総合的に必要となる。本講演では、バイオインフォマティクスの一例として、モデル植物であるシロイヌナズナ (*Arabidopsis thaliana*) の一塩基多型 (Single Nucleotide Polymorphism: SNP) のデータベース開発とその解析を取り上げる。

ゲノム塩基配列は、同種の生物個体間であっても互いに若干異なっており、この違いが個体間の表現型の差異に反映している。このような塩基配列の違いを一般的に多型 (polymorphism) と呼んでいる。遺伝子多型は、ゲノム情報に基づく疾患遺伝子の解明、予防的な遺伝子診断などに重要な役割を果たす。多型にはいくつかの種類があり、次の3種類に大別できる。(1) DNAの1つの塩基が他の塩基に置き換わったもの (= SNP)、(2) 数十から数千の塩基が欠失 (または挿入) しているもの、(3) 数個から数十個程度の単位の塩基配列が繰り返し存在する部位においてその繰り返し回数が異なるもの。この中でなぜ今 SNP が注目されているのだろうか。それは第一に、他の多型と比べて出現頻度が高く、多型マーカーとして優れているためである。次に、高速・大量の SNP タイピング技術が実現化されつつあることが挙げられる。最後に、情報処理の観点からは、多型の状態を(0,1)の信号に置き換えることができるため、情報処理操作が容易であるという利点がある。最近では、異なる遺伝子座間の対立遺伝子の非独立性を定式化した連鎖不平衡(linkage disequilibrium)という概念に基づいて、体系的・網羅的に SNP マッピングを行い、原因遺伝子の探索を行う手法の有用性が広く認められるようになってきている。

本プロジェクトは、シロイヌナズナにおける SNP を体系的にサーベイし、その結果をデータベースに格納して Web を通じて世界に公開するというものである。University of Southern California の M. Nordberg と University of Chicago の M. Kreitman および J. Bergelson がリーダーとなり NSF からファンドを得て行っており (<http://walnut.usc.edu/2010/>)、SNP を利用した網羅的・体系的遺伝子マッピングの先駆けである。このプロジェクトは、生物学者、数学者、コンピュータ科学者たちの連携のもとで遂行され、演者はそのデータベースの設計・実装およびデータ解析部分を担当した。

