

課題名 (タイトル) :

全ゲノムシーケンスデータ解析パイプラインのスーパーコンピューター・理研 RICC 上での
動作確認とチューニング

利用者氏名 : ○角田 達彦*, 藤本 明洋*, 阿部 哲雄*, 中村 英二*, **

所属 : * 横浜研究所 ゲノム医科学研究センター 統計解析・技術開発グループ 情報解析研究チーム

**株式会社ダイナコム 研究部

1. 本課題の研究の背景、目的、関係するプロジェクトとの関係

近年の著しいシーケンス技術の発展により、個人ゲノムシーケンスが可能となった。超並列シーケンサーの発展は著しく、現在では 600Gbp (ヒトゲノムの約 200 倍) の塩基配列データが約 2 週間で得られている。全ゲノムシーケンスは世界中で活発に行われており、昨年は 1000 人ゲノム計画の論文が出版されたほか、がんの全ゲノムシーケンスも数多く報告されている。今後もシーケンサーのデータ産出量は増大していくことは確実であり、全ゲノムシーケンスは次世代の疾患研究において、極めて重要な役割を担っていくと考えられる。しかしながら、現在のシーケンサーには、読み取り長 (リード長) が短い、エラー率が高いなどの問題があり、現在に至るまで解析手法が確立されているとはいえない。そこで、我々はシーケンサーからのデータの解析プログラムの開発を行った (Fujimoto et al Nat Genet 42: 931-936)。この解析パイプラインのさらなる高速化を行うため、RICC へ移植作業を行った。

2. 具体的な利用内容、計算方法

次世代シーケンサーから得られたリード配列を、標準ゲノム配列にマッピングし、確率計算に基づいて一塩基多様性の同定を行った。先行研究で解析したデータを用いて、RICC 上で解析を行った。RICC 上で開発したコードは、「京」へ移植し、チューニングを行っている。

3. 結果

先行研究のデータの一部を用いて、解析パイプラインが正常に動作することを確認した。

4. まとめ

我々は、次世代シーケンサーの解析データを高精度に解析するパイプラインを構築し、RICC への移植作業を行った。小規模なテストデータを用いて、解析パイプラインが正常に動作することを確認した。さらに、RICC でテストしたパイプラインを「京」へ移植を行った。現在は主に「京」でのチューニングを行っている。

5. 今後の計画・展望

シーケンスコストの低下にともない、ゲノムデータの産出量は爆発的に増加している。今後は、シーケンスデータの大量解析に向けて、高並列化を行う必要がある。