

ゲノムワイド関連解析ソフトウェア PLINK 次期バージョン 1.9.0 における性能評価

金井 仁弘¹, 山根 健治^{1,2}, 樋口 千洋¹, 田中 敏博^{1,3,4}, 岡田 随象^{1,5}

1. 東京医科歯科大学 大学院医歯学総合研究科 疾患多様性遺伝学分野 2. ソニー株式会社 メディカル事業ユニット 研究開発部門 LE 開発部 3. 東京医科歯科大学 疾患バイオリソースセンター
4. 理化学研究所 統合生命医科学研究センター 循環器疾患研究グループ 5. 理化学研究所 統合生命医科学研究センター 統計解析研究チーム

概要

米ハーバード大のPurcellらによって開発されたPLINKは、ゲノムワイド関連解析 (GWAS) において広く利用されているソフトウェアである。我々は大幅な性能改善が謳われる次期バージョン 1.9.0 に対し、ベータ版を用いて具体的な処理性能を評価した。まず PLINK のサンプルデータ作成機能を用い、7 段階のサンプル数 (1,000 ~ 100,000) 及び SNP 数 (10,000 ~ 1,000,000) からなる計 49 通りのジェノタイプデータを作成した。本データセットに対し異なるバージョンの PLINK (1.06, 1.07, 1.90b) を用いて、一般的な GWAS データに適応される QC (quality control) 処理を実施した。その結果、1.90b は 1.06, 1.07 に比べて大幅な性能改善が認められ、大規模ジェノタイプデータに対する適合性が分かった。また、PLINK 1.90 はソースコード共有ウェブサービス GitHub 上でオープンソースとして公開されている。我々はこの GitHub を通じて、ベータ版に存在していたバグ除去に貢献したのであわせて報告する。

1. 背景

米ハーバード大のPurcellらによって開発されたゲノムワイド関連解析ソフトウェア PLINK の次期バージョン 1.9.0 beta が登場した。

Version	Release Date	Developer
1.06	Apr. 24, 2009	S. Purcell
1.07	Oct. 10, 2010	S. Purcell
1.90b	Jul. 1, 2014	S. Purcell, C. Chang

計算機・次世代シーケンサーの性能が年々向上し、大規模データセットを扱う機会が増えた昨今、解析ソフトウェアの処理性能改善は研究時間の短縮に大きく寄与する。また適切な研究計画を練る上でも、各解析処理にどのくらいの時間を要するのか把握することが重要である。

2. 方法

一般的な GWAS データの QC プロセスで用いられる、以下 6 つのコマンドについて各バージョン (1.06, 1.07, 1.90b) での処理時間を測定した。同一データの下でのバージョン間の

1	コールレート	--missing
2	ヘテロ接合度	--het
3	連鎖不平衡	--indep-pairwise
4	IBS/IBD	--genome
5	Hardy-Weinberg 平衡	--hardy
6	アレル頻度	--freq

比較に加え、全測定データを用いてバージョン毎に処理時間とサンプル・SNP 数の関係を計算した。またプロファイリングツール gprof を用いて、Genome コマンドの処理を計測した。

測定に用いたデータセットは 1.90b のサンプルデータ作成機能を用いて生成されたサンプル数・SNP 数、各 7 種類の計 49 通りのジェノタイプデータである。測定に用いたデータ・計算機の仕様は以下。

サンプル数	1,000 2,000 5,000 10,000 20,000 50,000 100,000
SNP 数	10,000 20,000 50,000 100,000 200,000 500,000 1,000,000
CPU	Intel Xeon CPU E5-2450 v2 @ 2.50GHz × 16
Memory	96 GB

3. 結果 1.90b の処理時間は圧倒的に短縮されていた！

サンプル数 m : 5,000 · SNP 数 n : 100,000 のデータセットにおける各処理と QC プロセス全体での所要時間をそれぞれ図 1, 2 に示す。1.90b は 1.06, 1.07 に比べて圧倒的に処理時間が短いことが分かる。QC プロセス全体 (図 2 左) では、**1.90b は 1.07 に比べて約 2,680 倍速かった** (1.90b: 15 秒、1.07: 11 時間 43 分)。

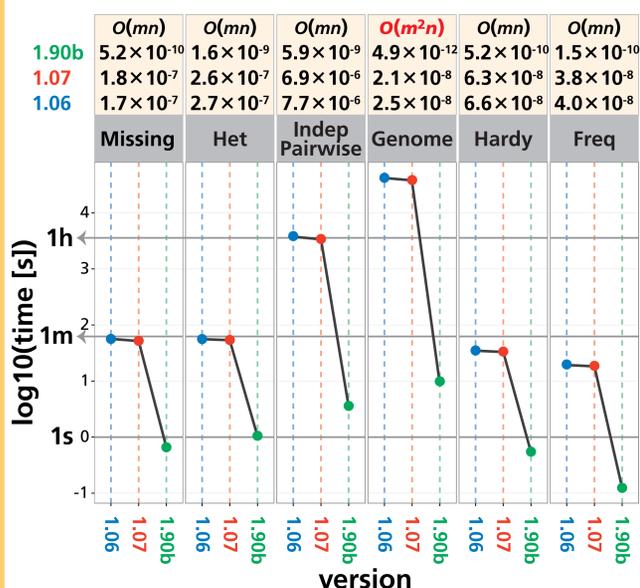


図 1. 各コマンドの処理時間の比較

上: 全測定データの基での比較 下: 同一データの基での比較

特に時間を要する Genome と IndepPairwise を除いた比較 (図 2 右) でも、1.90b は 1.07 に比べて約 66 倍速かった。また、全測定データを用いて処理時間 [s] とサンプル数 m · SNP 数 n の関係を計算したところ図 1 上表を得た。この値に mn をかけると概ねの処理時間を得ることが出来る (Genome のみ m^2n)。バージョン間での処理結果に丸め誤差以上の差異は認められなかった。

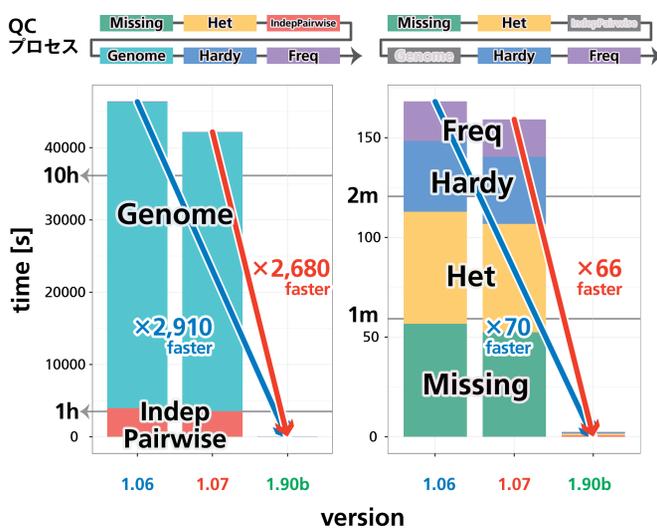


図 2. QC プロセスの所要時間の比較

左: 全コマンド 右: Genome · IndepPairwise 以外

gprof を用いて、PLINK の Genome コマンドの内部処理の様子を計測したところ図 3 を得た。PLINK 1.90b と 1.07 の間には抜本的な設計変更があるため単純な比較は出来ないが、やはり IBS/IBD の計算過程に大きな改善があることが分かる。また 1.07 はメモリ・文字列操作に関連する処理だけで 1.90b の処理時間に匹敵する時間を要した。

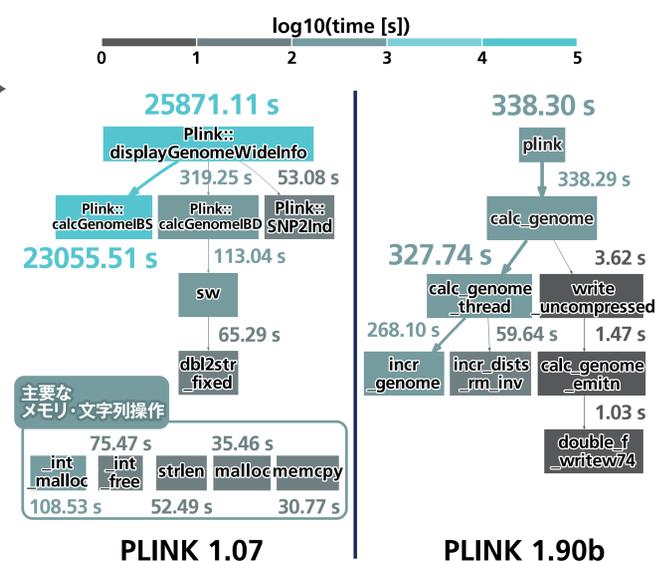


図 3. Genome コマンドの内部処理

4. ソースコードの解析

ソースコードレベルで確認すると、PLINK 1.90b で飛躍的に処理速度が向上した主要な要因として、以下の 3 つが挙げられる。

- 1 bit 演算や効率的なメモリアクセスといった抜本的な設計変更
- 2 アルゴリズムの改善・変更
- 3 並列計算への対応 (マルチスレッド、クラスター演算)

また gprof を用いた解析 (図 3) でも、これらの寄与が裏付けられた。

5. 結論

PLINK の次期バージョン 1.90b の性能評価を行ったところ、QC プロセス所要時間の比較において **現行バージョン (1.06, 1.07) に比べて約 2,500 ~ 3,000 倍速い**ことが分かった。また今回測定した処理においては、1.90b と現行バージョンの出力の間に明らかな差異は認められなかった。

ソースコードの解析や実行時プロファイリングの結果から、1.90b には抜本的な設計変更やアルゴリズムの改善が施されており、これらが高速化に大きく寄与していることが確認された。

開発版への貢献: バグ修正

開発中の PLINK 1.9 のソースコードが、GitHub を通じて公開されているため、開発者は自由にその設計を確認したり、機能の追加・修正を行ったりすることが出来る。我々も開発版の以下の機能のバグについて修正パッチを作成し、本体に取り込まれた。

- 1 VCF パーサーの half-missing call の取り扱い
- 2 PED パーサーの複数塩基・トリアレルの取り扱い
- 3 フェノタイプを含んだ共変量の出力

<https://github.com/chrchang/plink-ng>

連絡先: 金井 仁弘 mkanai.brc@tmd.ac.jp