# EM2: Enhanced Computational Algorithm for Haplotype-Based Association Analysis in Case-Control Studies

**Isaac Amundson, Kevin M. Bradley, Brian L. Yaspan, Jeffrey R. Smith**
**Departments of Computer Science, Medicine, and Cancer Biology**
**Vanderbilt University**

## Overview

Haplotype-based association analysis within a case-control study is broadly anticipated to offer a powerful approach to identify genetic variants causing common disease. The Expectation-Maximization (EM) algorithm is commonly employed to estimate haplotype frequency in unphased data, such as case-control genotypes. Accurate assignment of phase to enable tests of heterogeneity between case and control haplotype profiles is computationally very demanding. Use of permutation testing to estimate exact P values compounds this demand. This can limit the feasibility of large studies.

We modified algorithms implemented in the Fortran 77 program of Fallin and Schork (AJHG 67:947-959, 2000) to facilitate large-scale studies. Daniele Fallin kindly provided the original source code for these modifications. The original program limits each run to a maximum of 10 SNPs and 2000 samples, and does not accommodate multi-allelic markers such as STRs. In general these restrictions favor application speed. We altered the program to enable use of multi-allelic as well as bi-allelic markers, and to place no intrinsic limit on the number of markers or samples. "EM2" attains significant speed improvements despite the more generalized implementation. Our enhanced version is written in C++ and is developed for both Linux and Windows platforms. The program core incorporates object-oriented design pattern reuse and utilizes the C++ standard template library functionality, ensuring efficient algorithm usage. EM2 also employs parallel-processing. We are currently modifying permutation test code and anticipate further significant speed gains. Improved computational algorithms will be crucial for large-scale genome-wide association studies using a case-control design.

Several runs may be performed consecutively in batch mode, allowing the user to create a sliding window (subset of markers) over a region of interest. This is accomplished by creating a batch file specifying the markers in the window and the extent of marker overlap of the sliding window. The user specifies the number of processors for use, as well as other run parameters such as number of restarts, maximum number of iterations, and number of case/control permutations. An option also exists to only display haplotypes that exist within the population above a specified frequency. Output can be presented in HTML or text format, and/or displayed directly to the terminal. EM2 runs from the command line, and can be employed as a stand-alone application or incorporated into a larger application (such as an interactive website).

## Tests

| Program | Processors | Markers | |
|---|---|---|---|
| | | 10 htSNPs | 7 htSNPs & 3 STRs |
| EM2 | 8 | 5 hr, 1 min | 6 hr, 49 min |
| | 1 | 25 hr, 48 min | - |
| Fallin & Schork | 1 | 58 hr, 3 min | N/A |

As a test of relative speed we employed a laboratory dataset of 735 breast cancer cases and 735 unaffected controls at 10 htSNPs and 3 STRs spanning 38 kb at the CYP11A gene. There was no missing experimental data for these samples. Missing data for a given sample in a specified window of markers would result in the sample being dropped from the analysis. SNP minor allele frequency ranged from 5.1% to 43.4%. STR heterozygosity ranged from 0.52 to 0.79. Run parameters were: 15 restarts, 150 maximum iterations, 10,000 permutations, rare haplotypes (<5%) grouped for overall chi-square test of heterogeneity, and a single window of 10 specified markers. Tests were executed on an SGI Altix 3000 configured with eight 1.3 GHz Itanium 2 processors and 16 GB of DDR RAM, running 64 bit Red Hat Linux.

The convergence results, as part of the maximization step, show log-likelihood information associated with the overall, case, and control data sets. Estimated frequencies are provided for observed haplotypes for overall, case, and control data sets. Haplotypes are represented by SNP alleles and STR lengths in base pairs. Chi square and p value is calculated for a given haplotype in case v.s. control groups, and calculated for overall haplotype profile across the two groups ($\chi^2$ $n$df). Permutation test results are also presented for each haplotype, estimating exact p values. For each permutation, the case/control status of each sample in the entire dataset is randomized using a random number generator and random seed value, and the EM algorithm re-run. Fallin and Schork's omnibus likelihood ratio test is also presented, assessing the significance of overall haplotype profile difference between cases and controls.

## Output



Note that the highlighted haplotype of our example includes the same CYP11A allele originally associated with breast cancer risk in the full study population of 1193 cases and 1310 controls. In the published study, a lone STR marker comprised of a [(TAAAA)$_n$] repeat of the distal promoter region was analyzed by conditional logistic regression (Zheng, W., Gao, Y., Shu, X., Wen, W., Cai, Q., Dai, Q., and Smith, J. (2004) *Cancer Epidemiology, Biomarkers & Prevention* **13**, 709-714). In that study, allele 281 of STR CYP11A_91557-97 was in significant excess in cases (12.56% freq) vs. controls (8.46% freq), p <0.0001, OR 1.6 (1.3-1.9). This data illustrates that specific alleles of each STR marker and htSNP CYP11A_87195 each mark the disease-associated haplotype. Note that at this gene, alleles of the STRs as well as htSNPs efficiently tag the study population haplotypes. EM2 identifies ten htSNP-STR and six htSNP-only haplotypes of 1-5% frequency, which were grouped for analysis and not displayed per our run option.