Reliable and Fast Estimation of Recombination Rates by Convergence Diagnosis and Parallel Markov Chain Monte Carlo

Jing Guo, Ritika Jain, Peng Yang, Rui Fan, Chee Keong Kwoh, and Jie Zheng

Abstract—Genetic recombination is an essential event during the process of meiosis resulting in an exchange of segments between paired chromosomes. Estimating recombination rate is crucial for understanding the process of recombination. Experimental methods are normally difficult and limited to small scale estimations. Thus statistical methods using population genetics data are important for large-scale analysis. LDhat is an extensively used statistical method using rjMCMC algorithm to predict recombination rates. Due to the complexity of rjMCMC scheme, LDhat may take a long time for large SNP data sets. In addition, rjMCMC parameters should be manually defined in the original program which directly impact results. To address these issues, we designed an improved algorithm based on LDhat implementing MCMC convergence diagnostic algorithms to automatically predict values of parameters and monitor the mixing process. Then parallel computation methods were employed to further accelerate the new program. The new algorithms have been tested on ten samples from HapMap phase 2 data set. The results were compared with previous code and showed nearly identical output. However, our new methods achieved significant acceleration proving that they are more efficient and reliable for the estimation of recombination rates. The stand-alone package is freely available for download http://www.ntu.edu.sg/home/zhengjie/software/CPLDhat.

Index Terms—Recombination hotspot, reversible jump MCMC, convergence diagnosis, parallel computation, genome instability

1 INTRODUCTION

M EIOTIC recombination occurs in the pairing of homologous chromosomes in meiosis leading to the generation of novel gene combinations. The transfer of genes from parents into offspring by genetic recombination during meiosis is a major engine of genetic variation [1]. The meiotic recombination events break down the genealogical history within a genome which is critical for analyses of genetic variations [2]. The improper segregation of chromosomes can lead to aneuploidy, a significant risk factor for fetal loss and developmental disability in humans [3]. In addition, deleterious variations can be removed from the gene pool by recombination.

The rate and location of meiotic recombination have implications for understanding of recombination process and its evolution. They vary markedly between species and among individuals. The estimation of the rate at which recombination occurs can theoretically provide guidance for biologists to explore biological problems, e.g., gene targeting,

- J. Guo, R. Jain, R. Fan, and C.K. Kwoh are with the School of Computer Engineering, Nanyang Technological University, Nanyang Ave, Singapore 639798.
- P. Yang is with the School of Computer Engineering, Nanyang Technological University, Nanyang Ave, Singapore 639798, and with the Institute for Infocomm Research, Agency for Science, Technology & Research, 1 Fusionopolis Way, Singapore 138632.
- J. Zheng is with the School of Computer Engineering, Nanyang Technological University, Nanyang Ave, Singapore 639798, and with the Genome Institute of Singapore, Agency for Science, Technology, and Research, Biopolis, Singapore 138672. E-mail: zhengjie@ntu.edu.sg.

Manuscript received 9 Oct. 2013; accepted 14 Oct. 2013; date of publication 23 Oct. 2013; date of current version 7 May 2014.

For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below. Digital Object Identifier no. 10.1109/TCBB.2013.133 mutation mechanisms [4]. Tracking distance between two genes on a chromosome by recombination rate could detect the presence of certain disease-causing genes [5].

Obtaining accurate prediction of recombination rates could be challenging and prohibitively expensive through direct experimental methods. Sperm typing produces high-resolution estimates; however, this procedure is complex, only applicable for male [6], and limited to small scale prediction. Hence, indirect statistical methods are useful. Patterns of genetic variation among DNA sequences have been used to analyze recombination rate [7]. Hudson [8] proposed a composite-likelihood estimator of the population recombination rate that combines the coalescent likelihoods of all pairwise comparisons for segregating sites. McVean et al. [9] extended Hudson's method to allow for a finite-sites mutation model, and also introduced a likelihood permutation test. Later the heterogeneity implied by recombination hotspots is incorporated to improve the accuracy [10]. Li and Stephens [11] developed a method considering all loci simultaneously rather than pairwise comparisons based on an approximation to the conditional likelihood (implemented in PHASE). Instead of approximate likelihood method, Wang and Rannala [12] proposed a full-likelihood Markov chain Monte Carlo method (implemented in InferRho).

The algorithm of [10] has been implemented in the program LDhat package. It has been extensively used for detection and calculation of variable recombination rates in population genetic data via composite likelihood method. A typical change point scheme of reversible jump Markov chain Monte Carlo (rjMCMC) algorithm [13] is employed to

^{1545-5963 © 2013} IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications_standards/publications/rights/index.html for more information.

predict the interval constant rates. The final recombination rate is composed of the morphology of hotspots and background rate. Due to the complexity of rjMCMC scheme, LDhat is a time-consuming program that would take several hours to generate results. In addition, the accuracy of the outputs and execution time are determined by parameters of rjMCMC which can only be manually specified by users. The rjMCMC parameters consist of a set of numbers, including transition times, initial samples discarded called 'burnin' and sampling frequency. Insufficient running would cause unstable status of the Markov Chain. Conversely, over calculation would waste extra resources.

To address the above issues, we propose an improved algorithm for the prediction of recombination rates based on LDhat. First we evaluated the performances of LDhat identifying the bottleneck of running time and testified the impact of rjMCMC parameters on recombination profile and time complexity. Second, to avoid manually setting the parameters of rjMCMC, we incorporated algorithms for MCMC convergence assessment to automatically predefine those arguments. In addition, the chain convergent status is monitored during iteration process until it reaches the target distribution. Then we made use of parallel computation methods in order to further speed up the process of calculation.

In order to evaluate our new algorithm, we utilized 10 sets of test samples extracted from HapMap phase 2 data. We compared the recombination profiles, running time and iteration numbers of the original LDhat program and our improved methods. The result showed that our methods achieved significant speedup without affecting the accuracy of outputs. The parallel computation method resulted in even more significant reduction of execution time with identical outputs.

2 Метнор

In this section, we analyzed the LDhat program to identify the most time-consuming part. In addition, we evaluated the influence of parameters, i.e., iteration number and single nucleotide polymorphism (SNP) number, on output profiles. In allusion to rjMCMC scheme, we proposed an improved algorithm applying MCMC convergence diagnostic methods and parallel computation.

2.1 LDhat Program Analysis

LDhat (specifically, the *rhomap* program) employs the rjMCMC algorithm which incorporates genomic polymorphisms to estimate the pairwise constant rates by composite likelihood. Composite likelihood [8] is an approximation of the coalescent likelihood [14] which is more easily implemented and based on independent pairwise single nucleotide polymorphisms to estimate the recombination rate ρ . According to the composite likelihood estimator, the maximum-likelihood estimate of ρ can be obtained as the maximum product of conditional-likelihood functions of all independent pairs in *n* samples. These two-locus conditional-likelihoods of a fix *n* samples can be precalculated and stored for future researches.

However, the ad hoc estimator have underestimated the effects of mutations to genetic variation assuming infinite site model whereby mutation rate θ tends to be negligible. To address this problem, Auton and McVean extended Hudson's work to provide an improved estimation procedure incorporating mutation models [10]. Suppose that the population mutation rate θ is constant across the sequence, it is estimated by Watterson algorithm [15]. Then pairwise segregating sites with two alleles are classified into equivalent sets for further likelihood calculation. The execution burden depends on the number of segregating sites with an order of n_{seq}^3 where n_{seq} is the number of haplotypes. Assuming that pairs of SNPs are independent, given the number of haplotypes, all of the possible combinations of allelic states could be consulted from tabulated files which contain precalculated likelihoods. Then the likelihood of each pair of segregating sites is estimated over a grid extracted from those files.

In addition to the contribution of background rate, the morphology of hotspots that reveals the relationship between recombination and genome features [16], [17] is incorporated into the pseudoposterior distribution of recombination rate. The rjMCMC algorithm is implemented to determine the parameters of the mutation model, i.e., change-points of SNPs, background rate, hotspot locations, hotspot heat and hotspot scale.

The scheme of rjMCMC algorithm is a typical change point problem [13] (Appendix Table A.1, available in the online supplemental material). Set *L* as the position of the last SNP, and let *k* be the number of change-points drawn from a Poisson distribution. The locations of change points are s_i , where $0 < s_1 < s_2 < \cdots < s_k < L$. The recombination rate is given by a step function x(.) on [0, L].

In the algorithm of LDhat, the interval background rate h_j on the *j*th block $[s_j, s_{j+1}]$ is initialized with the prior as exponential distribution, denoted as $P(h_j) \sim Exp(\phi)$. And the prior on the *k*th hotspot rate is defined as a truncated double-exponential curve, presented as $f_k \propto \lambda Laplace(\mu, b)$, where *b* is the central position of hotspot. Two parameters λ and μ are defined for evaluating the heat and scale of hotspots. The priors on both of them are in gamma distribution, i.e., $\lambda \sim \Gamma(\alpha_1, \beta_1), \mu \sim \Gamma(\alpha_2, \beta_2)$. The hyperparameters α_1 , $\beta_1, \alpha_2, \beta_2$ were obtained by Maximum Likelihood estimation to fit a gamma distribution to empirical hotspot data sets [10]. The contribution of a recombination hotspot to the final recombination rate depends on its relative location to blocks.

The mutation model is designed as four independent random transitions for background model: (a) 'death'of a randomly chosen block, (b) 'birth'of a new block at a randomly chosen location in [0, L], (c) a change of the height of a randomly chosen block, and (d) a change of the position of a randomly chosen block. In addition, there are five transitions for hotspot model: (a) 'delete'of a randomly chosen hotspot, (b) 'insert'of a new hotspot at a randomly chosen location in [0, L], (c) a change of the heat of a randomly chosen hotspot, (d) a change of the scale of a randomly chosen hotspot, and (e) a change of the position of a randomly chosen hotspot.

According to the rjMCMC algorithm, mutations occur during each transition. To compute the Metropolis-Hastings acceptance ratio, the recombination map composed of constant rates and a pseudo-likelihood of the data in each transition have to be calculated which take most of the execution time. For N iterations, the complexity of LDhat program scales with an order of $N \times l_{seq}$, where l_{seq} is the number of SNPs. For large scale matrix, the running time of LDhat is prohibitively lengthy. Moreover, in order to get accurate recombination rate profiles, an appropriate setting of parameters, specifically the number of iterations, is critical to guarantee that a Markov chain reaches its equilibrium distribution.

Therefore, we evaluated the influences of parameters on LDhat. Two-sample Kolmogorov-Smirnov test [18] is employed to compare two outputs. The control object **x** is attained from the results of the same data set with 11 million iterations, 1,000,000 burn-in and 2,000 sample. The accuracy of the predicted recombination profile **x**' with *n* intervals is estimated by Kolmogorov-Smirnov statistic, defined as

$$KSZ = \sqrt{\frac{n}{2}} - \max_{i} |x_i - x'_i|, \quad where \ i = 0, 1, 2 \dots n - 1.$$
(1)

Given rejection level $\alpha = 0.05$, the reference $KSZ_{ref} = 1.36\sqrt{\frac{2}{lseq}}$. First, a test data set is applied to examine the impact of the number of iterations on execution time and recombination profile. Then we use a group of data sets with different sizes to analyse the correlation between data size and running time with the same number of iterations. The running time shows approximately linear correlations between the number of iterations and data sizes which is consistent with the analysis on LDhat complexity. In addition, the accuracy of outputs is highly correlated with the two parameters.

Thus in the original LDhat program, a major limitation is that the parameters have to be defined by users without references. Estimation of the parameters of MCMC, such as the iteration number and the number of discarded initial samples, is a critical issue for the application of LDhat. To address the issue, several algorithms have been developed to determine how many steps are needed to ensure the convergence of Markov chains. However, due to the complexity and specialization, direct prediction of parameters is only theoretically described and thus impractical [19]. Hence a variety of empirical tools for the diagnosis of MCMC convergence which are well designed and implemented are used, e.g., Gelman and Rubin diagnostic method [20], Brook and Giudici's method [21], Raftery and Lewis diagnostic algorithm [22].

In the next section, we presented our convergence diagnostic method based on the framework of Raftery and Lewis diagnostic algorithm [22], thereby solving the above major issue of LDhat, as well as speeding up large-scale estimation of recombination rates. For standard Markov Chain Monte Carlo algorithms, the dimension of the parameter vector is fixed, whilst in rjMCMC scheme it has varying dimensions. Normal convergence assessment algorithms cannot be applied directly to outputs from an rjMCMC sampler. Castelloe and Zimmerman's method extends the work of [20] by encompassing all of the parameter spaces and monitoring several parameters simultaneously [23] which is especially designed for rjMCMC situation. Thus we employed the method of Castelloe and Zimmerman in this paper to monitor the status of Markov chain.

2.2 Convergence Diagnostic Methods

Here we propose an improved algorithm for the prediction of meiotic recombination rates which makes use of convergence diagnostic methods. The original LDhat program not only takes large amount of time on calculation, but also requires users to specify the values of parameters which cannot ensure the convergence of Markov chains. Thus our main purpose is to control the process of rjMCMC iteration to monitor the Markov chain convergent status and supervise the adaptation of parameters in order to accelerate the mixing process and ensure the accuracy. To achieve these goals, the key point is to determine the appropriate number of iterations for the convergence of Markov chains.

Raftery and Lewis diagnostic algorithm and Castelloe and Zimmerman convergence assessment method are adopted in our program. The former is widely used to predict the number of iterations, burn-in and sample parameters in MCMC applications, and the latter cannot predict parameters, but it is dedicated for rjMCMC convergence diagnosis. In our program, the numbers of iterations and burn-in are determined by Raftery and Lewis diagnostic algorithm for a given level of precision. Castelloe and Zimmerman's convergence assessment method runs periodically to check if the chain has reached its target distribution.

First, a pilot chain is run with initial iterations. Using the output sample, Raftery and Lewis diagnostic algorithm will generate a new Markov chain to predict how many steps are needed for each parameter to get equilibrium status and how long the burn-in should be. A reliable factor I will be calculated as well. Values of I must be greater than 1, but when I > 5 it often indicates problems [22]. Since a bad starting value or high posterior correlations may cause unreliable results, the estimations are just used as references of initiate settings. The set of parameters satisfying the threshold with the maximum number of iterations will be chosen as the input arguments. Sometimes all of the reliable factors are larger than 5. In that case, none of the results are reliable. The minimal iteration number is selected as the initial values of parameters. Castelloe and Zimmerman convergence assessment method will check the status repeatedly to rectify the values.

Castelloe and Zimmerman's algorithm needs multiple testing chains running for certain steps. It is particularly designed for rjMCMC convergence diagnosis. Since the solution dimention of rjMCMC is not fixed, different models with variant sizes of parameter vectors are generated. Markov chains transit between these models. Thus the evaluation of the variations within each sample, testing chains, different models, could reflect the convergent status of whole Markov chains.

Let *C* be the number of chains required by Castelloe and Zimmerman's method, *T* be the number of sweeps in each chian and θ be a vector of parameters. *M* is the number of distinct models visited by any chain. R_{cm} stands for the number of times model *m* occurrs in chain *c*. The total variance \hat{V} is estimated by

$$\hat{V} = \frac{1}{CT - 1} \sum_{c=1}^{C} \sum_{m=1}^{M} \sum_{r=1}^{R_{cm}} (\boldsymbol{\theta}_{cm}^{r} - \overline{\boldsymbol{\theta}}_{..}) (\boldsymbol{\theta}_{cm}^{r} - \boldsymbol{\theta}_{..})', \qquad (2)$$

where $\overline{\theta}_{...}$ is the average of all samples. Variation within chains, variation within models and variation within models and chains are defined as

$$W_{c} = \frac{1}{C(T-1)} \sum_{c=1}^{C} \sum_{m=1}^{M} \sum_{r=1}^{R_{cm}} (\theta_{cm}^{r} - \bar{\theta}_{c.}) (\theta_{cm}^{r} - \bar{\theta}_{c.})', \quad (3)$$

$$W_m = \frac{1}{CT - M} \sum_{c=1}^{C} \sum_{m=1}^{M} \sum_{r=1}^{R_{cm}} (\boldsymbol{\theta}_{cm}^r - \overline{\boldsymbol{\theta}}_{.m}) (\boldsymbol{\theta}_{cm}^r - \overline{\boldsymbol{\theta}}_{.m})', \qquad (4)$$

$$W_m W_c = \frac{1}{C(T-M)} \sum_{c=1}^C \sum_{m=1}^M \sum_{r=1}^{R_{cm}} (\boldsymbol{\theta}_{cm}^r - \overline{\boldsymbol{\theta}}_{cm}) (\boldsymbol{\theta}_{cm}^r - \overline{\boldsymbol{\theta}}_{cm})'.$$
(5)

The above four factors reflect the convergence in different levels. The convergence assessment algorithm is used to check whether they reach the stable states. Four ratios in equations (6), (7), (8) and (9) are created to evaluate the chain mixing status. When $MPSRF_1$ and $MPSRF_2$ are settled close to 1, \hat{V} and W_c , W_m and W_mW_c are all settled approximately to a common value, indicating that it has achieved the desired distribution of convergence:

$$PSRF_1 = \frac{max \ eigen \ V}{max \ eigen \ W_c},\tag{6}$$

$$PSRF_2 = \frac{max \ eigen \ W_m}{max \ eigen \ W_m W_c},\tag{7}$$

$$MPSRF_1 = max \ eigen \ [W_c]^{-1}\hat{V},\tag{8}$$

$$MPSRF_2 = max \ eigen \ [W_m W_c]^{-1} W_m. \tag{9}$$

The convergence algorithm needs to calculate the inverse matrixes and eigenvalues. The complexities of them are about $O(n^3)$, where *n* is the dimension of the matrix, equal to the number of parameters. When the diagnosis algorithm is frequently called with a large number *n*, it may take considerable time. To accelerate the process, a new parameter *addon* is defined to control this process. It is initially set to 2 and will be added by 5 percent of the number of SNPs in each round of convergence diagnosis. Then the iteration number will be set to *addon* times the estimated value. Once the *C* chains are diagnosed as convergent, the final output is generated by combining the results of all chains.

In our improved program for estimating recombination rates, a convergence diagnostic model is invoked to estimate



Fig. 1. Workflow of convergence diagnosis model.

MCMC parameters and monitor the convergence process. The computational workflow is shown in Fig. 1.

Instead of manually setting values, an automatic definition process of parameters is initially run, then a convergence diagnosis procedure repeatedly to check the status. Finally, the results from each chain are combined to a final recombination profile. Since the rjMCMC assessment method requires multiple chains for diagnosis, the sequential scheme would make this new method even more timeconsuming than running a single chain when the iteration loop is fixed. However, in most cases the iteration number is unknown to users, thus we cannot directly compare our improved method with the original program with the same number of iterations.

2.3 Parallel Method

Due to the increasing availability of cheap computing power, parallel computing has received impetus. It has long been employed for scientific computing. Next, we are concerned with parallel implementation of MCMC in the context of accelerating our convergence diagnostic method.

Current algorithms for parallelizing MCMC can be classified into two main categories: one is parallelization of a single chain, and the other is parallel generation of multiple different chains [24]. Conceptually, parallel processing can be applied to almost any problem. However, MCMC is not easy to run in parallel owing to its serial nature. Due to the tight synchronization requirements of MCMC, the singlechain parallelization strategy requires considerable modification of the serial algorithm [25]. While the total iterations could be divided by multiple processors due to the independence nature of samples [25]. We concentrate on introducing a parallel algorithm that significantly decreases the execution time on multiple short Markov chains.

Assuming that each iteration takes roughly the same time to compute, an iteration may be used as a unit of time. Since the samples collected from MCMC chains are independent, it is possible to allocate the n required samples to N available processors, where the same program is run on each processor.

For a long chain the burn-in only happens once, whereas for several short chains, each must have a respective burn-in, resulting in many wasted samples [25]. With increasing numbers of processors, the performance of

Parallel Algorithm :
Step 1. Master program runs with <i>n</i> iterations, <i>b</i> burn-in, and passes $b + (n-b)/N$
iterations to each available processor.
Step 2. Each processor k:
(a) Simulates $b + (n-b)/N$ independent realizations of Markov chain
(b) Passes result M_k back to master program
Step 3. Master program combines the $M_0, M_1 \dots M_{N-1}$ to get final M .
Step 4. Master program returns M/N

Fig. 2. Parallel algorithm for a single Markov chain.

parallel computation becomes limited owing to redundant burn-in. Thus the issue of burn-in is of particular concern in a parallel computing environment. Here we make each process with the same burn-in phase identical with the sequential program.

We take advantage of parallel computation integrating convergence diagnosis model. The scheme of integrated method is nearly the same as convergence diagnosis program except that the convergent diagnostic tasks are divided by N processors. Each of the C diagnostic chains is run on N/C processors. The algorithm for parallel simulation of a single Markov chain can be described in Fig. 2.

The theoretical speedup is proportioned to the number of processors. However, due to the burn-in overhead, if N processors run one chain with a burn-in of b and n total iterations, then b + (n - b)/N iterations are allocated for each chain. When neglecting the communication time between processors and the handling time on file combination, it gives an optimal speedup of

$$SpeedUp_1(N) = \frac{n}{b + \frac{n-b}{N}}.$$
 (10)

Let n = 10b, then $SpeedUp_1(5) = 3.5714$, $SpeedUp_1(8) = 4.7059$. However, when using 10 processors, there is only fivefold speedup indicating that the effect of parallel computation on large clusters becomes limited. Theoretically, due to the burn-in, it could reach a maximum of 10 times speedup, when $N \rightarrow \infty$. Furthermore, effective utilization

of multiple processors is also limited due to the aggregation of communication time.

Based on a single chain parallel algorithm, the C diagnostic chains are divided into multiple sub-tasks. In Fig. 3, it shows the strategy of the parallel approach employed in our program. Each processor runs independent copies of the program with n/N iterations, and generates individual output files. The length of burn-in period keeps the same ratio with sequential execution. These numerous files are then compiled to obtain the final outputs of diagnostic chains for convergence evaluation. The implementation of this approach is done using OpenMPI programming language for communicating messages between multi-core processors.

Suppose that *C* parallel chains run on *N* processors with consistent burn-in of *b* and *n* total iterations, then n/C iterations and N/C processors are assigned to each chain. This gives a speedup of

$$SpeedUp_2(N) = \frac{n}{b + \frac{n/C-b}{N/C}}.$$
(11)

Empirically n = 10bC, C = 5 is a useful rule-of-thumb. We can get $SpeedUp_2(5) = 5$, $SpeedUp_2(10) = 9.09$, $SpeedUp_2(20) = 15.38$ with a maximum speedup of 50 theoretically.

Another issue we have addressed is the random number generator. The correlation among random number streams on separated processors should be reduced by assigning identical random number seeds to each machine [26]. The original LDhat program uses the default random number generator provided by the C language. It can only be used by one processor at a time. The other processors need to wait for their turn to obtain the random number losing the benefit of parallelization. Since each loop of the program requires random number generation over ten times, the over-all impact of improving the random number generator can be very significant. A sophisticated approach is to change the random number generator to SIMD-Oriented Fast Mersenne Twister (SFMT), which supports multicore



Fig. 3. Parallel algorithm for multiple chains.

TABLE 1 Test Data Sets

Datasets	Haplotypes	SNPs	Length(kb)
Test1	48	61	9.730
Test2	180	100	38.771
Test3	120	110	35.449
Test4	50	251	504.492
Test5	120	401	796.496
Test6	70	520	1102.571
Test7	70	610	983.183
Test8	60	790	1385.201
Test9	60	850	1437.376
Test10	50	1000	2643.03

parallel random number generation and has been shown suitable for use in Monte Carlo simulations [27]. Therefore we replace the random number generator in order to make it applicable for parallel computation.

3 RESULTS

To investigate the performance of the new method, we conducted two comparison studies. Since the program is used for fundamental genetics studies, it is imperative that the optimization techniques used do not affect the results. The new program should expedite the calculation process meanwhile retaining the accuracy. Not only the recombination rates but also the change positions should be predicted within the acceptable deviation. Hence we analyzed recombination profiles, running time and iteration numbers to evaluate the performance of the new method. We use *LDhat* to refer to the original LDhat implementation, *CLDhat* to refer to the convergence method, and *PLDhat* to refer to the parallel approach.

Ten sets of test data with equal iterations are used to evaluate the performance. They are drawn from human genomes with different numbers of haplotypes, SNPs and sequence lengths (Table 1). In the first study, we compared recombination profiles on outputs of 10 data sets by *LDhat*, *CLDhat* and *PLDhat*. In the second study, the execution time and the number of iterations are compared to show the efficiency of our improved programs.

The experiments were implemented on an IBM cluster of 24 quan-CPU 2.53 GHz Intel Xeon Linux Systems, connected to each other by 100 Mbps Ethernet connections.

3.1 LDhat Analysis

Before the comparison studies, we analyzed the influence of parameters, i.e., the numbers of SNPs and iterations, on execution time and recombination profiles. First, a test data set with 61 SNPs is used to examine the effect of iteration number on output profiles and execution time. The iteration numbers are set 3,000, 6,000, 10,000, 15,000 and 18,000 respectively. Assuming that each mutation transition consumes the same time, the iteration loop in Line 6 of LDhat pseudocode (Appendix Fig. A.1, available in the online supplemental material) constructs the main component of rjMCMC leading to a significant linear correlation between iterations and execution time. The comparison results in Fig. 4a shows an approximately linear relationship with r = 0.992. Thus controlling the loop length is important for the speedup



Fig. 4. Execution time (a) and KSZ (b) analysis on 61SNP test data set with different iterations. The dotted line in (b) represents the reference KSZ value for the data sets with an acceptance probability of 0.95.

of LDhat. With increasing iterations, the KSZ value decreases gradually (Fig. 4b) indicating that the Markov chain is close to the target distribution. However, when the iteration is set too small, e.g., less than 6,000, in this case, a deviation occurs in the output profile. On the contrary, over calculation with a large number of iterations would waste computational time without gain of additional information.

We analyzed the correlation between the number of SNPs and execution time, the number of SNPs and recombination profile with a fixed iteration number 10,000. As demonstrated above, the time complexity of LDhat has a linear correlation with the number of SNPs (Fig. 5a, r = 0.999). The red line in Fig. 5b shows the threshold of *KSZ* values for different SNPs with an acceptance value of 0.95. For small scale data sets, the setting of 10,000 iterations is enough for the convergence of Markov chains. When the number of SNPs exceeds 400, more training is required to make the chains reach the target distribution.

3.2 Comparison of Recombination Profiles

In the first study, we conduct experiments to compare recombination profiles of *LDhat*, *CLDhat* and *PLDhat* on the 10 data sets in Table 1. For all data sets, *LDhat* is running for 11 million iterations, and the initial 1,000,000 samples are discarded as burn-in. Samples of the chain are taken every 2,000 iterations after the burn-in. Then the output recombination rates are recorded as control groups to evaluate other methods. By contrast, we don't have to specify the numbers of iteration, burn-in and samples in the *CLDhat* method. By convention, five chains are generated to check the mixing status [28].

For *PLDhat*, the sequential procedure is divided into five parallel tasks making use of 15 processors. One processor



Fig. 5. Execution time (a) and KSZ (b) analysis on different size of test data sets with the same iterations. The dotted line in (b) represents the reference KSZ value for the individual data set with an acceptance probability of 0.95.

operates as the master running the Raftery and Lewis diagnostic algorithm to estimate the parameters and control the process of convergence assessment. Every three processors are applied to generate a single chain with parameters received from master processor.

Comparing the output graphs, the *CLDhat* and *PLDhat* methods got almost the same figures as the original program (Appendix Fig. A.1, available in the online supplemental material). Although the peak values are slightly changed in some points, the outputs showed high correlation coefficients among the three methods. The error is acceptable by KS test (see methods). Fig. 6 shows the *KSZ* of 10 data sets for our improved methods with reference values in red line. In most cases, the *KSZ* values of *CLDhat* are smaller than



Fig. 6. KSZ values on the 10 data sets by *CLDhat* and *PLDhat*. The solid line indicates reference values.



Fig. 7. Prediction of numbers and positions of hotspots on Test 2. (a) Posterior distribution of numbers of hotspots. (b) Posterior density estimates of positions of hotspots conditional on the number of hotspots k = 3 (solid curves) and k = 4 (broken curves).

PLDhat. This may be due to the loss of accuracy in frequent 'split-and-combine' process during parallel computation. But for small data sets, it converges more quickly with no need for frequently calling diagnosis and combination. Even so, both the outputs of *CLDhat* and *PLDhat* are under the threshold indicating the accuracy of our new methods.

Take Test 2 for instance. Fig. 7a shows the posterior distribution of the number of hotspots. More than 50 percent of models contain three hotspots. Conditional on the numbers of hotspots k = 3 and 4, Fig. 7b shows the posterior densities of the step positions. The positions of three hotspots are accurately identified. The density estimates are obtained using a Gaussian kernel with standard deviation.

3.3 Comparison of Running Time

In the second study, the total time consumed by *LDhat* and the improved methods on the 10 data sets are shown in Fig. 8a (details in Appendix Table A.2, available in the online supplemental material). The execution time was tremendously decreased when using our new methods. There are almost 80 times speedup in *CLDhat*. Using *PLDhat* on 15 processors, we got 622 times acceleration. In Fig. 8b, it shows the separate running time of each test data for *LDhat*, *CLDhat* and *PLDhat*. As the number of SNPs increases, our methods take a linear growth in time which is consistent with previous analysis in Section 2.1.

Unlike *LDhat* program, *CLDhat* method is a non-parameter approach under rjMCMC scheme. It is a more reliable and faster method. The mixing process is automatically monitored and checked periodically for convergence. So the



Fig. 8. Running time (in second and log scale) on the 10 data sets by LDhat, CLDhat and PLDhat. (a) Total time spent by each method. (b) Separate time on each test data of different methods.

MCMC chain could reach the equilibrium distribution rapidly in moderate iteration. In Table 2, the iteration numbers of test data sets by *CLDhat* are significantly decreased compared with *LDhat* leading to an expressively optimization of time efficiency. The *PLDhat* approach has successfully obtained more significant speedup than *CLDhat*.

We replace the original random number generator with SFMT for parallel computation. Since the program frequently requests for random number generation and SFMT is an efficient and faster random number generator, the replacement of original function reaches approximately three times speedup (data not shown). For large data sets, such as the calculation of Test 7-10, they take more iterations for convergence when the accelerating effect by parallel becomes more apparent.

The parameter *addon* controls the span length of each diagnosis round that correlates with the number of SNPs which makes large data sets mix faster. Conversely, this jumping scheme is suboptimal for small data sets.

4 DISCUSSION AND CONCLUSIONS

The main purpose of optimization of LDhat is to decrease the time complexity and increase the accuracy and reliability of output recombination profiles. Besides, there are no strategy to set the rjMCMC parameters in the original LDhat program, such as the iteration number, burn-in length and sample frequency. The bottleneck identified as the main loop in the original LDhat program is normally suggested to carried out a million iterations or more which may result in over calculation or insufficient running.

In this paper, we exploited MCMC convergence diagnostic algorithms and proposed two improved methods based

TABLE 2 Numbers of iterations by *CLDhat* and *PLDhat*

Datasets	CLDhat	PLDhat
Test1	156109	22561
Test2	185691	23768
Test3	293663	33746
Test4	297257	68746
Test5	390294	108746
Test6	1525275	144821
Test7	2395086	158746
Test8	3215746	297466
Test9	3695314	396121
Test10	4721811	447752

on LDhat. A major advantage of the new methods is significant acceleration compared with original program. In addition, the parameters are automatically estimated by our algorithms and only depend on input data. The mixing process is dynamic and monitored until the Markov chain reaches its target distribution. This could avoid unnecessary consumption of resources while also guarantees the accuracy of outputs.

Although the running time of the convergence method is tremendously decreased compared to the original program, it was further improved by implementation of parallel computation method due to the sequential scheme of the generation process of diagnostic chains. Hence we developed a parallel algorithm to allocate separate tasks to individual processors running a single chain in parallel. It achieves significant speedup.

The outputs of the above two methods were compared with the original LDhat program which showed similar output graphs. Since the results were generated through strict convergence assessment procedure, our methods achieved low values of KSZ (i.e., high accuracy) in much less iterations presenting extraordinarily similar recombination rate profiles.

Therefore our improved programs provide efficient and accurate methods for recombination rate prediction. Especially the parallel program provides a practicable, time saving and effective method. The improved methods, *CLDhat* and *PLDhat*, including the original LDhat (*rhomap*) program are implemented in a stand-alone package written in Java which is freely available for download at web site http://www.ntu.edu.sg/home/zhengjie/software/CPLDhat/. It could run in both Linux and Windows OS.

APPENDIX A

Tables A.1 and A.2 and Fig. A.1 can be found on the Computer Society Digital Library at http://doi.ieeecomputersociety.org/ 10.1109/TCBB.2013.133.

ACKNOWLEDGMENTS

This project was supported in part by Singapore Ministry of Education (MOE) AcRF Tier 1 Grant RG32/11. Jie Zheng is the corresponding author.

REFERENCES

- E. Shabanova, "Patterns of Genetic Recombination and Variation in the Human Genome," dissertation, Univ. zu Kln, 2009.
- [2] J.C. Avise, *Phylogeography: The History and Formation of Species*. Harvard Univ. Press, 2000.

- [3] B.L. Dumont and B.A. Payseur, "Genetic Analysis of Genome-Scale Recombination Rate Evolution in House Mice," *PLoS Genetics*, vol. 7, no. 6, pp. 847-856, 2011.
- [4] U. Mller, "Ten Years of Gene Targeting: Targeted Mouse Mutants, from Vector Design to Phenotype Analysis," *Mechanisms* of Development, vol. 82, no. 1, pp. 3-21, 1999.
- [5] N.J. Risch, "Searching for Genetic Determinants in the New Millennium," *Nature*, vol. 405, no. 6788, 2000.
- [6] R. Hubert, M. MacDonald, J. Gusella, and N. Arnheim, "High Resolution Localization of Recombination Hot Spots Using Sperm Typing," *Nature Genetics*, vol. 7, no. 3, pp. 420-424, 1994.
 [7] R.A. Gibbs, J.W. Belmont, P. Hardenbol, T.D. Willis, F. Yu, H.
- [7] R.A. Gibbs, J.W. Belmont, P. Hardenbol, T.D. Willis, F. Yu, H. Yang, L.-Y. Ch'ang, W. Huang, B. Liu, and Y. Shen, "The International HapMap Project," *Nature*, vol. 426, no. 6968, pp. 789-796, 2003.
- [8] R.R. Hudson, "Two-Locus Sampling Distributions and Their Application," *Genetics*, vol. 159, no. 4, pp. 1805-1817, 2001.
- [9] G. McVean, P. Awadalla, and P. Fearnhead, "A Coalescent-Based Method for Detecting and Estimating Recombination from Gene Sequences," *Genetics*, vol. 160, no. 3, pp. 1231-1241, 2002.
- [10] A. Auton and G. McVean, "Recombination Rate Estimation in the Presence of Hotspots," *Genome Research*, vol. 17, no. 8, pp. 1219-1227, 2007.
- [11] N. Li and M. Stephens, "Modeling Linkage Disequilibrium and Identifying Recombination Hotspots Using Single-Nucleotide Polymorphism Data," *Genetics*, vol. 165, no. 4, pp. 2213-2233, 2003.
 [12] Y. Wang and B. Rannala, "Bayesian Inference of Fine-Scale
- [12] Y. Wang and B. Rannala, "Bayesian Inference of Fine-Scale Recombination Rates Using Population Genomic Data," *Philosophical Trans. Royal Soc. B: Biological Sciences*, vol. 363, no. 1512, pp. 3921-3930, 2008.
- [13] P.J. Green, "Reversible Jump Markov Chain Monte Carlo Computation and Bayesian Model Determination," *Biometrika*, vol. 82, no. 4, pp. 711-732, 1995.
- [14] P. Beerli and J. Felsenstein, "Maximum-Likelihood Estimation of Migration Rates and Effective Population Numbers in two Populations using a Coalescent Approach," *Genetics*, vol. 152, no. 2, pp. 763-773, 1999.
- [15] G. Watterson, W.J. Ewens, T. Hall, and A. Morgan, "The Chromosome Inversion Problem," J. Theoretical Biology, vol. 99, no. 1, pp. 1-7, 1982.
- [16] A.J. Jeffreys, A. Ritchie, and R. Neumann, "High Resolution Analysis of Haplotype Diversity and Meiotic Crossover in the Human TAP2 Recombination Hotspot," *Human Molecular Genetics*, vol. 9, no. 5, pp. 725-733, 2000.
- [17] M.I. Jensen-Seaman, T.S. Furey, B.A. Payseur, Y. Lu, K.M. Roskin, C.-F. Chen, M.A. Thomas, D. Haussler, and H.J. Jacob, "Comparative Recombination Rates in the Rat, Mouse, and Human Genomes," *Genome Research*, vol. 14, no. 4, pp. 528-538, 2004.
- [18] H.W. Lilliefors, "On the Kolmogorov-Smirnov Test for Normality with Mean and Variance Unknown," J. Am. Statistical Assoc., vol. 62, no. 318, pp. 399-402, 1967.
- [19] M. Plummer, N. Best, K. Cowles, and K. Vines, "CODA: Convergence Diagnosis and Output Analysis for MCMC," *R News*, vol. 6, no. 1, pp. 7-11, 2006.
- [20] A. Gelman and D.B. Rubin, "Inference from Iterative Simulation Using Multiple Sequences," *Statistical Science*, pp. 457-472, 1992.
 [21] S. Brooks and P. Giudici, "Markov Chain Monte Carlo Conver-
- [21] S. Brooks and P. Giudici, "Markov Chain Monte Carlo Convergence Assessment via Two-Way Analysis of Variance," J. Computational and Graphical Statistics, vol. 9, no. 2, pp. 266-285, 2000.
- [22] A.E. Raftery and S. Lewis, "How Many Iterations in the Gibbs Sampler," *Bayesian Statistics*, vol. 4, no. 2, pp. 763-773, 1992.
- [23] J.M. Castelloe and D.L. Zimmerman, "Convergence Assessment for Reversible Jump MCMC Samplers," technical report, Dept. of Statistics and Actuarial Science, Univ. of Iowa, vol. 313, 2002.
- [24] J. Ye, A.M. Wallace, A. Al Zain, and J. Thompson, "Parallel Bayesian Inference of Range and Reflectance from LaDAR Profiles," J. Parallel and Distributed Computing, vol. 73, no. 4, pp. 383-399, 2012.
- [25] D.J. Wilkinson, "Parallel Bayesian Computation," Statistics: A Series of Textbooks and Monographs, pp. 477-508, Chapman and Hall/CRC, 2006.
- [26] A. Brockwell, "Parallel Markov Chain Monte Carlo Simulation by Pre-Fetching," J. Computational and Graphical Statistics, vol. 15, no. 1, pp. 246-261, 2006.
- [27] M. Saito and M. Matsumoto, "SIMD-Oriented Fast Mersenne Twister: A 128-Bit Pseudorandom Number Generator," Monte Carlo and Quasi-Monte Carlo Methods 2006, pp. 607-622, Springer, 2008.

[28] M.K. Cowles and B.P. Carlin, "Markov Chain Monte Carlo Convergence Diagnostics: A Comparative Review," J. Am. Statistical Assoc., vol. 91, no. 434, pp. 883-904, 1996.





degrees in biomedical engineering from Southeast University, Nanjing, China. She has been working toward the PhD degree at the School of Computer Engineering, Nanyang Technological University (NTU), Singapore, since 2013. She worked as a project officer in BioInformatics Research Centre, NTU, before the PhD program. She is now focusing on bioinformatics algorithm design for meiotic recombination hotspots and stem cell development modeling.

Jing Guo received the bachelor's and master's

Ritika Jain received the bachelor's of engineering degree in computer science from the School of Computer Engineering, Nanyang Technological University, Singapore, in 2012. She is currently working as a technical analyst in Barclays Capital, Singapore, as a C++ developer. She specialized in machine learning during her bachelor's degree. Her research interests include data mining, machine learning, and parallel processing.



Peng Yang received the bachelor's of computer science degree from the Beijing University of Posts and Telecommunications (BUPT) in 2008, and the PhD degree in computer science from the School of Computer Engineering, Nanyang Technological University (NTU) in 2010. His PhD supervisors are Professor Chee Keong Kwoh and Xiaoli Li. He is currently working as a scientist in the Data Analytics Department of Institute for Infocomm Research (I2R). Before joining I2R, he worked as a teaching assistant in Nanyang

Technological University from 2011 to 2012.



Rui Fan received the BSc degree in computer science and math from Caltech, and the PhD degree in computer science from MIT. He is an assistant professor in the School of Computer Engineering at Nanyang Technological University. His research interests include parallel and distributed computing, especially algorithmic, and theoretical aspects. His current work focuses on parallel computing in real-world architectures such as multicores, clouds and GPUs. He has served on the technical and organizational com-

mittees of a number of conferences, including ACM PODC, DISC, ICDCS, IPDPS, etc. His work has received multiple awards, including the MIT Sprowls Award for best computer science theses.



Chee Keong Kwoh received the bachelor's degree in electrical engineering (first class) and the master's degree in industrial system engineering from the National University of Singapore in 1987 and 1991, respectively. He received the PhD degree from the Imperial College of Science, Technology and Medicine, University of London, in 1995. He has been with the School of Computer Engineering, Nanyang Technological University (NTU), since 1993. He is the programme director, MSc (Bioinformatics), at NTU.

His research interests include data mining and soft computing and graph-based inference; applications areas include bioinformatics and biomedical engineering. He has done significant research work his research areas and published many quality international conferences and journal papers. He is an editorial board member of The International Journal of Data Mining and Bioinformatics, The Scientific World Journal; Network Modeling and Analysis in Health Informatics and Bioinformatics, Theoretical Biology Insights, and Bioinformation. He has been a guest editor for many journals such as Journal of Mechanics in Medicine and Biology, International Journal on Biomedical and Pharmaceutical Engineering, and others. He has been often invited as a organizing member or referee and reviewer for a number of premier conferences and journals, including GIW, IEEE BIBM, RECOMB, PRIB, BIBM, ICDM, and iCBBE just to name a few. He is a member of the Association for Medical and Bio-Informatics. Imperial College Alumni Association of Singapore. He has provided many service to professional bodies, and the Singapore and was conferred the Public Service Medal by the President of Singapore in 2008.



Jie Zheng received the BEng degree (first class honors) from Zhejiang University, China, in 2000, and the PhD degree from the University of California, Riverside, in 2006, both in computer science. He is a tenure-track assistant professor at the School of Computer Engineering, Nanyang Technological University (NTU), Singapore, and an adjunct faculty (senior research scientist) of Genome Institute of Singapore. Before joining NTU in 2011, he was a research scientist at the National Center for Biotechnology Information,

National Library of Medicine, National Institutes of Health, USA. His research interests include bioinformatics, genomics and systems biology, aiming to develop novel algorithms and in silico models to help answer biomedical questions (such as the mechanisms of cancer). While trained as a computer scientist, he maintains active and long-standing collaborations with life scientists.

▷ For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.