

# Simultaneous search for multiple QTL using the global optimization algorithm DIRECT

K. Ljungberg<sup>1</sup>, S. Holmgren

Information technology, Division of Scientific Computing  
Uppsala University, PO Box 337,751 05 Uppsala, Sweden

and Ö. Carlborg

Genetics and Biometry, Roslin Institute  
Roslin, Midlothian EH25 9PS, United Kingdom

<sup>1</sup>To whom correspondence should be addressed.

## ABSTRACT

**Motivation:** Epistatic interactions are important for quantitative traits. To maximize the power to detect epistatic quantitative trait loci (QTLs), a simultaneous search is necessary. The computational complexity demands that the traditional exhaustive search be replaced by a more efficient global optimization algorithm.

**Results:** We have adapted DIRECT, an algorithm presented in (Jones *et al.*, 1993), to the problem of simultaneous mapping of two and three QTL. We have compared DIRECT, in terms of accuracy and speed analyzing real data sets, with standard exhaustive search and a genetic algorithm previously used for QTL mapping in two dimensions. In all two- and three-QTL test cases, DIRECT accurately finds the global optimum two to four orders of magnitude faster than when using an exhaustive search, and one order of magnitude faster than when using the genetic algorithm. A search using a model with three fully interacting QTL is finished in six CPU minutes when using DIRECT, while an exhaustive search takes 142 CPU days. Thus three-QTL randomization testing for determining empirical significance thresholds is made feasible by the use of DIRECT. This opens the possibility to thoroughly investigate the power of simultaneous search to detect at least three interacting QTL.

**Availability:** The source code of the prototype implementation is available at

[http://www.tdb.uu.se/~kl/qtl\\_software.html](http://www.tdb.uu.se/~kl/qtl_software.html).

**Contact:** kl@tdb.uu.se

## INTRODUCTION

Rapid development in molecular genetics has led to the development of dense genetic maps, which are powerful tools for studying the molecular basis for quantitative genetic variation. One way to dissect the genetic architecture behind quantitative traits, i.e. traits affected by multiple genes and the environment, is to identify quantitative trait loci, QTL, in the genome. A QTL is a chromosomal region, locus, harboring one or several genes that affect the trait under study. The first methods used to locate, or map, QTL focused on detection of QTL by their marginal, i.e. additive and dominance, effects. These methods are presented in (Lander and Botstein, 1989; Haley and Knott, 1992). They are based on the concept of interval mapping, where the analyzed trait is modeled to depend on the genetic effects of a single QTL in the genome. A one-dimensional scan is performed using a dense grid covering the genome, and the

single QTL model is fitted at each grid-point. The most likely position of the QTL is taken to be the grid-point with the best model fit. In composite interval mapping (Zeng, 1993) and multiple QTL mapping (Jansen, 1992), a window of analysis is introduced in the one-dimensional scan. These schemes still search for the position of a single QTL, but markers outside the window of analysis are included as cofactors in the model. In this way the problem with variation caused by other QTL is reduced. A randomization test (Churchill and Doerge, 1994) is normally used to derive an empirical significance threshold for a statistical test of the putative QTL. During randomization testing normally 1000-10000 genome scans are performed on permuted datasets to obtain a stable distribution of the model fit under the null hypothesis of no QTL. A recent overview of current QTL mapping techniques is given in (Doerge, 2002).

Since a quantitative trait by definition is affected by multiple genes, it is desirable to simultaneously model the effects of these genes. Furthermore, simultaneous mapping is necessary for finding groups of interacting QTL where all loci involved lack significant marginal effects. Several methods have recently been proposed to simultaneously model the effects of multiple QTL and their interactions e.g. (Kao *et al.*, 1999; Wang *et al.*, 1999; Jannink and Jansen, 2001; Sen and Churchill, 2001; Carlborg and Andersson, 2002). A fundamental problem when using a multiple QTL model is that of computational complexity. For a model including  $n$  interacting QTL, the one-dimensional scan in a single QTL model is replaced by a  $n$ -dimensional search for the most likely positions of the interacting loci. When using randomization testing to derive significance thresholds for multiple QTL, the computations become very demanding even for models involving only two QTL.

To reduce the number of combinations of locations to evaluate, several approaches have been suggested. One suggestion (Kao and Zeng, 1997; Kao *et al.*, 1999; Zeng *et al.*, 1999, 2000) is that the computational complexity of the search is decreased by pre-selection of genomic regions with marginal effects. This potentially leads to a reduction in power since regions with primarily epistatic effects are disregarded. (Sen and Churchill, 2001) propose that a two-dimensional exhaustive search is performed on a sparse grid. This procedure reduces the resolution and would still be computationally burdensome in higher dimensions. To retain the true global search without introducing a prohibitive computational demand, the exhaustive search technique must be replaced by a more sophisticated algorithm for multi-dimensional global optimization. (Carlborg *et al.*, 2000) suggest that a ge-

netic optimization algorithm is used, and this type of algorithm was shown to be an efficient tool for mapping interacting QTL pairs in simulated data. Subsequently, a procedure for mapping and significance testing for epistatic QTL pairs was derived (Carlborg and Andersson, 2002). This method has recently been used to map QTL in experimental data, where multiple QTL pairs were detected in which neither of the QTL had significant marginal effects (Carlborg *et al.*, 2003). Similar results have been obtained using the method of (Sen and Churchill, 2001), e.g. in (Sugiyama *et al.*, 2001; Shimomura *et al.*, 2001).

To further investigate the evidence for higher order epistasis in experimental crosses, efficient numerical methods are needed for simultaneous mapping of QTL in two and higher dimensions. In this study we will explore the properties of a global optimization algorithm named DIRECT to perform QTL searches in two and three dimensions faster and more reliably than when using the genetic algorithm proposed in (Carlborg *et al.*, 2000). We will show that it is possible to perform simultaneous mapping, including randomization testing, of three fully interacting QTL, using a standard single-processor computer.

## SYSTEMS AND METHODS

### Computations in QTL mapping

There are two main elements in the computations when searching for QTL; the kernel problem and the global optimization problem. In general, any algorithm for the global optimization problem can be used together with any type of kernel algorithm.

The kernel problem consists of evaluating the objective function, i.e. calculating the model fit for one specific combination of putative QTL. Many different genetic models with or without interaction parameters can be used. The model parameters can be determined using e.g. ordinary linear regression (Haley and Knott, 1992; Haley *et al.*, 1994), or maximum likelihood estimation (Zeng, 1994). Both linear regression and maximum likelihood estimation, via the ECM algorithm (Meng and Rubin, 1993), involve solving a least squares problem, which is normally done using standard software library routines. The kernel problem was investigated in (Ljungberg *et al.*, 2002), where we presented efficient objective function evaluation algorithms based on updated QR fac-

torizations for both linear regression and maximum likelihood kernels.

The global problem consists of optimizing the objective function, i.e. out of all possible QTL combinations finding the one giving the best model fit. It appears in two flavors. When searching the original data, the goal is to find both the most likely positions of the QTL in the set *and* the corresponding value of the parameters and model fit. However, during randomization testing, only the optimal value of the model fit is needed. As long as the value found by the algorithm is sufficiently accurate, the significance thresholds will also be accurate. This is an important observation, since the problem of determining the position of the true global optimum is more difficult for the permuted data where the connection between genotype and phenotype is broken. In this case the optimization landscape will often have many smaller local optima, scattered over the search space, with almost the same value of the objective function.

### The global optimization problem

When performing simultaneous mapping of a set of  $n$  QTL, we search a point  $\bar{x}^{\text{opt}} = (x_1^{\text{opt}} \ x_2^{\text{opt}} \ \dots \ x_n^{\text{opt}})$  in the  $n$ -dimensional hypercube defined by  $0 \leq x_i \leq L$ . Here,  $L$  is the size of the genome in  $cM$  and  $x_i$  is the position of the  $i$ :th QTL in the set. The optimal value of the test statistics is independent of the ordering of the QTL in the set. Therefore, the optimization problem exhibits an  $n!$ -fold symmetry, equivalent to the  $n!$  possible orderings of the QTL. This represents a significant reduction of the search space. In QTL mapping, the search space can be divided into boxes where each edge corresponds to one chromosome. Such a chromosome combination box  $(c_1, c_2, \dots, c_n)$  encloses all points where QTL 1 is assumed to be on chromosome  $c_1$ , QTL 2 is assumed to be on chromosome  $c_2$  and so on. The search space symmetry is employed by restricting the search to chromosome combination boxes where  $c_1 \leq c_2 \leq \dots \leq c_n$ . Boxes where two or more QTL are located on the same chromosome are also affected by the symmetry, and only part of them need to be considered.

The most likely QTL position combination  $\bar{x}^{\text{opt}}$  minimizes an objective function which may be written as (Ljungberg *et al.*, 2002)

$$f(\bar{x}) = \min_b (y - Ab)^T G (y - Ab), \quad (1)$$

where  $y$  is the vector of trait values,  $b$  is a vector of re-

gression parameters and  $A$  is the matrix of regression indicator variables. The matrices  $G$  and  $A$  depend on the QTL mapping method being used. When using the linear regression method  $G = I$ , and the entries of  $A$  are either constants or continuous functions of  $\bar{x}$  within chromosomes. Hence, the objective function  $f(\bar{x})$  depends continuously on  $\bar{x}$  within every chromosome combination box. However, at the boundaries between chromosomes,  $f(\bar{x})$  is normally not continuous.

## Models

Name	Description
2:m	A two-QTL model including fixed effects and additive and dominance marginal effects.
3:m	The three-QTL version of 2:m.
2:m+p	The 2:m model with pairwise interaction effects added.
3:m+p	The three-QTL version of 2:m+p.
3:m+p+t	3:m+p adding the full three-way interaction.

**Table 1.** Abbreviations for the two- and three-QTL genetic models used in the study. **m** denotes marginal effects, **p** pairwise interaction effects and **t** three-way interaction effects.

Throughout this work we have used the Haley-Knott regression method for experimental crosses between outbred lines (Haley *et al.*, 1994). Table 1 describes the five genetic models that are used. For example, model 3:m+p+t is defined as

$$\begin{aligned}
 y = & \sum_{j=1}^{n_{fix}} b_j^{fix} a_j^{fix} \text{ (fixed effects including the mean)} \\
 & + \sum_{j=1}^3 (b_j^a a_j + b_j^d d_j) \\
 & \text{(marginal additive and dominance effects)} \\
 & + \sum_{j=1}^2 \sum_{k=j+1}^3 (b_{jk}^{aa} a_{jk} + b_{jk}^{ad} a_{jk} + b_{jk}^{da} d_{jk} + b_{jk}^{dd} d_{jk}) \\
 & \text{(pairwise interaction effects)} \\
 & + b^{aaa} aaa_{123} + b^{aad} aad_{123} + b^{ada} ada_{123} + b^{daa} daa_{123} \\
 & + b^{dda} dda_{123} + b^{dad} dad_{123} + b^{add} add_{123} + b^{ddd} ddd_{123} \\
 & \text{(three-way interaction effects) ,}
 \end{aligned}$$

where  $y$  is the phenotype,  $a_j^{fix}$  are the indicator regression variables for the fixed effects including the mean,  $a_j$  and  $d_j$  denote the regression indicator variables for the marginal additive and dominance effect of QTL  $j$  in an outbred line cross as described in (Haley *et al.*, 1994).  $aa_{jk}$ ,  $ad_{jk}$ ,  $da_{jk}$  and  $dd_{jk}$  are regression indicator variables for the interaction effects of QTL pair  $jk$ , obtained by multiplying the respective additive and dominance regression variables for QTL  $j$  and  $k$  (Haley and Knott, 1992).  $aaa_{123}$ ,  $daa_{123}$  and so on are regression indicator variables for the interaction effects of QTL triplet 123, obtained analogously to the QTL pair variables. The  $b$  values are the partial regression coefficients for the genetic parameters corresponding to the indicator regression variables.

In this paper we have not evaluated the power to detect epistatic QTL using the different models, nor have we investigated the best choice of model for the data sets used. This question will be addressed separately. The purpose of the current study is to compare the computational methods in terms of speed and their ability to find the global optimum of the objective function using real data, and the models were chosen with the intention to give a varied set of optimization landscapes.

## Data

We have tested the computational methods on data from two mapping populations. The first population is denoted WB/LW. It consists of 191 animals from an  $F_2$  intercross between European Wild Boar and Large White domestic pigs (Andersson *et al.*, 1994). The genome size is approximately 2300  $cM$  and we used phenotypic data for six growth-related traits. The second population, denoted JF/WL, consists of 852 animals from an  $F_2$  intercross between red jungle-fowl and White Leghorn chickens described in (Schütz *et al.*, 2002). The genome size is approximately 2500  $cM$ , and phenotypic data for nine different growth traits were used. We leave out further details about the phenotypes since we are not currently looking for new QTL.

In addition to optimizing the objective function for various models in the original data sets, it is relevant to compare empirical significance thresholds derived when using the three methods. For this purpose four sets of randomized data were generated, 1000 randomizations each of two JF/WL traits and two WB/LW traits.

## ALGORITHMS

### Exhaustive grid search

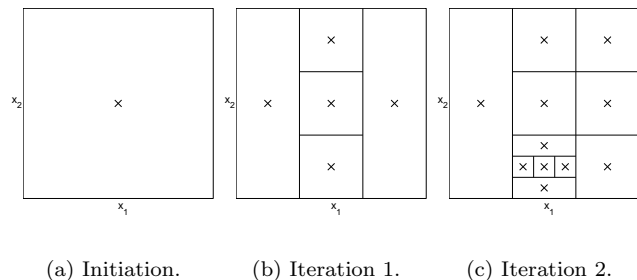
The standard method for solving the global optimization problem is to use an exhaustive grid search, evaluating the objective function for every possible QTL combination using steps of e.g. 1 cM. We have performed exhaustive two- and three-dimensional searches for all test cases. The symmetry of the search space was easily exploited. To make the computations feasible, the exhaustive searches were performed on a parallel computer. We measure the accuracy of DIRECT and GA as their ability to find the same optimum as the one found by exhaustive search, which is the global optimum.

### The DIRECT algorithm

The original DIRECT algorithm was presented in (Jones *et al.*, 1993). It searches for the global minimum  $\bar{x}^{\text{opt}}$  of multi-dimensional Lipschitz continuous functions  $f(\bar{x})$  with the same type of constant constraints as the QTL mapping problem described above. The practical interpretation of a function  $f(\bar{x})$  being Lipschitz continuous is that the slope of  $f(\bar{x})$  is limited by some constant  $K$  everywhere.

DIRECT systematically divides the search space into smaller and smaller boxes, see Figure 1. The Lipschitz continuity condition is used for deterministically determining which boxes to select for further division in each iteration. Suppose that the search space at iteration  $i$  has been divided into  $L$  boxes, and that  $f(\bar{x})$  has been computed at the center of each box. Given  $K$ , a lower bound on  $f(\bar{x})$  in each box could be computed, and the box with the lowest bound would be selected for further division. In practice  $K$  is unknown, so DIRECT divides all boxes where  $f(\bar{x})$  has the lowest bound for *any* value of  $K$  from zero to infinity. The center-point of each new box is sampled, and the selection procedure is repeated. The box selection step is very fast. It should be noted that the Lipschitz continuity condition is only used for bounding  $f(\bar{x})$  within each box, which is important for the application of the algorithm to QTL mapping problems.

In the original formulation of the algorithm, no box is ever discarded from the search. A box not considered potentially optimal in one iteration can be chosen for division in a later iteration. If the algorithm is run for sufficiently long time, it is possible to prove that the global



**Fig. 1.** Illustration of DIRECT search space division.

optimum will always be found (Jones *et al.*, 1993). In practice, the global optimum is normally found after a rather small number of iterations. However, a general problem for global optimization algorithms is how to determine when to stop the iterations. In the original paper (Jones *et al.*, 1993), it is suggested that a fixed number of function evaluations should be used.

The original algorithm has been modified to fit the QTL search problem. As observed above,  $f(\bar{x})$  is a continuous function of  $\bar{x}$  within every chromosome combination box. However, at the boundaries between chromosomes,  $f(\bar{x})$  is normally not continuous. To guarantee that the continuity condition of the algorithm is fulfilled, the search is initiated by sampling the center point of all chromosome combination boxes in the search space. In the original algorithm only the center-point of the complete search space is to be sampled at initiation. Also, we do not normalize the  $x_i$  coordinates as in the original algorithm, and do not divide boxes with edges smaller than 1 cM.

We present no proof that  $f(\bar{x})$  is Lipschitz continuous as well as continuous within the chromosome combination boxes. However, a simple argument along this line can be applied to the computations. Obviously,  $f(\bar{x})$  cannot exceed  $y^T G y$ , which is finite, nor be smaller than 0. When performing the search for the set of QTL, a resolution limit of typically 1 cM is used, and thus there exists a practical Lipschitz constant which is bounded by  $y^T G y$ .

The only parameter in DIRECT with a significant influence on performance is the number of function evaluations allowed. In two-dimensional searches we performed 6000 evaluations, and in three-dimensional searches we performed 46000 function evaluations plus 10000 in the intermediate refinement step. Using these settings we found the global optimum in all test cases using non-randomized data.

We have observed, in accordance with other authors e.g.

(Cox *et al.*, 2001; Bartholomew-Biggs *et al.*, 2002), that DIRECT quickly locates the region of the global optimum but that local convergence is rather slow. We therefore finish the search by performing a local exhaustive search,  $\pm 5cM$  in each dimension, around the best point. This is similar to the procedure suggested in (Cox *et al.*, 2001). In the three-dimensional searches we also use an intermediate refinement step. After a set number of iterations the chromosome combination box containing the best point is located, and a number of additional iterations are performed in this box only, before the final local exhaustive search.

## The genetic algorithm

We have compared DIRECT with a genetic algorithm, GA, from a library named PGAPack (Levine, 1996). The same GA was used in (Carlborg *et al.*, 2000), where a position in the search space is encoded as a string of  $2n$  real numbers representing the chromosomes and the chromosome positions of the  $n$  QTL. One QTL position string is called a GA-chromosome, and the fit of a chromosome is given by the objective function value at the corresponding position in the search space. A GA-population is a set of GA-chromosomes, and in each iteration new GA-chromosomes are created by mutation and crossover among the existing ones, selecting for best fit. The GA is thus partly related to forward selection in the sense that mutation and/or crossover on a good candidate GA-chromosome often results in keeping one QTL position fixed and changing the other. The symmetry of the search space is exploited by not allowing the algorithm to evaluate the reflection of a position already visited. After the GA is finished a local exhaustive search  $\pm 5cM$  is performed around the found optimum in the same way as for DIRECT.

Name	Number of populations	Iterations/population	Population size
GA(75k)	25	1500	20
GA(20k)	10	1000	20
GA(6k)	3	980	20
GA(1M)	25	2000	200

**Table 2.** Parameter settings for the genetic optimization algorithm. Parameters not explicitly described in the table are set as in (Carlborg *et al.*, 2000).

A significant effort was spent tuning the parameters to obtain the best possible accuracy for all test cases. Table 2 shows the different parameter settings chosen for this study. We refer to the settings chosen in (Carlborg

*et al.*, 2000) as GA(20k), the name reflecting the approximate number of function evaluations performed. The best parameter choice found was a modified version of GA(20k) which we call GA(75k). GA(6k) is the settings giving only the same number of function evaluations as DIRECT in two dimensions. The parameterization used in the three-dimensional searches is called GA(1M).

## IMPLEMENTATION

All objective function evaluations were done using the efficient kernel algorithm presented in (Ljungberg *et al.*, 2002). The experiments showed that, in practice, the only factors determining the CPU time for the three methods are the number of function evaluations performed and the time required for a single evaluation. The CPU time for one evaluation depends on the model and data set, but not on the optimization method since they use the same kernel algorithm.

All code is written in Fortran90, and the computations were done on SPARC UIII, 900MHz processors. The exhaustive searches were performed on a parallel computer using MPI, and the CPU times reported are the sums of the CPU times for each processor, not including overhead time for the parallelization.

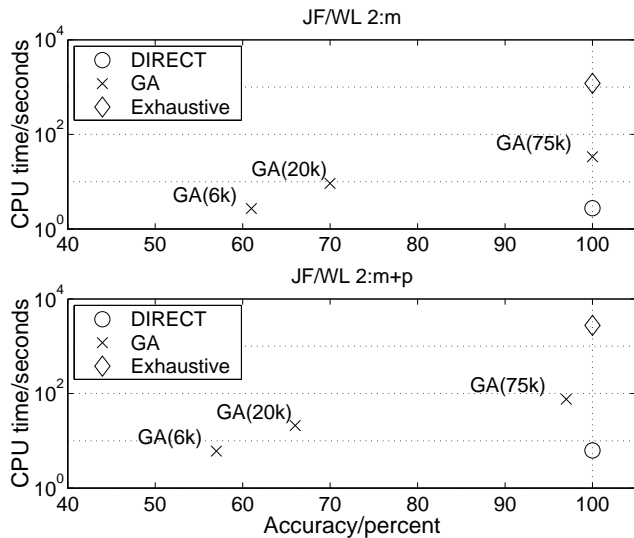
## RESULTS

### Original, non-randomized, data

The accuracy is reported as the percentage of successful localizations of the exact global optimum out of the total number of searches. Since the GA has a random element, the result will depend on the random seed. Therefore, using this method, each search was repeated 15 times to give a reasonable statistic. DIRECT is deterministic and gives the same result every time.

First we report the results for searches in two dimensions. We have tested the methods for the 2:m and 2:m+p models in Table 1 on all data sets described in the **Data** subsection, which gives a total of 30 tests.

Figure 2 shows the average CPU times and accuracy over the 9 phenotypes of the JF/WL data set using the 2:m and 2:m+p models. The results for WB/LW data and the same models were very similar. An exhaustive search with the 2:m model requires about 20 minutes, and 46 minutes with the 2:m+p model. DIRECT finds

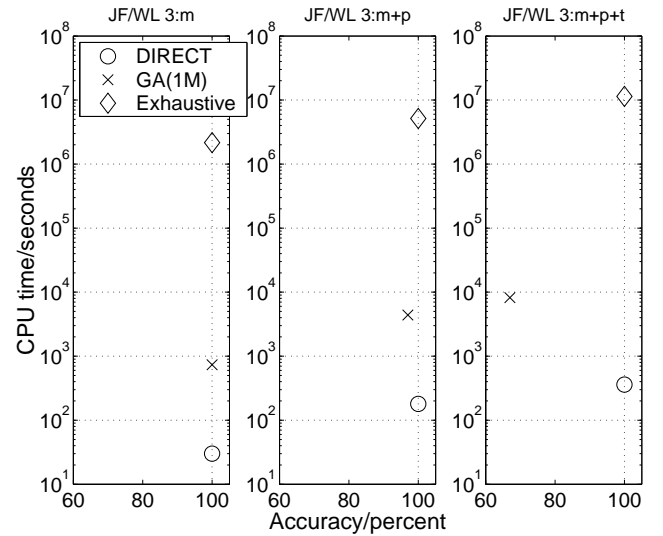


**Fig. 2.** CPU time for two-dimensional searches as a function of the percentage of successful localizations of the global optimum.

the global optimum in less than 3 and 7 seconds respectively. GA(75k) gives the global optimum in close to 100% of the runs, with CPU time 34 and 76 seconds. Using GA(6k), the genetic algorithm with the same number of function evaluations and thus practically the same CPU time as DIRECT, reduces the accuracy from close to 100% to around 60%. GA(20k), the settings of (Carlborg *et al.*, 2000) give intermediate results. The GA has more difficulties finding the global optimum when epistasis is included in the model. It was observed already in (Carlborg *et al.*, 2000) that the GA sometimes failed when a QTL pair lacked significant marginal effects. This can be explained by the forward selection property of the algorithm.

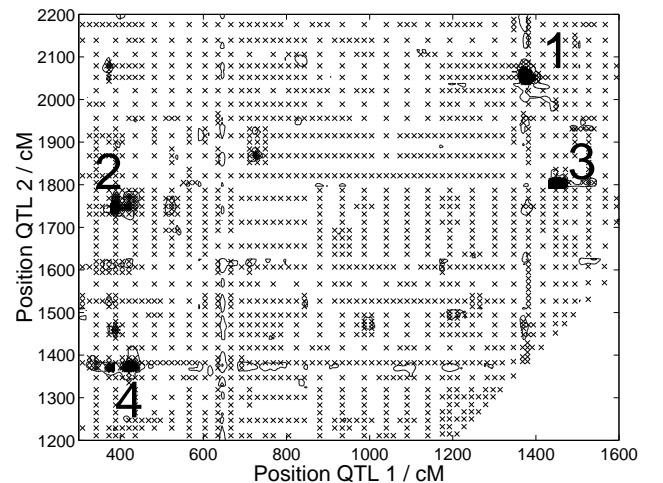
Now we turn to three-QTL results. We have used the 3:m and 3:m+p models combined with four JF/WL traits, one of which was also used with the 3:m+p+t model, giving nine tests in total.

Figure 3 shows the average CPU times and accuracy over 4 phenotypes of the JF/WL data set using the 3:m and 3:m+p models, and 1 phenotype using the 3:m+p+t model. The JF/WL 3:m, 3:m+p and 3:m+p+t exhaustive searches would take approximately 25, 60 and 142 days respectively, on a single processor computer. The gain in using DIRECT over exhaustive search is more than four orders of magnitude in speed, the searches taking 0.5, 3 and 6 minutes, while not losing accuracy. Using GA(1M) gives high accuracy for the 3:m and 3:m+p models but lower for 3:m+p+t and is over one order of



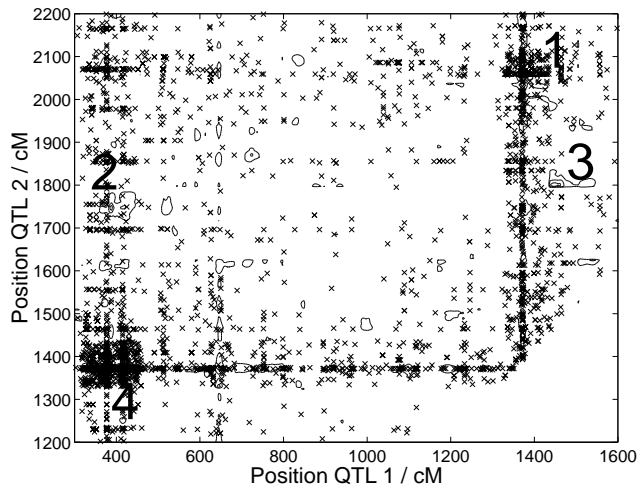
**Fig. 3.** CPU time for three-dimensional searches as a function of the percentage of successful localizations of the global optimum.

magnitude slower than DIRECT, the searches requiring 12, 73 and 136 minutes respectively.



**Fig. 4.** Search pattern after 6000 function evaluations with DIRECT in the region around the four largest peaks, numbered 1-4 according to their relative ranks.

Figures 4 and 5 illustrate the difference in search pattern between DIRECT and the genetic algorithm. Here we show results from model 2:m+p with WB/LW data. The two figures show the sampling pattern after a complete run, i.e. 6000 function evaluations, using DIRECT (Figure 4) and using GA(6k) (Figure 5). The locations where the objective function has been evaluated are marked with 'x' in contour plots of the objective function



**Fig. 5.** The GA sampling pattern after 6000 function evaluations, GA(6k), in the region around the four largest peaks.

around the four largest peaks. For clarity most contours for lower levels are not shown. DIRECT uses the function evaluations very efficiently. It gives even coverage of the search space with dense clusters of function evaluations around the largest peaks. This indicates that the algorithm can locate the global optimum for original data also in difficult cases when there are many local optima of similar magnitude. Using the same number of function evaluations the GA sometimes does not find the global optimum, even if the regions around all the four largest peaks are sampled. If many peaks are of similar height, the best position found so far when the local search is initiated might be at the wrong peak. Or the right peak might have been found, but local exhaustive search  $\pm 5cM$  is not a good enough method to localize the very best position on the peak. The GA samples the search space stochastically to a large extent.

## Randomized data

Finding the global optimum can be expected to be more difficult in a randomized data set, since the optimization landscape will be smoothed out and the peaks smaller for most of the randomizations when the connections between genotype and phenotype is broken.

We used the 2:m+p model for the four randomized data set. We determined the 1.0, 5.0, 10 and 20% genome-wide significance thresholds for 0 against 2 QTL using an exhaustive search. The thresholds were also calculated using DIRECT and GA on the same data. In Table 3 we

report the true levels (as given by the exhaustive search) of the thresholds derived using DIRECT, GA(20k) and GA(75k) intended to give the 1.0, 5.0, 10 and 20% significance levels. A number 5.6% in the 5.0% row means that in 5.0% of the randomizations a global optimum better than  $x$  ( $x$  not reported) was found when using the global optimization algorithm, i.e. the 5.0% significance threshold would be taken to be  $x$ , while in reality 5.6% of the true global optima, obtained using exhaustive search, were better than the same  $x$ . A threshold that is too low, i.e. at 5.6% instead of 5.0%, gives a slight increase in the type I error rate. This could in part explain the increased rate of type I errors in (Carlborg and Andersson, 2002) where the genetic algorithm is used.

Exhaustive search	JF/WL	2:m+p	
	DIRECT	GA(20k)	GA(75k)
1.0%	1.0%	1.0%	1.0%
5.0%	5.2%	5.6%	5.3%
10%	10%	12%	10%
20%	21%	24%	21%
32 days	1.7 hours	5.8 hours	21 hours

Exhaustive search	WB/LW	2:m+p	
	DIRECT	GA(20k)	GA(75k)
1.0%	1.0%	1.0%	1.0%
5.0%	5.1%	5.4%	5.0%
10%	10%	11%	10%
20%	21%	22%	20%
17 days	57 min	3.5 hours	14 hours

**Table 3.** Derived empirical thresholds and the corresponding CPU times using exhaustive search, DIRECT and two parameterizations of the genetic algorithm.

Looking at the individual runs it can be seen that DIRECT finds the wrong position in about 9% of the randomizations. The function values are however accurate enough to give nearly the same threshold values as exhaustive search, and they are calculated between two and three orders of magnitude faster. Using GA(20k) the wrong position is found in 23 – 35% of the cases. This is about the same error rate as was found with non-randomized data. However the computed thresholds are still sufficiently accurate for practical use. The error rate is about 1 – 14% when using GA(75k), which gives very accurate threshold values.

There is a tendency for the 1% and 5% computed thresholds to be more accurate than the 10% and 20%. This reflects that it is more easy for both algorithms to find large peaks, while the randomizations giving a “smeared”



landscape with many smaller peaks is more difficult from an optimization point of view.

## DISCUSSION

This study has shown that DIRECT is a fast and accurate algorithm for global optimization in QTL mapping. The exact optimum is found in real data sets, and searches in randomized data are accurate enough to give almost the same empirical significance thresholds as exhaustive search. Two-dimensional searches take a few seconds, and three-dimensional searches are finished in a few minutes. DIRECT makes randomization testing of two-QTL models faster, and randomization testing of three-QTL models fully feasible. This opens the possibility to thoroughly investigate the power of simultaneous search to detect triplets of interacting QTL, which will be done in future research. We will also implement DIRECT and the GA to simultaneously search for four and five interacting QTL, to further explore higher order epistatic interactions.

DIRECT is developed to find the optima of Lipschitz continuous functions, i.e. functions where the rate of change of the objective functions is everywhere limited by some constant  $K$ , where  $K$  is normally unknown. We gave a motivation for Lipschitz continuity of the QTL mapping objective function based on that  $0 \leq f(\bar{x}) \leq y^T G y$ , and that the resolution is limited. A more interesting observation is that genetic distance is a measure of change, a measure of recombination events. Recombinations are reflected by change in the indicator variable matrix  $A$  and consequently in  $f(\bar{x})$ . The magnitude of the change in  $f(\bar{x})$  depends on the phenotype values of the individuals who switch genotype between the flanking markers, but this still imposes a limit on the possible rate of change in  $f(\bar{x})$ . No such limit is assumed in the calculations, but we believe it is the explanation for the good performance of DIRECT.

When analyzing data from other types of experimental crosses, the optimization landscape will probably be different than in this study. However we believe that an  $F_2$  cross between outbred lines is one of the most difficult cases, since the objective function will contain more noise and less distinct peaks. Also we have used models with many parameters, and that will also make the peaks smaller and more difficult for the optimization algorithms to find. When adapting DIRECT to other experimental designs, an advantage is that the only parameter necessary to adjust is the number of function evalu-

ations allowed. We have used 6000 function evaluations in the two-dimensional searches and 56000 evaluations in the three-dimensional search, which corresponds to 0.2% and 0.002% of the total number of positions.

## ACKNOWLEDGEMENTS

Ö. Carlborg gratefully acknowledges financial support from The Knut and Alice Wallenberg foundation. We thank Leif Andersson for contributing the data sets. K. Ljungberg is affiliated with The Linnaeus Centre for Bioinformatics and The Graduate School in Mathematics and Computing, both at Uppsala University.

## References

- Andersson, L., Haley, C., Ellegren, H., Knott, S., Johansson, M., Andersson, K., Andersson-Eklund, L., Edfors-Lilja, I., Fredholm, M., and Hansson, I. (1994) Genetic mapping of quantitative trait loci for growth and fatness in pigs, *Science*, **263**, 1771–1774.
- Bartholomew-Biggs, M., Parkhurst, S., and Wilson, S. (2002) Using direct to solve an aircraft routing problem, *Comp Optim Appl*, **21**, 311–323.
- Carlborg, Ö. and Andersson, L. (2002) The use of randomization testing for detection of multiple epistatic QTL, *Genet Res*, **79**, 175–184.
- Carlborg, Ö., Andersson, L., and Kinghorn, B. (2000) The use of a genetic algorithm for simultaneous mapping of multiple interacting quantitative trait loci, *Genetics*, **155**, 2003–2010.
- Carlborg, Ö., Kerje, S., Schütz, K., Jacobsson, L., Jensen, P., and Andersson, L. (2003) A global search reveals epistatic interaction between QTL for early growth in the chicken, *Genome Res*, **13**, 413–421.
- Churchill, G. and Doerge, R. (1994) Empirical threshold values for quantitative trait mapping, *Genetics*, **138**, 963–971.
- Cox, S., Haftka, R., Baker, C., Grossman, B., Mason, W., and Watson, L. (2001) A comparison of global optimization methods for the design of a high-speed civil transport, *J Global Optim*, **21**, 415–433.
- Doerge, R. (2002) Mapping and analysis of quantitative trait loci in experimental populations, *Nature Rev Genet*, **3**, 43–52.

- Haley, C. and Knott, S. (1992) A simple regression method for mapping quantitative trait loci in line crosses using flanking markers, *Heredity*, **69**, 315–324.
- Haley, C., Knott, S., and Elsen, J.-M. (1994) Mapping quantitative trait loci in crosses between outbred lines using least squares, *Genetics*, **136**, 1195–1207.
- Jannink, J.-L. and Jansen, R. (2001) Mapping epistatic quantitative trait loci with one-dimensional genome searches, *Genetics*, **157**, 445–454.
- Jansen, R. (1992) A general mixture model for mapping quantitative trait loci by using molecular markers, *Theor Appl Genet*, **85**, 252–260.
- Jones, D., Perttunen, C., and Stuckman, B. (1993) Lipschitzian optimization without the lipschitz constant, *J Optimiz Theory App*, **79**, 157–181.
- Kao, C.-H. and Zeng, Z.-B. (1997) General formulae for obtaining the MLEs and the asymptotic variance-covariance matrix in mapping quantitative trait loci when using the EM algorithm, *Biometrics*, **53**, 653–665.
- Kao, C.-H., Zeng, Z.-B., and Teasdale, R. (1999) Multiple interval mapping for quantitative trait loci, *Genetics*, **152**, 1203–1216.
- Lander, E. and Botstein, D. (1989) Mapping mendelian factors underlying quantitative traits using RFLP linkage maps, *Genetics*, **121**, 185–199.
- Levine, D. (1996) *Users Guide to the PGAPack Parallel Genetic Algorithm Library*. Mathematics and Computer Science Division, Argonne National Laboratory, Illinois.
- Ljungberg, K., Holmgren, S., and Carlborg, Ö. (2002) Efficient algorithms for quantitative trait loci mapping problems, *J Comput Biol*, **9**, 793–804.
- Meng, X.-L. and Rubin, D. (1993) Maximum likelihood estimation via the ECM algorithm: A general framework, *Biometrika*, **80**, 267–278.
- Schütz, K., Kerje, S., Carlborg, Ö., Jensen, P., and Andersson, L. (2002) Analysis of a red junglefowl  $\times$  white leghorn intercross reveals trade-off in resource allocation between behaviour and production traits, *Behav Genet*, **32**, 423–433.
- Sen, S. and Churchill, G. A. (2001) A statistical framework for quantitative trait mapping, *Genetics*, **159**, 371–387.
- Shimomura, K., Low-Zeddies, S. S., King, D. P., Steeves, T. D., Whiteley, A., Kushla, J., Zemenides, P. D., Lin, A., Vitaterna, M. H., Churchill, G. A., and Takahashi, J. S. (2001) Genome-wide epistatic interaction analysis reveals complex genetic determinants of circadian behavior in mice, *Genome Res*, **11**, 959–980.
- Sugiyama, F., Churchill, G. A., Higgins, D. C., Johns, C., Makaritsis, K. P., Gavras, H., and Paigen, B. (2001) Concordance of murine quantitative trait loci for salt-induced hypertension with rat and human loci, *Genomics*, **71**, 70–77.
- Wang, D. L., Zhu, J., Li, Z. K. L., and Paterson, A. H. (1999) Mapping QTLs with epistatic effects and QTL  $\times$  environment interactions by mixed linear model approaches, *Theor Appl Genet*, **99**, 1255–1264.
- Zeng, Z.-B. (1993) Theoretical basis for separation of multiple linked gene effects in mapping quantitative trait loci, *Proc Nat Acad Sci USA*, **90**, 10972–10976.
- Zeng, Z.-B. (1994) Precision mapping of quantitative trait loci, *Genetics*, **136**, 1457–1468.
- Zeng, Z.-B., Kao, C.-H., and Basten, C. (1999) Estimating the genetic architecture of quantitative traits, *Genet Res*, **74**, 279–289.
- Zeng, Z.-B., Liu, J., Stam, L. F., Kao, C.-H., Mercer, J. M., and Laurie, C. C. (2000) Genetic architecture of a morphological difference between two drosophila species, *Genetics*, **154**, 299–310.