

# *Efficient Algorithms for Haplotype Phasing with RFLPs*

April 25, 2002

W. CASEY<sup>a</sup> AND B. MISHRA<sup>b,c, 1</sup>

a. Mathematics Department, Courant Institute of Mathematical Sciences, NYU.

b. Professor of Computer Science and Mathematics, Courant Inst., NYU

c. Professor, Watson School of Biological Sciences, Cold Spring Harbor Lab.

## Abstract

The determination of feature maps, such as STSs, SNPs or RFLPs maps, for each chromosome copy or haplotype in an individual has important potential applications to association studies. We present a method to recover RFLP feature maps for each haplotype starting from genotype data which is an ambiguous superposition of all haplotypes' data.

Our method is an inference method which is able to interpret data in two key ways: 1) We must be able to determine when a feature expresses polymorphic identity, thereby suggesting some differences on the different haplotypes. 2) Given any combination of polymorphic features we co-associate the alternatives of each feature thereby partitioning the alternatives to haplotypes.

We design an expectation maximization (EM) algorithm to detect the polymorphic markers. Secondly, we design an efficient algorithm to rapidly determine the co-associations of alternatives for each polymorphic feature: This process is called the phasing of polymorphisms.

The problem of SNP (single nucleotide polymorphism) phasing has been investigated in earlier literature and found NP-complete [10]. In contrast, using the RFLP (restriction fragment length polymorphism) markers, we show that our algorithm can produce marker phasing and hence haplotypes, when the genome-wide ordered restriction site data are produced by an available technology such as optical mapping [4].

A prior model of the data, comprising a set of restriction fragment lengths, allows us to analyze the proposed algorithm and provide a probabilistic guarantee for its correctness. Our algorithm can be suitably modified for a wide class of haplotyping problems, relying on unrelated markers and technologies. Independently, as a significant fraction of RFLP markers are directly caused by SNP's, the RFLP phasing may be an important tool for reducing the complexity of the SNP-phasing problem.

## 1 Introduction and Related Literature

A diploid organism contains two very similar copies of each chromosome, with the exception of sex-chromosomes. We call the pair of copies *haplotypes*, and refer to them individually as "Haplotype I" and "Haplotype II". Haplotypes are mainly contributed by each parent, but because of crossover events during meiosis, the haplotypes of the progeny may differ locally from either haplotype of either parent,

---

<sup>1</sup>Authors' Current Address: Courant Institute, New York University, 251 Mercer St, NYC, NY-10012. Cold Spring Harbor Laboratory, P.O. Box 100, 1 Bungtown Rd., Cold Spring Harbor, NY 11724. The research presented here was partly supported by a DOE Grant, NYU Research Challenge Grant and an NIH Grant. Will Casey: 212-998-3377 wcasey@cims.nyu.edu, Bud Mishra: 212-998-3464 mishra@nyu.edu. Fax to Will Casey 212-673-3242

thus exhibiting genetic diversity. As a result, understanding the statistics of a population of haplotypes is directly related to understanding the genetic structure of the population. Such information can be derived accurately only from genome-wide haplotype maps. In this paper, we discuss a restriction-enzyme-based technology and algorithms capable of producing such information.

A restriction enzyme (e.g., *Bam*H I) finds and cuts double-stranded DNA at specific recognition sites (e.g. ‘ggatcc’) with high specificity. Thus, slight positional variations of these restriction sites on the mostly similar copies of a chromosome can potentially separate and identify the haplotypes. For instance, a segmental insertion or deletion between two consecutive restriction sites on one copy of the chromosome will be observed as a difference in the two restriction fragment lengths and thus, can be used to locate those significant events in the development of haplotypes. This event represents one of two distinct causes that result in the underlying restriction site marker patterns on two haplotypes to differ—the other being an “in-del” or substitution of a nucleotide within the actual restriction site. Thus the second kinds of RFLPs represent a significant subset of SNP’s (single nucleotide polymorphisms). While RFLPs have not received as wide an attention as SNPs, they also hold the same kinds of promise as SNPs for association studies and classification of genetic diseases—most likely for a significantly lower cost.

More formally, quantities such as the base-pair length between two consecutive *restriction sites* are called *restriction fragment lengths* and are modeled as random variables, already extensively studied in [4]. When the homologous regions on the two haplotypes contain different lengths between consecutive restriction sites, they are said to result in a ‘restriction fragment length polymorphism’ (abbreviated, *RFLP*). These RFL’s are subject to measurement errors (e.g., sizing error, partial restriction digestion and false positive site errors) and locally confound length-based polymorphism detection. Exacerbating the problem further, the current technology is further limited in how many such features can be mapped over a small genomic segment. Consequently, making haplotype maps directly is much more difficult than making genotype maps. Nonetheless, if a group of ordered restriction fragments can be sampled from either haplotype, a large number of such samples allows first to identify RFLPs, and then determine how these RFLPs co-associate locally. Furthermore, with the increase in the number of such samples and the increase in expected number of markers in each sample, it is possible to improve the accuracy and resolution of the haplotype maps in spite of the statistical errors alluded to. Note that, at the end of this process, the result can be interpreted as two haplotype ordered restriction maps, further annotated with the location of the RFLPs.

Similar haplotyping or ‘the phasing problem’ has been investigated in [7, 10], but with a different data model. Of particular interest are the results of NP-hardness in [10] of a SNP phasing problem. First detailed statistical models of optical mapping process, precise computational complexity of the problem under various error models and a robust practical algorithm (the only one in use!) to produce ordered restriction maps are to be found in [4]. A related problem, “*K*-population” problem, has been studied extensively in [12, 13]. Also of related interest are [2, 3, 5].

The paper is organized as follow: In section 2, we review the terminology and capabilities of *Optical Mapping* and the Bayesian approaches developed for mapping restriction fragment sites on the genome. In section 3, we discuss the recognition of polymorphic sites from data. We develop an *EM-Algorithm* for this problem by modifying a well-known and previously-studied problem of estimating the parameters of a mixture distribution. The mixture distribution model also serves as our prior model for data. In section 4, we discuss the difference of genotype and haplotype data and define the concept of phasing mathematically. The problem of phasing using maximum likelihood estimates of probabilities for pairwise co-association of polymorphic features is obtained by multiplication on a multiplicative group. In section 5, we develop algorithms, capable of phasing the polymorphic sites in efficient time. In section 6,

we demonstrate via simulations on reasonably sized data sets how experimental parameters and other aspects of the problem interact with each other. The results of the simulation are highly promising. In section 7, we indicate our probabilistic analysis for the results of our algorithms. In the full paper we give an example of an adversarial event that is capable of causing an error in phasing and provide Chernoff bounds to show that the probability of such events decay exponentially with the sample size. Additionally, the experimental parameters for making such *false positive* events negligible are shown to be realistic. The full paper contains detailed mathematical proofs for the propositions whose proofs are omitted in this extended abstract.

## 2 Optical Mapping

Optical mapping is a physical mapping approach that provides an ordered enumeration of the restriction sites along with the estimated lengths of the restriction fragments between consecutive restriction sites. A *restriction site* is the location of a short specific nucleotide sequence (4-8 bp long) where a particular restriction enzyme cleaves the DNA by breaking a phosphodiester bond. The fragment of DNA generated by cleaving at two consecutive restriction sites is a restriction fragment.

The physico-chemical approach underlying optical mapping is based on immobilizing long single DNA molecules on an open glass surface, digesting the molecules on the surface and visualizing the gaps created with fluorescence microscopy. Thus the resulting image, in the absence of any error, would produce an ordered sequence of restriction fragments, whose masses can be measured via relative fluorescence intensity and interpreted as fragment lengths in bps. The corrupting effects of many independent sources of errors affect the accuracy of an optical map created from one single DNA molecule, but can be tamed statistically by combining the optical maps of many single molecules covering completely or partially the same genomic region and by exploiting accurate statistical models of the error sources. To a rough approximation the resolution and accuracy of an optical map can be arbitrarily improved by simply increasing the number of enzymes and number of molecules involved.

The success of optical maps for characterizing microbial genomes has been repeatedly proven over the last five years and the current research with the rice and human genomes indicates its scalability to larger, multi-chromosomal, complex and polyploid genomes [5, 11]. The ability to create haplotype maps using the algorithms described in this paper will only enhance the power of optical mapping.

### 2.1 Parameters of the experiment

We consider a set of  $M$  fragments of average length  $L$  which cover the genome of length  $G$  with coverage  $c = \frac{ML}{G}$ . On this genome we have a set of  $N$  restriction sites. Each molecule is a contiguous region from one of two haplotypes, and contained on the molecule are some restriction sites. Each molecule provides a local view of the ordered restriction sites taken from one haplotype. For each of the restriction sites found on a molecule we have data for position in the interval  $[1, G]$ . For the sake of simplicity, in this paper, we assume that non-digestion rates are negligible, and that distance data may be scaled to a consensus map so that positional data may be understood.

The consensus map may be represented by an ordered set of lengths. After the construction of a consensus map all of the observed data may be represented by an  $M \times N$  matrix  $D$ . Each row of  $D$  represents a molecule. Each column of  $D$  represents a restriction site found in the consensus map. The entry found at  $D_{ij}$  is the position of restriction site  $i$  (corresponding to consensus restriction site  $i$ ) found on molecule  $j$ , in the event that there is not a site  $i$  on molecule  $j$  then the entry  $D_{ij}$  is set to zero. The position may be specified by a metric amounting to a scaled distance in base-pairs from the 3' end of the chromosome.  $D$  is a large banded matrix whose band width is equal to the coverage  $c$  in expectation,

hence the expected sparsity of matrix  $D$  is  $\frac{c(N-1)}{MN} = \frac{c}{M}$  or  $\frac{L}{G}$ . We denote by  $c'$  be the smallest number which bounds the band width of this matrix.

### 3 EM-Algorithm for Detection of RFLP's

This section details the problem of detecting the *RFLP* event. An expectation maximization or EM-algorithm is given for the estimation of a mixed distribution model, and is followed by a criteria for deciding if data supports an RFLP event.

Consider the  $j$ th column of  $D$  and take the non-zero entries as a column vector denoted by  $a$  and consider it as an instance of random vector  $A$ , an  $n \times 1$  vector with  $\langle A_i \rangle_{i=1:n}$  i.i.d. random variables representing the position of restriction sites taken from a distribution with p.d.f. function:

$$f(A_i = x) = q_1 \frac{1}{\sqrt{2\pi\sigma^2}} \exp -\frac{(x - \mu_1)^2}{2\sigma^2} + q_2 \frac{1}{\sqrt{2\pi\sigma^2}} \exp -\frac{(x - \mu_2)^2}{2\sigma^2}. \quad (1)$$

We do not know the parameters yet, but without any loss of generality, we may assume that  $\mu_1 \leq \mu_2$ . Also,  $q_i (i = \{1, 2\})$  may be interpreted as a probability that point  $x$  is derived from the Gaussian with mean  $\mu_i$ . Further, we have  $\sum_i q_i = 1$ .

For each random column  $A$  of  $D$  there is an estimation problem: Namely, determine the values of  $\Theta := \langle \mu_1, \mu_2, \sigma, q_1 \rangle$ . Once this step is complete, we may detect RFLPs as events involving the distance between  $\mu_1$  and  $\mu_2$ , and also compute probabilities of pairwise events.

The typical approach to such problems is through maximum likelihood estimators (MLE), whereby one maximizes the probability that a particular parameter vector  $\Theta$  may produce the observed data.

$$L(\Theta) = P(A = a : \Theta) = \prod_i f_{\Theta}(A_i = a_i) \quad (2)$$

A necessary condition for  $L$  to attain a maxima at vector  $\Theta^*$  is that the gradient vanishes:

$$\frac{\partial L}{\partial \Theta_{\zeta}}(\Theta^*) = 0, \quad (3)$$

for all  $\zeta$  indexing parameters. If the Hessian or second variation is nonpositive then it is a sufficient condition for local maxima. This situation suggests a value  $\Theta^*$  for the argument of  $L$  which yields the best possible model parameters for the data. Note that a log likelihood function  $L' = \log L$  follows the same principle but has the feature that products of independent random variables are transformed to sums.

We choose to treat one of the parameters  $q_{1j}$  as the probability of a hidden random variable for each derivate  $a_i$ . Let  $Y_{1i}$  be a Bernoulli random variable whose  $p$ -value is equal to  $q_{1i}$  and represents the probability that the data point  $a_i$  is derived from the Gaussian with the leftmost mean.

We note that we have a distribution:

$$P(A_i = x, Y_{1i} = \nu | \Theta) = \mathbb{I}_{\nu=1} f_1(a_i) + \mathbb{I}_{\nu=2} f_2(a_i)$$

$$f_1(a_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp -\frac{(a_i - \mu_1)^2}{2\sigma^2} \quad \text{and} \quad f_2(a_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp -\frac{(a_i - \mu_2)^2}{2\sigma^2}$$

whose marginals are:

$$P(A_i = x | \Theta) = \mathbb{E}_{\nu} [\mathbb{I}_{\nu=1} f_1(a_i) + \mathbb{I}_{\nu=2} f_2(a_i)] = q_{1i} f_1(a_i) + q_{2i} f_2(a_i)$$

We now formulate the EM-algorithm by use of Jensen's inequality:

$$\begin{aligned}
L'(\Theta) &= \log(L(\Theta)) = \log \prod_{i=1:N} P(A_i = x | \Theta) \\
&= \sum_{i=1:N} \log \sum_{\nu=1:2} Q(Y_i = \nu) \frac{P(A_i = x, Y_i = \nu | \Theta)}{Q(Y_i = \nu)} \\
&\geq \sum_{i=1:N} \sum_{\nu=1:2} Q(Y_i = \nu) \log \left( \frac{P(A_i = x, Y_i = \nu | \Theta)}{Q(Y_i = \nu)} \right) \\
&= \sum_{i=1:N} \sum_{\nu=1:2} Q(Y_i = \nu) \log P(A_i = x, Y_i = \nu | \Theta) + H(Q_i).
\end{aligned}$$

Here  $Q$  is an arbitrary measure and  $H$  is the Entropy Function on probability vectors:

$$H(Q) = \sum_i Q_i \log\left(\frac{1}{Q_i}\right).$$

We may now define the function:

$$F(Q, \Theta) := \sum_{i=1:N} \sum_{\nu=1:2} Q(Y_i = \nu) \log P(A_i = x, Y_i = \nu | \Theta) + H(Q_i)$$

and devise an EM-algorithm as a process of increasing the  $F$  function value [14]. We note that a gradient ascent may be performed on the "likelihood surface" by alternately maximizing  $Q$  followed by  $\Theta$ .

**E-Step:**  $Q_{k+1} \leftarrow \{Q^* : \max_Q F(Q, \Theta_k) = F(Q^*, \Theta_k)\}$ .

**Lemma 1** *E-Step.*

Let  $Q$  be a vector  $\langle q_{i\nu} \rangle_{i=1:N, \nu=1:2}$  where  $N$  is the number of non-zero entries in the column of data. The Arg-Max can be solved for explicitly with :

$$Q_{k+1} = \langle q_{i1}^{k+1} \rangle_{i=1:N} \quad \text{and} \quad q_{i1}^{(k+1)} \leftarrow \left( \frac{1}{\exp\left(\frac{(a_i - \mu_1^{(k)})^2}{2\sigma^{(k)}} - \frac{(a_i - \mu_2^{(k)})^2}{2\sigma^{(k)}}\right) + 1} \right)$$

The proof is omitted.

**M-Step:**  $\Theta_{k+1} \leftarrow \{\Theta^* : \max_{\Theta} F(Q_{k+1}, \Theta) = F(Q_{k+1}, \Theta^*)\}$ .

**Lemma 2** *M-Step.*

The Arg-Max can be solved for explicitly with :

$$\begin{aligned}
\mu_1^{(k+1)} &\leftarrow \frac{\sum_{i=1:N} q_{i1}^{(k+1)} a_i}{\sum_{i=1:N} q_{i1}^{(k+1)}}, & \mu_2^{(k+1)} &\leftarrow \frac{\sum_{i=1:N} q_{i2}^{(k+1)} a_i}{\sum_{i=1:N} q_{i2}^{(k+1)}} \\
\sigma^{(k+1)} &\leftarrow \sqrt{\frac{1}{2N} \sum_{\nu=1:2} \sum_{i=1:N} q_{i\nu}^{(k+1)} (a_i - \mu_{\nu}^{(k)})^2}, & \Theta_{k+1} &= \langle \mu_1^{(k+1)}, \mu_2^{(k+1)}, \sigma^{(k+1)} \rangle
\end{aligned}$$

The proof is omitted.

**Lemma 3** *In the limit, the EM algorithm converges to a local maximum of the likelihood function in the parameter space.*

The proof is omitted.

With the lemmas we assume we have procedures called ESTEP and MSTEP. The EM-Algorithm is now:

**Algorithm 1**

```

EM( A )
  QPREV ← .5*ONES( MAX(SIZE(A)), 2)
  M ← MEAN(A )
  S ← STD(A )
  TPREV ← ( M( 1- S), M( 1 + S), S )
  QNEW ← INF
  TNEW ← INF
  WHILE( MAX( NORM( QPREV - QNEW ) , NORM( TPREV - TNEW ) ) > ε )
    QNEW ← ESTEP( QPREV, TPREV )
    TNEW ← MSTEP( QNEW, TPREV )
  ENDWHILE
  return ( QNEW, TNEW )

```

Denote the return values QNEW with matrix  $[\hat{q}_1, \hat{q}_2]_{M \times 2}$ , and TNEW with vector  $\langle \hat{\mu}_1, \hat{\mu}_2, \hat{\sigma} \rangle$ ; the “overscript hats” denoting that these quantities are estimates.

**3.1 Detection of RFLPs**

We define a detected RFLP as an outcome to our EM algorithm, it is an event such that  $|\mu_2 - \mu_1| > \delta$  for some positive  $\delta(c)$  as a function of local coverage  $c$ .

$$\text{detected RFLP} = (|\mu_2 - \mu_1| > \delta(c)) \tag{4}$$

**4 Mathematics of Phasing**

The *Phasing problem* consists of inferring two haplotype data sources which combine to provide the observed genotype data. The idea is straightforward but requires some explanation.

Each polymorphism gives rise to two alternatives for a diploid organism. We may enumerate the alternatives of polymorphism 1 by the set  $\{\uparrow_1, \downarrow_1\}$ , and likewise for polymorphism 2 the alternatives are a set  $\{\uparrow_2, \downarrow_2\}$ . Haplotype I must have one of the alternatives from each set while haplotype II must have the others. The potential co-association of alternatives for the haplotypes are shown below:

Haplotype I	$\uparrow_1$	$\uparrow_2$	(5)
Haplotype II	$\downarrow_1$	$\downarrow_2$	

Haplotype I	$\uparrow_1$	$\downarrow_2$	(6)
Haplotype II	$\downarrow_1$	$\uparrow_2$	

Haplotype I	$\downarrow_1$	$\uparrow_2$	(7)
Haplotype II	$\uparrow_1$	$\downarrow_2$	

Haplotype I	$\downarrow_1$	$\downarrow_2$	(8)
Haplotype II	$\uparrow_1$	$\uparrow_2$	

These four events ( 5 – 8 ) give all possible outcomes to pairwise events.

In practice, it will not be important to identify haplotypes, but rather which alternatives are co-associated on a haplotype. We may identify events ( 5 ) and event ( 8 ) as the event that  $\uparrow_1$ , and  $\uparrow_2$  are found on the same haplotype, and denote this *covariant* event by the symbol  $\uparrow\uparrow$ . Similarly we may identify events ( 6 ) and ( 7 ) as the event that  $\uparrow_1$ , and  $\downarrow_2$  are found on the same haplotype and denote this *contravariant* event by the symbol  $\uparrow\downarrow$ . One can compute the probabilities of these pairwise events  $\uparrow\uparrow$ , and  $\uparrow\downarrow$  with the use of a continuous multiplicative group, to be introduced below.

#### 4.1 Data maps to Group Elements, MLE homomorphism

The results of EM on each column  $A_j$  of  $D$  is value  $(Q(j), \Theta(j))$ . Consider a data point in the  $j$ th column  $d_{i'j}$ , as such it is derived from the distribution given by  $Q(j), \Theta(j)$  i.e.:

$$f(x) = q(x) \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp - \frac{(x - \mu_{j1})^2}{2\sigma_j^2} + (1 - q(x)) \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp - \frac{(x - \mu_{2j})^2}{2\sigma_j^2},$$

where  $q(x)$  is some random function giving p-values of the point at  $x$  being derived from the distribution with  $\mu_{j1}$ . Let  $q_i$  be the the probability that  $d_{i'j}$  derives from the left distribution (i.e., corresponding to the distribution  $N(\mu_{j1}, \sigma_j^2)$ ) as estimated by the EM algorithm of the preceding section.

Let  $p_i = 1 - q_i$  and identify the data point  $d_{i'j}$  with the  $2 \times 2$  matrix:

$$\begin{bmatrix} q_i & p_i \\ p_i & q_i \end{bmatrix}$$

We similarly define a map for each element in the  $j$ th column of data, and denote the dependence on column with an additional subscript  $j$ .

$$\Phi_j : d_{i'j} \rightarrow \begin{bmatrix} q_{ji} & p_{ji} \\ p_{ji} & q_{ji} \end{bmatrix}$$

The map is an injection into the set  $\mathcal{G}'$  the  $2 \times 2$  symmetric matrices with both row-sums and column-sums equal to one.  $\mathcal{G}'$  is a set with a natural group structure

Let us define a continuous Abelian group by its set members and its operation:

$$\mathcal{G} = \left\{ \begin{bmatrix} a & b \\ b & a \end{bmatrix} : a \neq b, a + b = 1 \right\}$$

$$* : \mathcal{G} \times \mathcal{G} \rightarrow \mathcal{G} : \begin{bmatrix} a & b \\ b & a \end{bmatrix} \begin{bmatrix} A & B \\ B & A \end{bmatrix} \rightarrow \begin{bmatrix} (aA + bB) & (aB + bA) \\ (aB + bA) & (aA + bB) \end{bmatrix}$$

We note that  $a \neq b$  implies that the matrix has full rank and excludes the case

$$\begin{bmatrix} .5 & .5 \\ .5 & .5 \end{bmatrix}$$

which corresponds to non-RFLPs. It has no inverse and acts as an idempotent under the operation of multiplication.

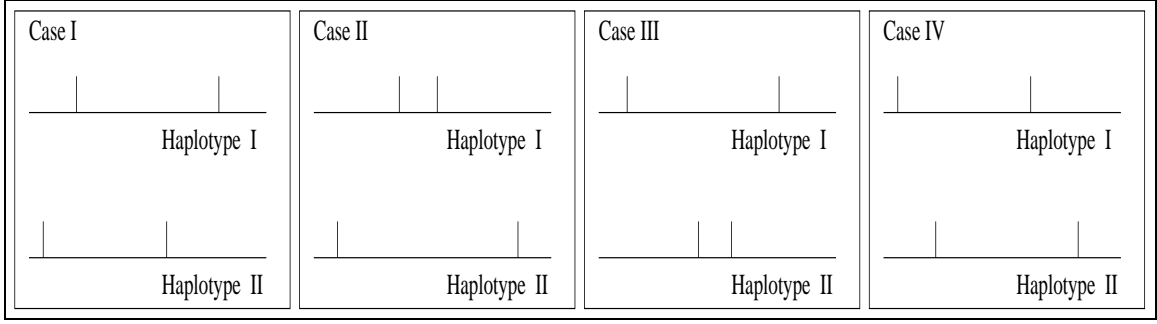


Figure 1: Case I, II, III, IV

## 4.2 Additional Operation and Closure

The set  $\mathcal{G}'$  is closed under any affine combination of elements, that is to say :  $g_1, g_2, \dots, g_n \in \mathcal{G}'$  then:

$$\sum_{k=1:n} w_k g_k \in \mathcal{G}', \quad \sum_{k=1:n} w_k = 1.$$

Here  $+$  is the usual matrix addition.

## 4.3 Computing Pairwise Events

Let us focus on two restriction sites: Site  $i$  and site  $j$  that are believed to be RFLPs. These sites have non-constant  $q_i(x), q_j(x)$  functions.

Consider molecules that span both RFLPs, these molecules contain data points  $x$  and  $y$  which were used to estimate both the functions  $q_i(x)$  and  $q_j(x)$ . Let us look at the possible sets of haplotypes:

- Case 1: Haplotype I contains the left most restriction site of  $i$  and the left most restriction site of  $j$  while haplotype II contains the right most restriction site of  $i$  and the right most restriction site of  $j$ . Denote this event as  $(i_{11}, j_{11}) \cap (i_{22}, j_{22})$  ( see figure 1 ).
- Case 2: Haplotype I contains the left most restriction site of  $i$  and the right most restriction site of  $j$  while haplotype II contains the right most restriction site of  $i$  and the left most restriction site of  $j$ . Denote this event as  $(i_{11}, j_{12}) \cap (i_{22}, j_{21})$  ( see figure 1 ).
- Case 3: Haplotype I contains the right most restriction site of  $i$  and the left most restriction site of  $j$  while haplotype II contains the left most restriction site of  $i$  and the right most restriction site of  $j$ . Denote this event as  $(i_{12}, j_{11}) \cap (i_{21}, j_{22})$  ( see figure 1 ).
- Case 4: Haplotype I contains the right most restriction site of  $i$  and the right most restriction site of  $j$  while haplotype II contains the left most restriction site of  $i$  and the left most restriction site of  $j$ . Denote this event as  $(i_{12}, j_{12}) \cap (i_{21}, j_{21})$  ( see figure 1 ).

Since we do not care to determine which haplotype the pairs are on, but rather just to determine which pairs are found together on a haplotype, the events of interest are:

$$E_1 = ((i_{11}, j_{11}) \cap (i_{22}, j_{22})) \cup ((i_{12}, j_{12}) \cap (i_{21}, j_{21})),$$

$$E_2 = ((i_{11}, j_{12}) \cap (i_{22}, j_{21})) \cup ((i_{12}, j_{11}) \cap (i_{21}, j_{22})).$$



Now we computing the probability that molecule  $\zeta$  with point  $x_\zeta$  and point  $y_\zeta$  support the event  $E_1$ :

$$\begin{aligned} P(E_1|\zeta) &= P((i_{11}, j_{11}, \zeta \in H_1) \cap (i_{22}, j_{22}, \zeta \in H_2)) \cup ((i_{12}, j_{12}, \zeta \in H_1) \cap (i_{21}, j_{21}, \zeta \in H_2)) \\ &= q_i(x_\zeta)q_j(x_\zeta) + p_i(y_\zeta)p_j(y_\zeta). \end{aligned}$$

Similarly

$$P(E_2) = q_i(x_\zeta)p_j(x_\zeta) + p_i(y_\zeta)q_j(y_\zeta).$$

But notice the connection with the group structure:  $P(E_1)$  is the entry on the diagonal while  $P(E_2)$  is the entry on the off diagonal of the product:

$$\begin{bmatrix} P(E_1|\zeta) & P(E_2|\zeta) \\ P(E_2|\zeta) & P(E_1|\zeta) \end{bmatrix} = \begin{bmatrix} q_i(x_\zeta) & p_i(x_\zeta) \\ p_i(x_\zeta) & q_i(x_\zeta) \end{bmatrix} *_G \begin{bmatrix} q_j(y_\zeta) & p_j(y_\zeta) \\ p_j(y_\zeta) & q_j(y_\zeta) \end{bmatrix}$$

Since sites on different molecules are independent, various probabilities of events ( $E_1$  and  $E_2$ ) are computed as follows:

$$\begin{aligned} \begin{bmatrix} P(E_1|\cup_{v=1:m} \zeta_v) & P(E_2|\cup_{v=1:m} \zeta_v) \\ P(E_2|\cup_{v=1:m} \zeta_v) & P(E_1|\cup_{v=1:m} \zeta_v) \end{bmatrix} &= \sum_{v=1:m} w_v \begin{bmatrix} q_i(x_{\zeta_v}) & p_i(x_{\zeta_v}) \\ p_i(x_{\zeta_v}) & q_i(x_{\zeta_v}) \end{bmatrix} *_G \begin{bmatrix} q_j(y_{\zeta_v}) & p_j(y_{\zeta_v}) \\ p_j(y_{\zeta_v}) & q_j(y_{\zeta_v}) \end{bmatrix} \\ &\text{with } \sum_{v=1:m} w_v = 1 \end{aligned}$$

When all molecules are equally “informative”:  $w_v = \frac{1}{m}$

Given two restriction sites  $\alpha$  and  $\beta$ , we define the *support* of the pair as:  $\text{Supp}(\alpha, \beta) = \{\zeta : d_{\zeta\alpha} \neq 0 \wedge d_{\zeta\beta} \neq 0\}$  or equivalently as the number of molecules indexed by  $\zeta$  that span both sites. The *phase* between two sites: RFLP  $\alpha$  and RFLP  $\beta$  may be defined as:

$$\phi(\alpha, \beta) = \frac{1}{|\text{Supp}(\alpha, \beta)|} \sum_{\zeta \in \text{Supp}(\alpha, \beta)} \Phi_\alpha(x_\zeta) *_G \Phi_\beta(y_\zeta)$$

We can also define the distance between two fragments as:

$$d_{\alpha, \beta} = \frac{1}{|\text{Supp}(\alpha, \beta)|}$$

Computing all pairwise spins can be done with a few sparse matrix multiplications:

**Algorithm 2**

```
PWS ( P )
  DIST ← ( P != 0 ) *( P != 0 )
  Q ← ONES( SIZE( P ) ) - P
  θ ← ( (P'*P) + (Q'*Q) ) ./ DIST
  return ( θ )
```

For use in large data sets we use a threshold to guard against a worst case, This idea is explained in the section on Chernoff bounds. We define a *dead state* as a spin

$$\begin{bmatrix} p & q \\ q & p \end{bmatrix},$$

where  $p$  is within an  $\epsilon$  ball of .5.

**Algorithm 3**

```

PWSDEADSTATE ( P )
  PWS ← PWS( P );
  PWS( ((PWS > .5 - ε) && (PWS < .5 + ε)) ) ← .5;
  return ( PWS )

```

## 5 Algorithms

We define the phasing problem as follows: *Given a sequence of polymorphisms (whose parameters and distributions have been estimated from a mixture model), use pairwise data to assign the polymorphisms to the haplotypes (i.e., a consistent phasing structure) such that the local assignments are consistent with the data, in the sense of maximum likelihood.*

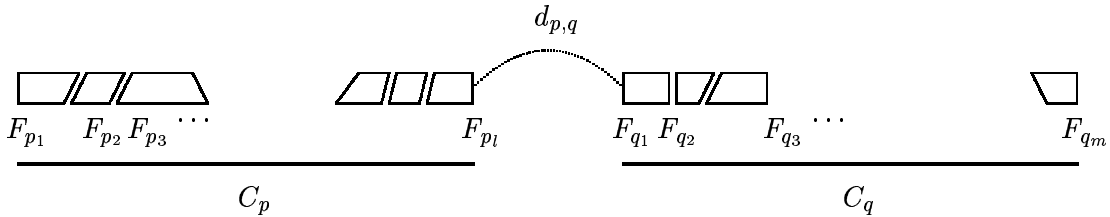
### 5.1 Weighted $k$ -Neighbor Phase-Contig Algorithm

We can define the phased contigs recursively as follows: The base case is a singleton contig:  $C_i = \{F_i\}$  shall be phased as follows:

$$\forall i, \Phi(C_i) = \begin{cases} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, & \text{if } F_i \text{ is a detected RFLP, (nontrivial contigs) ;} \\ \begin{bmatrix} .5 & .5 \\ .5 & .5 \end{bmatrix}, & \text{otherwise (trivial contigs).} \end{cases}$$

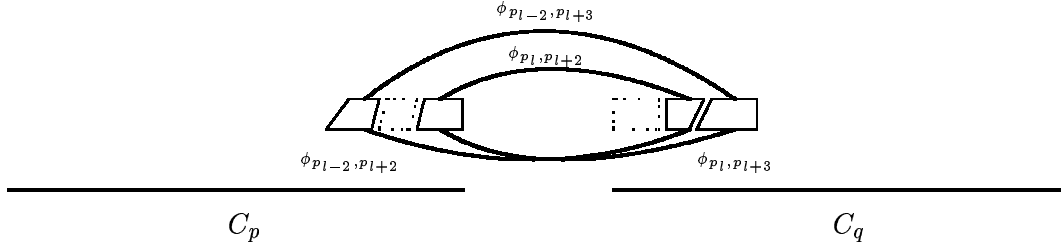
Induction cases: If  $C_p = \{F_{p_1}, F_{p_2}, \dots, F_{p_l}\}$  and  $C_q = \{F_{q_1}, F_{q_2}, \dots, F_{q_m}\}$  are phased contigs with well defined *phasing*, then the union  $C_p \cup C_q$  may be phased by a *phase-join* operation.

We define the *distance* between two phased contigs as the minimum distance between two fragments within contigs.

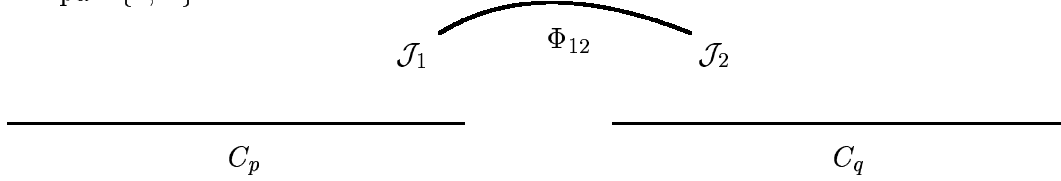


The *phase-join* operation may be performed on  $C_p$  and  $C_q$  if and only if there is a molecule  $\zeta$  which contains a data point  $x_\zeta$  from a restriction site  $F_{p_l}$  found in contig  $C_p$ , and a data point  $y_\zeta$  from a restriction site  $F_{q_1}$  found in contig  $C_q$ , as otherwise the distance is undefined.

For every pair  $F_\alpha \in C_p$  and  $F_\beta \in C_q$  there are pairwise “phasing” variables to consider in the phase-join. These pairwise phasings tell us how to orient the phased-contig  $C_q$  relative to the phased-contig  $C_p$ : we will consider a weighted combination of this information, where weights depend on distance between fragments, confidence in RFLP assignment etc.



To attempt a join of  $C_p$  to  $C_q$  we compute a *mean group action* which is a ‘least squares’ rotation to be applied similarly to all variables in the right contig to make a “fit” for all pairwise spins in the union of  $C_p \cup C_q$ . To compute the *group action* for a pair of RFLPs one in each of the phased-contigs with spin assignment  $\mathcal{J}_1$  and  $\mathcal{J}_2$ , and pairwise spin  $\Phi_{12}$ , we derive the chain of computations, let  $k_{12}$  be the *group action* for pair  $\{1, 2\}$ .



$$k_{12}\mathcal{J}_2 = \Phi_{12}\mathcal{J}_1, \quad (9)$$

$$k_{12} = \mathcal{J}_2^{-1}\Phi_{12}\mathcal{J}_1. \quad (10)$$

Solving for  $k_{12}$  we find the best rotation for these pairs in that after we update the phasing  $\mathcal{J}_2 \leftarrow k_{12}\mathcal{J}_2$  the variables would be in a state which satisfies the pairwise spin data. Thus in our algorithm the pair 1,2 casts a “weighted vote” of  $k_{12} = \mathcal{J}_2^{-1}\mathcal{J}_1\Phi_{12}$  as the *mean group element* needed to phase contig  $C_q$  correctly. In summary, the contigs are phased by the *mean group action*  $\Phi_{MGA}$ :

$$\Phi_{MGA} = \sum_{F_\alpha \in C_p, F_\beta \in C_q, \beta - \alpha < k+1} w_{\alpha\beta} k_{\alpha\beta}, \text{ where } \sum w_{\alpha\beta} = 1.$$

Now if the resulting *mean group action*

$$\Phi_{MGA} = \begin{bmatrix} p & q \\ q & p \end{bmatrix}$$

is not a dead state, the phase-join operation is successfully executed, and a parent contig  $C_p \wedge C_q$  is assigned the value  $\Phi(C_p \wedge C_q) \leftarrow \Phi_{MGA}$ . We omit the discussion of the special attention that needs to be given to the situation when either of the constituent contigs is trivial.

In addition, the phasing can also be locally improved by a *phase-adjust* operation. Phase-Adjust constitutes computing the *mean group action* for all pairwise fragments in the  $k$ -neighborhood of  $F_\alpha$ . If  $C$  denotes the contig containing  $F_\alpha$ , then the *adjust* operation involves the following steps:

$$\begin{aligned} \Phi_{MGAnew} &= \sum_{F_\beta \in C: \beta - \alpha < k+1} w_{\alpha\beta} k_{\alpha\beta}, \text{ where } k_{\alpha\beta} = \mathcal{J}_\alpha^{-1}\Phi_{\alpha\beta}\mathcal{J}_\beta \text{ and } \sum w_{\alpha\beta} = 1 \\ \mathcal{J}_\alpha &\leftarrow \Phi_{MGAnew}\mathcal{J}_\alpha \end{aligned}$$

We omit all the details of an efficient implementation, based on the “contig-tree data structure” which is constructed with the Union-Find structure [15], operating on disjoint sets of polymorphic markers

(phased contigs). Assuming that the maximum molecular coverage at any region is  $c_{max}$ , the worst case complexity of the phasing algorithm is bounded by  $O(c_{max}^2 N \gamma(c_{max}^2 N))$ . In practice, as the parameters  $c_{max}$  and  $k$  are likely to be small constants, the algorithm performs almost linearly in the number of polymorphic markers  $N$ .

## 6 Simulations and Examples

We demonstrate our algorithm on two simulated data sets. The views below are broken up into bands, the simulated haplotypes are in the bottom-most band of the layout. Above that is the haplotype molecule map for a diploid organism, these molecule maps are available to the algorithm as mixed data, the segmentation shown is unknown to the algorithm. The third band indicates estimate values and here we can see what features the EM algorithm for mixed Gaussian chooses as RFLPs. Mistakes occur with the lack of a deep library. The fourth band in the layout indicates the history of contig-operations and from this tree one can view: 1) the developing  $k$ -neighborhoods used to compute mean group action, and 2) the distinct phased contigs. The top band in the layout gives the algorithmic output to this problem, complete with phasing in subsets that span the distance indicated by the bars. Areas where phase structure overlaps but cannot extend indicate regions that are of interest to target with more specific sequences to extend the phasing.

Parameters of the simulations are summarized in the table:

Parameter	Symbol	Data Set 1	Data Set 2
number of molecules	M	80	150
number of fragments RFLP and non RFLP	F	20	100
size of the genome	G	12000	50000
expected molecule size	EMS	2000	2000
variance in molecule size	VMS	50	500
variance in fragment length size	VFS	1	20
P-value that any given Fragment is an RFLP	P-BIMODE	.5	.3
Expected separation of means for RFLP	ERFLPSEP	10	50
Variance in the separation of means for RFLP	VRFLPSEP	.01	6

Any parameter with both an expectation and variance is generated with a normal distribution. From these parameters one can compute some additional symbols that we use in the paper  $L = EMS$  and  $c = \frac{LM}{G}$ .

For the first simulation on data set I seen in figure 5.1 a relatively small set is chosen so that one can view the action of the algorithm, here the neighborhood size is set to  $k = 5$  and there is no  $\epsilon$  guard of the dead state, still things work pretty well, and one can see that any mistakes are due to the low coverage library.

In the second simulation on data set II seen in figure 5.1 we illustrate that similar results may be achieved on large data sets.

## 7 Future Research

In the full paper we provide analysis of the EM algorithm. We found that a false negative event ( a true polymorphisms mis-classified by our EM algorithm) generally do not pose any threat to the phasing results. False positives ( non-polymorphic sites that appear to be polymorphic) may introduce trouble, but we were able to attain Chernoff bounds, eliminating this treat to the phasing with reasonable experimental parameter values.



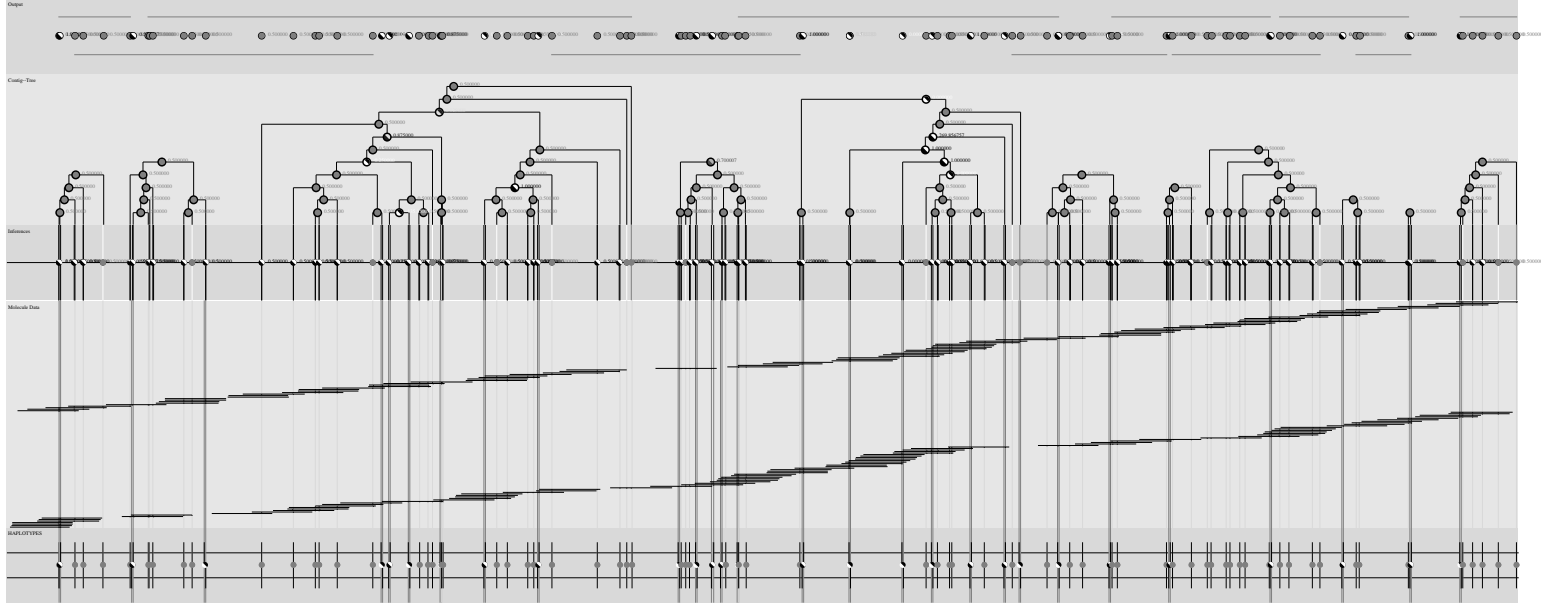


Figure 3: Data set II

We have investigated the algorithms ability to reduce the large search space of phase configurations, and we have results on the expected combinatorial reduction, as it relates to the number of join operations.

Currently we are investigating how the successful phasing of chromosome copies haplotypes may be extended to population studies, where in a large population of  $N$  individuals there exists some  $K < N$  actual genomic variations.

## References

- [1] A. AHO, J. HOPCROFT AND J. ULLMAN. **The Design and Analysis of Computer Algorithms**, In *Addison-Wesley Series in Computer Science and Information Processing*, Addison-Wesley Publishing Company, Reading Massachusetts, 1974.
- [2] T.S. ANANTHARAMAN AND B. MISHRA. "Genomics via Optical Mapping I: Probabilistic Analysis of Optical Mapping Models," Courant Technical Report, No. TR # 1998-770, August, 1998.
- [3] T.S. ANANTHARAMAN AND B. MISHRA. "A Probabilistic Analysis of False Positives in Optical Map Alignment and Validation," *Algorithms in Bioinformatics*, First International Workshop, WABI 2001 Proceedings, LNCS 2149:27–40, Springer-Verlag, 2001.
- [4] T.S. ANANTHARAMAN, B. MISHRA AND D.C. SCHWARTZ. "Genomics via Optical Mapping II: Ordered Restriction Maps," **Journal of Computational Biology**, 4(2):91–118, 1997.
- [5] T.S. ANANTHARAMAN, B. MISHRA AND D.C. SCHWARTZ. "Genomics via Optical Mapping III: Contiging Genomic DNA," *Intelligent Systems for Molecular Biology: ISMB '99*, (Heidelberg, Germany), 7:18–27, AAAI Press, 1999.
- [6] W. CASEY, B. MISHRA, AND M. WIGLER. "Placing Probes on the Genome with Pairwise Distance Data," **Algorithms in Bioinformatics: first international workshop: proceedings WABO 2001**, 52–68, Springer, New York, 2001
- [7] A. CLARK. "Inference of Haplotypes from PCR-Amplified Samples of Diploid Populations," **Mol. Biol. Evol.** 7:111–122, 1990.
- [8] A. CLARK, K. WEISS, AND D. NICKERSON. "Haplotype Structure and Population Genetic Inferences from Nucleotide-Sequence Variation in Human Lipoprotein Lipase," **Am. J. Human Genetics**, 63:595–612, 1998.
- [9] H. CHERNOFF. "A Measure of Asymptotic Efficiency for Tests of a Hypothesis Based on the Sum of Observations," **Annals of Mathematical Statistics**, 23:483–509, 1952.
- [10] D. GUSFIELD. "Inference of Haplotypes from Samples of Diploid Populations: Complexity and Algorithms," **Journal of Computational Biology**, 8-3:305–323, 2001.
- [11] J. LIN, R. QI, C. ASTON, J. JING, T.S. ANANTHARAMAN, B. MISHRA, O. WHITE, M.J. DALY, K.W. MINTON, J.C. VENTER, ET AL. 1999. "Whole-Genome Shotgun Optical Mapping of *Deinococcus radiodurans*," **Science**, 287:675–680
- [12] L. PARIDA, AND B. MISHRA. "Partiioning K clones: Hardness results and practical algorithms for the K-populations problem," **RECOMB98**, 192–201, ACM press, 1998.
- [13] L. PARIDA, AND B. MISHRA. "Partitioning Single-Molecule Maps into Multiple Populations: Algorithms And Probabilistic Analysis," *Discrete Applied Mathematics*, (The Computational Molecular Biology Series), 104(1-3):203–227, August, 2000.
- [14] S. ROWEIS, AND Z. GHAHRAMANI. "A Unifying Review of Linear Gaussian Models," **Neural Computation**, 11(2):305–345, 1999
- [15] R. E. TARJAN. **Data Structures and Network Algorithms**, CBMS 44, SIAM, Philadelphia, 1983.