

# Necessary and Sufficient Conditions for a Phenotype System to Have a Factor-Union Representation

GEORGE MARKOWSKY

*Computer Sciences Department, IBM T. J. Watson Research Center,  
Yorktown Heights, New York 10598*

*Received 8 June 1982; revised 22 April 1983*

---

## ABSTRACT

This paper presents an algorithm for determining whether a phenotype system has a factor-union representation, i.e., whether sets of properties (factors) can be assigned to alleles so that the phenotype of each genotype is determined by the union of the factor sets of the alleles. If such a representation is possible, the algorithm can construct one using  $p$  or fewer factors, where  $p$  is the number of phenotypes. Refinements to the algorithm are presented that often construct representations using few factors. An example is presented to show that some phenotype systems can be represented only by using almost as many factors as there are phenotypes in the system. Finally, it is shown that the problem of finding factor-union representations for multilocus systems reduces to the problem for 1-locus systems.

---

## 1. INTRODUCTION

Cotterman [1] investigated the existence of factor-union representations for phenotype systems and derived some partial answers. This paper completely solves the problem and presents an algorithm which always determines whether a phenotype system has a factor-union representation. This algorithm is efficient enough to be used manually on some large examples.

If a factor-union representation is possible for a phenotype system, the algorithm can construct one using the phenotypes themselves as factors. In general, this may use many more factors than necessary, since it is possible to represent  $2^k$  distinct phenotypes with  $k$  factors. Refinements to the algorithm can significantly reduce the number of factors required. As an example, two representations are constructed for an example left unresolved in [1]. The first one uses 12 factors, while the second one uses only 4.

An example in this paper shows that some phenotype systems require almost as many factors to represent them as they have phenotypes. This

example is also of interest because the number of factors required grows quadratically in the number of alleles, which shows that simple representations are not always possible.

The representations constructed by the algorithm need not be based on actual physical factors. Determining a physical basis for a representation is a problem that cannot be solved purely mathematically, but the mathematical results can guide the search for suitable physical factors. Of course, if the algorithm shows that no representation is possible, much valuable time can be saved.

In Section 5, techniques are developed to translate multilocus systems into 1-locus systems so that finding a factor-union representation for one system finds it for the other. Thus, this algorithm can be used with any phenotype system.

## 2. DEFINITIONS AND NOTATION

Listed below are several simple definitions and concepts required for understanding this paper. Many of them come from [1], which should be consulted for general background on this topic. A few come from the theory of partial orders.

- (1) Lowercase letters denote alleles, Greek letters denote genotypes, and uppercase letters denote phenotypes. Pairs of lowercase letters will be used when it is desirable to show the alleles contained in a genotype. Bracketed expressions such as  $\langle ab, cd \rangle$  are used to show the genotypes in a phenotype.
- (2) If the two alleles in a genotype are the same, it is called a *homozygote*. If they are different, it is called a *heterozygote*.
- (3) The symbol  $+$  denotes the union of individual sets,  $\Sigma$  the union of a family of sets, and  $\Pi$  the intersection of a family of sets. The symbols  $\leq$  and  $<$  denote set inclusion and proper set inclusion respectively.
- (4) A *factor-union representation* of a phenotype system  $\Phi$  consists of a set  $X$  of factors, called the *factor set*, and an assignment  $f$  of subsets of  $X$  to the alleles of  $\Phi$  with the following properties:
  - (a) If  $ab$  and  $cd$  are two genotypes in the same phenotype,

$$f(a) + f(b) = f(c) + f(d).$$

The set  $f(a)$  is called the *factor set* of  $a$ . If  $ab$  is in  $P$  then

$$f(a) + f(b) = \sum_{x \text{ in } P} f(x).$$

The quantity on the right side is abbreviated  $f(P)$  and called the *factor set* of  $P$ . It is the union of the factor sets of all alleles in  $P$ .

- (b) If  $P$  and  $Q$  are distinct phenotypes,  $f(P) \neq f(Q)$ . This condition states that phenotypes are determined by their factors.
- (5) A *partially ordered set (poset)* is a set  $X$  together with a partial order  $<$  defined on its elements. A *partial order* is any binary relation that satisfies the following two rules:
- (a) *Reflexive property.* For all  $x$  in  $X$ ,  $x < x$  never holds.
- (b) *Transitive property.* For all  $x, y, z$  in  $X$ , if  $x < y$  and  $y < z$ , then  $x < z$ .

The string  $x \leq y$  means that either  $x < y$  or  $x = y$ . By convention,  $y > x$  and  $y \geq x$  mean the same as  $x < y$  and  $x \leq y$ . The same symbols are used for partial order and set inclusion. This should cause no confusion, since set inclusion is a form of partial order.

- (6) Let  $(X, <)$  be a poset and  $S$  a subset of  $X$ . An element  $z$  of  $X$  is called an *upper bound for  $S$*  if it is greater than or equal to every element of  $S$ . An element is the *sup* or *least upper bound* of  $S$  if it is an upper bound that is less than or equal to every upper bound of  $S$ . In general, a subset of a poset need not have a sup. A poset in which every finite nonempty subset has a sup is called a *semilattice*.
- (7) If  $X$  is a set,  $\text{Power}(X)$  denotes the *power set* of  $X$ , which is the set of all subsets of  $X$ . Set inclusion is a partial order on  $\text{Power}(X)$ . With this order,  $\text{Power}(x)$  is a semilattice in which the sup of a family of sets is its union.

### 3. THE ALGORITHM

The algorithm attempts to construct a partial order on the set of phenotypes using necessary inclusion relations among the phenotypes. If the algorithm fails to construct a partial order, Theorem 2 shows that two distinct phenotypes will have identical factor sets in any factor-union representation. Thus no factor-union representation is possible. On the other hand, if the algorithm constructs a partial order, a factor-union representation can be constructed using Theorem 1.

The following describes the part of the algorithm that determines whether a factor-union system is possible. Theorem 2 proves that the algorithm is correct and shows how to construct a factor-union representation.

#### FACTOR-UNION REPRESENTATION ALGORITHM (Existence Part)

1. Represent each phenotype of  $\Phi$  by a point.
2. For each ordered pair of distinct phenotypes  $(P_i, P_j)$ , draw an arrow from  $P_i$  to  $P_j$  if and only if some genotype contained in  $P_i$  has all its alleles

belonging to some genotype in  $P_j$ . This step creates a directed graph having the phenotypes as its points and arrows as its edges. This graph is called  $\text{Graph}(\Phi)$ . If there is an arrow from  $P_i$  to  $P_j$ , then in any factor-union representation, the factor set of  $P_i$  is a subset of the factor set of  $P_j$ .

3. If  $\text{Graph}(\Phi)$  has a *cycle*, i.e., a sequence of arrows from a point to itself,  $\Phi$  does not have a factor-union representation. A cycle in  $\text{Graph}(\Phi)$  implies that all the factor sets connected by the arrows must be equal.

4. If  $\text{Graph}(\Phi)$  does not have a cycle, define a partial order  $<$  on  $\text{Graph}(\Phi)$  as follows:  $P_i < P_j$  if and only if there is a sequence of arrows from  $P_i$  to  $P_j$ . A sequence of arrows from one point to another is called a *path*. [The condition that there are no cycles in  $\text{Graph}(\Phi)$  implies that  $<$  is reflexive. Clearly,  $<$  is also transitive and so a partial order. If  $P_i \leq P_j$ , the factor set of  $P_i$  is a subset of the factor set of  $P_j$  in any factor-union representation.]

5. The next stage of the algorithm adds arrows to  $\text{Graph}(\Phi)$  and extends the partial order  $<$  until either a cycle appears in  $\text{Graph}(\Phi)$  or there is no further way to carry out the following rule.

#### ARROW-ADDING RULE

Let  $aa$  be in  $P_1$ ,  $bb$  in  $P_2$ ,  $ab$  in  $P_3$ , and  $P_4$  a phenotype such that  $P_1 \leq P_4$  and  $P_2 \leq P_4$  in  $\text{Graph}(\Phi)$ . Then add an arrow from  $P_3$  to  $P_4$  unless there already is such an arrow.

As arrows are added to  $\text{Graph}(\Phi)$ ,  $<$  must be updated. The partial order is always defined in the same way relative to  $\text{Graph}(\Phi)$ , i.e.,  $P < P^*$  if and only if there exists a path in  $\text{Graph}(\Phi)$  from  $P$  to  $P^*$ .

To use the Arrow-Adding Rule, choose a heterozygote  $ab$ . Then find the phenotypes  $P_1$  and  $P_2$  that contain the homozygotes  $aa$  and  $bb$ . Next, find the upper bounds of  $P_1$  and  $P_2$ . Finally, add the arrows from  $P_3$  to the upper bounds. *The algorithm always halts*, since  $\text{Graph}(P)$  can have at most  $p(p-1)/2$  arrows without having a cycle.

Instead of drawing all arrows in  $\text{Graph}(P)$ , I will only draw an arrow from  $P$  to  $Q$  if  $P < Q$  and there does not exist an  $R$  such that  $P < R < Q$ . In this case  $Q$  is said to *cover*  $P$ . This representation, called the *Hasse diagram*, contains all the information of the original graph, but is easier to draw and interpret. The original order can be reconstructed from this reduced order by transitivity. To further simplify the drawings, most arrow heads will be omitted. Line segments will point from lower to higher points.

*Example 1 (System  $\Phi_1$ ).* This example, from Cotterman [1, p. 8], has 6 alleles and 12 phenotypes. Cotterman was unable to determine whether this system had a factor-union representation.

Figure 1 shows the Hasse diagram of  $\text{Graph}(\Phi_1)$  after the algorithm has finished with stage 2. There are no cycles in it, so that the algorithm proceeds

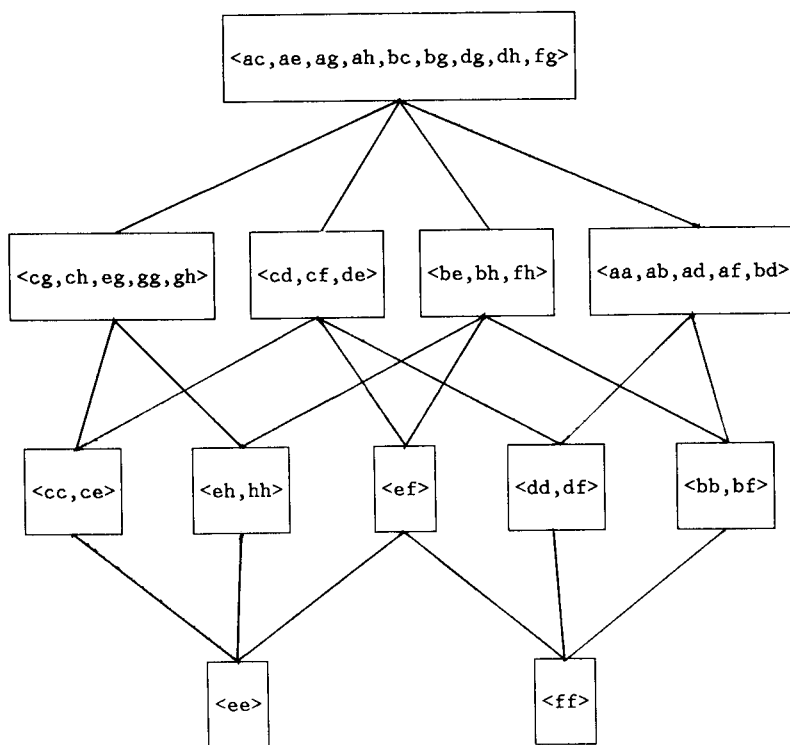


FIG. 1. The Hasse diagram of  $\text{Graph}(\Phi_1)$  after being processed by the algorithm.

to stage 4. Stage 5 lasts a short time, since the Arrow-Adding Rule cannot be applied. Thus  $\Phi_1$  has a factor-union representation.

A factor-union representation with 12 factors can be constructed for  $\Phi_1$  using Theorem 2. This representation will not be displayed since a representation using only 4 factors is produced in Section 4.

*Example 2 (System  $\Phi_2$ ).* This example also has 6 alleles. The phenotypes are shown in Figure 2, which displays the results after the first stage of the algorithm is completed.  $\text{Graph}(\Phi_2)$  does not have any cycles, and there are many ways to apply the Arrow-Adding Rule. Theorem 2 shows that the Arrow-Adding Rule can be applied to heterozygotes in any order. If, however, a factor-union representation does not exist, some sequences of heterozygotes will produce a cycle faster than others. For this example, the Arrow-Adding Rule will be applied to the heterozygotes  $gx$ ,  $hx$ ,  $hu$ ,  $uy$ , and  $ty$  in order. The intermediate steps are detailed below and the results shown in Figure 3.

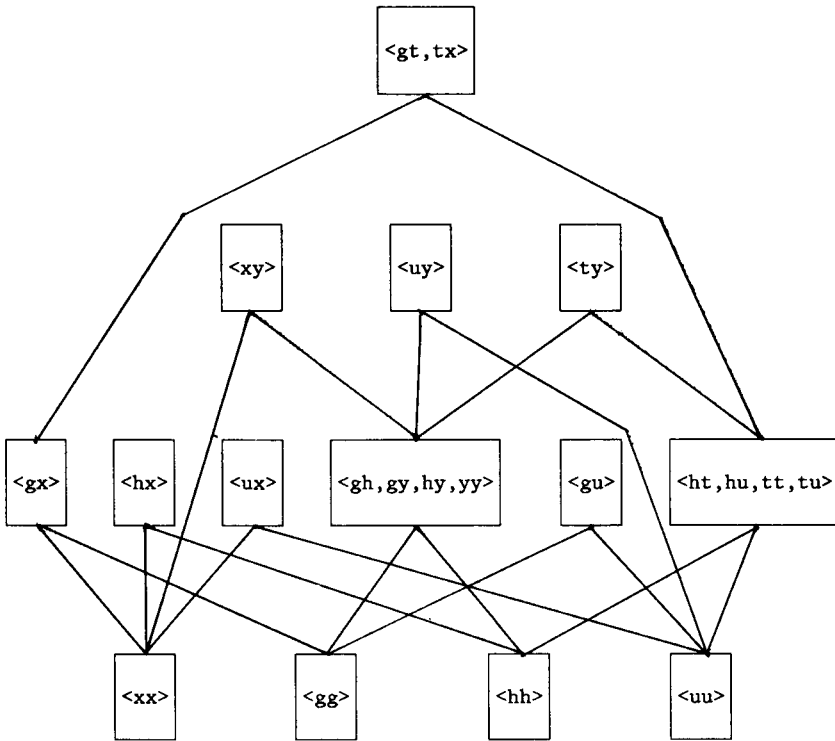


FIG. 2. The Hasse diagram of  $\text{Graph}(\Phi_2)$  after the first stage of the algorithm.

Since  $\langle xx \rangle$  and  $\langle gg \rangle$  are both less than  $\langle xy \rangle$ , an arrow must be added from  $\langle gx \rangle$  to  $\langle xy \rangle$ . The arrow from  $\langle xx \rangle$  to  $\langle xy \rangle$  can now be removed, since there is a path from  $\langle xx \rangle$  to  $\langle xy \rangle$  through  $\langle gx \rangle$ .

Since  $\langle xx \rangle$  and  $\langle hh \rangle$  are both less than  $\langle xy \rangle$  and  $\langle gt, tx \rangle$ , arrows must be added from  $\langle hx \rangle$  to  $\langle xy \rangle$  and  $\langle gt, tx \rangle$ .

Since  $\langle hh \rangle$  and  $\langle uu \rangle$  are both less than  $\langle uy \rangle$ , add an arrow from  $\langle ht, hu, tt, tu \rangle$  to  $\langle uy \rangle$ . Now the arrow from  $\langle uu \rangle$  to  $\langle uy \rangle$  can be removed, since there is a path from  $\langle uu \rangle$  to  $\langle uy \rangle$  through  $\langle ht, hu, tt, tu \rangle$ .

Again,  $\langle ty \rangle$  is greater than  $\langle uu \rangle$  and  $\langle gh, gy, hy, yy \rangle$ , so an arrow is added from  $\langle uy \rangle$  to  $\langle ty \rangle$ . Since  $\langle uy \rangle$  is greater than  $\langle ht, hu, tt, tu \rangle$  and  $\langle gh, gy, hy, yy \rangle$ , we add an arrow from  $\langle ty \rangle$  to  $\langle uy \rangle$ . This arrow combined with the previous one creates a cycle involving the two phenotypes  $\langle ty \rangle$  and  $\langle uy \rangle$ . Thus  $\Phi_2$  does not have a factor-union representation.

Based on the results just obtained, there is a direct way to show that  $\Phi_2$  cannot have a factor-union representation. Suppose, to the contrary, that it did have a factor-union representation  $f$ . Then the factors associated with

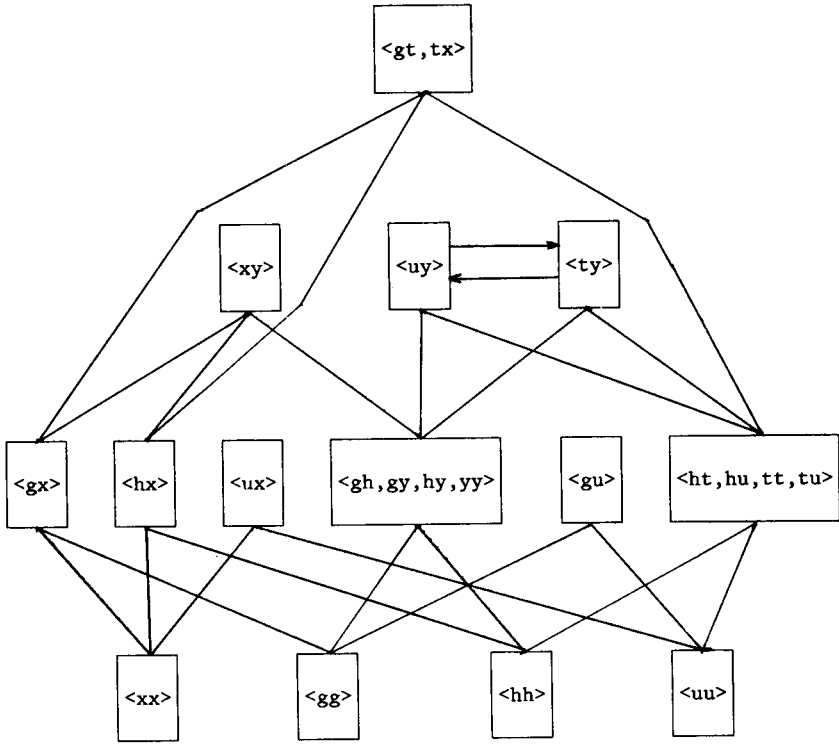


FIG. 3. The Hasse diagram of  $\text{Graph}(\Phi_2)$  when a cycle is first detected.

$\langle uy \rangle$  would be  $f(u) + f(y)$ . Similarly, the factors associated with  $\langle ty \rangle$  would be  $f(t) + f(y)$ . Since  $yy$  is in the same phenotype as  $gh$ , it would follow that  $f(y) = f(y) + f(y) = f(g) + f(h)$ . Similarly, since  $tt$  is in the same phenotype as  $hu$ , it follows that  $f(t) = f(h) + f(u)$ . Thus  $f(t) + f(y) = f(h) + f(u) + f(g) + f(h) = f(u) + f(g) + f(h) = f(u) + f(y)$ . This is impossible, since distinct phenotypes must have distinct factor sets associated with them.

Theorem 1 is a version of Theorem 2.1 in Markowsky [1973], and contains the basic facts about the representation of posets needed to prove that the algorithm is correct. Theorem 2 gives the proof that the algorithm is correct. For additional references and a more complete treatment of these ideas, see Markowsky [1980].

**THEOREM 1**

Let  $\Delta$  be a poset and  $\theta: \Delta \rightarrow \text{Power}(\Delta)$  a function given by the following rule:

$$\theta(x) = \{y \text{ in } \Delta \mid y \text{ is not } \geq x\}$$

The function  $\theta$  has the following properties:

- (1) For all  $x$  and  $y$  in  $\Delta$ ,

$$x \leq y \text{ if and only if } \theta(x) \leq \theta(y).$$

Thus distinct elements are mapped into distinct elements.

- (2) Let  $S$  be a subset of  $\Delta$ , and assume that  $S$  has a sup,  $s$ . Then

$$\theta(s) = \sum_{t \text{ in } S} \theta(t).$$

*Proof.* If  $x \leq y$ , and  $t$  is not  $\geq x$ , then  $t$  is not  $\geq y$ , since  $\geq$  is transitive. Thus  $\theta(x)$  is a subset of  $\theta(y)$ . Suppose that  $\theta(x)$  is a subset of  $\theta(y)$ . Further, suppose  $y$  is not  $\geq x$ . Then  $y$  must be an element of  $\theta(x)$  and hence of  $\theta(y)$ . But if  $y$  belongs to  $\theta(y)$ , then  $y$  is not  $\geq y$ , which is impossible. Thus  $y \geq x$ .

If  $s$  is the sup of  $S$ , then  $s \geq t$  for all  $t$  in  $S$ . From the first part of this theorem, it follows that  $\theta(s)$  contains the union of the  $\theta(t)$ 's, since it contains each individual  $\theta(t)$ . Now suppose that there exists an  $x$  which belongs to  $\theta(s)$  but not to any  $\theta(t)$  where  $t$  is an element of  $S$ . Thus  $x$  is not  $\geq s$ , but  $x \geq t$  for all  $t$  in  $S$ . Since  $s$  is the sup of  $S$ , it follows that  $x \geq s$ . This contradicts the fact that  $x$  is not in  $\theta(s)$ . ■

#### THEOREM 2

If the algorithm described above halts without finding a cycle,  $\Phi$  has a factor-union representation. If a cycle is found,  $\Phi$  cannot have a factor-union representation. Furthermore, the heterozygotes in stage 5 of the algorithm can be chosen arbitrarily.

*Proof.* If  $\Phi$  has a factor-union representation  $f$ , it is impossible for the algorithm to construct a cycle when processing  $\text{Graph}(\Phi)$ , since all the arrows introduced by the algorithm correspond to set inclusions that hold in any a factor-union representation.

If a cycle exists in  $\text{Graph}(\Phi)$ , at least two distinct phenotypes have the same factor sets, contradicting the assumption that  $f$  was a factor-union representation.

If the algorithm halts without finding a cycle, the relation  $<$  is a partial order on  $\text{Graph}(\Phi)$ . Theorem 1 produces a representation  $\theta$  assigning subsets of  $\text{Graph}(\Phi)$  to each point in  $\text{Graph}(\Phi)$ .

For each allele  $a$  in  $\Phi$ , let  $f(a) = \theta(\langle aa \rangle)$ . I claim that  $f$  is a factor-union representation of  $\Phi$ . Note that  $\theta(\langle ab \rangle) = f(a) + f(b)$ , since the Arrow-Adding Rule ensures that  $\langle ab \rangle$  is the sup of  $\langle aa \rangle$  and  $\langle bb \rangle$  in  $\text{Graph}(\Phi)$ , and the representation of Theorem 1 transforms sups into unions. If  $\alpha$  is a genotype, define  $f(\alpha)$  to be the union of the factor sets of the alleles in the genotype. It follows that  $f(\alpha) = \theta(P)$ , where  $P$  is the phenotype containing  $\alpha$ .



Since  $\theta$  maps distinct phenotypes into distinct subsets,  $f$  is a factor-union representation of  $\Phi$ .

The order in which heterozygotes are considered in stage 4 is immaterial, because either a phenotype system has a factor-union representation or it doesn't. If it does, it is impossible for the algorithm to find a cycle, since such a cycle would imply the equality of two distinct phenotypes. Similarly, if a system doesn't have a factor-union representation, the algorithm must eventually discover a cycle, else a factor-union representation could be constructed for the system. ■

#### 4. THE SIZE OF A REPRESENTATION

*Example 3 (System  $\Phi_3$ ).* This example has  $2k$  alleles:  $a_1, a_2, \dots, a_k$  and  $x_1, x_2, \dots, x_k$ , and  $k^2 + 2k + 1$  phenotypes:

- (1)  $\langle a_i a_i \rangle$  for  $i = 1, 2, \dots$  ( $k$  of this type).
- (2)  $\langle x_i x_i \rangle$  for  $i = 1, 2, \dots$  ( $k$  of this type).
- (3)  $\langle a_i x_j \rangle$  for  $i, j = 1, 2, \dots$  ( $k^2$  of this type).
- (4)  $\langle a_i a_j, x_i x_j | i \neq j \text{ and } i, j = 1, 2, \dots, k \rangle$  (1 of this type).

I claim that any factor-union representation of  $\Phi_3$  must have at least  $k^2$  factors in it. To see this, assume a factor-union representation  $f$  for  $\Phi_3$  has been given. Let  $S$  be factor set of the phenotype  $\langle a_1 a_2, \dots \rangle$ , and  $S_{i,j}$  the factor set of the phenotype  $\langle a_i x_j \rangle$ . Here  $i$  and  $j$  both range over  $1, \dots, k$ . Clearly,  $S_{i,j}$  is a subset of  $S$  for all  $i$  and  $j$ . Let  $C_{i,j}$  be  $S - S_{i,j}$ . Each  $C_{i,j}$  is not empty, since each  $S_{i,j}$  must be a proper subset of  $S$ .

I claim that the intersection of two distinct  $C_{i,j}$ 's must be empty. If  $(i, j) \neq (i', j')$ , then either  $i \neq i'$  or  $j \neq j'$ . In the first case,  $S_{i,j} + S_{i',j'}$  contains  $S$ , since  $S = f(a_i a_{i'})$ . The same result holds in the second case. Since the union of  $S_{i,j}$  and  $S_{i',j'}$  is  $S$ , the intersection of  $C_{i,j}$  and  $C_{i',j'}$  must be empty. Thus the  $C_{i,j}$ 's are a family of  $k^2$  mutually disjoint, nonempty subsets of  $S$ . This implies that  $S$  has at least  $k^2$  factors.

$\Phi_3$  has  $k^2 + 2k + 1$  phenotypes. As  $k$  increases, the minimum number of factors necessary to represent it asymptotically approaches the number of phenotypes, and is quadratic in the number of alleles.

The following theorem is useful for finding minimal factor-union representations when the poset turns out to be a semilattice. It is an immediate consequence of Proposition 2.6 in [2]. For additional details on this result see [3].

#### THEOREM 3

*Let  $\Delta$  be a finite semilattice and  $\chi: \Delta \rightarrow \text{Power}(X)$  a function which is injective (1-to-1), preserves order, and preserves sups. Then the number of elements in  $X$  is  $\geq$  the number of elements of  $\Delta$  that have exactly one element*

covering them in the Hasse diagram of  $\Delta$ . These elements are called meet-irreducibles. Furthermore, if  $M$  is the set of all meet-irreducible elements in  $P$ , then  $\theta: \Delta \rightarrow \text{Power}(M)$  given by

$$\theta(x) = \{y \text{ in } M \mid y \text{ is not } \geq x\}$$

maps distinct elements into distinct elements, preserves order, and preserves sups.

If  $\text{Graph}(\Phi)$  is a semilattice such that in any factor-union representation all sups correspond to unions, Theorem 3 shows how to easily find a factor-union representation using the smallest possible number of factors. Checking whether  $\text{Graph}(\Phi)$  is a semilattice only requires that you check that every pair of elements has an upper bound. Checking whether unions would always correspond to sups is more involved, and is illustrated in the following examples.

It is easy to check that  $\text{Graph}(\Phi_1)$  is a semilattice, with sups being unions. For example,  $P = \langle cd, cf, de \rangle$  and  $Q = \langle be, bh, fh \rangle$  have  $R = \langle ac, ae, \dots, fg \rangle$  as a sup in the partial order constructed. In any factor-union representation  $f$  of  $\Phi_1$ ,  $f(R)$  must contain the union of  $f(P)$  and  $f(Q)$ . To apply Theorem 3, it must be the union of  $f(P)$  and  $f(Q)$ . This follows because  $f(R) = f(b) + f(c)$ , but  $f(P) + f(Q) = f(c) + f(d) + f(b) + f(e)$ . The argument for all other sups is similar. Thus any factor-union representation of  $\Phi_1$  must be a mapping of the type described in Theorem 3.

$\text{Graph}(\Phi_1)$  has 4 meet-irreducible elements:  $\langle aa, \dots \rangle$ ,  $\langle cg, \dots \rangle$ ,  $\langle be, \dots \rangle$  and  $\langle cd, \dots \rangle$ . Let these be represented by the symbols  $w$ ,  $x$ ,  $y$ , and  $z$  for brevity. The 4-factor representation of  $\Phi_1$  derived from Theorem 3 is illustrated in Figure 4.

The phenotypes corresponding to the meet-irreducible elements have both a label (one of  $w$ ,  $x$ ,  $y$ , or  $z$ ) and a factor set of labels. It is important to distinguish between the two. The phenotype  $\langle cd, cf, de \rangle$  has the label  $z$ , but receives the factor set  $\{w, x, y\}$  because the other three meet-irreducible elements,  $w$ ,  $x$ , and  $y$ , are not greater than or equal to it. The phenotypes that are not meet-irreducibles only have a factor set associated with them, so there is no confusion.

Example 3 is analyzed similarly. The meet-irreducible elements turn out to be the  $\langle a_i x_j \rangle$ 's where  $i, j = 1, 2, \dots, k$ . Here Theorem 3 shows that  $k^2$  is a possible number of factors for a factor-union representation and that it is also the smallest possible number.

## 5. MULTILOCUS SYSTEMS

Finding factor-union representations for multilocus systems reduces in a straightforward way to finding factor-union representations for 1-locus sys-

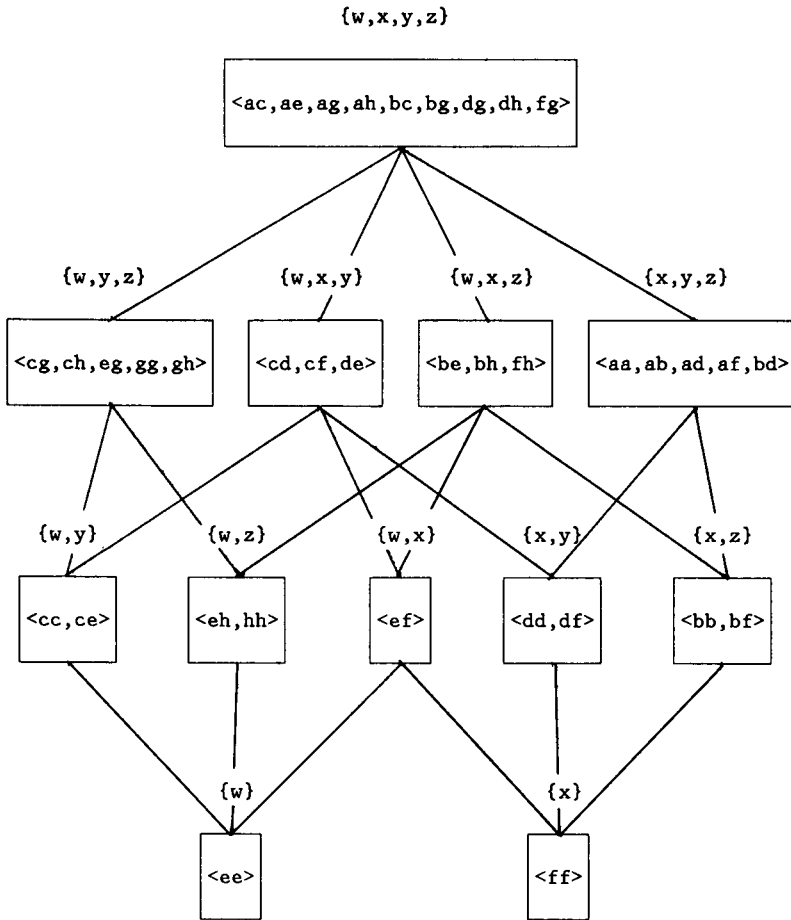


FIG. 4. A 4-factor representation of  $\Phi_1$ .

tems. A computer is probably required for determining whether even small multilocus systems have a factor-union representation.

**DEFINITIONS AND NOTATION**

- (1) The conventions introduced earlier for alleles and phenotypes are still in effect. A new notation, however, is required to represent genotypes. Thus *aa-bc-dd* represents an individual homozygous at the first and third loci and heterozygous at the second locus. For simplicity, it will be assumed that an allele can occur only at a particular locus. Because of this assumption, distinct genotypes must have different sets of alleles.
- (2) The genotype *ab* can be considered the union of *aa* and *bb*. Similarly,

- $ab-cd-ee$  can be considered the union of  $aa-cd-ee$  and  $bb-cd-ee$ . In general, the union of two or more genotypes is possible if and only if there are no more than two different alleles at each locus.
- (3) In the remainder of this paper, the term *homozygote* will refer to a genotype which is homozygous at each locus. The term *heterozygote* will refer to a genotype which is heterozygous at some locus.
  - (4) Given a multilocus system  $\Phi$ , a 1-locus system, called the *1-locus model* of  $\Phi$  and denoted by  $\text{Trans}(\Phi)$ , can be derived from it as follows:
    - (a) The alleles of  $\text{Trans}(\Phi)$  are derived from the homozygotes of  $\Phi$ , i.e., for each homozygote of  $\Phi$ , use a new symbol to represent an allele in  $\text{Trans}(\Phi)$ . If  $\alpha$  is a homozygote in  $\Phi$ , let  $A(\alpha)$  represent the corresponding allele of  $\text{Trans}(\Phi)$ .
    - (b) Since each genotype  $\alpha$  in  $\Phi$  can be written as the union of two homozygotes and the union of any two homozygotes is possible, every genotype in  $\Phi$  corresponds to some genotype in  $\text{Trans}(\Phi)$ . Any such genotype is denoted by  $\text{Trans}(\alpha)$ .
    - (c) Define the phenotypes as follows.  $\text{Trans}(\alpha)$  and  $\text{Trans}(\beta)$  are in the same phenotype if and only if  $\alpha$  and  $\beta$  are in the same phenotype in  $\Phi$ . For each phenotype  $P$  in  $\Phi$ ,  $\text{Trans}(P)$  denotes the corresponding phenotype in  $\text{Trans}(\Phi)$ . [In general,  $\text{Trans}(\Phi)$  has more genotypes than  $\Phi$ , but the same number of phenotypes.]
  - (5) A *factor-union representation* for a multilocus system  $\Phi$  consists of a set of factors and an assignment of subsets of factors to the alleles of  $\Phi$ , such that the factor sets associated with all the genotypes within a given phenotype are equal, and distinct phenotypes have distinct factor sets. The factor set of a genotype is computed by taking the union of factor sets associated with the alleles at each locus. This definition is essentially the same as the definition for 1-locus systems.

*Example 4.* Let  $\Phi$  be a 2-locus system with two alleles at each locus. Suppose the alleles at the first locus are  $w$  and  $x$ , while the alleles at the second locus are  $y$  and  $z$ .  $\Phi$  has nine genotypes: four homozygotes and five heterozygotes. Further, suppose each genotype of  $\Phi$  forms a distinct phenotype. To construct  $\text{Trans}(\Phi)$ , pick four new symbols corresponding to the homozygotes of  $\Phi$ . I will use the following correspondence:

$$\langle ww-yy \rangle \rightarrow a, \quad \langle ww-zz \rangle \rightarrow b, \quad \langle xx-yy \rangle \rightarrow c, \quad \langle xx-zz \rangle \rightarrow d.$$

$\text{Trans}(\Phi)$  has ten genotypes and nine phenotypes. Each phenotype, except for  $\langle ad, bc \rangle$ , contains just a single genotype. The phenotype  $\langle ad, bc \rangle$  corresponds to the phenotype  $\langle wx-yz \rangle$  in  $\Phi$ .

#### THEOREM 4

*Let  $\Phi$  be a multilocus system, and  $\text{Trans}(\Phi)$  the 1-locus model of  $\Phi$ . Then  $\Phi$  has a factor-union representation if and only if  $\text{Trans}(\Phi)$  does.*

*Proof.* Suppose that  $\Phi$  has a factor-union representation. Let  $f$  be the function that assigns sets of factors from a set  $S$  to the alleles in  $\Phi$ . The factor-union representation of  $\text{Trans}(\Phi)$  will use the same set  $S$  and a function  $g$  similar to  $f$ . For each allele  $A(\alpha)$  in  $\text{Trans}(\Phi)$ , define

$$g(A(\alpha)) = f(\alpha).$$

Since all unions in  $\text{Trans}(\Phi)$  are derived from unions in  $\Phi$ , it is easy to see that we have constructed a factor-union representation for  $\text{Trans}(\Phi)$ .

Deriving a factor-union representation of  $\Phi$  from one of  $\text{Trans}(\Phi)$  is more involved, since none of the entities of  $\text{Trans}(\Phi)$  correspond directly to the alleles of  $\Phi$ .

Some additional concepts are needed to conclude the proof. If  $X$  is a set of alleles of  $\Phi$ , let  $\text{Hom}(X)$  be the set of homozygotes that contain all the elements in  $X$ . If  $X$  has more than one allele from the same locus,  $\text{Hom}(X)$  will be empty. Now define

$$F(X) = \prod_{\alpha \text{ in } \text{Hom}(X)} g(A(\alpha)).$$

I claim that if  $X$ ,  $Y$ , and  $Z$  are sets of alleles of  $\Phi$  such that  $Z = X + Y$  and  $\text{Hom}(Z)$  is nonempty,

$$F(Z) = F(X) + F(Y) \tag{*}$$

Since  $\text{Hom}(Z)$  is a subset of  $\text{Hom}(X)$  and  $\text{Hom}(Y)$ , both of them are also nonempty. Pick arbitrary  $\chi$  in  $X$  and  $\nu$  in  $Y$ . There is a  $\zeta$  in  $\text{Hom}(Z)$  and another homozygote  $\lambda$  such that  $\zeta + \lambda$  are the same genotype as  $\chi + \nu$ . For brevity, use  $\mu$  to denote this genotype. Because  $g$  is a factor-union representation,

$$g(A(\zeta)) \leq g(A(\mu)) = g(A(\chi)) + g(A(\nu))$$

Furthermore,

$$F(Z) \leq g(A(\zeta)).$$

Thus

$$F(Z) \leq g(A(\chi)) + g(A(\nu))$$

for all  $\chi$  in  $X$  and  $\nu$  in  $Y$ . By fixing  $\chi$ , letting  $\nu$  vary through all of  $Y$ , and taking the intersection of all terms, we get

$$F(Z) \leq g(A(\chi)) + F(Y).$$

Now by letting  $\chi$  vary through all of  $X$  and taking the intersection of all

terms produced, we get

$$F(Z) \leq F(X) + F(Y)$$

Since  $F(Z)$  is  $\geq F(X)$  and  $F(Y)$ , we have  $F(Z) \geq F(X) + F(Y)$ . Consequently, (\*) follows.

For each allele  $a$  in  $\Phi$  define

$$f(a) = F(\{a\}).$$

Let  $\text{Alleles}(\alpha)$  be the set of all alleles in the homozygote  $\alpha$ . By induction and (\*) it follows that

$$F(\text{Alleles}(\alpha)) = \sum_{a \text{ in } \alpha} f(a).$$

But

$$\text{Hom}(\text{Alleles}(\alpha)) = \{\alpha\},$$

so

$$F(\text{Alleles}(\alpha)) = g(A(\alpha)).$$

This means that the "factor sets" assigned to the alleles by  $f$  are such that when the appropriate unions are taken, each homozygote receives the same set of factors  $g$  assigns to its image in  $\text{Trans}(\Phi)$ . Because heterozygotes are unions of homozygotes, it follows that  $f$  is a factor-union representation of  $\Phi$ . ■

*I gratefully thank Dr. Kenneth Lange for introducing me to this problem, supplying me with various references, and some very helpful comments on an early draft of this paper. I would also like to thank Dr. Joshua Lederburg for an interesting discussion that led to my cleaning up the section on multilocus systems. At a seminar where I presented these results he indicated that my reduction of 2-locus systems should also hold for multilocus systems. I thought the matter through and produced the proof that this is correct. Finally, I would like to thank the referee for making some suggestions that improved the exposition of the paper.*

#### REFERENCES

- 1 C. W. Cotterman Factor-union phenotype systems, in *Computer Applications in Genetics* (N. E. Morton, Ed.), Univ. of Hawaii Press 1969, 1-19.
- 2 G. Markowsky, Combinatorial aspects of lattice theory with applications to the enumeration of free distributive lattices, Ph.D thesis, Harvard Univ., 1973.
- 3 G. Markowsky, The representation of posets and lattices by sets, *Algebra Universalis* 11:173-192 (1980).