# Relationship between heterosis and heterozygosity at marker loci: a theoretical computation

A. Charcosset, M. Lefort-Buson and A. Gallais

C.N.R.S. – I.N.R.A. – U.P.S Station de Génétique Végétale, Ferme du Moulon, F-91190 Gif/Yvette, France

**Summary.** In this paper we have studied the linear correlation between a genetic distance index between two parent lines (based on marker loci information) and the heterosis observed in the $F_1$ hybrid from the two lines, for a quantitative character (determined by several loci, or QTL). Theoretical computations of the correlation coefficient ($\varrho$) between the distance index and the heterosis were made, assuming the biallelic model (defined by Fisher). When the alleles at both marker loci and QTL are equally distributed among the whole population of considered lines, the coefficient $\varrho$ is a function of the squares of linkage disequilibria between alleles at marker loci and alleles at QTL. The QTL that are not marked by marker loci and marker loci that do not mark any QTL play symmetrical roles and can decrease $\varrho$ greatly. We conclude that the prediction of $F_1$ hybrid heterosis based on marker loci would be more efficient if these markers were selected for their relationship to the alleles implicated in the heterotic traits considered.

**Key words:** Heterosis – Heterozygosity – Markers

## Introduction

Maize breeders continuously face the choice of inbred combinations to be tested. Many studies have discussed the interest of distance indexes between parent lines in order to achieve this choice, and particularly the use of enzymatic indexes [see Brunel (1985) for a review].

With maize, numerous authors have obtained negative results when trying to relate enzymatic diversity between lines and either heterosis or $F_1$ value, for yield or other agronomic characteristics (Hunter and Kannenberg 1971; Heidrich-Sobrinho and Cordeiro 1975; Hadji-

nov et al. 1982; Price et al. 1986; Lamkey et al. 1987; A. Charcosset, M. Lefort-Buson and M. Grenèche, unpublished results). These negative results could be partially explained by the low number of loci that were considered in each study. Recent experimental results using a large number of RFLP markers showed much better correlations between the parental distance and the hybrid trait (Fidgore 1987; Walton and Helentjaris 1987; Smith and Smith 1989; Lee et al. 1989).

However, several authors have found a positive relationship between the enzymatic distance between parents and $F_1$ yield value, even when few enzymatic loci were considered (Frei et al. 1986; Edwards et al. 1987). In both cases, there was a strong linkage disequilibrium among the population of potential parents. Therefore, the lack of a linkage disequilibrium between marker loci and loci involved in the variation of a quantitative character is probably another explanation of the numerous negative results that have been mentioned.

The aim of this paper is to present a theoretical approach to the relationship between a distance index based on markers and the observed heterosis for a polygenic character, using a simple genetic model. This relationship is discussed according to the strength of the linkage between marker loci and loci involved in the variation of the considered character. The need for selecting "efficient" markers to build up the distance index is also discussed.

## Basis of the model

We shall consider a population of inbred lines and assume that phenotypes of lines and hybrids follow the biallelic genetical model defined by Fisher. If we consider Hayman's notation, the genotypes at the $l^{th}$ locus (LL, Ll, ll) are represented by the variable $\theta_l$, which takes the

572

values $+1, 0, -1$, respectively. When the locus $l$ controls a quantitative trait (i.e., a QTL), the phenotypes of individuals can be written: $c + a_l \theta_l + d_l(1 - \theta_l^2)$ ($c$ is a constant, $a_l$ is half the difference between homozygotes LL and ll phenotypes, $d_l$ is the difference between heterozygote phenotype and the mean of homozygotes phenotypes, i.e., the dominance effect). When the trait is controlled by $n_l$ loci acting independently, the phenotype is: $\sum_{l=1}^{l=nl} a_l \theta_l + d_l(1 - \theta_l^2)$. The per se value of an inbred line $i$ is:

$$L(i) = Y_{ii} = \sum_l a_l \theta_l^i.$$

The value of the single cross $i \times j$ is:

$$Y_{ij} = \sum_l a_l(\theta_l^i + \theta_l^j)/2 + d_l(1 - \theta_l^i \theta_l^j)/2.$$

The value of heterosis, computed as the difference between hybrid and midparental values, is:

$$H_{ij} = \sum_l d_l(1 - \theta_l^i \theta_l^j)/2.$$

If we consider $n_k$ marker loci, the number of loci for which the two lines $i$ and $j$ differ (i.e., the heterozygosity of the hybrid) is:

$$\Delta_{ij} = \sum_k (1 - \theta_k^i \theta_k^j)/2.$$

This number will be considered as the distance between lines $i$ and $j$.

When $N$ lines are considered for given loci $l$ and $k$, to determine the population's characteristics we establish:

$$w_l = \left(\sum_{i=1}^{i=N} \theta_l^i\right)/N \quad \text{and} \quad W_{lk} = \left(\sum_{i=1}^{i=N} \theta_l^i \theta_k^i\right)/N.$$

The frequency of the allele $L$ in the population of inbred lines is $f_L = (1 + w_l)/2$. The linkage disequilibrium between alleles $L$ and $K$ is $D_{LK} = (W_{lk} - w_l w_k)/4$.

## Relationship between distance and heterosis

### Correlation calculation

When we consider all potential hybrids that can be obtained when crossing the lines of the previous "population" (considering $i \times i$ as a "hybrid"), the correlation between heterozygosity at marker loci (i.e., the number of loci for which the parent lines are different) and heterosis is the ratio of the covariance between those two parameters to the square root of the product of variances among heterosis values and among distances values.

### Variance among distance indexes ($\Delta$).

$$\text{Var}(\Delta_{ij}) = 1/4 \, \text{Var}\left(\sum_k (1 - \theta_k^i \theta_k^j)\right)$$

$$= 1/4 \sum_k \sum_{k'} \text{Cov}(\theta_k^i \theta_k^j; \theta_{k'}^i \theta_{k'}^j)$$

$$\text{Cov}(\theta_l^i \theta_l^j, \theta_k^i \theta_k^j) = E(\theta_l^i \theta_l^j \theta_k^i \theta_k^j) - E(\theta_l^i \theta_l^j) E(\theta_k^i \theta_k^j)$$

$$= E^2(\theta_l^i \theta_k^i) - E^2(\theta_l^i) E^2(\theta_k^i)$$

$$= W_{lk}^2 - w_l^2 w_k^2$$

and so:

$$\text{Var}(\Delta_{ij}) = 1/4 \sum_k (1 - w_k^4) + 1/2 \sum_k \sum_{k' < k} (W_{kk'}^2 - w_k^2 w_{k'}^2)$$

### Variance among heterosis values.

$$\text{Var}(H_{ij}) = 1/4 \, \text{Var}\left(\sum_l d_l(1 - \theta_l^i \theta_l^j)\right)$$

$$= 1/4 \sum_l d_l^2(1 - w_l^4) + 1/2 \sum_l \sum_{l' < l} d_l d_{l'}(W_{ll'}^2 - w_l^2 w_{l'}^2)$$

### Covariance between distance and heterosis.

$$\text{Cov}(H_{ij}, \Delta_{ij}) = 1/4 \, \text{Cov}\left(\sum_l d_l(1 - \theta_l^i \theta_l^j); \sum_k (1 - \theta_k^i \theta_k^j)\right)$$

$$= 1/4 \sum_l \sum_k d_l(W_{kl}^2 - w_l^2 w_k^2)$$

### Correlation coefficient.

$$\varrho(H_{ij}, \Delta_{ij}) = \frac{1/4 \sum_l \sum_k d_l(W_{kl}^2 - w_l^2 w_k^2)}{\sqrt{1/4 \sum_l d_l^2(1 - w_l^4) + 1/2 \sum_l \sum_{l' < l} d_l d_{l'}(W_{ll'}^2 - w_l^2 w_{l'}^2)} \sqrt{1/4 \sum_k (1 - w_k^4) + 1/2 \sum_k \sum_{k' < k} (W_{kk'}^2 - w_k^2 w_{k'}^2)}}$$

### Influence of linkage disequilibria on the relationship

If we consider the cases for which the dominance parameter ($d_l$) is constant ($d$) for all loci ($d_l = d$, $\forall l$), $\varrho$ is the correlation between mean heterozygosities computed for two sets of loci ($k$ and $l$) [see Mitton and Pierce (1980) and Chakraborty (1981) for a discussion of the case where one group of loci is a subset of the other group].

The quantity $(W_{kl}^2 - w_l^2 w_k^2)$ has been studied for different sets of allelic frequencies at loci $k$ and $l$ (see Appendix). Negative values of this quantity can be found unless ($w_l = 0$ or $w_k = 0$), so that one can imagine situations such as: $\text{Cov}(H_{ij}; \Delta_{ij}) < 0$; however, in these cases absolute values are low. When the frequency of an allele at one locus is low ($f = 0.1$) and of equal probability ($f = 0.5$) at the other locus, the maximum value of ($W_{kl}^2 - w_l^2 w_k^2$) is 0.04, so the correlation between heterozygosity at both loci is extremely low (0.052).

When $w_l = w_k = 0$ for all loci (i.e., alleles are equally distributed in the whole population), the covariance is reduced to: $\frac{1}{4} \sum_l \sum_k d_l W_{lk}^2$. Since $W_{lk}^2$ equals $16 D_{lk}^2$ in this

particular case, the covariance between heterosis and distance is a function of the squares of the linkage disequilibria between marker alleles and alleles involved in heterosis variation. Thus:

$$\varrho(H_{ij}, \Delta_{ij}) =$$
$$= \frac{16\sum_l \sum_k d_l D_{kl}^2}{\sqrt{K + 32\sum_k \sum_{k'<k} D_{kk'}^2} \sqrt{\sum_l d_l^2 + 32\sum_l \sum_{l'<l} d_l d_{l'} D_{ll'}^2}}.$$

If we consider that: (1) $(d_l = d, \forall l)$, (2) all $k$ loci are independent, and (3) all $l$ loci are independent, then:

$$\varrho(H_{ij}, \Delta_{ij}) = \frac{16\sum_l \sum_k D_{kl}^2}{\sqrt{n_k} \sqrt{n_l}}.$$

If we then consider a situation for which: (1) each locus $l$ is associated to a single marker locus $k_l (D_{kl} = 0$ if $k \neq k_l)$, and (2) $(\forall l, D_{k_l l} = D)$, then:

$$\varrho(H_{ij}, \Delta_{ij}) = 16 D^2.$$

**Table 1.** Evolution of correlation coefficient between heterosis and heterozygosity at marker loci

| $X(l, k)$ Generation | 1 | 5 | 10 | 20 | 50 |
|---|---|---|---|---|---|
| 2 | 0.961 | 0.819 | 0.670 | 0.449 | 0.135 |
| 3 | 0.942 | 0.743 | 0.554 | 0.313 | 0.063 |
| 4 | 0.923 | 0.674 | 0.458 | 0.219 | 0.030 |
| 5 | 0.905 | 0.611 | 0.379 | 0.152 | 0.014 |
| 7 | 0.870 | 0.503 | 0.259 | 0.074 | 0.003 |
| 10 | 0.819 | 0.375 | 0.147 | 0.025 | 0.000 |
| 15 | 0.742 | 0.230 | 0.057 | 0.004 | 0.000 |
| 20 | 0.672 | 0.142 | 0.022 | 0.001 | 0.000 |
| 30 | 0.550 | 0.053 | 0.003 | 0.000 | 0.000 |
| 40 | 0.451 | 0.020 | 0.000 | 0.000 | 0.000 |
| 50 | 0.370 | 0.008 | 0.000 | 0.000 | 0.000 |

$X(l, k)$: distance between QTL and marker locus (cMorgans)

Now let us consider $n_{l'}$ additional loci and $n_{k'}$ additional marker loci, assuming that these loci are not linked to previous $l$ loci and $k$ marker loci and that there is no linkage between these additional loci, in which case:

$$\varrho(H_{ij}, \Delta_{ij}) = \frac{16 n l D^2}{\sqrt{n_l + n_{l'}} \sqrt{n_k + n_{k'}}}$$
$$= \frac{16 D^2}{\sqrt{1 + \frac{n_{l'}}{n_l}} \sqrt{1 + \frac{n_{k'}}{n_l}}} < 16 D^2.$$

The correlation coefficient decreases as $n_{l'}$ and $n_{k'}$ increase; nonmarked loci and "nonmarking" marker loci play symmetrical roles.

*Application within a random-mating population*

The relationship between heterosis and heterozygosity at marker loci depends on the germ plasm considered. For example, let us consider a random mating population of infinite size. The loci that are involved in heterosis are assumed to segregate independently (to be on different chromosomes) and to all have the same impact on character variation $(d_l = d, \forall l)$. Each allele is also assumed to be linked to a marker allele with a recombination fraction of $r$.

If we consider that the population has been founded by the self-fertilization of one single individual $(w_l = 0$ for each locus), linkage disequilibrium for gametes from this individual (generation 1) would be represented as:

$$D_1 = (1 - 2r)/4,$$

and, at a given generation $n$, the linkage disequilibrium $Dn$ between a marker and the marked allele would be represented as:

$$D_n = \tfrac{1}{4}(1 - 2r)(1 - r)^{n-1}.$$

**Table 2.** Evolution of correlation between heterosis and heterozygosity at marker loci. The first marker, $k1$, is 5 cMorgans from the QTL. An additional marker, $k2$, is $X(l, k2)$ from QTL, to either side of the QTL

| $k2$ rel. to QTL $X(l, k2)$ Generation | Beyond $k1$ | | | | Opposite side from $k1$ | | | |
|---|---|---|---|---|---|---|---|---|
| | 5 | 10 | 20 | 50 | 5 | 10 | 20 | 50 |
| 2 | 0.819 | 0.781 | 0.720 | 0.625 | 0.896 | 0.846 | 0.767 | 0.640 |
| 3 | 0.743 | 0.695 | 0.628 | 0.548 | 0.842 | 0.771 | 0.671 | 0.556 |
| 4 | 0.674 | 0.619 | 0.550 | 0.488 | 0.789 | 0.698 | 0.588 | 0.492 |
| 5 | 0.611 | 0.552 | 0.485 | 0.438 | 0.736 | 0.629 | 0.515 | 0.440 |
| 7 | 0.503 | 0.440 | 0.383 | 0.357 | 0.634 | 0.505 | 0.400 | 0.357 |
| 10 | 0.375 | 0.315 | 0.275 | 0.266 | 0.496 | 0.359 | 0.282 | 0.266 |
| 15 | 0.230 | 0.183 | 0.165 | 0.163 | 0.317 | 0.202 | 0.166 | 0.163 |
| 20 | 0.142 | 0.108 | 0.100 | 0.100 | 0.198 | 0.115 | 0.101 | 0.100 |
| 30 | 0.053 | 0.039 | 0.038 | 0.075 | 0.040 | 0.038 | 0.038 | 0.038 |
| 40 | 0.020 | 0.014 | 0.014 | 0.014 | 0.028 | 0.015 | 0.014 | 0.014 |
| 50 | 0.008 | 0.005 | 0.005 | 0.005 | 0.011 | 0.005 | 0.005 | 0.005 |

$X(l, k2)$: distance between QTL and marker locus ($k2$), in cMorgans

The evolution of the correlation coefficient according to the generation number and the distance $x$ between the marker and the locus involved in the quantitative trait is given in Table 1. This assumes that $r = 1/2(1 - e^{-2x})$, according to a Poisson process along the chromosome, which corresponds to Haldane's mapping function.

The correlation decreases quickly with the evolution of generation number, unless the linkage is very tight ($x \le 0.01$). The introduction of additional markers, which are not tightly linked to alleles involved in the variation of the character, into the distance computation may lead to a dilution of the relationship between heterosis and heterozygosity at marker loci (Table 2). When a first marker is 5 cMorgans away from the marked allele, the introduction of a second marker at 20 cMorgans from the marked allele and 15 cMorgans from first marker leads to a drop in the correlation coefficient from 0.611 to 0.485, at the 5$^{th}$ generation.

## Discussion and conclusions

### Linkage disequilibria in a set of 50 public lines

When lines of several origins are considered as a population, linkage disequilibrium is the result of a very complex series of events. On the one hand, the generation number is very large; on the other hand, the germ plasm is generally structured into partially isolated groups. Such groups with little genetic exchange between them maintain the linkage disequilibrium between group-specific alleles. A group of 50 publicly available lines, which are (or have been) commonly used in Northern France, has been analyzed for 18 isoenzymatic loci. A significant association was found between alleles of the loci *Mdh5* and *Pgm2*, which in maize are known to be about 15 cMorgans apart; no association was found between alleles of the loci *Idh1* and *Mdh1* (1 cMorgan apart) or between alleles of the loci *Idh2* and *Mdh2* (also 1 cMorgan apart). Thus, despite their greater proximity, these latter were found to have no association.

In the first case, we should mention that alleles *Mdh5-15* and *Pgm2-1* are both specific to the early European flint group, which has been partly protected from contamination by other groups. In the second case, two explanations can be given. Firstly, the accumulation of recombination events may have led to this situation. Secondly, because of the poor resolution of electrophoresis technique, a given isozymic class may include several genotypes, which may have very different origins. Use of RFLP markers should eliminate this problem, since the probability that a mutation will lead to a previously existing type is extremely low.

### Selection of "efficient" markers

Prediction of heterosis on the basis of parental markers would require selection of which markers should be in-troduced into the heterozygosity computation. Development of RFLP techniques suggests that the number of potential markers will become almost infinite; however, strategies for the selection of appropriate markers remain to be defined. As an illustration, when considering the germ plasm that can be used under French conditions, we found that 7 lines (out of 48), which had the allele *Mdh5-15*, also had the same parental origin (the *French line* $F_2$), shared some physiological characteristics (earliness, relatively short height when crossed to testers), and led to low-yielding hybrids when intercrossed. However, recent attempts have been made to intercross this group with American later-growing types (M. Derieux (1990), personal communication), suggesting that what we had known about the allele *Mdh5-15* may no longer be sufficient when predicting the outcome of crossings with such new material. Studies of segregating populations ($F_2$ or recombinant inbred lines) involving European and other parents should help to answer this question and provide information about the genetic basis of the relationship between parental marker distance and hybrid heterosis.

## Appendix

*Study of the variation of* $Cov(\theta_l^i \theta_l^j; \theta_k^i \theta_k^j)$
*for given allelic frequencies at loci $l$ and $k$*

$$Cov(\theta_l^i \theta_l^j; \theta_k^i \theta_k^j) = W_{kl}^2 - w_l^2 w_k^2$$
$$= 16D^2 + 8 w_l w_k D.$$

The value of the covariance depends on the value of $D$ but also on the allelic frequencies at both loci. For a given set of allelic frequencies, the values that $D$ may have fall within the following interval:

$$D \in \left[ -\min\left(\frac{(1+w_l)(1+w_k)}{4}, \frac{(1-w_l)(1-w_k)}{4}\right); \right.$$
$$\left. \min\left(\frac{(1+w_l)(1-w_k)}{4}, \frac{(1-w_l)(1+w_k)}{4}\right)\right].$$

For fixed allelic frequencies, the covariance is a function of $D$, with a minimum of $-\frac{w_l w_k}{4}$.

Therefore, the maximum and minimum values of $Cov(\theta_l^i \theta_l^j; \theta_k^i \theta_k^j)$ will depend on the relative positions of $-\min\left(\frac{(1+w_l)(1+w_k)}{4}, \frac{(1-w_l)(1-w_k)}{4}\right)$, $\min\left(\frac{(1+w_l)(1-w_k)}{4}, \frac{(1-w_l)(1+w_k)}{4}\right)$, and $-\frac{1}{4}w_l w_k$. If we call those three values, respectively, $a$, $b$, and $c$:

| Situation | Cov$^{min}$ | Cov$^{max}$ |
| --- | --- | --- |
| $c \le a \le b$ | Cov($a$) | Cov($b$) |
| $a \le c \le b$ | Cov($c$) | max(Cov($a$), Cov($b$)) |
| $a \le b \le c$ | Cov($b$) | Cov($a$) |

And so:

| Minimum covariances | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 0.90 | −0.050 | −0.070 | −0.062 | −0.026 | 0.000 | −0.026 | −0.062 | −0.070 | −0.050 |
| 0.80 | −0.070 | −0.090 | −0.058 | −0.014 | 0.000 | −0.014 | −0.058 | −0.090 | −0.070 |
| 0.70 | −0.062 | −0.058 | −0.026 | −0.006 | 0.000 | −0.006 | −0.026 | −0.058 | −0.062 |
| 0.60 | −0.026 | −0.014 | −0.006 | −0.002 | 0.000 | −0.002 | −0.006 | −0.014 | −0.026 |
| 0.50 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 0.40 | −0.026 | −0.014 | −0.006 | −0.002 | 0.000 | −0.002 | −0.006 | −0.014 | −0.026 |
| 0.30 | −0.062 | −0.048 | −0.026 | −0.006 | 0.000 | −0.006 | −0.026 | −0.058 | −0.062 |
| 0.20 | −0.070 | −0.090 | −0.058 | −0.014 | 0.000 | −0.014 | −0.058 | −0.090 | −0.070 |
| 0.10 | −0.050 | −0.070 | −0.062 | −0.026 | 0.000 | −0.026 | −0.062 | −0.070 | −0.050 |
| $\dfrac{1+w_l}{2} \Big/ \dfrac{1+w_k}{2}$ | 0.100 | 0.200 | 0.300 | 0.400 | 0.500 | 0.600 | 0.700 | 0.800 | 0.900 |

| Maximum covariances | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 0.90 | 0.590 | 0.410 | 0.258 | 0.134 | 0.040 | 0.134 | 0.258 | 0.410 | 0.590 |
| 0.80 | 0.410 | 0.870 | 0.582 | 0.346 | 0.160 | 0.346 | 0.582 | 0.870 | 0.410 |
| 0.70 | 0.258 | 0.582 | 0.974 | 0.634 | 0.360 | 0.634 | 0.974 | 0.582 | 0.258 |
| 0.60 | 0.134 | 0.346 | 0.634 | 0.998 | 0.640 | 0.998 | 0.634 | 0.346 | 0.134 |
| 0.50 | 0.040 | 0.160 | 0.360 | 0.640 | 1.000 | 0.640 | 0.360 | 0.160 | 0.040 |
| 0.40 | 0.134 | 0.346 | 0.634 | 0.998 | 0.640 | 0.998 | 0.634 | 0.346 | 0.134 |
| 0.30 | 0.258 | 0.582 | 0.974 | 0.634 | 0.360 | 0.634 | 0.974 | 0.582 | 0.258 |
| 0.20 | 0.410 | 0.870 | 0.582 | 0.346 | 0.160 | 0.346 | 0.582 | 0.870 | 0.410 |
| 0.10 | 0.590 | 0.410 | 0.258 | 0.134 | 0.040 | 0.134 | 0.258 | 0.410 | 0.590 |
| $\dfrac{1+w_l}{2} \Big/ \dfrac{1+w_k}{2}$ | 0.100 | 0.200 | 0.300 | 0.400 | 0.500 | 0.600 | 0.700 | 0.800 | 0.900 |

## References

Brunel D (1985) Distance génétique et hétérosis. Utilisation des marqueurs moléculaires. In: Les distances génétiques. INRA Paris, pp 159–168

Chakraborty R (1981) The distribution of the number of heterozygous loci in an individual in natural populations. Genetics 98:461–466

Edwards MD, Stuber CW, Wendel JF (1987) Molecular-facilitated investigations of quantitative trait loci in maize. I. Numbers, genomic distribution, and types of gene action. Genetics 116:113–125

Fidgore S (1987) Application of restriction fragment length polymorphisms to quantitative genetics in maize. 2nd Int Conf Quant Genetics, Raleigh NC

Frei OM, Stuber CW, Goodman MM (1986) Use of allozymes as genetic markers for predicting performance in maize single-cross hybrids. Crop Sci 26:37–42

Hadjinov MI, Scherbak VS, Benko NI, Gusen VP, Sukhorzhevskaya TP, Vorona LP (1982) Interrelationships between isozymic diversity and combining ability in maize lines. Maydica 27:135–149

Heidrich-Sobrinho E, Cordeiro AR (1975) Codominant isoenzymic alleles as markers of genetic diversity correlated with heterosis in maize (Zea mays). Theor Appl Genet 46:197–199

Hunter RB, Kannenberg LW (1971) Isozyme characterization of corn (Zea mays) inbreds and its relationship to single-cross hybrid performance. Can J Genet Cytol 13:649–655

Lamkey KR, Hallauer AR, Kahler AL (1987) Allelic differences at enzyme loci and hybrid performance in maize. J Hered 78:231–234

Lee M, Godshalk EB, Lamkey KR, Woodman WW (1989) Association of restriction fragment length polymorphisms among maize inbred lines with agronomic performance of their crosses. Crop Sci 29:1067–1071

Mitton JB, Pierce BA (1980) The distribution of individual heterozygosity in natural populations. Genetics 95:1043–1054

Price SC, Kahler AL, Hallauer AR, Chermley P, Giegel DA (1986) Relationship between performance and multilocus heterozygosity at enzyme loci in single-cross hybrids of maize. J Hered 77:341–344

Smith JSC, Smith OS (1989) The use of morphological, biochemical, and genetic characteristics to measure distance and to test for minimum distance between inbred lines of maize (Zea mays L.). UPOV Workshop, Versailles, France, October 1989

Walton M, Helentjaris T (1987) Application of restriction fragment length polymorphism (RFLP) technology to maize breeding. 42nd Ann Corn Res Conf 48–75