T. J. L. van Hintum · D. Haalman

# Pedigree analysis for composing a core collection of modern cultivars, with examples from barley (*Hordeum vulgare* s. lat.)

**Abstract** A method for analyzing the pedigrees of cultivars is developed that allows for the calculation of the 'effective number of origin lines' $(n_{OL})$. The $n_{OL}$ is defined as the average number of alleles, not identical by descent, per locus in a set of lines. Its relationship with the commonly used 'coefficient of parentage' is clarified. A related quantity, the 'effective overlap of origin lines' $(r_{OL})$ is defined as the average number of alleles, not identical by descent, per locus common in two sets of individuals. A set of 85 modern barley cultivars is used to illustrate the application of $n_{OL}$ and $r_{OL}$. This set originated from 153 mutually unrelated ancestors. The degree to which each ancestor contributed was quantified, and the result was a $n_{OL}$ of only 43.1. In the set were 51 spring and 34 winter cultivars, with a $n_{OL}$ of 25.0 and 21.0, respectively. Consequently, the $r_{OL}$ of these two groups was 2.9, indicating that the two groups can be considered to be nearly distinct genetically since they have only 2.9 origin lines in common. How the effective number of origin lines can be used to create a core collection of cultivars with known pedigrees by maximizing the $n_{OL}$ in a set of cultivars of given size is also discussed.

**Key words** Pedigree analysis · Genetic diversity · Core collections · Genetic resources · *Hordeum vulgare*

## Introduction

To improve the accessibility of increasingly large germ plasm collections, core collections should be developed. A core collection represents the genetic diversity of a crop and its relatives with a minimum of repetitiveness (Frankel and Brown 1984), and although it will not replace an existing collection, it will make the latter more accessible: an accession in a core collection represents accessions in a reserve collec-

tion (Brown 1989). A core collection has two major advantages: (1) due to its limited size, it can be more comprehensively documented than an ordinary germ plasm collection, thus allowing a more effective choice of material for utilization; (2) due to its well-defined structure, a core collection allows for optimization of the genetic diversity within material selected for utilization from the core or another collection.

The availability of information on the material usually determines the methods that can be applied in constructing such core collections. The pedigrees of cultivars are sometimes available, especially in well-studied crops like wheat (Cox et al. 1986) and maize (Smith et al. 1990). This type of information is very useful in diversity studies, as has been shown in comparative methodological studies (Smith et al. 1990; Souza and Sorrells 1991). In the present article we apply and extend the theory of pedigree analysis for the construction of a core collection of cultivars. This will be illustrated with examples from barley.

## Core collection

When constructing a core collection several criteria have to be established:

- *Diversity to be covered.* A core collection can cover any gene pool, varying from modern Dutch barley cultivars to the entire *Hordeum* gene pool.
- *Material to choose from.* In the case of cultivars, it can be very effective to include parental lines in the core collection, but these do not always exist or are not always available. This can especially be a problem in the case of breeding lines and old cultivars.
- *Number of core accessions.* Though the methodology of choosing accessions is usually independent of the number to be chosen, this obviously is a very important parameter.
- *Preconditions related to the objectives.* The objectives of the core collections can imply such preconditions as, for example, an a-priori set of accessions that should be included or the prerequisite that representatives of all historic phases in barley breeding and of all utilization types should be included.

T. J. L. van Hintum (✉) · D. Haalman
Centre for Plant Breeding and Reproduction Research (CPRO-DLO), Centre for Genetic Resources, the Netherlands, P.O. Box 16, 6700 AA Wageningen, The Netherlands

Bearing these criteria in mind and using the methods and data available, we should be able to assemble accessions that will maximize the number of alleles in the core collection.
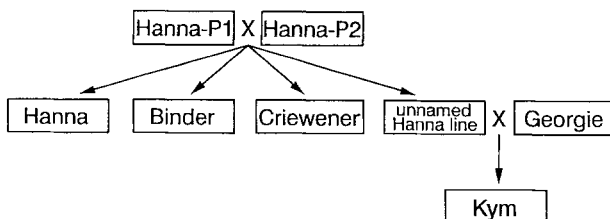
## Pedigree analysis

In pedigree analysis the ancestral relationship between two individuals is studied. The degree of coancestry is usually quantified with the coefficient of parentage ($r$) as defined by Kempthorne (1969). The $r$ between two individuals is the probability that a random allele at a random locus in one individual is identical by descent to a random allele at the same locus in the other individual (Cox et al. 1985). In order to apply this $r$ to the pedigree analysis of cultivars, the following assumptions have to be made (Martin et al. 1991): (1) a cultivar receives half of its genes from each parent; (2) parents in crosses are homozygous and homogeneous; (3) ancestors for which no pedigree information is available are unrelated; (4) the $r$ between a cultivar and a selection from that cultivar is 0.75.

All of these assumptions are to varying extents disputable, and some are even contradictory. To avoid the clear contradiction between the second and fourth assumption the latter should be replaced by: if selections are made from a cultivar, this cultivar is assumed to be the variable offspring of the cross of two unrelated lines. A selection from the cultivar is one of the offspring lines; if the cultivar itself is used as a parent in a cross, one of the offspring lines is considered to be used.

The barley landrace 'Hanna' may serve as an example (Fig. 1). 'Binder' and 'Criewener' are selections from 'Hanna' (Linde-Laursen et al. 1982; Baum et al. 1985). The modern cultivar 'Kym' is an offspring of a cross between 'Georgie' and 'Hanna' (NIAB 1989). If we assume that, 'Hanna' consists of the offspring of the cross between the unrelated lines 'Hanna-P1' and 'Hanna-P2', 'Binder' and 'Criewener' can be considered to have 'Hanna-P1 × Hanna-P2', and 'Kym' to have 'Georgie × (Hanna-P1 × Hanna-P2)', as their pedigrees. This avoids contradiction between the 'homogeneity-' and 'selection-assumption', and facilitates calculations. It also implies that the $r$ between a cultivar and a selection from that cultivar equals 0.50 as opposed to the 0.75 obtained in the regular system. The $r$ between two selections from the same cultivar also equals 0.50 as opposed to the 0.56 obtained in the regular system.

**Fig. 1** Example of the proposed representation of a landrace ('Hanna'), selections from it ('Binder' and 'Criewener') and a cross with it ('Hanna' × 'Georgie' → 'Kym')



Pedigree analysis has been used to describe the genetic basis of crops (Knauft and Gorbet 1989) and development in time (Cox et al. 1986; Souza and Sorrells 1989), to predict hybrid performance (Smith et al. 1990), and to compare different measures of genetic similarity (Cox et al. 1985). It can also be used to select material for inclusion in a core collection, but for that purpose it is necessary to extend the theory of pedigree analysis.

## Effective number of origin lines

Effective number of origin lines ($n_{OL}$) can be defined as the average number of alleles, not identical by descent, per locus in a set of lines. The $n_{OL}$ of a set of three unrelated lines will be 3.0; the $n_{OL}$ of an offspring line and one of the two unrelated parents will be 1.5, since all alleles of the included parent plus half of the alleles of the parent not included (via the offspring line) will be in the set. The $n_{OL}$ of a set of two unrelated parents and any number of offspring lines will be 2.0, since only alleles of the two parents can be found in the set.

Another quantity has to be defined: the effective overlap of origin lines ($r_{OL}$) of two sets of individuals is the average number of alleles, not identical by descent, per locus common to both sets. So the maximum value of $r_{OL}(A, B)$ will equal the smallest of the two effective numbers of origin lines $n_{OL}(A)$ and $n_{OL}(B)$.

There is a simple relationship between effective overlap of origin lines ($r_{OL}$) and the effective number of origin lines ($n_{OL}$):

$$r_{OL}(A, B) = n_{OL}(A) + n_{OL}(B) - n_{OL}(A \cup B) \tag{1}$$

where $r_{OL}(A, B)$ is the effective overlap of origin lines of sets A and B, $n_{OL}(A)$ is the effective number of origin lines of the lines in set A, $n_{OL}(B)$ is effective number of origin lines of the lines in set B, and $n_{OL}(A \cup B)$ is the effective number of origin lines of the lines in the combined sets A and B.

The coefficient of parentage ($r$) is a special case of the effective overlap of origin lines, i.e., the case that both sets have only one element. So if '$a$' and '$b$' are the only elements of sets 'A' and 'B', respectively, relation (1) becomes:

$$r(a, b) = 2 - n_{OL}(A \cup B) \tag{1a}$$

A big advantage of $n_{OL}$ is that it is exactly the quantity that should be optimized when selecting accessions for a core collection.

The effective number of origin lines of a set of lines can be seen as the sum of contributions of the mutually unrelated 'origin lines'. These contributions equal the effective overlap of origin lines of these lines individually with the target set:

$$n_{OL}(A) = \Sigma r_{OL}(A, l_0) \tag{2}$$

where $n_{OL}(A)$ is the effective number of origin lines of the lines in set A, and $\Sigma r_{OL}(A, l_0)$ is sum of the effective overlaps of origin lines of set A and all origin lines (ancestors for which no pedigree information is available).

The individual $r_{OL}$ of an origin line and the target set indicates the contribution of that origin line to the target set.
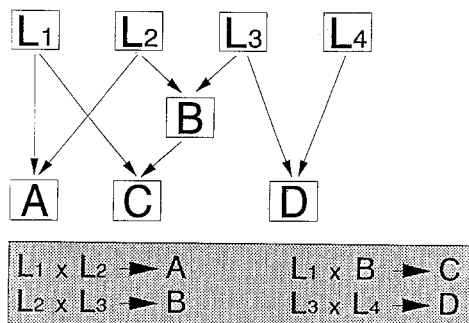
**Fig. 2** Example of a crossing scheme

## Example

As an example, to elucidate the concept, the crossing scheme presented in Fig. 2 will be used. All lines are homozygous, as in the case of barley cultivars. There are four origin lines: $L_1, L_2, L_3$ and $L_4$, and four offspring lines: $A$, $B$, $C$, and $D$. In Table 1 the effective overlaps of the four origin lines with several combinations of offspring lines are given, together with the effective number of origin lines, which equals the sum of the preceding values.

The $n_{OL}$ of a single line is always one, since for each locus there will be one allele. The $n_{OL}$ of a set of the two unrelated lines $A$ and $D$ equals the number of lines, which is two, since for each locus there will be two alleles not identical by descent.

It can be seen from $r_{OL}$ that 0.5 of the alleles of origin line $L_1$ are in line $A$ and 0.5 are in line $C$. Since these halfs are independent $r_{OL}(L_1, AC)$ equals 0.75. The combination $AC$ has 0.625 of the alleles of $L_2$: line $A$ has 0.5 of the $L_2$ alleles and 0.25 of the rest will be in line $C$.

The set BC has 0.5 of the alleles of origin line $L_1$ (in $C$) and 0.5 of line $L_3$ (in $B$). It also has 0.5 of the alleles of $L_2$ since all alleles of $L_2$ in $C$ will also be in $B$; therefore $C$ does not contribute any 'new' $L_2$ alleles.

**Table 1** The effective overlap $(r_{OL})$ and the effective number of origin lines $(n_{OL})$ of several combinations of lines as presented in Fig. 2

|  | $r_{OL}$ | | | | $n_{OL}$ |
| --- | --- | --- | --- | --- | --- |
|  | $L_1$ | $L_2$ | $L_3$ | $L_4$ | |
| A | 0.5 | 0.5 | 0.0 | 0.0 | 1.0 |
| B | 0.0 | 0.5 | 0.5 | 0.0 | 1.0 |
| C | 0.5 | 0.25 | 0.25 | 0.0 | 1.0 |
| D | 0.0 | 0.0 | 0.5 | 0.5 | 1.0 |
| AB | 0.5 | 0.75 | 0.5 | 0.0 | 1.75 |
| AC | 0.75 | 0.625 | 0.25 | 0.0 | 1.625 |
| AD | 0.5 | 0.5 | 0.5 | 0.5 | 2.0 |
| BC | 0.5 | 0.5 | 0.5 | 0.0 | 1.5 |
| BD | 0.0 | 0.5 | 0.75 | 0.5 | 1.75 |
| CD | 0.5 | 0.25 | 0.625 | 0.5 | 1.875 |
| ABC | 0.75 | 0.75 | 0.5 | 0.0 | 2.0 |
| ABD | 0.5 | 0.75 | 0.75 | 0.5 | 2.5 |
| ACD | 0.75 | 0.625 | 0.625 | 0.5 | 2.5 |
| BCD | 0.5 | 0.5 | 0.75 | 0.5 | 2.25 |
| ABCD | 0.75 | 0.75 | 0.75 | 0.5 | 2.75 |

The set of all offspring lines $A$, $B$, $C$, and $D$ has a $n_{OL}$ of 2.75. The $n_{OL}$ of the complete set including origin lines will be equal to the number of origin lines, which is four. From Table 1 it can be seen that 0.25 of the alleles of $L_1$, $L_2$, and $L_3$ are not in the set ABCD and that 0.5 of the alleles of $L_4$ are missing.

## Calculation methods

If the relevant pedigree information is available it is always possible to calculate the effective number of origin lines of a set of lines. Basic to the calculation is the assumption that the chance of an off-spring line getting an allele at a certain locus from its mother is equal to the chance of getting it from its father and that these events are totally interdependent: i.e., if a line gets the allele for a certain locus from one parent it will not get it from the other. This interdependence implies that when calculating the effective number of origin lines all inter-linked pedigrees have to be studied simultaneously. It can be shown that this will result in an exponential algorithm which, in turn, implies that it is not possible to calculate exactly the effective number of origin lines for groups of considerable size.

It is possible to calculate an approximation by using a linear algorithm. This algorithm falsely assumes that the chances of an allele coming from one parent or the other are independent. By starting at the target set and working in the direction of the origin lines it is possible to calculate for each line the chance of its alleles ending up in the target set by combining the chances for its offspring lines. The effective number of origin lines can be calculated using relationship 2, i.e., by summing up the chances of origin lines reaching the target set. This algorithm underestimates the actual value (if both parents have the same allele the algorithm assumes the chance of an offspring line getting this allele at 0.75 instead of 1.0). The error will be largest in sets with much inbreeding. In the sets used in the calculations presented in the next section the error appeared to be around 1%. Since the algorithm is very fast it is very suitable to be used in optimization procedures.

A third alternative uses the Monte Carlo simulation of the flow of alleles from the origin lines to the target set. In this way the effective overlap of origin lines of each origin line and the target set can be calculated by determining the chance of an allele of the origin line reaching the target set. If the effective number of origin lines of a set has to be known with a given reliability in a situation where the exponential algorithm would take too much time, this Monte Carlo simulation has proven to be a practical and relatively quick way of calculating it.

A computer program written in BASIC, showing the 'linear' and the 'Monte Carlo' algorithms, is available from the authors.

## Modern European barley cultivars

An analysis was made of barley cultivars grown in 1990 over an area exceeding 20000 ha in countries participating in the European Brewery Convention (EBC 1991). The pedigrees of

**Table 2** Cultivars included in the analysis

| Cultivar | Year[a] | RT[b] | Acreage[c] | Countries[d] | Cultivar | Year[a] | RT[b] | Acreage[c] | Countries[d] |
|---|---|---|---|---|---|---|---|---|---|
| Alexis | 1986 | 2S | 420 013 | D, DK, GB, I | Kalle | | 6S | 38 816 | SF |
| Alis | | 2S | 61 520 | DK | Kira | | 2W | 56 672 | GB |
| Alpaca | 1987 | 6W | 54 946 | D, NL | Koru | 1979 | 2S | 50 438 | E |
| Alpha | | 2W | 279 655 | E, F | Kustaa | | 6S | 77 632 | SF |
| Andrea | 1984 | 6W | 108 800 | D, DK | Kym | 1980 | 2S | 277 409 | E |
| Apex | 1983 | 2S | 128 566 | A, D, NL | Kymppi | | 2S | 77 632 | SF |
| Aramir | 1972 | 2S | 34 237 | A, F, I | Lina | 1985 | 2S | 68 884 | S |
| Arra | | 6S | 77 632 | SF | Magie | 1986 | 2W | 220 131 | D, F, GB, IRL |
| Atem | 1980 | 2S | 69 226 | A, GB, F | Mammut | 1978 | 6W | 66 978 | CH, D |
| Aura | 1975 | 2S | 52 860 | D, I | Marinka | 1985 | 2W | 367 268 | D, DK, GB, NL |
| Barbarossa | | 6W | 247 792 | E, F | Mars | 1967 | 6S | 25 500 | H |
| Baronesse | | 2S | 24 320 | D | Menuet | 1955 | 2S | 42 582 | F |
| Beka | 1954 | 2S | 580 037 | E | Mette | 1984 | 2S | 54 123 | S |
| Blenheim | 1985 | 2S | 318 090 | DK, GB, IRL, NL | Mogador | | 2W | 103 132 | E, F |
| Bruenhild | 1986 | 6W | 4 364 | D | Moulon | 1966 | 6W | 55 122 | E |
| Camargue | | 2S | 58 702 | GB, IRL | Natasha | | 2S | 111 890 | DK, F, GB |
| Cameo | | 2S | 23 227 | F | Pallas | 1958 | 2S | 75 657 | E |
| Canor | 1985 | 2S | 23 070 | DK | Panda | 1983 | 2W | 43 667 | F, GB |
| Carina | 1971 | 2S | 48 269 | A, I, P | Pastoral | | 2W | 201 212 | DK, F, GB, IRL |
| Carmen | 1970 | 2S | 37 259 | A | Pernilla | | 2S | 54 123 | S |
| Catinka | 1983 | 6W | 76 370 | D | Pipkin | 1959 | 2W | 80 960 | GB |
| Cheri | 1987 | 2S | 42 560 | D | Plaisant | | 6W | 649 517 | E, F, GB, I |
| Corona | 1980 | 6W | 43 739 | D, NL | Pohto | 1987 | 6S | 58 224 | SF |
| Danilo | 1984 | 2W | 21 820 | D | Pokko | 1980 | 6S | 24 260 | SF |
| Defra | 1987 | 2S | 24 320 | D | Prisma | 1985 | 2S | 69 191 | F, GB, I, NL |
| Digger | 1986 | 2S | 161 151 | DK, GB, IRL | Puffin | | 2W | 93 876 | GB, IRL |
| Escort | 1986 | 2S | 32 910 | DK, IRL | Rachel | 1979 | 6W | 36 378 | A, H |
| Express | | 6W | 219 802 | F | Regatta | 1987 | 2S | 80 450 | DK, GB |
| Flamenco | | 2W | 74 823 | DK, F, NL | Reinette | | 2W | 126 095 | E |
| Formula | 1987 | 2S | 30 881 | DK, I, S | Ribeka | | 2S | 23 320 | P |
| Franka | 1980 | 6W | 21 820 | D | Roland | 1981 | 2S | 21 643 | A, S |
| Frolic | | 2W | 33 088 | GB | Sewa | 1983 | 2S | 69 210 | DK |
| Gaulois | | 6W | 41 213 | F | Sherpa | | 2S | 22 317 | GB |
| Gimpel | 1979 | 2S | 32 490 | D, I, P | Sonja | 1974 | 2W | 94 360 | A, D, F |
| Golf | 1982 | 2S | 238 596 | CH, D, DK, F, GB, NL, S | Steffi | | 2S | 42 582 | D |
| | | | | | Steptoe | 1973 | 6W | 110 244 | E |
| Grit | 1983 | 2S | 75 000 | DK, IRL | T. Union | | 2S | 100 876 | E |
| Halcyon | 1968 | 2W | 82 368 | GB | Tapir | 1980 | 6W | 98 190 | D |
| Hart | | 2S | 24 583 | GB | Torrent | | 2W | 41 184 | GB |
| Hassan | 1971 | 2S | 100 876 | E | Triumph | | 2S | 131 582 | F, GB, IRL, NL |
| Ida | 1980 | 2S | 39 158 | S, SF | Trixi | 1987 | 2W | 161 493 | A, D, DK |
| Igri | 1976 | 2W | 220 564 | A, E, F, GB | Tyne | | 2S | 76 841 | GB |
| | | | | | Union | 1955 | 2S | 50 438 | E |

[a] Year of release
[b] Row number followed by annuality (S, spring; W, winter)
[c] Acreage (harvest 1990) in mentioned countries
[d] Countries where the cultivar is grown (A, Austria; B, Belgium; D, Germany; GB, Great Britain; CH, Switzerland; DK, Denmark, E, Spain; F, France; H, Hungary; I, Italy; IRL, Ireland; NL, Netherlands; P, Portugal; S, Sweden, SF, Finland)

85 of these 97 cultivars were traced using several sources, including Arias et al. (1983), Baum et al. (1985), Baumer and Goppel (1988), EBC (1991), Linde-Laursen (1982), NIAB (1989) and personal communication with breeders. These 85 cultivars (Table 2) were analyzed.

The total effective number of origin lines calculated via Monte Carlo simulation of the 85 cultivars was 43.1 (see Fig. 3). The set of 85 cultivars could be shown to originate from 153 mutually unrelated 'origin lines'. These lines contributed in varying extents to the modern gene pool; the line 'Lyallpur' contributed via the cvs 'Escort', 'Frolic', and 'Regatta' only 0.6 percent of its alleles, while the Swedish selection 'Gull' contributed via 58 cultivars 97% of its alleles. The contribution of 'Gull' varied from an average of 2.3% for 'Cheri' to 53.5% for 'Pernilla'.

There are 51 spring and 34 winter cultivars, with effective numbers of origin lines ($n_{OL}$) of 25.0 and 21.0, respectively. This means that the effective overlap of origin lines ($r_{OL}$) between these groups is only 2.9. Therefore, the two groups can be considered to be nearly distinct genetically.

The set of 85 modern cultivars was also used to look at possibilities of optimization of the $n_{OL}$ in samples for a core collection. In this case, $n_{OL}$ was calculated with the linear algorithm. Via an optimization procedure a sample of given size was selected from cultivars that contained the highest possible $n_{OL}$. In addition to the optimization procedure, a
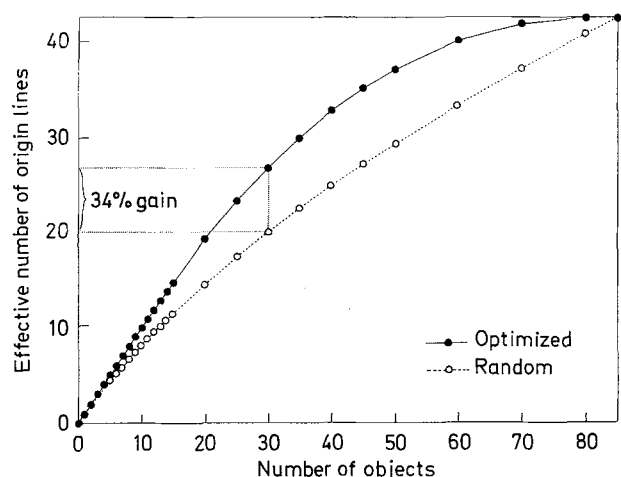
**Fig. 3** Effective number of origin lines in an optimized sample and in a random sample of 85 modern European barley cultivars

Monte Carlo simulation of 10 000 runs was performed to determine the average $n_{OL}$ of a random sample of given size. The results are presented in Fig. 3. It can be seen that up to 34% gain (with 30 objects out of 85) can be made via optimization. This percentage may increase if, for example, material from different historic phases in barley breeding is being selected, since such a selection is more likely to contain ancestor-offspring combinations than a selecton of modern cultivars.

A core collection of barley cultivars grown in 1990 throughout member countries of the European Barley Convention, consisting of 10 accessions, should include the following cultivars: 'Baronesse', 'Beka', 'Danilo', 'Defra', 'Kalle', 'Mammut', 'Moulon', 'Steptoe', 'Tapir', and 'Trixi'. This core collection has an effective number of origin lines of 9.99. If a prerequisite was set, for example, that both the spring and winter cultivar with largest acreage, 'Beka' and 'Plaisant', respectively, should be included, the others would be 'Baronesse', 'Danilo', 'Gaulois', 'H. Grignon', 'Kalle', 'Steptoe', 'T. Union', and 'Trixi'. This set would have an effective number of origin lines of 9.97. A random set of 10 cultivars would have an average effective number of origin lines of 8.22; if it would have to include 'Beka' and 'Plaisant', the average effective number of origin lines would be 7.53.

## Conclusion

Besides being a useful tool for numerous genetic analytical purposes, the effective number of origin lines can be used to select a core collection if pedigree information is available.

The exact calculation of the effective number of origin lines is time consuming if the target set includes more than a few

lines. There is a simple alternative algorithm slightly underestimating the exact value. This algorithm can easily be applied to optimization procedures for the selection of a core collection. Monte Carlo simulation is a third alternative.

## References

Arias G, Reiner L, Penger A, Mangstl A (1983) Directory of barley cultivars and lines. Technical University of Munich, Freising-Weihenstephan

Baum BR, Bailey LG, Thompson BK (1985) Barley register. Publ no. 1783/B, Agriculture Canada, Ottawa

Baumer M, Goppel W (1988) Gerste. Sorten, Züchter, Ursprungsland, Zulassungsjahr, Abstammung. Bayerische Landesanstalt fur Bodenkultur und Pflanzenbau, Freising-Weihenstephan, March 1988

Brown AHD (1989) Core collections: a practical approach to genetic resources management. Genome 31:818–824

Cox TS, Kiang YT, Gorman MB, Rodgers DM (1985) Relationship between coefficient of parentage and genetic similarity indices in the soybean. Crop Sci 25:529–532

Cox TS, Murphy JP, Rodgers DM (1986) Changes in the genetic diversity in the red winter wheat region of the United States. Proc Natl Acad Sci USA 83:5583–5586

EBC (1991) Advances in malting barley (harvest 1990) SECR European Brewery Convention, Zoeterwoude, The Netherlands

Frankel OH, Brown AHD (1984) Plant genetic resources today: a critical appraisal. In: Holden JHW, Williams JT (eds) Crop genetic resources: conservation and evaluation. George Allen and Unwin, London, pp 249–257

Kempthorne O (1969) An introduction to genetic statistics. Iowa State University Press, Ames

Knauft DA, Gorbet DW (1989) Genetic diversity among peanut cultivars. Crop Sci 29:1417–1422

Linde-Laursen I, Doll H, Nielsen G (1982) Giemsa C-Banding patterns and some biochemical markers in a pedigree of European barley. Z Pflanzenzuecht 88:191–219

Martin JM, Blake TK, Hockett EA (1991) Diversity among North American spring barley cultivars based on coefficients of parentage. Crop Sci 31:1131–1137

NIAB (1989) Detailed descriptions of varieties of wheat, barley, oats, rye and triticale. National Institute of Agricultural Botany, Cambridge

Smith OS, Smith JSC, Bowen SL, Tenborg RA, Wall SJ (1990) Similarities among a group of elite maize inbreds as measured by pedigree, $F_1$ grain yield, grain yield, heterosis, and RFLPs. Theor Appl Genet 80:833–840

Souza E, Sorrells ME (1989) Pedigree analysis of North American oat cultivars released from 1951 to 1985. Crop Sci 29:595–601

Souza E, Sorrells ME (1991) Relationships among 70 North American oat germplasms: 2. cluster analysis using qualitative characters. Crop Sci 31:605–612