

Genomic approaches to selection in outcrossing perennials: focus on essential oil crops

David Kainer¹ · Robert Lanfear² · William J. Foley¹ · Carsten Külheim¹

Received: 15 November 2014 / Accepted: 23 July 2015
© Springer-Verlag Berlin Heidelberg 2015

Abstract The yield of essential oil in commercially harvested perennial species (e.g. ‘Oil Mallee’ eucalypts, Tea Trees and Hop) is dependent on complex quantitative traits such as foliar oil concentration, biomass and adaptability. These often show large natural variation and some are highly heritable, which has enabled significant gains in oil yield via traditional phenotypic recurrent selection. Analysis of transcript abundance and allelic diversity has revealed that essential oil yield is likely to be controlled by large numbers of quantitative trait loci that range from a few of medium/large effect to many of small effect. Molecular breeding techniques that exploit this information could increase gains per unit time and address complications of traditional breeding such as genetic correlations between key traits and the lower heritability of biomass. Genomic selection (GS) is a technique that uses the information from markers genotyped across the whole genome in order to predict the phenotype of progeny well before they reach maturity, allowing selection at an earlier age. In this review, we investigate the feasibility of genomic selection (GS) for the improvement of essential oil yield. We explore the challenges facing breeders selecting for oil yield, and how GS might deal with them. We then assess the factors that affect the accuracy of genomic estimated breeding values, such as linkage disequilibrium (LD), heritability, relatedness and

the genetic architecture of desirable traits. We conclude that GS has the potential to significantly improve the efficiency of selection for essential oil yield.

Introduction

Essential oils are a diverse group of around 3000 natural plant products, of which about 300 are traded commercially for purposes such as flavourings, cosmetics, pharmaceuticals, aromatherapy and solvents. They are typically composed of a mix of volatiles (mostly terpenoids) and aromatics, often dominated by one or two major compounds. A wide variety of oil-bearing plant species, ranging from herbs and grasses to trees, are cultivated in plantations or harvested from wild stands in order to obtain essential oils for trade. Although a few essential oils have been extracted since the Middle Ages (Bakkali et al. 2008), until recently many cultivated species had undergone little selection and improvement for oil yield, especially when compared to major agricultural crops such as maize, wheat and fruits. Commercially important essential oil-bearing species include Orange, Cornmint, Lemon, Eucalyptus, Tea Tree, Peppermint, Citronella and Hop. Pharmaceutical-grade *Eucalyptus* oil, the 4th largest essential oil by annual tonnage (CBI Ministry of Foreign Affairs 2012), has only been distilled commercially since the 1850s (Pearson 1993). The market for *Eucalyptus* oil became globally competitive during the 20th century due to the large-scale extraction of leaf oil as a by-product of wood and pulp production in China, South Africa and Brazil. Given the highly competitive nature of the essential oils market, in which uses for oils shift regularly and demand and supply can fluctuate rapidly, improvements in oil yield can be of great benefit to producers.

Communicated by R. K. Varshney.

✉ David Kainer
david.kainer@anu.edu.au

¹ Research School of Biology, The Australian National University, Canberra, ACT 2601, Australia

² Department of Biological Sciences, Macquarie University, Sydney, NSW 2109, Australia

For many selection and improvement programs, the primary goal is to increase the yield per unit area of the harvested product in a cost-effective manner. The expense of the breeding technique and labour must be at least offset by the longer-term gain in revenues (Luby and Shaw 2001; Heffner et al. 2010). This economic equation has been balanced in recent decades with techniques such as mass selection and recurrent phenotypic selection on relatively small breeding populations. Studies and trials using molecular techniques such as Marker-Assisted Selection (MAS) for essential oil traits are few. Byrne (2007) surveyed over 150 perennial fruit and ornamental breeding programs from around the world to examine if and how they were making use of molecular markers. Only 14 % of the trials were using MAS for research, and only 3 % were actively using markers to aid selection. The small scale of many breeding programs, lack of available markers and poor cost per unit gain relative to phenotypic selection were cited as the primary impediments to the use of molecular markers. In particular, Byrne (2007) noted that many of the crops included in his study had been recently domesticated and consequently had high genetic variability—a situation in which phenotypic selection can be the quickest and least expensive route to develop new cultivars. Although many essential oil crops are also likely to have high genetic variability, the cost effectiveness equation has steadily shifted further in favour of genetic markers since 2007 (Bernardo 2008), through rapid improvements in genetic technologies.

In addition to the perception that MAS is more expensive than other methods, the use of MAS in selecting for complex quantitative traits in plants has some well-documented problems (Holland 2004; Hospital 2009). In short, MAS combines phenotypic and pedigree information with a priori knowledge of markers for specific genes, or quantitative trait loci (QTLs), associated with the trait of interest. Individuals with the most favourable breeding values are selected using phenotypic data supported by genotype data for those key markers. It is the goal of the breeder, through crossing, to produce a new generation in which at least some individuals will contain the majority, or even all, of the favourable QTL alleles. The more QTLs that are included in the MAS process, the more progeny are required to ensure that at least some of those progeny will contain the majority of the favourable alleles. In order to keep the scale realistic and avoid the inclusion of false-positive associations, only QTLs that are deemed to be highly significant (e.g. $P < 0.0001$) are used and the rest are culled. This has been shown to upwardly bias estimates of the effects of the chosen QTLs (Beavis 1994) and to cause breeders to miss out on the cumulative effects of many minor QTLs. In practice, most markers identified in candidate gene association studies in forest trees explain less than 5 % of the total variation of the trait, so for complex

traits that are influenced by many QTLs of small effect MAS is often not particularly useful or cost-effective (Luby and Shaw 2001; Hospital 2009; Thavamanikumar et al. 2013).

Recent advances in the theory of Genomic Selection (GS) have generated renewed interest in using molecular markers for plant breeding. Genomic Selection involves the selection of favourable individuals based solely on the predictive value of genetic markers (Meuwissen et al. 2001). The process involves two main stages. First, a training population (TP) is phenotyped and genotyped across the whole genome to develop a model of breeding value. Cross-validation techniques are often applied, where a subset of the training population is excluded from the process of estimating parameters so their phenotypic values can be used to verify the model's predictive accuracy. Second, a separate breeding population (BP) is genotyped and the model derived from stage 1 is applied to estimate each individual's Genomic Estimated Breeding Value (GEBV) which is used for selection. The model of breeding value in the first stage is developed by simultaneously estimating the additive effect on the phenotype of every chromosomal segment of the genome that is bounded by the genotyped markers. GS enables selection to be applied before the mature phenotype is measurable, and the unit of selection is the allele rather than the line (Lorenz et al. 2011). By avoiding the need to wait for plants to mature before selection, GS can considerably shorten the selection cycle, decrease labour costs and increase the gain per unit time (Wong and Bernardo 2008; Heffner et al. 2010). Also, by estimating effects for all available markers, GS can capture the effects of many small-effect QTLs, thus avoiding the problems of missing trait variance and biased QTL effects inherent in MAS. This aspect of GS is particularly powerful for the breeder—in a scientific context, the majority of marker effects would be rejected as statistically insignificant, but GS for breeding purposes presents no such restrictions.

When it was first proposed by Meuwissen et al. (2001), the feasibility of GS was questionable since the concept hinges on the ability to genotype many markers across the whole genome to ensure that all QTLs are in association with at least one proximate marker. The advent of high-throughput SNP genotyping technologies, e.g. SNP chips, Genotyping-by-Sequencing (GBS) and whole-genome re-sequencing, has since lowered the barrier to high density, low cost genotyping. As a consequence, a variety of simulated and empirical GS studies have been performed in plants since 2007, with accuracies and genetic gains usually exceeding both phenotypic selection and MAS. The majority of plant-based GS studies have taken place in highly inbred crops with large-scale breeding programs; maize (Zhao et al. 2012; Massman et al. 2013), wheat (Heffner et al. 2010), barley (Lorenzana and Bernardo 2009; Crossa

et al. 2010), cassava (Oliveira et al. 2012), apples (Kumar et al. 2012), sugarcane (Gouy et al. 2013) and sugar beet (Würschum et al. 2013). Commercially important forest tree species such as *Eucalyptus grandis* (Resende et al. 2012a; Denis and Bouvet 2013), *Picea glauca* (Beaulieu et al. 2014) and *Pinus taeda* (Resende et al. 2012b) have also received attention to improve wood and growth traits. Genomic Selection in plants has been the subject of several reviews in the past few years in both forest tree breeding (Isik 2014) and more generally in plant breeding (Jannink et al. 2010; Lorenz et al. 2011; Nakaya and Isoe 2012).

Here we review the feasibility of genomic selection for the improvement of essential oil yield. We explore the challenges facing breeders when selecting for oil yield with traditional means and how GS might deal with them. We then assess the factors that affect the accuracy of genomic estimated breeding values (GEBVs) such as Linkage Disequilibrium (LD), heritability, relatedness between the training and breeding populations and the genetic architecture of desirable traits in order to determine if GS is a viable technique for increasing oil yield in certain essential oil species, with a focus on out-crossing perennials such as *Eucalyptus*, Tea Tree (*Melaleuca* sp.) and Hop (*Humulus lupulus* L.).

Selecting for essential oil yield

Essential oil yield is complex and comprises multiple quantitative traits (Doran et al. 2002) that should be accounted for during a selective breeding process. These traits include: (1) oil concentration per leaf; (2) biomass (leaf mass for some species; flowers, bark, wood or seeds for others); (3) broad adaptability to variable environments; and (4) resistance to pests and diseases. The first two traits form the basis of oil yield ‘per plant’, which combined with the other two traits forms the basis for overall yield per unit area of plantation. Additionally, the composition, or quality, of the oil is often critical to the selection process in order to maintain levels of certain compounds at industry requirements. For *Eucalyptus* oil, at least 70 % (v/v) of the monoterpene 1,8-cineole is required for the oil to be classed as pharmaceutical grade (BP) along with a

negligible amount of undesirables such as α -phellandrene (Coppen 2002). Tea tree oil quality is more complex as there are multiple known chemotypes, each with their own compound profile (Butcher et al. 1996; Keszei et al. 2010) but commercially valuable oil must contain >40 % (v/v) of terpinen-4-ol and <4 % (v/v) of 1,8-cineole. In hop, the essential oil accumulated in flower cones is used to impart flavour and aroma in beer, so hop cultivars are developed with varied oil concentration and profile in order to meet the requirements of the brewing industry. Finally, for those species that are continually harvested through coppicing (e.g. various *Eucalyptus* “oil mallees” and Tea Tree plants), the ability to regenerate rapidly after being harvested, and to produce consistent oil yield at the time of the next harvest is also critically important.

Despite its complexity, certain factors combine to present a strong case for the potential for improving oil yield. Firstly, the lack of long-term selection or domestication in many oil-bearing species means that populations show great phenotypic variation in oil traits and contain a vast array of allelic diversity (Thumma 2005; Külheim et al. 2009; Goodger and Woodrow 2012; Webb et al. 2013). For example, the oil concentration in *Eucalyptus polybractea* (Blue Mallee) can range from 0.7 to 13 % of leaf dry weight (King et al. 2006), while in *Melaleuca alternifolia* (Medicinal Tea Tree) it ranges from 2.5 to 14.5 % of dry weight (Homer et al. 2000). Secondly, much of the observed variation in foliar oil concentration and composition has been shown to be moderately to highly heritable in a variety of species: *Eucalyptus* (Doran and Matheson 1994; Grant 1997; King et al. 2004; Goodger and Woodrow 2012), Tea Tree (Butcher et al. 1996; Doran et al. 2002), Fennel (Izadi-Darbandi et al. 2013) and Peppermint (Kumar et al. 2014) (Table 1). High heritability leads to increased accuracy of selection since much of the observed variation is due to genetic rather than environmental effects. Under these conditions, recurrent phenotypic selection has the power to generate large gains per selection cycle. Indeed this has been the case for various essential oil crops over the past decades. For example, five cycles of recurrent selection in *Cymbopogon flexuosus* (Lemongrass) increased mean oil concentration from 0.7 to 1.7 % (Kulkarni et al. 2003),

Table 1 The narrow sense heritability (h^2) of essential oil concentration (oil conc) and of biomass in a range of commercial crops

| Species | Common name | h^2 (oil conc) | h^2 (biomass) | References |
|-------------------------|---------------|------------------|-----------------|---------------------------------------|
| <i>M. alternifolia</i> | Tea Tree | 0.67 | 0.25 | Butcher et al. (1996) |
| <i>M. alternifolia</i> | Tea Tree | 0.51–0.93 | 0.25 | Doran et al. (2002) |
| <i>M. piperita</i> | Peppermint | 0.54 | – | Kumar et al. (2014) |
| <i>E. polybractea</i> | Blue Mallee | 0.36 | 0.05 | Grant (1997) |
| <i>E. camaldulensis</i> | River Red Gum | 0.54 | – | Doran and Matheson (1994) |
| <i>E. kochii</i> | Oil Mallee | 0.83 | – | Barton et al. (1991) |
| <i>H. lupulus</i> | Hop | 0.37 | 0.03 | Murakami (1999), McAdam et al. (2014) |

while in *Carum carvi* (Annual Caraway) mean oil concentration increased from 3.4 to 7.4 % over 20 years of recurrent selection (Pank 2010). The Australian Tea Tree breeding program has doubled commercial Tea Tree oil yield from 150 to 300 kg ha⁻¹ since 1993 through selection based on a weighted multi-trait index (Baker et al. 2014). Estimated gains from one cycle of selection for oil concentration in *Eucalyptus* species *E. camaldulensis* (Doran and Matheson 1994) and *E. polybractea* (Grant 1997) are around 30 %, though Goodger and Woodrow (2008) noted that in practice, trial plantations of *E. polybractea* often failed to achieve such gains due partly to large variation in open-pollinated half-sibling progeny.

Limitations of phenotypic selection for oil yield

Although phenotypic selection often performs well for quantitative trait improvement, it has its limitations. Notably long cycle times in perennial crops, large and costly progeny trials and difficulty selecting for multiple traits simultaneously can limit the gain per unit time and cost.

Long cycle times

The usual cycle time for selection in *E. polybractea* is 3–5 years, in Tea Tree it is 3 years, while in *E. camaldulensis* the time to first flowering averages around 14 years making the selection gain per unit time far smaller than is achievable in many annuals. For example, the significant oil yield gains made by the Australian Tea Tree breeding program (see above), operating since 1993, must be considered in the light of the commercial release of only three improved cultivars to date (Baker et al. 2014). The long time to maturity also adds large costs to breeding programs for such species since a great number of trees must be nurtured, consuming resources and labour, only to later be culled at the point of selection.

Genetic correlations

To get the most benefit out of an essential oil breeding program, it is desirable to select for oil concentration, biomass, oil composition, coppice ability and plant adaptability simultaneously. Genetic correlations, r_g , can affect the accuracy and size of the gains that can be made for multiple traits with artificial selection. A negative correlation between two traits means that selection for one is likely to result in deterioration in the other. Estimates of genetic correlations are often imprecise due to large sampling errors, and they are strongly influenced by allele frequencies and so may differ between populations (Falconer et al. 1996). Nevertheless, various examples provide guidance on how selection

gains in oil yield can be affected. In predictive studies of Tea Tree, Butcher et al. (1996) estimated $r_g = -0.42$ for oil concentration and dry biomass, though recent results from two related seedling orchards (Baker et al. 2014) show wide variation in the genetic correlation between oil concentration and leafiness ($r_g = 0.624$ at one site and $r_g = -0.246$ at the other). Recurrent selection for oil concentration in this population might eventually lead to a reduction in total oil yield due to loss in biomass. Doran and Matheson (1994) also found a negative correlation for oil concentration and growth traits such as height ($r_g = -0.481$) in *E. camaldulensis*, though with a large standard error. In an *E. polybractea* progeny test, Grant (1997) found a small negative correlation between oil concentration per leaf and leaf biomass of $r_g = -0.174$. In hop, overall cone yield and essential oil concentration are highly important traits to breeders, but selection for cone yield may negatively affect total oil content due to significant negative correlation (Henning et al. 1997) and therefore make the development of certain high yield cultivars difficult.

Negative correlation between oil concentration and biomass could occur if increased biosynthesis and accumulation of terpenes has a high cost to the plant, leading to fewer resources being allocated to growth. On the other hand, increased biosynthesis and/or accumulation of terpenes may improve the plant's defences against herbivores (Farmer 2014), or be an indicator of natural selection for factors other than growth. For example, King et al. (2006) found that the accumulation of foliar oil was actually associated with better growth in *E. polybractea*, but no evidence was found to suggest a mechanism of herbivory defence. It should be noted that in this latter study the correlation was measured in seedlings. It is possible that any positive correlation between oil concentration and growth disappears by maturity—the point at which phenotypic selection for oil content is most accurate.

Complex traits

Traits such as oil concentration and biomass are often controlled by large numbers of genetic loci of small effect. Different individuals can exhibit similar phenotypes despite possessing very different sets of alleles at those loci. Producing and detecting crossed progeny that possess favourable alleles across all loci is extremely difficult and many controlled crosses are needed, resulting in greater population sizes and lower gain per unit cost.

Phenotyping

The process of phenotyping presents its own unique set of challenges that scale with the size of the breeding population. Assessment of oil concentration and composition per individual plant using methods such as steam distillation

or solvent extraction followed by gas chromatography is costly and time-consuming. Estimating biomass based on growth traits and foliar measurements may be simpler but still requires significant labour per plant, while truly measuring biomass (rather than making estimates) often requires the destruction of the plant itself.

Phenotypic changes during growth

Phenotypic changes during growth can limit attempts to reduce cycle times and/or breeding population sizes through early selection. Oil composition and concentration often change dramatically as a plant matures, making it hard to accurately select or cull progeny based on immature phenotypes (Coppen 2002). In some species, certain desired chemotypes may not even be detectable until plants reach a certain age. Doran and Bell (1994) studied the yield of monoterpenes in *E. camaldulensis* under glasshouse conditions and found that leaves from 26 month old trees had 42 % greater average cineole content than the same trees at 7 months of age, although ranking of the best and worst trees did remain consistent in this case. Barton et al. (1991) estimated that the narrow sense heritability of oil concentration in *E. kochii* (Oil Mallee) was $h^2 = 0.83$ for mature trees, but only $h^2 = 0.19$ for 1-year-old juveniles highlighting the difficulty in estimating the true performance of progeny at early stages using purely phenotypic measurements. Similarly, in *E. polybractea*, maternal oil concentration and oil concentration in young half-sib progeny are only weakly correlated, due to the large variation within the half-sib families (King et al. 2006). These findings caution against early phenotypic selection for oil concentration and composition as it may compromise final gain.

Improving selection efficiency with genomic selection

Marker-assisted selection techniques such as genomic selection (GS) are designed to tackle the issues discussed above by selecting individuals based on genotypic values rather than phenotypic values. GS has been shown, both in simulations and empirically, to provide improved selection efficiency compared to phenotypic selection (PS) and MAS (although this is not always the case—see Jannink et al. 2010). For poorly heritable traits in particular, GS has been shown to produce equal or larger gain than PS and MAS due to the greater predictive accuracy of GEBVs (Heffner et al. 2010; Resende et al. 2012a). On the other hand, several studies have indicated that a single cycle of PS often outperforms a single cycle of GS. For example, in a simulation for breeding in cassava, Oliveira et al. (2012) estimated that PS would produce gains 13–30 % greater than

GS for various traits over a single 4 year cycle. Similarly, in an empirical study for the improvement of an index of yield-related traits in maize, Massman et al. (2013) showed that GS outperformed MAS, but produced lower gains for a single cycle than PS.

Despite some limitations in single cycle selection, GS consistently outperforms other methods in recurrent (multiple cycle) selection. Cycle times can be dramatically reduced with GS because markers can be genotyped from very young plants, so selection based on GEBVs can be performed without waiting for mature phenotype. By inducing early flowering in selected individuals, the breeding cycle can be truncated (Grattapaglia and Resende 2011). The rate-limiting factor for reducing cycle time with GS is therefore the ability for early propagation, and achieving this is not necessarily straightforward in all essential oil-bearing crops. In some *Eucalyptus* species, e.g. *E. globulus*, chemically induced early flowering has successfully reduced cycle time by up to 50 % (Hasan and Reid 1995). In other *Eucalyptus* species, it is possible to graft juvenile cuttings onto established rootstock, triggering earlier flowering in the juvenile genotype. In *Melaleuca*, there has been limited success with chemical methods (Doran et al. 2002), however, large variation in flowering time exists due to abiotic stresses (such as low winter temperatures). This effect can be exploited to reduce flowering time from 42 months to just 14 months (Baskorowati et al. 2010).

Although the actual gain per cycle may sometimes be lower with GS, the increased frequency of cycles serves as a multiplier that makes the GS approach more efficient per unit time than PS (see Fig. 1). This is particularly effective for perennial crops because of their long generation times (and hence long PS cycle times). In the earlier cassava example, a reduction in cycle time from 4 to 2 years through the use of GS results in a predicted efficiency gain of 39–74 % for various traits compared to PS. For wood growth traits in various *Eucalyptus* species, it was predicted that reducing the breeding cycle length by 50 % would result in efficiency gains of 50–100 %, while reducing cycle length by 75 % (if possible) could see efficiency gains of up to 300 % (Resende et al. 2012a). Wong and Bernardo (2008) predicted that genomic selection can shorten cycle time in oil palm from 19 years to 6. In *Malus × domestica* (Apple), cycle time was reduced from 7 to 4 years resulting in over 100 % improvement in gain per unit time compared to conventional phenotypic selection methods (Kumar et al. 2012).

Factors affecting GS accuracy in essential oil species

GS aims to use the information provided by genome-wide markers to model the additive genetic variance of a trait.

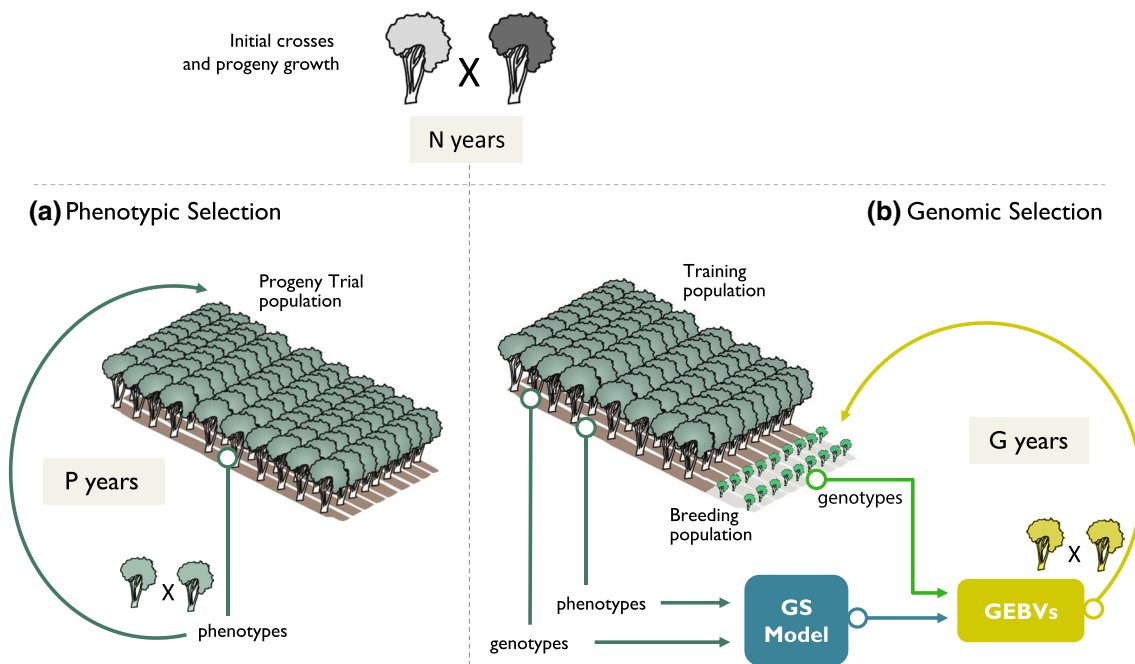


Fig. 1 A schematic representation of breeding approaches based on either phenotypic (PS) or genomic selection (GS). Both **a** PS and **b** GS start with a cross between parental lines or natural populations, requiring N years to reach maturity. After that, each cycle of PS requires P years in which to select, cross and grow the next generation to maturity. Each cycle of GS requires G years, but G is often

much smaller than P since the breeding population can be genotyped, have GEBVs calculated, and be selected and crossed at a young age. Over multiple cycles C , the time expended for PS is $N + CP$, while the time for GS is $N + CG$. Assuming similar gain per cycle from both methods, the gain from PS can be achieved in a much shorter time with GS

The markers carry two main forms of information that can improve predictive accuracy over traditional pedigree-based methods such as Best Linear Unbiased Prediction (BLUP). Firstly, the additive genetic effects of markers that are in LD with QTLs can be used to build a model of the trait variance based on the genetic architecture of the trait itself. Secondly, the markers provide an accurate measure of relatedness between individuals in the training and breeding populations based on identity-by-state or identity-by-descent of genotypes (Yang et al. 2010; de Los Campos et al. 2013). For example, in a pedigree two full-sibs are assumed to possess 50 % of common parental genetic material, however, due to random segregation of chromosomes during meiosis the real percentage may be significantly lower or higher. Accurately capturing this Mendelian sampling effect results in a finer grained measure of just how related two individuals are (Habier et al. 2007, 2013). While information about relatedness breaks down rapidly with each generation beyond the training population, LD information can persist and is more effective for predictions in individuals that are relatively unrelated to the training population (Habier et al. 2007).

The genome-wide scale of GS presents a modelling issue known as “large p , small n ” (Jannink et al. 2010), where the number of markers (p) for which effects are to

be estimated far exceeds the number of individuals (n) for which there are data. This results in over-fitting of the data, redundancy and multicollinearity between many markers, and the inability to model the marker effects using multiple regression by ordinary least squares. Aggressively culling the markers to a smaller subset containing only those with the largest effects often reduces the situation to that of MAS, forfeiting the inherent advantages of GS (Meuwissen et al. 2001; Moser et al. 2009). As a consequence, a range of modelling techniques have been designed to keep the advantage of including all or most marker effects while avoiding the ‘large p , small n ’ problem (de Los Campos et al. 2013). Detailed comparisons of various genomic selection models, both simulated and empirical, are available at Gianola (2013), Heslot et al. (2012), Lorenz et al. (2011) and Ogutu et al. (2012). They can broadly be categorized into two main strategies (Daetwyler et al. 2010): (1) BLUP-based methods (e.g. G-BLUP, RR-BLUP) that assume an infinitesimal model of genetic architecture, where all markers have effects drawn from a common normal distribution, though marker effects may be equally shrunken towards zero; (2) variable selection methods (e.g., Bayesian linear regression, LASSO, Elastic Net, machine-learning methods) that relax the assumption of a common distribution of marker effects across the genome, so

that portions of markers have significantly larger effects, smaller effects or are not included in the model at all. Both strategies model the additive genetic variance of the trait as described by a population's relatedness and LD (Habier et al. 2007, 2013; Zhong et al. 2009). However, their accuracies differ according the prevalence of each type of information, which in turn are affected by a range of factors: (1) the genetic architecture of the trait in question, (2) extent of LD in the populations, (3) degree of relatedness between the training and breeding populations, (4) the size of the training population, and (5) the density of markers used for genotyping.

One measure of accuracy is defined by Daetwyler et al. (2010) as the expected correlation between marker-predicted genotypic value and true genotypic value ($r_{g\hat{g}}$), which can be estimated by the equation:

$$r_{g\hat{g}} = \sqrt{Nh^2 / (Nh^2 + M_e)} \quad (1)$$

where N = training population size, h^2 = heritability of selected trait, M_e = the number of independent chromosomal regions, or QTLs, underlying the trait in the population. Equation 1 suggests that the accuracy of prediction improves with a larger training population, higher heritability and fewer QTLs. These predictions were mostly borne out in a recent study of five populations of maize, wheat and barley (Combs and Bernardo 2013). Likewise as M_e decreases, which occurs with increasing relatedness between individuals, accuracy improves (Daetwyler et al. 2013).

Below we examine how these factors might impact a genomic selection program for improving essential oil yield in perennial crops.

Genetic architecture

In GS, the additive effect of every genotyped marker on phenotypic variation is considered. The choice and accuracy of the GS model depends somewhat on the distribution of marker effects, which is ultimately tied to the number of QTLs underpinning the trait(s) and the distribution of QTL effects (Daetwyler et al. 2010). Understanding the genetic architecture of the traits under selection is highly important to the success of Genomic Selection.

Our understanding of the biosynthetic pathways that underlie terpene production is well-developed, and often a significant amount of variation in oil profile and concentration can be explained by the genes in those pathways (Fig. 2). QTL analysis in *E. nitens* identified 45 loci that were significantly associated with a range of monoterpene and sesquiterpene traits, each explaining from 3 to 16 % of variance (Henery et al. 2007). The authors noted that terpene concentration in eucalypts may therefore be affected

by relatively few loci of relatively large effect. Additionally, QTLs for several phenotypically correlated monoterpene traits were clustered together, pointing to putative genes with impact on the monoterpene precursor compound geranyl diphosphate, or perhaps regulatory factors for terpene synthase genes. QTL analysis also identified 13 widely spread QTL regions associated with the foliar concentration of terpenes in *E. globulus* explaining up to 71 % of trait variance (O'Reilly-Wapstra et al. 2011). In *Humulus lupulus* (Hop), linkage mapping and QTL analyses (Cerenak et al. 2009; McAdam et al. 2013) have revealed several large genomic regions of significance for total oil content, terpene concentrations (e.g. humulene) and biomass (e.g. cone weight). Certain putative QTLs clustered together within a linkage group and were associated with multiple oil traits, possibly reflecting the presence of gene families from terpene synthesis pathways. Other QTLs, however, showed large and isolated effects on individual terpene compounds, suggesting the presence of regulatory factors involved in the latter stages of biosynthesis. The small sample sizes and low number of genotyped markers used in these studies suggests that estimated QTL effects, such as 20 % of total oil content variation explained, are probably exaggerated. Additionally the narrow phenotypic and genotypic diversity present in the mapping populations limited the range of potential QTLs to be discovered. Finally, the common practice of using the same population to both detect QTLs and estimate their effect size has been shown to cause upward bias on estimates of the effects of QTLs (Utz et al. 2000). The majority of heritable phenotypic variation, not surprisingly, remains unexplained and would require genome-wide investigation in larger, more diverse populations.

QTL analyses have provided a low-resolution estimate of the location and effect size of major QTLs for oil traits. As a result, association studies in populations of the Myrtaceae family (which includes *Eucalyptus* and Tea Tree) have since focused on the specific genes involved in the synthesis of terpenoids and their effect on quantitative variation in oil content and composition. Külheim et al. (2011) investigated genetic associations between SNPs in 24 candidate genes from biosynthetic pathways and quantitative variation in plant secondary metabolites in *E. globulus*. The study revealed 37 significant associations in 11 genes, each explaining between 2 and 6 % of phenotypic variation in 19 oil traits. It should be noted that this study used a low density of markers so probably missed many QTLs, while the use of candidate genes and significance thresholds probably resulted in over-estimates of the effect sizes of associated QTLs. A candidate gene approach was also used by Webb et al. (2013) to investigate pathways of genetic control of terpene concentration in a small wild population of *M. alternifolia* (Tea Tree). This study revealed that, in

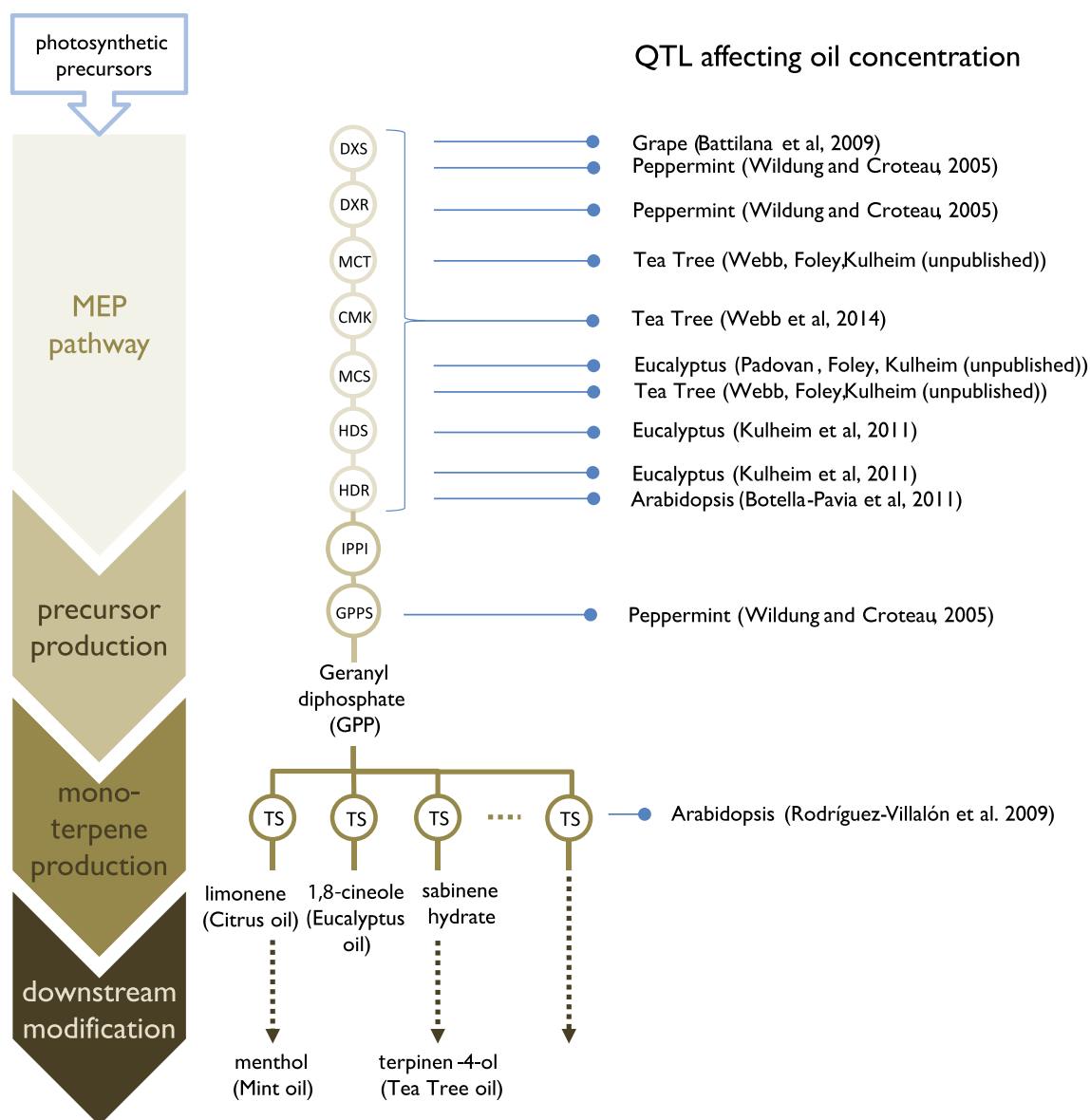


Fig. 2 An overview of the enzymes involved in the monoterpene biosynthesis pathway and QTL that are associated with terpene concentration. Total oil concentration appears to be influenced by the overall availability of photosynthetic precursors as input to the MEP path-

way, plus allelic variation at several points within the pathway. Variation at later stages, e.g., the terpene synthases (TS), mostly affects the ratio between individual terpenoids rather than overall concentration

addition to the relevance of individual genes within the terpene synthase pathway (see Fig. 2), the coordinated regulation of the precursor MEP pathway showed a strong and significant correlation with the concentration of the commercially-important terpinene-4-ol ($R^2 = 0.87$) in that species. The strength of this result, however, must be considered in light of the small sample size ($N = 48$).

Teasing out the more elaborate or precise genetic architecture of oil traits requires going beyond QTL approaches to genome-wide association studies (GWAS), though difficulties persist. QTLs in plants have been shown to

have varying estimated effect sizes from large ($>10\%$) to extremely small ($\ll 1\%$), with a skew towards smaller effect sizes (see Ingvarsson and Street 2011). The power of GWAS to detect a QTL is a function of effect size (a^2) and LD (R^2), so the smaller the effect of a QTL the harder it is to detect it (Hill 2012). When a trait is affected by a multitude of small-effect QTLs in a study population with short LD, then much of the genetic variation underpinning that trait may still remain unexplained—part of the classic ‘missing heritability’ in GWAS and QTL mapping studies (Myles et al. 2009). Additionally, few association studies

in forest trees have detected QTLs that explain greater than 5 % of trait variation (Grattapaglia et al. 2012), though rare alleles which explain a greater percentage of the total trait variance may exist but go undetected due to the lack of power when the study population size is small.

Little is known about the genome-wide architecture of essential oil yield in natural populations (Webb et al. 2014) as the rapid decay of LD in many outcrossed perennial species has made GWAS unfeasible until very recently. Zhu et al. (2008) and Hall et al. (2010) both presented lists of contemporary GWAS studies in plants, though none directly involved essential oil producing species, let alone any traits associated with essential oil production. Indeed most studies of the genetic architecture of oil concentration and biomass pertain to major commercial crops. Nevertheless, these studies provide insight into the complexity of these traits in plants in general. For example, kernel oil concentration analysed in a large maize population is under control of at least fifty QTLs of estimated small and mostly positive effect, that account for ~50 % of genetic variance (Laurie et al. 2004; Li et al. 2013).

Variation in essential oil concentration is most probably controlled by several key QTLs within and near to terpene synthesis pathway genes with large effect (Fig. 2), plus a greater number of QTLs of small effect throughout the genome which are likely regulatory elements. For the estimation of GEBVs, it may be prudent to consider modelling methods that distinguish these few well-characterized loci of larger effect from the many other unknown loci across the genome. A recent model, W-BLUP (weighted best linear unbiased prediction), was proposed by Zhao et al. (2014) with the intent to treat specific markers of large effect known from prior association studies differently, while still simultaneously modelling the many minor unknown effects. W-BLUP aims to bridge the gap between MAS and GS and could be appropriate for GS for essential oil yield due to a priori knowledge of important QTLs in the terpene biosynthesis pathway. Another recent model, MultiBLUP (Speed and Balding 2014), clusters markers, or genomic regions, into partitions based on effect size, with each partition being treated as a different random effect. Since significant oil trait QTLs have been mapped in clusters within linkage groups, this may be an effective approach worth exploring further. Models that assume constant marker-effect variance across the genome, such as RR-BLUP, are probably more appropriate for biomass traits where the infinitesimal model is realistic. In reviewing a wide range of GS models, de Los Campos et al. (2013) noted that in empirical studies model choice often makes little difference to accuracy, but also noted that few studies to date have used natural populations with short LD in which case model choice is likely to carry more weight.

Linkage disequilibrium (LD) and marker density

The resolution of QTL discovery is a function of LD decay, and therefore LD is at the heart of marker-based breeding techniques such as GS. Linkage disequilibrium refers to non-random association between pairs of loci, e.g., between two markers, between two QTLs, or between a QTL and a marker (Gupta et al. 2005). The intensity of LD between two loci is typically a function of the physical distance between them on a chromosome and the frequency of recombination in that region. Loci that are closer together and/or in a low recombination region have higher LD, since historical recombination events are less likely to have ‘shuffled’ the common stretch of DNA that links them. It is recombination events that cause LD to decay over time within a population (Fig. 3).

When a marker is associated with a phenotype, it acts as a predictor for the surrounding chromosomal region that is in LD with that marker—we can infer that a causative QTL probably lies somewhere within that linked region. When LD decays quickly, the linked chromosomal region surrounding any given marker is short, and so many uniformly distributed markers are required to ensure that every segment of the genome is linked with at least one nearby marker. Therefore, the average genomic distance over which LD decays determines the density of markers that will be required in a genomic selection program in order to adequately model marker-QTL associations.

Strong LD between two loci is commonly considered to be $R^2 > 0.1$ (Nakaya and Isobe 2012), though 0.2 or even 0.3 are also commonly used (see Table 1). Calus et al. (2008) demonstrated through simulation that the accuracy of GEBVs increased as the average LD between adjacent markers increased from $R^2 = 0.1$ to $R^2 = 0.2$, so for genomic selection it has been suggested that adjacent markers have LD of at least $R^2 > 0.1$ or 0.2 (Massman et al. 2013). The reasoning is well described by Ersoz et al. (2008). A large effect QTL may explain, for example, 15 % of the phenotypic variation. A marker in LD with that QTL at intensity of $R^2 = 0.1$ explains 10 % of the variation in the QTL, which in turn means that the marker itself only explains 1.5 % of the phenotypic variation. Therefore, the power to detect a QTL is a function of the effect size of the QTL and the strength of LD between the QTL and a nearby marker. Accordingly, GS accuracy increases with increasing marker density until it eventually reaches a plateau when the genome is ‘saturated’ with markers that are in strong LD with all QTLs (Meuwissen and Goddard 2010; Combs and Bernardo 2013).

For the reasons above, the first step in an association study design is to assess the extent of LD in the study population (Myles et al. 2009) in order to determine how many markers are required. Much of the research on the extent

Fig. 3 A schematic depiction of the decay in linkage disequilibrium (LD) in outcrossed populations over time. The decay is particularly rapid when there is a large effective population size (N_e) as the effect of genetic drift in reducing allelic variation is diminished. LD can be lengthened through breeding with a small effective population or inbreeding

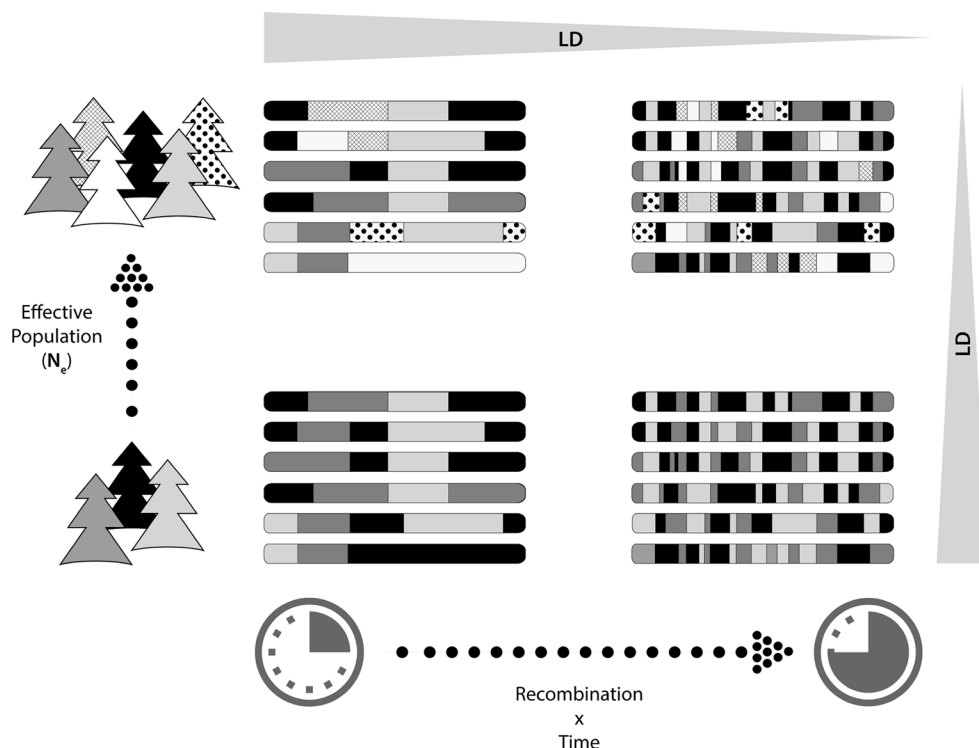


Table 2 The extent of significant linkage disequilibrium (LD) in various perennial species including *Pinus*, *Eucalyptus*, *Melaleuca* and *Vitis*

| Species | Significant r^2 | Distance (bp) ^a | References |
|-------------------------------|-------------------|----------------------------|---|
| <i>Picea. abies</i> | Unspecified | 100–200 | Rafalski and Morgante (2004), Heuertz et al. (2006) |
| <i>Pinus taeda</i> | <0.2 | 1500 | Neale and Savolainen (2004) |
| <i>Pinus nigra</i> | <0.2 | 300 | Chu et al. (2009) |
| <i>Pseudotsuga menziesii</i> | <0.2 | 300 | Krutovsky and Neale (2005) |
| <i>Populus balsamifera</i> | Unspecified | >750 | Olson et al. (2010) |
| <i>Eucalyptus globulus</i> | <0.2 | 200–500 | Thavamanikumar et al. (2011), Külheim et al. (2011) |
| <i>Eucalyptus nitens</i> | Unspecified | Low | Thumma (2005) |
| <i>Melaleuca alternifolia</i> | <0.3 | 500–1000 | Keszei et al. (2010), Webb (2015) (unpub) |
| <i>Vitis vinifera</i> | <0.2 | 50–200 | Lijavetzky et al. (2007), Myles et al. (2010) |

^a LD is typically short and its decay can be determined by the mean distance (in bp) at which the pairwise correlation (r^2) between markers drops below a threshold of significance

and distribution of LD has been reported in humans, animals and annual crop species, but there are examples in outcrossing perennial plants (Table 2).

In undomesticated outcrossing species, the LD between any two polymorphic markers typically decays rapidly with increasing genomic distance due to many generations of effective historical recombination in a large effective population (Fig. 3). This is certainly the case in *Eucalyptus* and *Melaleuca*, which are often highly outcrossing in the wild (Grattapaglia and Kirst 2008; Myburg et al. 2014) and have large effective population sizes. The very short range of LD in essential oil-bearing species such as *E. polybractea* and *M. alternifolia* implies that GS for oil yield in progeny

derived from naturally sourced progenitors would require a very high density of markers across the whole genome, possibly to a density whereby the causative SNPs themselves are genotyped. *Eucalyptus polybractea* has an estimated genome size of 550 Mbp. Linkage disequilibrium likely decays within a similar distance to that observed in *E. nitens* and *E. globulus* (i.e., 100 bp) as the three species share similarly small geographical distributions and probably similar historical effective population sizes. Therefore, at least 5.5 million genome-wide markers would be required to ensure adequate coverage across all regions of LD in the genome, and preferably more to increase power. Considering that the SNP density in *E. globulus* is about 1

every 31 bp (Külheim et al. 2009), obtaining 5.5 m genotyped markers is biologically and technically feasible using current whole genome re-sequencing technology (though this says nothing of the cost of doing so in many individuals!), and could result in virtually the entire additive component of the genetic variance being accounted for by the markers (Daetwyler et al. 2010).

The benefits of using whole genome SNP data for estimating genetic breeding values, as opposed to less dense genotyping, were demonstrated in a simulation study by Meuwissen and Goddard (2010). Firstly, the accuracy of prediction doubled as marker density increased from 1000 per morgan to 33,000 per morgan, irrespective of whether many or few QTLs were simulated for the trait. Secondly, the accuracy of GEBVs is likely to hold for many more generations since the markers for which effects are estimated are so close to, if not actually, the causative SNPs for the trait. Thirdly, while reduced representation sequencing techniques such as Genotyping-by-Sequencing (GBS) can still generate large numbers of SNPs, there is a risk of missing major QTLs, especially if LD is short. For example, Romay et al. (2013) used GBS for a GWAS of flowering time in maize and found only one marker significantly associated with the most important gene associated with flowering time (*ZmCCT*). In other words, the GBS markers almost failed to detect a known major QTL, even with 680 k SNPs genotyped in inbred lines. The rapid LD decay in the region surrounding *ZmCCT* was cited as a reason for the near failure to detect it, and many other unknown QTLs would have undoubtedly gone undetected. Similarly, Myles et al. (2010) used reduced representation genotyping to characterize the *Vitis vinifera* (Grape) genome and came to the conclusion that due to the presence of very short LD, progress towards GWAS and GS in grape would require whole genome sequencing to ensure association with most functional QTLs.

Relatedness and training population size

When the training and validation/breeding populations are closely related, much of the accuracy achieved with GS can come from the relatedness information carried by markers. The G-BLUP model, which uses markers to define a genomic relationship matrix to replace the pedigree matrix used in standard phenotypic BLUP, is often highly effective in this scenario (de Los Campos et al. 2013), and can be efficiently implemented with relatively low marker density and small training population size. Indeed this may be a straightforward approach for GS in Hop due to its long history of domestication. However, many other essential oil crops are largely undomesticated and little genetic relatedness exists in individuals sourced from natural populations. Here, information due to LD becomes the dominant

component of GS accuracy (Habier et al. 2007), assuming a model that effectively estimates marker effects of varying size is used, thus compensating for the lack of relationship information (Meuwissen and Goddard 2010). Consequently a higher density of markers is needed to ensure all relevant QTLs are detected, particularly in populations with short LD [see “Linkage disequilibrium (LD) and marker density” for more detail]. As marker density increases, a larger training population is required in order to accurately estimate additional marker effects (especially those of relatively small effect). In general, a larger training population results in increased accuracy of prediction (Zhong et al. 2009; Grattapaglia and Resende 2011; Lorenz et al. 2011).

Genotyping a very high density of markers has been a limitation for practical implementation of GS in outcrossing, undomesticated tree populations (Nakaya and Isobe 2012). Beaulieu et al. (2014) were one of the first to assess the accuracy of GS in a large, diverse, undomesticated population of outcrossing trees (White spruce *Picea abies*). Training and predictions were made both within and between half-sib families, with accuracies being significantly lower in the latter as expected, but still higher than that of pedigree-based models. They recommended that for the time being, for most tree species, GS models should be trained and used within related populations in order to obtain high accuracies with limited marker density. For undomesticated species this issue can be addressed in the short term by increasing the relatedness within the study population through an initial breeding phase, which reduces the effective population size and lengthens LD (see Fig. 3), as demonstrated in *Pinus taeda* (Resende et al. 2012b) and *Eucalyptus* (Resende et al. 2012a; Denis and Bouvet 2013). These studies resulted in good prediction accuracy with only sparse marker coverage but the models are unlikely to work well in future breeding populations because relatedness to the training population declines rapidly per generation. With the decreasing cost of genotyping, GS may in future be performed with higher accuracy in undomesticated populations with greater allelic diversity.

Heritability (h^2)

The accuracy of genomic selection is lower for traits with lower h^2 , though this can be improved if the training population size is increased, thereby keeping the Nh^2 term of Eq. (1) constant (Combs and Bernardo 2013). Nevertheless, for traits with low heritability, GS has been shown to produce equal or larger gain than PS and MAS due to the greater predictive accuracy of GEBVs (Heffner et al. 2010; Resende et al. 2012a). Thus, GS is likely to be the best method for artificial selection on essential oil yield, for which the all-important biomass traits are often of low to moderate heritability.

The method used for the estimation of the heritability of a trait may also have an effect on the estimated accuracy of GS. Downwardly biased estimates of h^2 may occur if genotypes are assumed to be independent when, in reality, they are correlated (Estaghvirou et al. 2013).

Selection for multiple traits with GS

Selecting for oil yield is, in reality, selecting for multiple complex traits, or a selection index formed from those traits. For example, in breeding for pharmaceutical grade *Eucalyptus* oil an index comprising total oil concentration, leaf biomass, % cineole, % undesirable compounds, family survival rate and other traits could be used.

Bernardo and Yu (2007) speculated that GS would outperform other methods for improving a selection index in maize comprising multiple traits, as there would be a large number of QTLs involved, many of which would be associated with traits of low heritability. This prediction was borne out in a yield-based index of traits in maize (Massman et al. 2013) which resulted in significantly increased grain yield per hectare despite little improvement in each of the component traits within the index.

Three approaches may be taken for genomic selection of multiple traits: (1) estimate marker effects for each individual trait and then form a selection index based on the weighted GEBVs of each trait (Resende et al. 2012a); (2) estimate marker effects for the index as a trait itself. (3) Use a multiple-trait genomic selection model (MT-GS) when a trait with low heritability is correlated with another trait of high heritability (Calus and Veerkamp 2011). A full comparison of these three approaches to selecting for essential oil yield requires further investigation.

Conclusion

Selection for complex quantitative traits has presented challenges to breeders that do not arise with more simple Mendelian traits. In plants, molecular assisted selection using small numbers of significant QTL has not proven particularly effective, especially in outcrossing species with little prior domestication. Genomic Selection, on the other hand, has shown great promise and could improve the breeding process in essential oil bearing crops. The highly complicated genetic architecture involved in oil yield traits may be most adequately detected and accounted for using whole genome re-sequencing and genotyping. Coupled with advanced modelling techniques, the gain per unit time using genomic selection could well outstrip traditional breeding practises, especially in perennials such as *Eucalyptus*, Tea Tree and Hop where the reduction in cycle time has the greatest impact on overall gain.

Acknowledgments We would like to acknowledge funding from the Australian Research Council Linkage Programme (LP110100184) to WJF, and from the Rural Industries Research and Development Corporation (RIRDC), Australia. We are grateful to Richard Davis of GR Davis Pty Ltd for providing firsthand insight into the realities of commercial breeding for essential oils. Finally, thanks to John Henning (USDA-ARS-Forage Seed Research Center Unit, USA) for providing answers to questions on Hop breeding.

Compliance with ethical standards

Conflict of interest None.

References

- Baker GR, Doran JC, Williams E, Morris G (2014) Highly improved tea tree varieties to maximise profit, rural industries research and development corporation. PRJ-003689, Barton, ACT
- Bakkali F, Averbeck S, Averbeck D, Idaomar M (2008) Biological effects of essential oils—a review. *Food Chem Toxicol* 46:446–475
- Barton AF, Cotterill PP, Brooker MI (1991) Heritability of cineole yield in *Eucalyptus kochii*. *Silvae Genet* 40:37–38
- Baskorowati L, Moncur MW, Doran JC, Kanowski PJ (2010) Reproductive biology of *Melaleuca alternifolia* (Myrtaceae) 1. Floral biology. *Aust J Bot* 58:373–383
- Beaulieu J, Doerksen T, Clément S, MacKay J, Bousquet J (2014) Accuracy of genomic selection models in a large population of open-pollinated families in white spruce. *Heredity* 113:343–352
- Beavis WD (1994) The power and deceit of QTL experiments: lessons from comparative QTL studies. In: Proceedings of the Forty-Ninth Annual Corn and Sorghum Industry Research Conference, pp 250–266
- Bernardo R (2008) Molecular markers and selection for complex traits in plants: learning from the last 20 years. *Crop Sci* 48:1649–1664
- Bernardo R, Yu J (2007) Prospects for genomewide selection for quantitative traits in Maize. *Crop Sci* 47:1082
- Butcher PA, Matheson A, Slee MU (1996) Potential for genetic improvement of oil production in *Melaleuca alternifolia* and *M. linariifolia*. *New Forest* 11:31–51
- Byrne D (2007) Molecular marker use in perennial plant breeding. *Acta Hort* 751:163–167
- Calus MPL, Veerkamp RF (2011) Accuracy of multi-trait genomic selection using different methods. *Genet Sel Evol* 43:1–14
- Calus MPL, Meuwissen T, de Roos APW, Veerkamp RF (2008) Accuracy of genomic selection using different methods to define haplotypes. *Genetics* 178:553–561
- CBI Ministry of Foreign Affairs (2012) Promising EU export markets for essential oils. <http://www.cbi.eu/marketintel/Essential-oils-for-cosmetics-promising-EU-export-markets/164831>. Accessed 14 May 2014
- Cerenak A, Satovic Z, Jakse J, Luthar Z, Carovic-Stanko K, Javornik B (2009) Identification of QTLs for alpha acid content and yield in hop (*Humulus Lupulus* L.). *Euphytica* 170:141–154
- Chu Y, Su X, Huang Q, Zhang X (2009) Patterns of DNA sequence variation at candidate gene loci in black poplar (*Populus nigra* L.) as revealed by single nucleotide polymorphisms. *Genetica* 137:141–150
- Combs E, Bernardo R (2013) Accuracy of genomewide selection for different traits with constant population size, heritability, and number of markers. *Plant Genome* 6
- Coppen JJW (2002) *Eucalyptus: the genus Eucalyptus*. Taylor & Francis, London

- Crossa J, de Los Campos G, Perez P, Gianola D, Burgueno J, Araus JL, Makumbi D, Singh RP, Dreisigacker S, Yan J et al (2010) Prediction of genetic values of quantitative traits in plant breeding using pedigree and molecular markers. *Genetics* 186:713–724
- Daetwyler HD, Pong-Wong R, Villanueva B, Woolliams JA (2010) The impact of genetic architecture on genome-wide evaluation methods. *Genetics* 185:1021–1031
- Daetwyler HD, Calus MPL, Pong-Wong R, de Los Campos G, Hickey JM (2013) Genomic prediction in animals and plants: simulation of data, validation, reporting, and benchmarking. *Genetics* 193:347–365
- de Los Campos G, Hickey JM, Pong-Wong R, Daetwyler HD, Calus MPL (2013) Whole-genome regression and prediction methods applied to plant and animal breeding. *Genetics* 193:327–345
- Denis M, Bouvet J (2013) Efficiency of genomic selection with models including dominance effect in the context of *Eucalyptus* breeding. *Tree Genet Genomes* 9:37–51
- Doran JC, Bell R (1994) Influence of non-genetic factors on yield of monoterpenes in leaf oils of *Eucalyptus camaldulensis*. *New Forest* 8:363–379
- Doran JC, Matheson A (1994) Genetic parameters and expected gains from selection for monoterpene yields in Petford *Eucalyptus camaldulensis*. *New Forest* 8:155–167
- Doran JC, Baker GR, Williams E, Southwell I (2002) Improving Australian Tea Tree through selection and breeding (1996–2001). Rural Industries and Research Development Corporation, ACT
- Ersoz ES, Yu J, Buckler ES (2008) Applications of linkage disequilibrium and association mapping in maize. In: Kriz AL, Larkins BA (eds) *Molecular genetic approaches to maize improvement*. Springer Science & Business Media, pp 173–195
- Estaghirou SBO, Ogutu JO, Schulz-Streeck T, Knaak C, Ouzunova M, Gordillo A, Piepho H (2013) Evaluation of approaches for estimating the accuracy of genomic prediction in plant breeding. *BMC Genom* 14:860
- Falconer DS, Mackay TFC, Frankham R (1996) *Introduction to quantitative genetics* (4th edn). Trends Genet 12:280
- Farmer EE (2014) *Leaf defense*. Oxford University Press, New York
- Gianola D (2013) Priors in whole-genome regression: the bayesian alphabet returns. *Genetics* 194:573–596
- Goodger JQ, Woodrow IE (2008) Selection gains for essential oil traits using micropropagation of *Eucalyptus polybractea*. *Forest Ecol Manag* 255:3652–3658
- Goodger JQ, Woodrow IE (2012) Genetic determinants of oil yield in *Eucalyptus polybractea* R.T. Baker. *Trees* 26:1951–1956
- Gouy M, Roussele Y, Bastianelli D, Lecomte P, Bonnal L, Roques D, Efile J, Rocher S, Daugrois J, Toubi L et al (2013) Experimental assessment of the accuracy of genomic selection in sugarcane. *Theor Appl Genet* 126:2575–2586
- Grant G (1997) Genetic variation in *Eucalyptus polybractea* and the potential for improving leaf oil production. Thesis. Australian National University, Canberra
- Grattapaglia D, Kirst M (2008) *Eucalyptus* applied genomics: from gene sequences to breeding tools. *New Phytol* 179:911–929
- Grattapaglia D, Resende MD (2011) Genomic selection in forest tree breeding. *Tree Genet Genomes* 7:241–255
- Grattapaglia D, Vaillancourt RE, Shepherd M, Thumma BR, Foley WJ, K ulheim C, Potts BM, Myburg AA (2012) Progress in Myrtaceae genetics and genomics: *Eucalyptus* as the pivotal genus. *Tree Genet Genomes* 8:463–508
- Gupta PK, Rustgi S, Kulwal PL (2005) Linkage disequilibrium and association studies in higher plants: present status and future prospects. *Plant Mol Biol* 57:461–485
- Habier D, Fernando RL, Dekkers JCM (2007) The impact of genetic relationship information on genome-assisted breeding values. *Genetics* 177:2389–2397
- Habier D, Fernando RL, Garrick DJ (2013) Genomic BLUP decoded: a look into the black box of genomic prediction. *Genetics* 194:597–607
- Hall D, Tegstrom C, Ingvarsson PK (2010) Using association mapping to dissect the genetic basis of complex traits in plants. *Briefings Funct Genom* 9:157–165
- Hasan O, Reid JB (1995) Reduction of generation time in *Eucalyptus globulus*. *Plant Growth Regul* 17:53–60
- Heffner EL, Lorenz AJ, Jannink J, Sorrells ME (2010) Plant breeding with genomic selection: gain per unit time and cost. *Crop Sci* 50:1681–1690
- Henery ML, Moran GF, Wallis IR, Foley WJ (2007) Identification of quantitative trait loci influencing foliar concentrations of terpenes and formylated phloroglucinol compounds in *Eucalyptus nitens*. *New Phytol* 176:82–95
- Henning J, Haunold A, Nickerson G, Gampert U (1997) Estimates of heritability and genetic correlation for five traits in female hop accessions. *J Am Soc Brew Chem* 55:161–165
- Heslot N, Yang H, Sorrells ME, Jannink J (2012) Genomic selection in plant breeding: a comparison of models. *Crop Sci* 52:146–160
- Heuertz M, de Paoli E, K allman T, Larsson H, Jurman I, Morgante M, Lascoux M, Gyllenstrand N (2006) Multilocus patterns of nucleotide diversity, linkage disequilibrium and demographic history of Norway spruce [*Picea abies* (L.) Karst]. *Genetics* 174:2095–2105
- Hill WG (2012) Quantitative genetics in the genomics era. *Curr Genomics* 13:196–206
- Holland JB (2004) Implementation of molecular markers for quantitative traits in breeding programs—challenges and opportunities. In: *New directions for a diverse planet. Proceedings of the 4th International Crop Science Congress*, pp 1–13
- Homer LE, Leach DN, Lea D, Slade Lee L, Henry RJ, Baverstock PR (2000) Natural variation in the essential oil content of *Melaleuca alternifolia* Cheel (Myrtaceae). *Biochem Syst Ecol* 28:367–382
- Hospital F (2009) Challenges for effective marker-assisted selection in plants. *Genetica* 136:303–310
- Ingvarsson PK, Street NR (2011) Association genetics of complex traits in plants. *New Phytol* 189:909–922
- Isik F (2014) Genomic selection in forest tree breeding: the concept and an outlook to the future. *New Forest* 45:379–401
- Izadi-Darbandi A, Bahmani K, Ramshini A, Moradi N (2013) Heritability estimates of agronomic traits and essential oil content in Iranian fennels. *J Agric Sci Technol* 15:1275–1283
- Jannink J, Lorenz AJ, Iwata H (2010) Genomic selection in plant breeding: from theory to practice. *Brief Funct Genom* 9:166–177
- Keszei A, Hassan Y, Foley WJ (2010) A biochemical interpretation of terpene chemotypes in *Melaleuca alternifolia*. *J Chem Ecol* 36:652–661
- King DJ, Gleadow RM, Woodrow IE (2004) Terpene deployment in *Eucalyptus polybractea*; relationships with leaf structure, environmental stresses, and growth. *Funct Plant Biol* 31:451–460
- King DJ, Gleadow RM, Woodrow IE (2006) The accumulation of terpenoid oils does not incur a growth cost in *Eucalyptus polybractea* seedlings. *Funct Plant Biol* 33:497–505
- Krutovsky KV, Neale DB (2005) Nucleotide diversity and linkage disequilibrium in cold-hardiness-and wood quality-related candidate genes in Douglas fir. *Genetics* 171:2029–2041
- K ulheim C, Yeoh SH, Maintz J, Foley WJ, Moran GF (2009) Comparative SNP diversity among four *Eucalyptus* species for genes from secondary metabolite biosynthetic pathways. *BMC Genom* 10:452
- K ulheim C, Yeoh SH, Wallis IR, Laffan S, Moran GF, Foley WJ (2011) The molecular basis of quantitative variation in foliar

- secondary metabolites in *Eucalyptus globulus*. *New Phytol* 191:1041–1053
- Kulkarni R, Baskaran K, Ramesh S (2003) Five cycles of recurrent selection for increased essential oil content in East Indian lemongrass: response to selection, and effects on heritabilities of traits and intertrait correlations. *Plant Breeding* 122:131–135
- Kumar S, Chagné D, Bink MCAM, Volz RK, Whitworth C, Carlisle C, Zhang T (2012) Genomic selection for fruit quality traits in apple (*Malus × domestica* Borkh.). *PLoS One* 7:e36674
- Kumar B, Mali H, Gupta E (2014) Genetic variability, character association, and path analysis for economic traits in menthofuran rich half-sib seed progeny of *Mentha piperita* L. *Biomed Res Int* 2014:1–7. Article ID 150830. doi:10.1155/2014/150830
- Laurie CC, Chasalow SD, LeDeaux JR, McCarroll R, Bush D, Hauge B, Lai C, Clark D, Rocheford TR, Dudley JW (2004) The genetic architecture of response to long-term artificial selection for oil concentration in the maize kernel. *Genetics* 168:2141–2155
- Li H, Peng Z, Yang X, Wang W, Fu J, Wang J, Han Y, Chai Y, Guo T, Yang N et al (2013) Genome-wide association study dissects the genetic architecture of oil biosynthesis in maize kernels. *Nat Genet* 45:43–50
- Lijavetzky D, Cabezas JA, Ibáñez A, Rodríguez V, Martínez-Zapater JM (2007) High throughput SNP discovery and genotyping in grapevine (*Vitis vinifera* L.) by combining a re-sequencing approach and SNPlex technology. *BMC Genom* 8:424
- Lorenz AJ, Chao S, Asoro FG, Heffner EL, Hayashi T, Iwata H, Smith KP, Sorrells ME, Jannink J (2011) Genomic selection in plant breeding: knowledge and prospects. In: Sparks DL (ed) *Advances in agronomy*. Academic Press, San Diego, pp 77–123
- Lorenzana RE, Bernardo R (2009) Accuracy of genotypic value predictions for marker-based selection in biparental plant populations. *Theor Appl Genet* 120:151–161
- Luby JJ, Shaw DV (2001) Does marker-assisted selection make dollars and sense in a fruit breeding program? *HortScience* 36:872–879
- Massman JM, Jung HG, Bernardo R (2013) Genomewide selection versus marker-assisted recurrent selection to improve grain yield and stover-quality traits for cellulosic ethanol in maize. *Crop Sci* 53:58–66
- McAdam EL, Freeman JS, Whittock SP, Buck EJ, Jakse J, Cerenak A, Javornik B, Kilian A, Wang C, Andersen D et al (2013) Quantitative trait loci in hop (*Humulus lupulus* L.) reveal complex genetic architecture underlying variation in sex, yield and cone chemistry. *BMC Genom* 14:360
- McAdam EL, Vaillancourt RE, Koutoulis A, Whittock SP (2014) Quantitative genetic parameters for yield, plant growth and cone chemical traits in hop (*Humulus lupulus* L.). *BMC Genet* 15:22
- Meuwissen T, Goddard ME (2010) Accurate prediction of genetic values for complex traits by whole-genome resequencing. *Genetics* 185:623–631
- Meuwissen T, Hayes BJ, Goddard ME (2001) Prediction of total genetic value using genome-wide dense marker maps. *Genetics* 157:1819–1829
- Moser G, Tier B, Crump RE, Khatkar MS, Raadsma HW (2009) A comparison of five methods to predict genomic breeding values of dairy bulls from genome-wide SNP markers. *Genet Sel Evol* 41:56
- Murakami A (1999) Inheritance of major chemical components in Hops. *J Brew* 105:107–111
- Myburg AA, Grattapaglia D, Tuskan GA, Hellsten U, Hayes RD, Grimwood J, Jenkins J, Lindquist E, Tice H, Bauer D et al (2014) The genome of *Eucalyptus grandis*. *Nature* 510:356–362
- Myles S, Peiffer JA, Brown PJ, Ersoz ES, Zhang Z, Costich DE, Buckler ES (2009) Association mapping: critical considerations shift from genotyping to experimental design. *Plant Cell Online* 21:2194–2202
- Myles S, Chia J, Hurwitz B, Simon C, Zhong GY, Buckler ES, Ware D (2010) Rapid genomic characterization of the genus *Vitis*. *PLoS One* 5:e8219
- Nakaya A, Isobe SN (2012) Will genomic selection be a practical method for plant breeding? *Ann Bot Lond* 110:1303–1316
- Neale DB, Savolainen O (2004) Association genetics of complex traits in conifers. *Trends Plant Sci* 9:325–330
- O'Reilly-Wapstra JM, Freeman JS, Davies NW, Vaillancourt RE, Fitzgerald H, Potts BM (2011) Quantitative trait loci for foliar terpenes in a global eucalypt species. *Tree Genet Genomes* 7:485–498
- Ogutu JO, Schulz-Streeck T, Piepho H (2012) Genomic selection using regularized linear regression models: ridge regression, lasso, elastic net and their extensions. *BMC Proc* 6(Suppl 2):S10
- Oliveira EJ, Resende MD, Silva Santos V, Ferreira CF, Oliveira GAF, Silva MS, Oliveira LA, Aguilar-Vildoso CI (2012) Genome-wide selection in cassava. *Euphytica* 187:263–276
- Olson MS, Robertson AL, Takebayashi N, Silim S, Schroeder WR, Tiffin P (2010) Nucleotide diversity and linkage disequilibrium in balsam poplar (*Populus balsamifera*). *New Phytol* 186:526–536
- Pank F (2010) Aims and results of breeding research on eight medicinal and aromatic plants—a survey. *Isr J Plant Sci* 58:241–249
- Pearson M (1993) The Good Oil. Eucalyptus oil distilleries in Australia. *Australas Hist Archaeol* 11:99–107
- Rafalski A, Morgante M (2004) Corn and humans: recombination and linkage disequilibrium in two genomes of similar size. *Trends Genet* 20:103–111
- Resende MD, Resende MFR, Sansaloni CP, Petrolini CD, Missiaggia AA, Aguiar AM, Abad JM, Takahashi EK, Rosado AM, Faria DA et al (2012a) Genomic selection for growth and wood quality in Eucalyptus: capturing the missing heritability and accelerating breeding for complex traits in forest trees. *New Phytol* 194:116–128
- Resende MFR, Munoz P, Resende MD, Garrick DJ, Fernando RL, Davis JM, Jokela EJ, Martin TA, Peter GF, Kirst M (2012b) Accuracy of genomic selection methods in a standard data set of Loblolly pine (*Pinus taeda* L.). *Genetics* 190:1503–1510
- Romay MC, Millard MJ, Glaubitz JC, Peiffer JA, Swarts KL, Castevens TM, Elshire RJ, Acharya CB, Mitchell SE, Flint-Garcia SA et al (2013) Comprehensive genotyping of the USA national maize inbred seed bank. *Genome Biol* 14:R55
- Speed D, Balding DJ (2014) MultiBLUP: improved SNP-based prediction for complex traits. *Genome Res* 24:1550–1557
- Thavamanikumar S, McManus LJ, Tibbits JF, Bossinger G (2011) The significance of single nucleotide polymorphisms (SNPs) in *Eucalyptus globulus* breeding programs. *Austral For* 74:23–29
- Thavamanikumar S, Southerton SG, Bossinger G, Thumma BR (2013) Dissection of complex traits in forest trees—opportunities for marker-assisted selection. *Tree Genet Genomes* 9:627–639
- Thumma BR (2005) Polymorphisms in *Cinnamoyl CoA Reductase* (CCR) are associated with variation in microfibril angle in *Eucalyptus* spp. *Genetics* 171:1257–1265
- Utz HF, Melchinger AE, Schön CC (2000) Bias and sampling error of the estimated proportion of genotypic variance explained by quantitative trait loci determined from experimental data in maize using cross validation and validation with independent samples. *Genetics* 154:1839–1849
- Webb H, Lanfear R, Hamill J, Foley WJ, Külheim C (2013) The yield of essential oils in *Melaleuca alternifolia* (Myrtaceae) is regulated through transcript abundance of genes in the MEP pathway. *PLoS One* 8:e60631

- Webb H, Foley WJ, Külheim C (2014) The genetic basis of foliar terpene yield: implications for breeding and profitability of Australian essential oil crops. *Plant Biotechnol* 31:363–376
- Wong CK, Bernardo R (2008) Genomewide selection in oil palm: increasing selection gain per unit time and cost with small populations. *Theor Appl Genet* 116:815–824
- Würschum T, Reif JC, Kraft T, Janssen G, Zhao Y (2013) Genomic selection in sugar beet breeding populations. *BMC Genet* 14:85
- Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, Nyholt DR, Madden PA, Heath AC, Martin NG, Montgomery GW et al (2010) Common SNPs explain a large proportion of the heritability for human height. *Nat Genet* 42:565–569
- Zhao Y, Gowda M, Liu W, Würschum T, Maurer HP, Longin FH, Ranc N, Reif JC (2012) Accuracy of genomic selection in European maize elite breeding populations. *Theor Appl Genet* 124:769–776
- Zhao Y, Mette MF, Gowda M, Longin FH, Reif JC (2014) Bridging the gap between marker-assisted and genomic selection of heading time and plant height in hybrid wheat. *Heredity* 112:638–645
- Zhong S, Dekkers JCM, Fernando RL, Jannink J (2009) Factors affecting accuracy from genomic selection in populations derived from multiple inbred lines: a Barley case study. *Genetics* 182:355–364
- Zhu C, Gore M, Buckler ES, Yu J (2008) Status and prospects of association mapping in plants. *Plant Genome* 1:5–20