

Document downloaded from:

<http://hdl.handle.net/10251/63024>

This paper must be cited as:

Pérez De Castro, AM.; Vilanova Navarro, S.; Cañizares Sales, J.; Pascual Bañuls, L.; Blanca Postigo, JM.; Díez Niclós, MJTDJ.; Prohens Tomás, J.... (2012). Application of genomic tools in plant breeding. *Current Genomics*. 13(3):179-195.  
doi:10.2174/138920212800543084.



The final publication is available at

<http://dx.doi.org/10.2174/138920212800543084>

Copyright Bentham Science Publishers

Additional Information

1 **Application of genomic tools in plant breeding**

2

3 A.M. Pérez de Castro, S. Vilanova, J. Cañizares, L. Pascual, J.M. Blanca, M.J. Díez, J. Prohens\* and B.

4 Picó

5

6 *Instituto de Conservación y Mejora de la Agrodiversidad Valenciana, Universitat Politècnica de*

7 *València, Camino de Vera 14, 46022 Valencia, Spain*

8

9 Running title: Genomic tools in plant breeding

10

11 \*Address correspondence to this author at the Instituto de Conservación y Mejora de la Agrodiversidad

12 Valenciana, Universitat Politècnica de València, Camino de Vera 14, 46022 Valencia, Spain; Tel:

13 +34963879424; Fax: +34963879422; E-mail: [jprohens@btc.upv.es](mailto:jprohens@btc.upv.es)

14

15 **Abstract:** Plant breeding has been very successful in developing improved varieties using conventional  
16 tools and methodologies. Nowadays, the availability of genomic tools and resources is leading to a new  
17 revolution of plant breeding, as they facilitate the study of the genotype and its relationship with the  
18 phenotype, in particular for complex traits. Next Generation Sequencing (NGS) technologies are allowing  
19 the mass sequencing of genomes and transcriptomes, which is producing a vast array of genomic  
20 information. The analysis of NGS data by means of bioinformatics developments allows discovering new  
21 genes and regulatory sequences and their positions, and makes available large collections of molecular  
22 markers. Genome-wide expression studies provide breeders with an understanding of the molecular basis  
23 of complex traits. Genomic approaches include TILLING and EcoTILLING, which make possible to  
24 screen mutant and germplasm collections for allelic variants in target genes. Re-sequencing of genomes is  
25 very useful for the genome-wide discovery of markers amenable for high-throughput genotyping  
26 platforms, like SSRs and SNPs, or the construction of high density genetic maps. All these tools and  
27 resources facilitate studying the genetic diversity, which is important for germplasm management,  
28 enhancement and use. Also, they allow the identification of markers linked to genes and QTLs, using a  
29 diversity of techniques like bulked segregant analysis (BSA), fine genetic mapping, or association  
30 mapping. These new markers are used for marker assisted selection, including marker assisted backcross  
31 selection, 'breeding by design', or new strategies, like genomic selection. In conclusion, advances in  
32 genomics are providing breeders with new tools and methodologies that allow a great leap forward in  
33 plant breeding, including the 'superdomestication' of crops and the genetic dissection and breeding for  
34 complex traits.

35

36 **Key Words:** bioinformatics, complex traits, genetic maps, marker assisted selection, molecular markers,  
37 next-generation-sequencing, quantitative trait loci

38

39

## 40 **INTRODUCTION**

41

42 Ever since the beginnings of the domestication of plants, some 10,000 years ago, plant breeding  
43 has been extremely successful in developing crops and varieties that have contributed to the development  
44 of modern societies, and have successively beaten (neo-)Malthusian predictions [1]. Application of

45 conventional pre-genomics scientific breeding methodologies has led to the development of modern  
46 cultivars, which have contributed to the dramatic improvement of yield of most major crops since the  
47 middle of the 20<sup>th</sup> century. The success of plant breeding in the last century has relied in the utilization of  
48 natural and mutant induced genetic variation and in the efficient selection, by using suitable breeding  
49 methods, of the favorable genetic combinations. In this respect, the evaluation and identification of  
50 genetic variants of interest as well as the selection methodologies used have largely been based in the  
51 phenotypic evaluation.

52           Nowadays, genomics provides breeders with a new set of tools and techniques that allow the  
53 study of the whole genome, and which represents a paradigm shift, by facilitating the direct study of the  
54 genotype and its relationship with the phenotype [2]. While classical genetics revolutionized plant  
55 breeding at the beginning of the 20<sup>th</sup> century, genomics is leading to a new revolution in plant breeding at  
56 the beginning of the 21<sup>th</sup> century.

57           The field of genomics and its application to plant breeding are developing very quickly. The  
58 combination of conventional breeding techniques with genomic tools and approaches is leading to a new  
59 genomics-based plant breeding. In this new plant breeding context, genomics will be essential to develop  
60 more efficient plant cultivars, which are necessary, according to FAO, for the new 'greener revolution'  
61 needed to feed the world's growing population while preserving natural resources.

62           One of the main pillars of genomic breeding is the development of high-throughput DNA  
63 sequencing technologies, collectively known as next generation sequencing (NGS) methods. These and  
64 other technical revolutions provide genome-wide molecular tools for breeders (large collections of  
65 markers, high-throughput genotyping strategies, high density genetic maps, new experimental  
66 populations, etc.) that can be incorporated into existing breeding methods [2, 3, 4, 5]. Recent advances in  
67 genomics are producing new plant breeding methodologies, improving and accelerating the breeding  
68 process in many ways (e.g., association mapping, marker assisted selection, 'breeding by design', gene  
69 pyramiding, genomic selection, etc.) [5, 6, 7].

70           Genomics approaches are particularly useful when dealing with complex traits, as these traits  
71 usually have a multi-genic nature and an important environmental influence. Thanks to these  
72 technological improvements it is now feasible for a small laboratory to generate in a short time span (e.g.,  
73 several months) enough molecular data to obtain a set of mapped quantitative trait loci (QTLs), even in a  
74 species lacking any previous genomic information [8]. Genomic tools are thus facilitating the detection of

75 QTLs and the identification of existing favorable alleles of small effect, which have frequently remained  
76 unnoticed and have not been included in the gene pool used for breeding [9, 10].

77 In this review, we present and discuss the most relevant advances in the development and  
78 application of genomic tools for plant breeding, in particular for complex traits. Firstly, we introduce the  
79 most relevant genomic tools and secondly we provide examples of the application of these tools to plant  
80 breeding. The objective is to provide modern breeders with an updated synthetic view of how genomics  
81 can effectively improve the efficiency of breeding programs.

82

## 83 **GENOMIC TOOLS AND RESOURCES FOR PLANT BREEDING**

84

### 85 **Genome and transcriptome sequencing**

86

87 The availability of the whole genome sequence of a crop is of great utility for plant breeding, but  
88 until recently it has been an unachievable goal for most crops. This privilege was restricted to a reduced  
89 number of model species with small genomes and to projects with an important investment in time and  
90 resources, but now has extended to an increasing number of crops. However, it is also true that for  
91 important cultivated species with large and complex genomes such as wheat, sugarcane, or coffee, the  
92 sequencing of the whole genome is very challenging and may take several years before a draft is  
93 completed. Given the high cost of whole genome sequencing, transcriptome sequencing has been a  
94 cheaper alternative. The cDNA sequences (expressed sequence tags, ESTs) provide relevant information  
95 about the genes expressed in a certain tissue or organ, at a given stage of development and under  
96 particular environmental conditions. ESTs sequencing projects do not provide information about non-  
97 coding sequences and, even using diverse libraries, it is difficult to identify all genes and transcripts  
98 variants. Despite these limitations, ESTs collections have been very useful for breeders.

99 The Sanger technology has been the predominant sequencing method for the past thirty years. It  
100 has been used to sequence several genomes as well as many transcriptomes. The first international  
101 collaborative project resulted in the whole genome sequence of the model plant *Arabidopsis thaliana*  
102 [11]. Since then, reference genomes of selected genotypes were completed in a limited number of crops  
103 such as rice [12], maize [13], sorghum [14], populus [15], grapevine [16], papaya [17], or soybean [18].  
104 The transcriptomes of most major crops, to a greater or lesser extent, were also sequenced. A global view

105 of the genomes and transcriptomes sequenced can be obtained from the Gene Index Project  
106 (<http://compbio.dfci.harvard.edu/tgi/plant.html>) or in the NCBI Unigene database  
107 (<http://www.ncbi.nlm.nih.gov/unigene>).

108 The field of genomics has changed with the arrival of NGS technologies [19]. These new  
109 technologies have reduced the cost of sequencing by more than one thousand times compared to Sanger  
110 technology, by avoiding time-consuming and tedious traditional cloning steps and making possible to  
111 perform millions of sequencing reactions in parallel (Table 1). Among the “second generation”  
112 technologies, the 454 (Roche, <http://www.454.com>) and Illumina (Illumina, <http://www.illumina.com>)  
113 platforms are already profusely used to sequence crop species. Others, like Solid (Applied Biosystems,  
114 <http://www.appliedbiosystems.com/technologies>), have been less exploited in plants. By using these  
115 NGS technologies, the sequencing capacity of laboratories is continuously increasing. For instance, one  
116 High-Seq 2000 Illumina Sequencer is able to generate 55 Gb per day, which is roughly eighteen times the  
117 size of the human genome. Moreover, new, “third generation” platforms are being developed and  
118 incorporated to sequencing projects, such as PacBio RS (Pacific Biosciences,  
119 <http://www.pacificbiosciences.com>), Helicos (Helicos, <http://www.helicosbio.com>), or Ion Torrent (Life  
120 Technologies, <http://www.iontorrent.com>). The sequence obtained by NGS are generally deposited in the  
121 NCBI Sequence Read Archive (<http://www.ncbi.nlm.nih.gov/unigene>).

122 Nowadays, it is feasible to sequence most crop genomes (excluding those with a very large and  
123 complex genome) with a relatively modest budget, by combining Sanger with NGS technologies.  
124 Sequencing projects for 135 plant genomes, most of them corresponding to cultivated species or wild  
125 relatives, have been completely sequenced (3), are being assembled (27) or are in progress (105), as  
126 reported at the NCBI database (<http://www.ncbi.nlm.nih.gov/genomes/leuks.cgi>). Other databases  
127 including plant genome sequences are Plant GDB (<http://www.plantgdb.org>) and Phytozome  
128 (<http://www.phytozome.net>). A fully sequenced and well annotated genome provides useful tools for the  
129 breeders, as it allows the discovery of genes, determining their position and function, as well as the  
130 development of large marker collections and high resolution maps. In the cases where only a draft  
131 sequence is available, its usefulness depends on the quality of the assembly. For instance on many  
132 occasions thousands of scaffolds are obtained but they are not anchored to the genetic map, which  
133 difficults the use of the information. Many transcriptomes have also been sequenced, a number of them in  
134 species for which no previous sequence information was available. Sweet potato [20], squash [21],

135 pigeonpea [22], or buckwheat [23] represent some examples published in the last months. These assays  
136 are showing the great complexity of plant transcriptomes, allowing the identification of rare transcript  
137 variants that are being used to improve gene annotation and our knowledge of gene function and  
138 regulation.

139

## 140 **Bioinformatics**

141

142 NGS technologies are facilitating sequencing projects, but have brought new challenges, as  
143 millions of short DNA reads have to be analysed and assembled [19]. Also, genetic maps, genotypes, or  
144 expression information at a genomic scale have to be processed in order to obtain the relevant biological  
145 information. Therefore, it is necessary to develop new bioinformatics tools (algorithms and software),  
146 which allow the analyses of huge amounts of genome-wide data, but it is also necessary to change the  
147 approaches used to understand complex biological traits [25, 26].

148 The field of sequence analysis has a long tradition and has enabled the assembly of many  
149 genome sequences obtained by Sanger sequencing. Nowadays, the huge amount of sequence reads  
150 generated by NGS and the low quality of individual reads requires new software tools and algorithms that  
151 allow dealing with these data in an efficient way [27]. We consider this to be a limiting factor for this kind  
152 of analyses right now. Although in the last years great advancements in the sequence processing  
153 algorithms have been achieved, the software currently available requires improvements in many cases.

154 Two of the most common analyses carried out on these NGS reads are genome assembly and  
155 annotation and mapping. Genome assembly is a complex task requiring powerful computers and skilled  
156 bioinformaticians [25]. In particular, the RAM memory requirements of the computers used to assemble  
157 eukariotic genomes could hinder the application of this technique by small laboratories. Some of the most  
158 commonly used assemblers are Roche's 454 Gsassembler, Celera Assembler, and Mira. Once a reference  
159 genome is available in the species it is common to study its variation [19]. To do this, a mapper software  
160 is commonly used instead of an assembler software. A mapper tries to align every read against the  
161 reference genome. This process is much simpler and faster than the assembly. In this case the computer  
162 requirements are usually less demanding and the limiting factor could be the storage capacity. Some  
163 commonly used mappers are Bowtie, BWA, and TopHat. Once the reads are aligned, single nucleotide  
164 polymorphism (SNPs) can be detected by using the SAMtools or the GigaBayes SNP callers [28].

165           The open source software mainly used by the bioinformaticians is cumbersome for users not well  
166           versed in the Unix command line operating system. Some commercial proprietary solutions easier to use  
167           have been developed (e.g., LaserGene or CLC Genomics Workbench), but they have not been widely  
168           embraced by the breeders. Galaxy is an open source project devoted to create an easy to use web interface  
169           to the open source CLI based applications used in this area.

170           An important amount of work has been devoted in this field to the creation of standard and open  
171           file formats capable of storing information regarding sequence alignment and modelling (SAM) [24],  
172           SNP calls using variant call format (VCF; [http://1000genomes.org/wiki/doku.php?id=1000\\_genomes:](http://1000genomes.org/wiki/doku.php?id=1000_genomes:analysis:vcf4.0)  
173           analysis:vcf4.0), genomic regions with browser extensible data (BED, [http://genome.ucsc.edu/FAQ/](http://genome.ucsc.edu/FAQ/FAQformat#format1)  
174           FAQformat# format1) and genomic annotations using the general feature format (GFF;  
175           <http://www.sequenceontology.org/resources/gff3.html>). These open and standard formats allow the  
176           interoperability of the different software tools that are being actively developed and used. In addition, the  
177           computer requirements might be strong as some analyses require a large amount of RAM memory or  
178           storage capability.

179           The algorithms and methods used to store and process raw genomic data generated by the  
180           different technological platforms will depend on the type of data being processed and on the result  
181           expected. In any case, once the relevant information is obtained by the bioinformaticians, results have to  
182           be made available to the breeders by using an interface as easy and friendly as possible [25]. To provide  
183           access to this information, the generation of an easily browseable web site is a common and usually  
184           successful approach. Several general purpose web databases exist to make the relevant biological  
185           information available to the researchers and breeders (Table 2), like GenBank  
186           (<http://www.ncbi.nlm.nih.gov/genbank/>), EBML (<http://www.ebi.ac.uk/embl/>), DDBJ  
187           (<http://www.ddbj.nig.ac.jp/>) and Swiss-prot (<http://expasy.org/sprot/>). These latter databases are devoted  
188           to store information about any species, but several other more specific databases focused on species of  
189           interest to the breeders also exist, like the Sgn (<http://solgenomics.net/>), Phytozome  
190           (<http://www.phytozome.net/>), Gramene (<http://www.gramene.org/>) or CropNet (<http://ukcrop.net/>), which  
191           hold information that could have more specific use for breeding programs.

192

193   **Expression studies, from microarrays to RNA-seq**

194



195 New genomic tools are also of interest to expand and accelerate gene expression studies. The  
196 analysis of gene expression has provided a rich source of biological information, which allows breeders to  
197 understand the molecular basis of complex plant processes, leading to the identification of new targets for  
198 manipulating these processes.

199 Gene expression studies were at first based on the classical Northern blot method that only  
200 allowed the quantification of tens of genes simultaneously. The QRT-PCR is a more affordable and  
201 quantitative technique; but the number of genes analyzed by experiment is also limited [29]. Other  
202 approaches allowing the study of thousands of genes were differential display and cDNA amplified  
203 fragment length polymorphisms (cDNA-AFLPs) [30]. However, these methods are not really quantitative  
204 and are limited by the ability of the developed libraries to capture low-abundance transcripts. Other  
205 methods that overcome part of these problems are the serial analysis of gene expression (SAGE) [31] and  
206 massively parallel signature sequencing (MPSS) [32]. Nevertheless, the most employed methods at  
207 present to analyze transcript profiling are the hybridization-based platforms or microarrays [33].

208 Expression arrays have several advantages when compared with other methods. They can measure tens of  
209 thousands of different transcripts in the same reaction, they are semi-quantitative and sensitive to low-  
210 abundance transcripts if those are represented in a given array.

211 There are several web resources that facilitate microarray data analysis (e.g.,  
212 <http://babelomics.bioinfo.cipf.es/>) [34] or even web pages where the breeder can download experiments  
213 already performed and analyzed. There are also software packages specialized in microarrays analysis  
214 as the Bioconductor (<http://www.bioconductor.org/help/workflows/oligo-arrays/>) or MeV  
215 (<http://www.tm4.org/mev/>) [35]. Probably one of the most useful database is Genevestigator  
216 ([https://www.genevestigator.com/gv/doc/plant\\_biotech.jsp](https://www.genevestigator.com/gv/doc/plant_biotech.jsp)) [36], which contains microarray data from  
217 different species. The most extensive data are from the model species *A. thaliana* [37], but an increasing  
218 number of studies in crops like maize, wheat, rice, barley, or soybean are already available. All published  
219 expression data are public and disposables in databases as GEO (<http://www.ncbi.nlm.nih.gov/geo/>) [38],  
220 ArrayExpress (<http://www.ebi.ac.uk/arrayexpress/>) [39] or species specific databases. These data can be  
221 really useful when analyzing gene expression in these species or other crops [40].

222 Microarrays make use of the existing EST collections and genome sequence data. The vast  
223 increase provided by NGS in the number of sequences opens the possibilities of expression studies in a  
224 large number of species lacking previous sequence information. Also, deep NGS sequencing of RNA

225 transcripts (RNA-seq) is emerging as an alternative to microarray studies to quantify gene expression [41,  
226 42]. RNA-seq does not depend on genome annotation or on the probes contained in the array platform.  
227 This technology is also very useful to improve genome annotation, improving the detection of rare  
228 transcripts and splicing variants and the mapping of exon/intron boundaries. Moreover, RNA-seq avoids  
229 bias introduced during hybridization of microarrays and saturation level problems, has a greater  
230 sensibility, and shows high reproducibility [41, 43]. This approach has been already used in different  
231 crops with different breeding objectives, leading to the identification of genes involved in several  
232 metabolic pathways, disease response, fruit development, etc. [44, 45, 46, 47]. All these studies show the  
233 potential of RNA-seq for complex traits breeding. However, RNA-seq is an emerging technology and the  
234 methods used to analyze this kind of data are still being developed.

235

### 236 **Mutant and germplasm collections in the genomics era: TILLING and EcoTILLING**

237

238 Plant breeding requires genetic variability to be selected in order to increase the frequencies of  
239 favourable alleles and genetic combinations. Sources of natural genetic variability can be found within the  
240 crop, mostly in the form of landraces, and also in the wild relatives. Although many landraces have been  
241 substituted by modern and uniform cultivars and genetic erosion has taken place in wild materials, gene  
242 banks preserve many of these materials, which constitute an important reservoir of genetic variation  
243 useful for breeding [48].

244 Another important source of variability corresponds to the artificial mutant collections developed  
245 in several crop species. These artificial collections are created by radiation, chemical mutagenesis, or  
246 transgenic and insertion methods. Artificial mutations, mostly obtained by radiation and chemical  
247 methods, were used in plant breeding in the pre-genomics era, but new technologies are allowing the  
248 development of other types of collections [49]. For instance, the transferred DNA tagged lines and  
249 transposon tagged lines have been used to develop mutant collections in several species such as the model  
250 plant *Arabidopsis* (The Arabidopsis Information Resource; <http://www.arabidopsis.org>) or rice  
251 (International Rice Functional Consortium; <http://irfgc.irri.org>). Gene silencing technologies, using RNA  
252 interference, have also been used to create gene specific mutant collections in *Arabidopsis*, like the  
253 AGRİKOLA project (<http://www.agrikola.org>). The artificial mutant collections frequently contain

254 variability not present in the natural collections, and so are also very useful for the study and development  
255 of new traits.

256 In order to facilitate the identification of the accessions of interest in these collections, a genetic  
257 reverse approach has been used. Targeting Induced Local Lesions in Genomes (TILLING) [50] is able to  
258 identify all allelic variants of a DNA region present in an artificial mutant collection. A similar procedure  
259 called ecotype TILLING (EcoTILLING) [51] can be used to identify allelic variants for targeting genes in  
260 natural collections. These two methods are based on the use of endonucleases, such as CEL I or Endo I,  
261 that recognize and cut mismatches in the double helix of DNA [52, 53]. Since the TILLING and  
262 EcoTILLING techniques identify all allelic variants for a certain genomic region, the phenotypic  
263 characterization effort can be concentrated in a reduced number of accessions with different variants.  
264 Obviously, the success of the identification of variation useful for breeding programmes will depend on  
265 the right selection of target genes. The availability of sequences coming from NGS sequencing projects  
266 and the information provided by gene expression studies is significantly increasing the number and  
267 quality of candidates for TILLING and EcoTILLING studies

268 These procedures have been successfully used in many crops [54]. For example, TILLING has  
269 been applied to *Arabidopsis* [55], *Lotus* [56], barley [57], maize [58], pea [59], and melon [60]. Rice was  
270 the first crop for which EcoTILLING was applied [61]. Subsequently, EcoTILLING has been used in  
271 other crops and wild relatives, like barley [62], wheat [63], or the wild peanut *Arachis duranensis* [64],  
272 using both genebank collections and natural populations [65]. These assays used gene targets involved in  
273 different processes. Many studies have been focused on detecting allelic variants in genes most related to  
274 organoleptic quality [66, 67] or disease resistance [68, 69].

275

## 276 **Re-sequencing for SNPs discovery and use in genotyping platforms**

277

278 One of the most interesting applications of NGS for plant breeders is the discovery of genetic  
279 variation. Now it is possible to sequence rapidly multiple individuals within a species with limited  
280 technical expertise and at minimal cost. The parallel development of computational pipeline tools is  
281 greatly accelerating the accurate mining of these sequences for genetic variants that can be converted into  
282 genetic markers, mainly microsatellites or simple sequence repeats (SSRs) and SNPs [70]. SSRs and  
283 SNPs are now the predominant markers in plant genetic analysis. SNPs are more abundant, stable,

284 amenable to automation, and increasingly cost-effective, thus are fast becoming the marker system of  
285 choice in modern genomics research [71].

286         The genome-wide SNPs discovery by massive re-sequencing has been performed in model  
287 species with small genomes, such as *Arabidopsis thaliana*, where the 1001 Genomes project  
288 (<http://www.1001genomes.org>) [72] is attempting to unveil the whole-genome sequence variation in this  
289 reference plant. Similar re-sequencing efforts are becoming possible in rice, maize, grape, soybean,  
290 poplar etc. by sequencing sets of related genotypes, individually or pooled, within each species (elite  
291 cultivars, breeding lines, ecotypes, landraces, and/or weedy and wild relatives of a crop) [73, 74, 75, 76].  
292 The higher complexity in architecture and repeat content of these genomes has made necessary the use of  
293 approaches for genomic complexity reduction that also reduce sequencing cost [77]. Identification of  
294 SNPs is also very challenging in species with high levels of heterozygosity and/or with complex ploidy  
295 levels, as pseudo-SNPs are identified by misassembly of paralogs [78, 79, 80].

296         Both Roche 454 and Illumina GA have been mostly used for genome re-sequencing. The  
297 alignment difficulties often associated to the use of short Illumina GA reads (Table 1) are less problematic  
298 in species for which available reference genomes facilitates SNPs calling and genome positioning of  
299 genetic variation [81]. For most of these species, limited collections of SSRs and SNPs were available  
300 from early re-sequencing efforts, previous to the advent of NGS, but new genome-wide re-sequencing is  
301 enlarging the SNP pools and making them more representative of the range of natural variation.

302         For an increasing number of species with high societal or economic value NGS genome re-  
303 sequencing is providing the first SSRs and SNPs resources. Examples are the grain amaranths  
304 (*Amaranthus* sp.), important pseudocereals, appreciated for their nutritional quality and tolerance to  
305 abiotic stress [82], for which no prior genome information existed. In these cases the combination of  
306 several sequencing techniques, and the use of paired-end sequencing facilitates SNP discovery. Roche  
307 454 and Illumina GA were combined for high-throughput SNP discovery in common bean [83] and also  
308 Solid was used to sequence diploid wheat species, which are donors of the subgenomes of modern  
309 hexaploid bread wheat [84].

310         Most of the effort in species lacking genomic resources is being made through transcriptome re-  
311 sequencing, as an alternative way of genome complexity reduction. Targeted amplicon re-sequencing is  
312 another strategy for discovering SNPs in specific genes [78].

313 One of the first examples of deep transcriptome sequencing was a study with two maize inbred  
314 lines [85]. This first study was followed by a large and rapidly increasing number of projects using non-  
315 model crops, some of them with large, complex, polyploid, and uncharacterized genomes, including  
316 forest trees, like *Eucalyptus* [86], oak [87], several polyploid crops, like oilseed rape [79], oats [80],  
317 coffee [88], and sweet potato [89], non-model grain legumes as chickpea and chickling pea [90], tomato  
318 [91], or several cucurbits, including *Cucurbita* spp. [21], cucumber [92], and melon [93].

319 These studies employ normalized/non-normalized cDNA libraries generated from single or  
320 multi-tissues samples, and derived from single or pooled genotypes, combined or not with multiplex  
321 identifier barcodes that allow allele origin identification. Sequencing is often focused on selected  
322 genotypes subjected to specific conditions, to detect SNPs in candidate genes involved in complex  
323 biological processes of interest to plant breeders. For example, the transcriptomes of two resistant and one  
324 susceptible lines of water yam, a major staple crop in Africa, under anthracnose infection, were  
325 successfully sequenced detecting SNPs in genes putatively involved in pathogen response [94]. Also, two  
326 alfalfa genotypes contrasting for cellulose and lignin content were sequenced, which allowed selecting  
327 SNPs useful to improve alfalfa as a forage crop and cellulosic feedstock [95].

328 The use of genome and transcriptome sequencing for SNP discovery has resulted in large SNPs  
329 collections in most of the crops. These large collections are being validated and applied for different  
330 purposes such as map construction, map saturation, genome-wide diversity studies, association mapping  
331 etc (Table 3). Some of the most important achievements will be described in later sections.

332 There are many SNPs genotyping techniques, which are more or less appropriate for different  
333 scales of individuals/SNPs to be genotyped [107]. The implementation of marker-assisted breeding  
334 strategies often requires the generation of thousands of genotypes per population. Thus, one practical way  
335 of optimizing the use of these large SNPs collections is using them with cost-effective platforms for  
336 medium to high density genotyping. A large number of commercial platforms are available for  
337 semiautomated or fully automated SNP genotyping [108, 109]. Genotyping assays usually require a  
338 previous process of selection of a set of SNPs, among the hundreds/tens of thousands detected, that are  
339 appropriate for the assay objectives.

340 The Illumina GoldenGate assays have been the most widely used for mid-throughput  
341 applications. SNPs platforms with 384, 768, or 1536 SNPs are available for a number of species (Table  
342 3). Popularity is also increasing for the Sequenom Mass array and the KASPar genotyping chemistry [82,

343 110]. Expanded arrays with tens of thousands SNPs for high-throughput applications have been also  
344 developed with the Infinium technology in maize, grape, tomato, pine, and poplar and are under  
345 development in soybean and several *Rosaceae* crops (apple, peach, and cherry) [74, 111].

346

### 347 **Construction of high density genetic maps**

348

349 One of the main applications of genomic advances is the development of high density genetic  
350 maps. The high-density map construction involves the location of hundreds or even thousand markers in  
351 the different linkage groups. In these maps the coverage should be very high and no large gaps must be  
352 present. NGS technologies and high-throughput genotyping platforms have allowed the improvement of  
353 genetic maps by increasing markers density. Several works include the integration of new markers,  
354 basically SNPs derived from re-sequencing studies, into previously developed genetic maps, both in  
355 diploid and polyploidy species [80, 112]. Golden Gate has been the most widely used platform. It has  
356 been estimated that this genotyping platform is 100-fold faster than gel-based methods for increasing 2-3  
357 times maize map density [101]. Also Sequenom-based SNP-typing assays are starting to be applied. In a  
358 recent study, a total of 1.016 SNPs, identified *via* comparative next-generation transcriptomic sequencing,  
359 were successfully mapped by genotyping 297 maize recombinant inbred lines (RILs) [113]. Other  
360 genotyping strategies based on arrays hybridization, such as the single-feature polymorphisms (SFP),  
361 variants detected by a single probe in oligonucleotide arrays, are speeding up genetic map construction.  
362 This technique has been used for the construction of a high-density linkage map in species poorly  
363 characterized, like *Eucalyptus* [114]. The newly developed maps, enriched in sequence-based markers are  
364 facilitating comparative mapping. Recent examples are high density SNPs maps of barley compared with  
365 wheat and rice [98, 115].

366 The decrease of sequencing costs is also allowing the detection of new types of genetic markers  
367 useful for increasing the density of genetic maps. In this respect, restriction-site associated DNA (RAD) is  
368 a kind of marker which detects genetic variation adjacent to restriction enzyme cleavage sites across a  
369 target genome. These markers are produced after NGS sequencing of genomic libraries obtained after  
370 digestion with different restrictases. As an example of the utility of this technique, a total of 445 RAD  
371 markers distributed across all seven barley chromosomes were located, which was very useful for linkage  
372 map construction in this crop [116].

373           The markers derived from NGS can also be useful to position sequence scaffolds onto physical  
374 maps. In this respect, a new method combining deep sequencing and the bin mapping strategy has been  
375 described [117]. The SNPs identified by re-sequencing genomic libraries from selected progeny  
376 individuals, derived from a cross between two closely related diploid strawberry species, were used to  
377 anchor 92.8% of the *Fragaria* genome to the genetic map. Results highlighted the potential of this  
378 methodology to obtain a robust framework for the anchoring of the genome sequence without the  
379 requirement of a high density physical mapping or a well-saturated genetic map.

380           Whole-genome re-sequencing at different coverage levels is being increasingly applied for map  
381 construction using different strategies. As an example, a genetic map for rice has been constructed using  
382 whole genome re-sequencing of 150 RILs [118]. These authors concluded that the sequencing-based  
383 method was approximately 35 times more precise in recombination breakpoint determination than PCR-  
384 based markers maps. Also, the whole genome of 128 chromosome segment substitution lines (CSSLs) of  
385 rice was re-sequenced and used it for the construction of an ultra-high quality physical map in this crop  
386 [119]. Based on low coverage re-sequencing, a new mapping strategy that allows inferring the parental  
387 genotypes of the assayed RILs population has been proposed [120]. An ultra-high density linkage map  
388 was obtained with this method and the quality of the map was evaluated by using it to identify a QTL  
389 controlling grain width. Further applications of new sequence-based denser genetic maps to QTL  
390 discovery and marker assisted selection (MAS) will be discussed later.

391

## 392 **TOWARDS A GENOMICS-BASED PLANT BREEDING**

393

### 394 **Genome-wide genetic diversity studies**

395

396           One of the main challenges in agricultural genetics is to access and use the tremendous genetic  
397 variation present in germplasm collections and in the wild relatives. A significant part of this variation  
398 remains untapped because of the difficulties in effectively identifying genetic differences in large  
399 collections. Some traits, with high heritability and of simple characterization, are easy to select for.  
400 However, desirable allelic variants and genetic combinations for complex traits are difficult to identify.  
401 Recent advances in genotyping are enabling genome-wide diversity studies capable of better capturing the  
402 spectrum of variability in natural and breeding populations.

403 Most of the mid- to high-throughput genotyping platforms described above are being used for  
404 studies on diversity and population structure in the corresponding crops (Table 3). By using representative  
405 diversity panels, polymorphism rates for individual SNP markers, minor allele frequencies (MAFs), etc.  
406 are estimated, facilitating the selection of those SNPs with biological interest and highly polymorphic in  
407 the different groups. For example, the Infinium arrays developed in some of these crops are being used to  
408 create haplotype maps for vast germplasm collections, such as the 18,000 accessions of the USDA  
409 soybean germplasm collection [121].

410 Haplotype maps (hapmap) of entire collections are useful to identify rare, potentially valuable,  
411 alleles. Hapmap projects are undergoing in a number of species such as the “*rice diversity project*”  
412 (<http://www.ricehapmap.org/index.aspx>) aimed to develop a 10,000 SNP chip for rice and create a  
413 haplotype map to document the differences in allelic variation within and between the different  
414 subpopulations of *O. sativa* and its progenitor *O. rufipogon*. Large-scale genetic diversity studies have  
415 also been accomplished in maize. Gore et al. [122] identified and genotyped several million sequence  
416 polymorphisms among 27 diverse maize inbred lines. This study allowed the discovery of regions with  
417 highly suppressed recombination that appear to have influenced the effectiveness of selection during  
418 maize inbred development and may be a major component of heterosis. Also, highly differentiated  
419 regions were found that probably contained *loci* that are key to geographic adaptation. Also in legumes,  
420 the Medicago HapMap Project, that consist in the sequencing the whole-genomes of 30 inbred lines, will  
421 explore the genetic basis of symbiosis, creating a robust platform for genome-scale association mapping  
422 [123].

423 The diversity panels can include representatives of close or more distantly related species to  
424 check if these sequence-based SNP assays also work in related species [74, 82]. Sometimes sets of SNPs  
425 specifically developed for detecting genetic diversity among closely related cultivars are used in  
426 genotyping platforms. For example, despite the large amounts of SNPs available in rice obtained from the  
427 comparison of the two reference genomic sequences (one *japonica* and one *indica* variety) [124],  
428 extremely low levels of DNA polymorphism were detected within *japonica* cultivars. A whole-genome  
429 sequencing of an elite Japanese rice cultivar, closely related to the reference *japonica* variety, has been  
430 conducted and the SNP information obtained by comparison of the two *japonica* sequences was applied  
431 to develop a high-throughput genotyping array used for genotyping a set of representative Japanese  
432 cultivars [125]. These experiments are useful for understanding the role of selection and breeding in the



433 distribution of genetic variation across the crop genome. In fact, this assay led to the identification of  
434 several haplotype blocks which are inherited from traditional landraces to current improved varieties.  
435 Moreover, it was found that, as predicted, modern breeding practices have generally decreased genetic  
436 diversity. On the practical level, the distribution of genetic diversity in modern cultivars plays an  
437 important role in the choice of specific mapping and crop improvement strategies.

438         Genome-wide survey of genetic diversity is useful to elucidate the causative genetic differences  
439 that give rise to observed phenotypic variation providing a foundation for dissecting complex traits  
440 through genome-wide association studies. However, its utility is limited if phenotypic data are not  
441 available. Not just genomics and transcriptomics, but the other 'omics' disciplines, like proteomics and  
442 metabolomics, are useful to understand how the changes in the genotype lead to differences in the final  
443 phenotype. Phenomics, which uses high-throughput technologies to characterize germplasm, is being  
444 developed and will help to deal with this issue [126].

445

#### 446 **Identification of molecular markers linked to single genes and QTLs**

447

448         NGS and high-resolution maps have led to a significant improvement in the identification of  
449 molecular markers linked to specific genes and to QTLs. The most important advantage comes from the  
450 dense genome coverage, which allows the identification of markers closely linked to any target genomic  
451 region, with the advantages that this tight linkage provides.

452         Methods already used in the pre-genomics era to facilitate the identification of markers linked to  
453 single loci, such as bulked segregant analysis (BSA), are now optimized. For example, a GoldenGate  
454 assay has been combined with BSA to significantly accelerate mapping of the dominant resistance *locus*  
455 to soybean rust *Rpp3* [127]. In this respect, there is an increasing number of reports on exploitation of  
456 NGS technologies to identify molecular markers tightly linked to major genes. For example, a fine  
457 genetic mapping of the single dominant scab resistance gene (*Ccu*) in RILs of cucumber (*Cucumis*  
458 *sativus*) has been conducted [128]. The resistant cucumber genome was sequenced with Solexa/Illumina  
459 NGS and compared with the susceptible cucumber genome. As a result, three insertion/deletion (indel)  
460 markers closely linked to the *Ccu locus* were obtained. A detailed study of the region delimited by  
461 markers revealed four resistance gene analogs as possible candidates for *Ccu*.

462 QTL detection has traditionally been conducted by linkage mapping. NGS technologies are  
463 significantly contributing to increase accuracy in detection of QTLs. They allow increases in many orders  
464 of magnitude of the number of markers mapped, ensuring high mapping resolution, and also aid in the  
465 development of mapping populations, such as RILs, near isogenic lines (NILs), and CSSLs, more  
466 appropriated for QTLs detection. These populations have conventionally been constructed and genotyped  
467 using a limited number of markers.

468 There are increasing reports describing accurate QTLs mapping with different NGS or high-  
469 throughput genotyping strategies. For example, a high density rice map constructed by whole-genome re-  
470 sequencing of a RILs population, was used to identify four QTLs controlling plant height [90]. On a  
471 different study [129] an ultra-high density genetic map based on SNPs, obtained with Illumina GA, was  
472 compared with a linkage map based on RFLPs/SSRs in rice. The positions of several cloned genes, two  
473 major QTLs for grain length and grain width, and a QTL for pigmentation were evaluated in a RIL  
474 population, arising the expected result that the SNPs map detected more QTLs and more accurately than a  
475 RFLPs/SSRs based linkage map.

476 QTL detection based on the linkage analysis method has the disadvantage that the number of  
477 recombination events is limited to the generations needed to develop the mapping population. Association  
478 mapping or linkage disequilibrium (LD) mapping is a new powerful approach to map complex traits. This  
479 method identifies genetic *loci* associated with phenotypic trait variation in a collection of individuals.  
480 Association mapping uses the natural diversity, which represents many more recombination events  
481 occurred in the history of the population, providing better resolution. Nowadays, two association mapping  
482 methodologies are in use: candidate gene association, where a good understanding of the biochemistry  
483 and genetics of the trait is needed, and whole genome scan, also called genome-wide association (GWA)  
484 studies. New genomic advances are providing the higher density of genetic markers required to ensure  
485 enough coverage to detect linkage between markers and a causal *locus*. Also the decrease of sequencing  
486 costs (Table 1) has allowed the use of whole genome sequencing in current studies [130].

487 Nevertheless, association mapping is just rising in model species and major crops. Maize is the  
488 most widely studied crop regarding association analysis. Many candidate genes have been successfully  
489 associated to morphological or quality traits. As an example, candidate genes *Dwarf8*, *Vgt1* and  
490 *ZmRap2.7* were successfully associated to flowering time [131]. Other candidate genes have been  
491 associated, among others, to forage quality, carotenoid content, oil content and kernel quality [132, 133,

492 134, 135]. GWA studies have been more limited, probably due to the large genome of maize (2300 Mbp)  
493 and the great number of markers needed to cover it. The first study identified a fatty acid desaturase gene  
494 (*fad2*) associated with increased oleic acid levels [99]. More recently, other authors found 32 QTLs  
495 associated with southern leaf blight disease resistance [100].

496         Examples of association mapping approaches in other crops are more limited. Studies based on  
497 the candidate gene approach have been reported in some crops, like grape, or conifers [102, 106].  
498 However, GWA studies have only been developed either in the model species *A. thaliana* [136] or in  
499 major crops such as rice [96], barley [97], or wheat [104]. Some articles also describe successful mapping  
500 processes combining classical linkage mapping and association mapping [137]. Although genetic  
501 association mapping is in its early steps, it is a promising tool for the dissection of complex traits in crop  
502 plants.

503

#### 504 **Marker assisted selection**

505

##### 506 *Marker assisted backcross selection*

507

508         Marker assisted selection (MAS) is an indirect process where selection is carried out on the basis  
509 of a marker instead of the trait itself. The successful application of MAS relies on the tight association  
510 between the marker and the major gene or QTL responsible for the trait. As we have described before, the  
511 new genomic tools accelerate the identification of markers tightly linked to target genomic regions. On  
512 the other hand, the new dense genotyping platforms available today accelerate the genotyping of large  
513 amounts of samples during the MAS process in a rapid and economically feasible manner. MAS can take  
514 benefit from these technologies, speeding up the release of new varieties.

515         In spite of the close linkage between the marker and the gene, the possibility of recombination  
516 limits the use of MAS. The use of intragenic markers, also called functional markers, can help to  
517 overcome this limitation [138]. NGS sequencing projects produce large collections of functional markers.  
518 These markers enhance real gene assisted breeding, reducing the possibility of losing the desirable trait  
519 due to recombination. This is today feasible in many crop species in which NGS cDNA sequencing is  
520 being conducted. Some of these studies perform expression profiling, identifying candidates and  
521 associated gene targeted markers.

522 MAS is also frequently applied to perform background selection in the context of backcrossing  
523 programmes. Background selection consists in the identification of plants with lower contents in donor  
524 genome to continue the breeding scheme, in order to achieve the recovery of the recipient genome. The  
525 use of background markers facilitates the quick recovery of the recurrent parent genome [139].  
526 Background selection is being used extensively in rice breeding. High-density genome maps are being  
527 effectively used in background analysis. For example, background selection integrated with foreground  
528 selection of bacterial blight resistance (*xa13* and *Xa21* genes), amylose content (*waxy* gene) and fertility  
529 restorer gene has been performed in order to identify superior lines with maximum recovery of Basmati  
530 rice genome along with the quality traits and minimum non-targeted genomic introgressions of the donor  
531 chromosomes [140].

532 In some cases, the problem of recovering the genetic background of the recurrent parent arises  
533 because of the linkage drag, that is, the introgression of chromosome regions with deleterious effects  
534 which are tightly linked to the gene or QTL of interest. The detection of QTLs responsible of the negative  
535 effects and the localization of molecular markers tightly associated to them can be an efficient way to  
536 break the genetic drag. A well known example concerns canola (rapeseed) breeding, which began with  
537 the discovery of germplasm with low erucic acid content in seeds of a spring forage cultivar in the 1950's.  
538 The problem arose because a high association between low erucic acid content and low seed oil content  
539 exists. The recent availability of high-density molecular maps has allowed the detection of several QTLs  
540 associated to both traits. Moreover, the identification of molecular markers very tightly linked to the  
541 QTLs made possible to break the linkage drag between the low oil content and erucic acid concentration  
542 in seeds in the process for breeding new high seed oil content canola cultivars [141].

543 Frequently, current breeding programmes involve the introgression of more than one gene or  
544 QTL controlling traits of interest into one genetic background, in a process that is called pyramiding. The  
545 most useful application of MAS in the process of pyramiding is related to the introgression of different  
546 genes or QTLs whose effect on the phenotype is undistinguishable. The accumulation of genes from  
547 different sources which confer resistance against the same disease is an example, and is indeed one of the  
548 most widespread applications of gene pyramiding [142]. The main advantages of recent advances in plant  
549 genomics incorporated into gene pyramiding will be related to two different aspects. On one hand, the  
550 number of plants to be analyzed in a gene pyramiding programme must be increased as the number of *loci*  
551 of interest is higher, to ensure with a reasonable likelihood that the genotype combining favorable alleles

552 is present in the population [143]. In this sense, the availability of genotyping platforms will provide the  
553 possibility to screen larger generations. On the other hand, the efficiency of the process strongly depends  
554 on the tightness of the linkage between markers used and the target genes or QTLs. Again, identification  
555 of functional markers will circumvent this limitation.

556

557 *'Breeding by design'*

558

559 The possibility to predict the outcome of a set of crosses on the basis of molecular markers  
560 information is known as 'breeding by design' [6]. The process includes three steps: mapping *loci*  
561 involved in all agronomically relevant traits, assessment of the allelic variation at those *loci*, and, finally,  
562 breeding by design. In the method as initially described by Peleman and van der Voort [6], the first step  
563 was proposed to be completed by either using mapping populations segregating for the trait of interest or  
564 based on a candidate gene approach (mainly exploiting information from model plant species and  
565 increasing understanding of gene function). Also linkage disequilibrium (LD) mapping was suggested,  
566 focused on the region previously identified as related to the trait ('targeted LD mapping'). Currently, as  
567 previously discussed, other possibilities such as GWA studies allow a more efficient way to accomplish  
568 this first step, avoiding limitations of biparental populations. The second step of the process consists in  
569 the identification of allelic variation for the *locus* of interest and the assignation of the phenotypic value to  
570 each of them. This step cannot be based on biparental populations, given that only two alleles per *locus*  
571 are segregating in this case. The analysis should then include plant materials representing the variability  
572 of the species. Genotypic and phenotypic data for each plant are required.

573

574 As previously stated, high level of saturation with markers is not the limiting factor in most  
575 cases, and so currently the restrictions mostly come from the phenotyping step. Strictly speaking,  
576 'breeding by design' exploits information obtained in the previous steps: once the *loci* of interest have  
577 been mapped, and the contribution of each allelic variant has been determined, crosses can be established  
578 to generate superior genotypes which combine all favourable alleles. Application of this breeding strategy  
579 has been used for different crops and with different objectives, such as breeding for heading date in rice  
580 [144] or seed length in soybean [145]. This procedure has also been used in patent applications; as an  
581 example, 'breeding by design' has been reported as part of the development of higher quality maize  
varieties. However, the most effective application of the 'breeding by design' approach will come from

582 the incorporation of the most advanced genomic tools into the process, which will allow the improvement  
583 of the predictions.

584

585 *Genomic selection*

586

587 MAS strategies described so far require the identification of markers associated to the traits of  
588 interest. This represents one of the weaknesses of traditional MAS approaches [146]. Nevertheless, MAS  
589 can also be applied eluding this step, using an approach known as genomic (or genome-wide) selection  
590 (Figure 1). The method was first described in 2001 [147], as an attempt to exploit information generated  
591 from emerging genotyping technologies. Genomic selection is based on simultaneous estimation of  
592 effects on phenotype of all *loci*, haplotypes, and markers available. The difference with other MAS  
593 methods relies on the fact that no previous selection of markers with effects on phenotype is developed  
594 [148]. Genomic selection requires the availability of phenotypic and genotypic data for the reference  
595 population. This data set will allow estimating the parameters for the model, so that the differences at the  
596 phenotype level are explained by the markers analysed. Once the model is established, application to  
597 breeding populations makes possible to determine the genomic value of each individual, i.e., the expected  
598 phenotype based on the genotypic data. The requirement is the availability of enough molecular markers  
599 to provide good genome coverage [5, 146].

600 Simulation studies carried out using maize proved the usefulness of genomic selection applied to  
601 an initial cross between an adapted line and exotic germplasm. With 512 markers and a reference  
602 population of 288 F<sub>2</sub> plants evaluated in six different environments, it was possible to obtain good  
603 selection response after 7-8 generations. [149]. Also with maize, simulations showed that response to  
604 selection was 18 to 43% larger for genomic selection than for marker assisted recurrent selection [150].  
605 Response obtained when using genomic selection can be lower than response by phenotypic selection.  
606 However, the reduction in cycle length due to early MAS results in an increase of gain per time unit. This  
607 reduction is even more accused for species with a long generation interval, such as tree species [148].

608 The availability of phenotypic databases for different crops has allowed the comparison of  
609 predictions about the genotypic value obtained using genomic selection with the true genotypic value as  
610 shown by the phenotypic manifestation of the trait. In a study developed with phenotypic and genotypic  
611 data from *Arabidopsis*, maize and barley, results obtained were more accurate when genome-wide

612 selection was carried out, if compared with results derived of previous selection of markers with effects  
613 on the phenotype [151].

614           Although when applying genomic selection there is no need to previously identify QTLs  
615 controlling a certain trait, the utilization of this approach allows detecting the chromosome regions  
616 involved in a given trait, as markers with greater effect on the phenotype will indicate the presence of a  
617 QTL for this trait [152]. Some studies go one step farther and propose the application of MAS prior to  
618 phenotyping. This approach involves the use of *prior indices*, i.e., marker selection indices which have  
619 been constructed from a given phenotyped and genotyped population and are applied to different  
620 populations which have not been phenotyped [129]. The decrease in the costs of genotyping provides the  
621 appropriate scenario for this strategy to become cost-effective.

622           In any case, even the identification of the QTLs responsible for a certain trait does not imply the  
623 identification of the gene or genes controlling the trait itself or the understanding of the mode of action.  
624 Models applied in genomic selection are useful to predict breeding values and, in some cases, detect  
625 regions associated to a trait, but further work is necessary from this point to identify the gene or genes  
626 responsible for the phenotypic variability observed. From plant breeders' perspective, the availability of  
627 molecular markers which allow MAS to be applied is generally sufficient. However, development of the  
628 new high throughput -omics technologies has provided breeders with new strategies to search for  
629 candidate genes, mainly based on microarray for differential gene expression, being the possibility to  
630 explore more genes the most important advantage. Future exploitation of these strategies could facilitate  
631 the identification of candidate genes underlying the traits of interest and make MAS more efficient.

632

## 633 **CONCLUSIONS**

634

635           For some major crops the pace experimented for genetic gains in yield and other complex traits  
636 in the 20<sup>th</sup> century will be difficult to be maintained if only existing pre-genomics technologies are used  
637 [153]. However, plant breeding is a dynamic science and, fortunately, genomics resources and tools are  
638 already available and are helping to give another quantitative leap in plant breeding. In this respect many  
639 advances are already taking place, and the superdomestication, i.e., “the processes that lead to a  
640 domesticate with dramatically increased yield that could not be selected in natural environments from  
641 naturally occurring variation without recourse to new technologies” [10], will require the combination of

642 conventional breeding with crop genomics. Also, genomic tools and approaches will help conventional  
643 breeding in achieving important advances in the breeding of crops that from the point of view of genetic  
644 improvement have remained either orphaned or neglected [8]. Therefore, while conventional pre-  
645 genomics plant breeding has been, is, and will be successful at improving our crops, the application of  
646 genomic tools and resources to practical plant breeding will push forward the genetic gains obtained by  
647 breeding programmes. New genomic advances, many of which are already being developed, will make  
648 easier for breeders to obtain new cultivars with improved characteristics, either by facilitating selection or  
649 by improving the variation available for breeders by using precision breeding approaches. In particular,  
650 the present and new genomics tools are of great value for the genetic dissection and breeding of complex  
651 traits.

652

### 653 **ACKNOWLEDGEMENTS**

654

655 This contribution has been partially funded by the Ministerio de Ciencia y Tecnología  
656 (AGL2008-05114, AGL2009-07257, and PIM2010PKB-00691), and Instituto Nacional de Investigación  
657 y Tecnología Agraria y Agroalimentaria (RTA2008-00035-C02-02).

658

### 659 **ABBREVIATIONS**

660

661 AFLP = Amplified Fragment Length Polymorphism

662 BED = Browser Extensible Data

663 BSA = Bulked Segregant Analysis

664 CSSL = Chromosome Segment Substitution Line

665 DArT = Diversity Arrays Technology

666 EcoTILLING = Ecotype TILLING

667 EST = Expressed Sequence Tag

668 GFF = General Feature Format

669 GWA = Genome Wide Association

670 LD = Linkage Disequilibrium

671 MAF = Minor Allele Frequency



672 MAS = Marker Assisted Selection  
673 MPSS = Massively Parallel Signature Sequencing  
674 NGS = Next Generation Sequencing  
675 NIL = Near Isogenic Line  
676 QTL = Quantitative Trait Locus  
677 RAD = Restriction-Site Associated DNA  
678 RIL = Recombinant Inbred Line  
679 RNA-seq = Sequencing of RNA Transcripts  
680 SAGE = Serial Analysis of Gene Expression  
681 SAM = Sequence Alignment and Modelling  
682 SPF = Single-Feature Polymorphism  
683 SNP = Single Nucleotide Polymorphism  
684 SSR = Simple Sequence Repeat  
685 TILLING = Targeting Induced Local Lesions on Genomes  
686 VCF = Variant Call Format

687

## 688 REFERENCES

689

- 690 [1] Fedoroff, N.V. The past, present and future of crop genetic modification. *New Biotechnol.* **2010**,  
691 27, 461-465.
- 692 [2] Tester, M.; Langridge, P. Breeding technologies to increase crop production in a changing world.  
693 *Science* **2010**, 327, 818-822.
- 694 [3] Varshney, R.K.; Tuberosa, R. *Genomics-assisted crop improvement vol. 1: Genomics approaches*  
695 *and platforms*. Springer; New York. **2007a**.
- 696 [4] Varshney, R.K.; Tuberosa, R. *Genomics-assisted crop improvement vol. 2: Genomics applications*  
697 *in crops*. Springer, New York. **2007b**.
- 698 [5] Lorenz, A.J.; Chao, S.; Asoro, F.G.; Heffner, E.L.; Hayashi, T.; Iwata, H.; Smith, K.P.; Sorrells,  
699 M.K.; Jannink, J.L. Genomic selection in plant breeding: knowledge and prospects. *Adv. Agron.*  
700 **2011**, 110, 77-123.
- 701 [6] Peleman J.D.; van der Voort, J.R. Breeding by design. *Trends Plant Sci.* **2003**, 8, 330-334.

- 702 [7] Collard, B.C.Y.; Mackill, D.J. Marker-assisted selection: an approach for precision plant breeding  
703 in the twenty-first century. *Phil. Transact. Royal Soc. B* **2008**, *363*, 557-572.
- 704 [8] Varshney, R.K.; Glaszmann, J.C.; Leung, H.; Ribaut, J.M. More genomic resources for less-  
705 studied crops. *Trends Biotechnol.* **2010**, *28*, 452-460.
- 706 [9] Morgante, M.; Salamini, F. From plant genomics to breeding practice. *Curr. Opinion Biotechnol.*  
707 **2003**, *14*, 214-219.
- 708 [10] Vaughan, D.A.; Balász, E.; Heslop-Harrison, J.S. From crop domestication to super-  
709 domestication. *Ann. Bot.* **2007**, *100*, 893-901.
- 710 [11] The Arabidopsis Genome Initiative. Analysis of the genome sequence of the flowering plant  
711 *Arabidopsis thaliana*. *Nature* **2010**, *408*, 796-815.
- 712 [12] International Rice Genome Sequencing Project. The map-based sequence of the rice genome.  
713 *Nature* **2005**, *436*, 793-800.
- 714 [13] Schnable, P.S.; Ware, D.; Fulton, R.S.; Stein, J.C.; Wei, F.; Pasternak, S.; Liang, C.; Zhang, J.;  
715 Fulton, L.; Graves, T.A.; Minx, P.; Reily, A.D.; Courtney, L.; Kruchowski, S.S.; Tomlinson, C.;  
716 Strong, C.; Delehaunty, K.; Fronick, C.; Courtney, B.; Rock, S.M.; Belter, E.; Du, F.; Kim, K.;  
717 Abbott, R.M.; Cotton, M.; Levy, A.; Marchetto, P.; Ochoa, K.; Jackson, S.M.; Gillam, B.; Chen,  
718 W.; Yan, L.; Higginbotham, J.; Cardenas, M.; Waligorski, J.; Applebaum, E.; Phelps, L.; Falcone,  
719 J.; Kanchi, K.; Thane, T.; Scimone, A.; Thane, N.; Henke, J.; Wang, T.; Ruppert, J.; Shah, N.;  
720 Rotter, K.; Hodges, J.; Ingenthron, E.; Cordes, M.; Kohlberg, S.; Sgro, J.; Delgado, B.; Mead, K.;  
721 Chinwalla, A.; Leonard, S.; Crouse, K.; Collura, K.; Kudrna, D.; Currie, J.; He, R.; Angelova, A.;  
722 Rajasekar, S.; Mueller, T.; Lomeli, R.; Scara, G.; Ko, A.; Delaney, K.; Wissotski, M.; Lopez, G.;  
723 Campos, D.; Braidotti, M.; Ashley, E.; Golser, W.; Kim, H.; Lee, S.; Lin, J.; Dujmic, Z.; Kim, W.;  
724 Talag, J.; Zuccolo, A.; Fan, C.; Sebastian, A.; Kramer, M.; Spiegel, L.; Nascimento, L.; Zutavern,  
725 T.; Miller, B.; Ambroise, C.; Muller, S.; Spooner, W.; Narechania, A.; Ren, L.; Wei, S.; Kumari,  
726 S.; Faga, B.; Levy, M.J.; McMahan, L.; Van Buren, P.; Vaughn, M.W.; Ying, K.; Yeh, C.T.;  
727 Emrich, S.J.; Jia, Y.; Kalyanaraman, A.; Hsia, A.P.; Barbazuk, W.B.; Baucom, R.S.; Brutnell,  
728 T.P.; Carpita, N.C.; Chaparro, C.; Chia, J.M.; Deragon, J.M.; Estill, J.C.; Fu, Y.; Jeddelloh, J.A.;  
729 Han, Y.; Lee, H.; Li, P.; Lisch, D.R.; Liu, S.; Liu, Z.; Nagel, D.H.; McCann, M.C.; SanMiguel, P.;  
730 Myers, A.M.; Nettleton, D.; Nguyen, J.; Penning, B.W.; Ponnala, L.; Schneider, K.L.; Schwartz,  
731 D.C.; Sharma, A.; Soderlund, C.; Springer, N.M.; Sun, Q.; Wang, H.; Waterman, M.; Westerman,

732 R.; Wolfgruber, T.K.; Yang, L.; Yu, Y.; Zhang, L.; Zhou, S.; Zhu, Q.; Bennetzen, J.L.; Dawe,  
733 R.K.; Jiang, J.; Jiang, N.; Presting, G.G.; Wessler, S.R.; Aluru, S.; Martienssen, R.A.; Clifton,  
734 S.W.; McCombie, W.R.; Wing, R.A.; Wilson, R.K. The B73 maize genome: complexity, diversity,  
735 and dynamics. *Science* **2009**, *326*, 1112-1115.

736 [14] Paterson, A.H.; Bowers, J.E.; Bruggmann, R.; Dubchak, I.; Grimwood, J.; Gundlach, H.; Haberer,  
737 G.; Hellsten, U.; Mitros, T.; Poliakov, A.; Schmutz, J.; Spannagl, M.; Tang, H.; Wang, X.; Wicker,  
738 T.; Bharti, A.K.; Chapman, J.; Feltus, F.A.; Gowik, U.; Grigoriev, I.V.; Lyons, E.; Maher, C.A.;  
739 Martis, M.; Narechania, A.; Otiillar, R.P.; Penning, B.W.; Salamov, A.A.; Wang, Y.; Zhang, L.;  
740 Carpita, N.C.; Freeling, M.; Gingle, A.R.; Hash, C.T.; Keller, B.; Klein, P.; Kresovich, S.;  
741 McCann, M.C.; Ming, R.; Peterson, D.G. Rahman, M.; Ware, D.; Westhoff, P.; Mayer, K.F.X.;  
742 Messing, J.; Rokhsar, D.S. The *Sorghum bicolor* genome and the diversification of grasses. *Nature*  
743 **2009**, *457*, 551-556.

744 [15] Tuskan, G.A.; Difazio, S.; Jansson, S.; Bohlmann, J.; Grigoriev, I.; Hellsten, U.; Putnam, N.;  
745 Ralph, S.; Rombauts, S.; Salamov, A.; Schein, J.; Sterck, L.; Aerts, A.; Bhalerao, R.R.; Bhalerao,  
746 R.P.; Blaudez, D.; Boerjan, W.; Brun, A.; Brunner, A.; Busov, V.; Campbell, M.; Carlson, J.;  
747 Chalot, M.; Chapman, J.; Chen, G.L. Cooper, D.; Coutinho, P.M.; Couturier, J.; Covert, S.; Cronk,  
748 Q.; Cunningham, R.; Davis, J.; Degroeve, S.; Déjardin, A.; Depamphilis, C.; Detter, J.; Dirks, B.;  
749 Dubchak, I.; Duplessis, S.; Ehlting, J.; Ellis, B.; Gendler, K.; Goodstein, D.; Gribskov, M.;  
750 Grimwood, J.; Groover, A.; Gunter, L.; Hamberger, B.; Heinze, B.; Helariutta, Y.; Henrissat, B.;  
751 Holligan, D.; Holt, R.; Huang, W.; Islam-Faridi, N.; Jones, S.; Jones-Rhoades, M.; Jorgensen, R.;  
752 Joshi, C.; Kangasjärvi, J.; Karlsson, J.; Kelleher, C.; Kirkpatrick, R.; Kirst, M.; Kohler, A.;  
753 Kalluri, U.; Larimer, F.; Leebens-Mack, J.; Leplé, J.C.; Locascio, P.; Lou, Y.; Lucas, S.; Martin,  
754 F.; Montanini, B.; Napoli, C.; Nelson, D.R.; Nelson, C.; Nieminen, K.; Nilsson, O.; Pereda, V.;  
755 Peter, G.; Philippe, R.; Pilate, G.; Poliakov, A.; Razumovskaya, J.; Richardson, P.; Rinaldi, C.;  
756 Ritland, K.; Rouzé, P.; Ryaboy, D.; Schmutz, J.; Schrader, J.; Segerman, B.; Shin, H.; Siddiqui,  
757 A.; Sterky, F.; Terry, A.; Tsai, C.J.; Uberbacher, E.; Unneberg, P.; Vahala, J.; Wall, K.; Wessler,  
758 S.; Yang, G.; Yin, T.; Douglas, C.; Marra, M.; Sandberg, G.; Van de Peer, Y.; Rokhsar, D. The  
759 genome of black cottonwood, *Populus trichocarpa* (Torr. & Gray). *Science* **2006**, *313*, 1596-604.

760 [16] Jaillon, O.; Aury, J.M.; Noel, B.; Policriti A.; Clepet, C.; Casagrande, A.; Choisne, N.; Aubourg,  
761 S.; Vitulo, N.; Jubin, C.; Vezzi, A.; Legeai, F.; Huguency, P.; Dasilva, C.; Horner, D.; Mica, E.;

762 Jublot, D.; Poulain, J.; Bruyère, C.; Billault, A.; Segurens, B.; Gouyvenoux, M.; Ugarte, E.;  
763 Cattonaro, F.; Anthouard, V.; Vico, V.; Del Fabbro, C.; Alaux, M.; Di Gaspero, G.; Dumas, V.;  
764 Felice, N.; Paillard, S.; Juman, I.; Moroldo, M.; Scalabrin, S.; Canaguier, A.; Le Clainche, I.;  
765 Malacrida, G.; Durand, E.; Pesole, G.; Laucou, V.; Chatelet, P.; Merdinoglu, D.; Delledonne, M.;  
766 Pezzotti, M.; Lecharny, A.; Scarpelli, C.; Artiguenave, F.; Pè, M.E.; Valle, G.; Morgante, M.;  
767 Caboche, M.; Adam-Blondon, A.F.; Weissenbach J.; Quétier, F.; Wincker, P.; French-Italian  
768 Public Consortium for Grapevine Genome Characterization. The grapevine genome sequence  
769 suggests ancestral hexaploidization in major angiosperm phyla. *Nature* **2007**, *449*, 463-7.

770 [17] Ming, R.; Hou, S.; Feng, Y.; Yu, O.; Dionne-Laporte, A.; Saw, J.H.; Senin, P.; Wang, W.; Ly,  
771 B.V.; Lewis, K.L.T.; Salzberg, S.L.; Feng, L.; Jones, M.R.; Skelton, R.L.; Murray, J.E.; Chen, C.;  
772 Qian, W.; Shen, J.; Du, P.; Moriah Eustice, M.E.; Tong, E.; Tang, H.; Lyons, E.; Paull, R.E.;  
773 Michael, T.P.; Wall, K.; Rice, D.W.; Albert, H.; Wang, M.L.; Zhu, Y.J.; Schatz, M.; Nagarajan,  
774 N.; Acob, R.A.; Guan, P.; Blas, A.; Wai, C.M.; Ackerman, C.M.; Ren, Y.; Liu, C.; Wang, J.;  
775 Wang, J.; Na, J.K.; Shakirov, E.V.; Brian Haas, B.; Thimmapuram, J.; Nelson, D.; Wang, X.;  
776 Bowers, J.E.; Gschwend, A.R.; Delcher, A.L.; Singh, R.; Suzuki, J.Y.; Tripathi, S.; Neupane, K.;  
777 Wei, H.; Irikura, B.; Paidi, M.; Jiang, N.; Zhang, W.; Presting, G.; Windsor, A.; Navajas-Pérez, R.;  
778 Torres, M.J.; Feltus, F.A.; Porter, B.; Li, Y.; Burroughs, A.M.; Luo, M.C.; Liu, L.; Christopher,  
779 D.A.; Mount, S.M.; Moore, P.H.; Sugimura, T.; Jiang, J.; Schuler, M.A.; Friedman, V.; Mitchell-  
780 Olds, T.; Shippen, D.E.; dePamphilis, C.W.; Palmer, J.D.; Freeling, M.; Paterson, A.H.;  
781 Gonsalves, D.; Wang, L.; Alam, M. The draft genome of the transgenic tropical fruit tree papaya  
782 (*Carica papaya* Linnaeus). *Nature* **2008**, *452*, 991-996

783 [18] Schmutz, J.; Cannon, S.B.; Schlueter, J.; Ma, J.; Mitros, T.; Nelson, W.; Hyten, D.L.; Song, Q.;  
784 Thelen, J.J.; Cheng, J.; Xu, D.; Hellsten, U.; May, G.D.; Yu, Y.; Sakurai, T.; Umezawa, T.;  
785 Bhattacharyya, M.K.; Sandhu, D.; Valliyodan, B.; Lindquist, E.; Peto, M.; Grant, D.; Shu, S.;  
786 Goodstein, D.; Barry, K.; Futrell-Griggs, M.; Abernathy, B.; Du, J.; Tian, Z.; Zhu, L.; Gill, N.;  
787 Joshi, T.; Libault, M.; Sethuraman, A.; Zhang, X.C.; Shinozaki, K.; Nguyen, H.T.; Wing, R.A.;  
788 Cregan, P.; Specht, J.; Grimwood, J.; Rokhsar, D.; Stacey, G.; Shoemaker, R.C.; Jackson, S.A.  
789 Genome sequence of the palaeopolyploid soybean *Nature* **2010**, *463*, 178-183.

790 [19] Metzker, M. Sequencing technologies-the next generation. *Nature Rev. Genet.* **2010**, *11*, 31.

- 791 [20] Wang, Z.; Fang, B.; Jingyi, C.; Zhang,X.; Luo, X.; Huang,L.; Chen, X.; Li, Y. De novo assembly  
792 and characterization of root transcriptome using Illumina paired-end sequencing and development  
793 of cSSR markers in sweetpotato (*Ipomoea batatas*). *BMC Genoics* **2011**, *11*, 726.
- 794 [21] Blanca, J.; Cañizares, J.; Roig, C.; Ziarsolo, P.; Nuez, F.; Picó, B. Transcriptome characterization  
795 and high throughput SSRs and SNPs discovery in *Cucurbita pepo* (Cucurbitaceae). *BMC*  
796 *Genomics* **2011**, *12*, 104.
- 797 [22] Dutta, S.; Kumawat, G.; Singh, B.P.; Gupta, D.K.; Singh, S.; Dogra, V.; Gaikwad, K.; Sharma,  
798 T.R.; Raje, R.S.; Bandhopadhya, T.K.; Datta, S.; Singh, M.N.; Bashasab, F.; Kulwal, P.; Wanjari,  
799 K.B.; Varshney, R.K.; Cook D.R.; Singh, N.K. Development of genic-SSR markers by deep  
800 transcriptome sequencing in pigeonpea [*Cajanus cajan* (L.) Millspaugh]. *BMC Plant Biol.* **2011**,  
801 *11*, 17.
- 802 [23] Logacheva, M.D.; Kasianov, A.S.; Vinogradov, D.V.; Samigullin, T.H.; Gelfand, M.S.; Makeev,  
803 V.J.; Penin, A.A. *De novo* sequencing and characterization of floral transcriptome in two species  
804 of buckwheat (*Fagopyrum*). *BMC Genomics* **2011**, *12*, 30.
- 805 [24] Li, H.; Handsaker, B.; Wysoker, A.; Fennell, T.; Ruan, J.; Homer, N.; Marth, G.; Abecasis, G.;  
806 Durbin, R. Genome Project Data P: The Sequence Alignment/Map format and SAMtools.  
807 *Bioinformatics* **2009**, *25*, 2078- 2079.
- 808 [25] Pop, M.; Salzberg, S.L. Bioinformatics challenges of new sequencing technology. *Trends Genet.*  
809 **2008**, *24*, 142-149.
- 810 [26] Horner, D.S.; Pavesi, G.; Castrignanò, T.; Meo, P.D.O.D.; Liuni, S.; Sammeth, M.; Picardi, E.;  
811 Presole, G. Bioinformatics approaches for genomics and post genomics applications of next-  
812 generation sequencing. *Briefings Bioinformatics* **2009**, *2*, 181-197.
- 813 [27] Li, H.; Horner, D.S. A survey of sequence alignment algorithms for next-generation sequencing.  
814 *Briefings Bioinformatics* **2010**, *2*, 473-483.
- 815 [28] Wang, L.; Li, P.; Brutnell, T.P. Exploring plant transcriptomes using ultra high-throughput  
816 sequencing. *Briefings Functional Genomics* **2010**, *9*, 118-128.
- 817 [29] VanGuilder, H.D.; Vrana, K.E.; Freeman, W.M. Twenty-five years of quantitative PCR for gene  
818 expression analysis. *Biotechniques* **2008**, *44*, 619-626.

- 819 [30] Bachem, C.W.; van der Hoeven, R.S.; de Bruijn, S.M.; Vreugdenhil, D.; Zabeau, M.; Visser, R.G.  
820 Visualization of differential gene expression using a novel method of RNA fingerprinting based on  
821 AFLP: analysis of gene expression during potato tuber development. *Plant J.* **1996**, *9*, 745-753.
- 822 [31] Anisimov SV. Serial Analysis of Gene Expression (SAGE): 13 years of application in research.  
823 *Curr. Pharm. Biotechnol.* **2008**, *9*, 338-350.
- 824 [32] Reinartz, J.; Bruyns, E.; Lin, J.Z.; Burcham, T.; Brenner, S.; Bowen, B.; Kramer, M.; Woychik,  
825 R. Massively parallel signature sequencing (MPSS) as a tool for in-depth quantitative gene  
826 expression profiling in all organisms. *Briefings Functional Genomics Proteomics*, **2002**, *1*:95-104.
- 827 [33] Schena, M.; Shalon, D.; Davis, R.W.; Brown, P.O. Quantitative monitoring of gene expression  
828 patterns with a complementary DNA microarray. *Science* **1995**, *270*, 467-470.
- 829 [34] Al-Shahrour, F.; Minguez, P.; Vaquerizas, J.M.; Conde, L.; Dopazo, J. BABELOMICS: a suite of  
830 web tools for functional annotation and analysis of groups of genes in high-throughput  
831 experiments. *Nucleic Acids Res.* **2005**, *33*, 460-464.
- 832 [35] Saeed, A.I.; Sharov, V.; White, J.; Li, J.; Liang, W.; Bhagabati, N.; Braisted, J.; Klapa, M.;  
833 Currier, T.; Thiagarajan, M.; Sturn, A.; Snuffin, M.; Rezantsev, A.; Popov, D.; Ryltsov, A.;  
834 Kostukovich, E.; Borisovsky, I.; Liu, Z.; Vinsavich, A.; Trush, V.; Quackenbush, J. TM4: a free,  
835 open-source system for microarray data management and analysis. *Biotechniques*, **2003**, *34*, 374-  
836 378.
- 837 [36] Zimmermann, P.; Laule, O.; Schmitz, J.; Hruz, T.; Bleuler, S.; Gruissem, W. Genevestigator  
838 transcriptome meta-analysis and biomarker search using rice and barley gene expression  
839 databases. *Mol. Plant.*, **2008**, *1*, 851-857.
- 840 [37] Schmid, M.; Davison, T.S.; Henz, S.R.; Pape UJ.; Demar M.; Vingron M.; Scholkopf B.; Weigel  
841 D.; Lohmann JU. A gene expression map of Arabidopsis thaliana development. *Nature Genetics*  
842 **2005**, *37*, 501-506.
- 843 [38] Barrett, T.; Troup, D.B.; Wilhite, S.E.; Ledoux, P.; Evangelista, C.; Kim, I.F.; Tomashevsky, M.;  
844 Marshall, K.A.; Phillippy, K.H.; Sherman, P.M.; Muerlter, R.N.; Holko, M.; Ayanbule, O.;  
845 Yefanov, A.; Soboleva, A. NCBI GEO: archive for functional genomics data sets -10 years on.  
846 *Nucleic Acids Res.*, **2011**, *39*, 1005-1010.
- 847 [39] Parkinson, H.; Sarkans, U.; Kolesnikov, N.; Abeygunawardena, N.; Burdett, T.; Dylag, M.; Emam,  
848 I.; Farne, A.; Hastings, E.; Holloway, E.; Kurbatova, N.; Lukk, M.; Malone, J.; Mani, R.;

849 Pilicheva, E.; Rustici, G.; Sharma, A.; Williams, E.; Adamusiak, T.; Brandizi, M.; Sklyar, N.;  
850 Brazma A. ArrayExpress update - an archive of microarray and high-throughput sequencing-based  
851 functional genomics experiments. *Nucleic Acids Res.*, **2011**, *39*, 1002-1004.

852 [40] Pascual, L.; Blanca, JM.; Canizares, J.; Nuez, F. Analysis of gene expression during the fruit set of  
853 tomato: A comparative approach. *Plant Sci.* **2007**, *173*, 609-620.

854 [41] Marioni, J.C.; Mason, C.E.; Mane, S.M.; Stephens, M.; Gilad, Y. RNA-seq: an assessment of  
855 technical reproducibility and comparison with gene expression arrays. *Genome Res.*, **2008**, *18*,  
856 1509-1517.

857 [42] Stiglic, G.; Bajgot, M.; Kokol, P. Gene set enrichment meta-learning analysis: next-generation  
858 sequencing versus microarrays. *BMC Bioinformatics* **2010**, *11*, 176.

859 [43] Cloonan, N.; Forrest, A.R.; Kolle, G.; Gardiner, B.B.; Faulkner, G.J.; Brown, M.K.; Taylor, D.F.;  
860 Steptoe, A.L.; Wani, S.; Bethel, G.; Robertson, A.J.; Perkins, A.C.; Bruce, S.J.; Lee, C.C.; Ranade,  
861 S.S.; Peckham, H.E.; Manning, J.M.; McKernan, K.J.; Grimmond, S.M. Stem cell transcriptome  
862 profiling via massive-scale mRNA sequencing. *Nat. Methods.*, **2008**, *5*, 613-619.

863 [44] Alagna, F.; D'Agostino, N.; Torchia, L.; Servili, M.; Rao, R.; Pietrella, M.; Giuliano, G.;  
864 Chiusano, M.L.; Baldoni, L.; Perrotta, G. Comparative 454 pyrosequencing of transcripts from  
865 two olive genotypes during fruit development. *BMC Genomics*, **2009**, *10*, 399.

866 [45] Zenoni, S.; Ferrarini, A.; Giacomelli, E.; Xumerle, E.; Fasoli, M.; Malerba, G.; Bellin, D.;  
867 Pezzotti, M.; Delledonne, M. Characterization of transcriptional complexity during berry  
868 development in *Vitis vinifera* using RNA-Seq. *Plant Physiol.* **2010**, *152*, 1787-1795.

869 [46] Wang Y.; Zhang H.; Li H.; Miao X. Second-generation sequencing supply an effective way to  
870 screen RNAi targets in large scale for potential application in pest insect control. *PLoS One* **2011**,  
871 *6*(4), e18644.

872 [47] Shi, C.Y.; Yang, H.; Wei, C.L.; Yu, O.; Zhang, Z.Z.; Jiang, C.J.; Sun, J.; Li, Y.Y.; Chen, Q.; Xia,  
873 T.; Wan, X.C. Deep sequencing of the *Camellia sinensis* transcriptome revealed candidate genes  
874 for major metabolic pathways of tea-specific compounds. *BMC Genomics* **2011**, *12*, 131.

875 [48] Gepts, P. Plant genetic resources conservation and utilization: the accomplishments and future of a  
876 societal insurance policy. *Crop Sci.* **2006**, *46*, 2278-2292.

- 877 [49] Parry, M.A.J.; Madgwick, P.J.; Bayon, C.; Tearall, K.; Hernandez-Lopez, A.; Baudo, M.;  
878 Rakszegi, M.; Hamada, W.; Al-Yassin, A.; Ouabbou, H.; Labhilili, M.; Phillips, A.L. Mutation  
879 discovery for crop improvement. *J. Exp. Bot.* **2009**, *60*, 2817-2825.
- 880 [50] Till, B.J.; Reynolds, S.H.; Greene, E.A.; Codomo, C.A.; Enns, L.C.; Johnson, J.E.; Burtner, C.;  
881 Odden, A.R.; Young, K.; Taylor, N.E.; Henikoff, J.G.; Comai, L.; Henikoff, S. Large-scale  
882 discovery of induced point mutations with high- throughput TILLING. *Genome Res.* **2003**, *13*,  
883 524-530.
- 884 [51] Comai, L.; Young, K.; Till, B.J.; Reynolds, S.H.; Greene, E.A.; Codomo, C.A.; Enns, L.C.;  
885 Johnson, J.E.; Burtner, C.; Odden, A.R.; Henikoff, S. Efficient discovery of DNA polymorphisms  
886 in natural populations by EcoTILLING. *Plant J.* **2004**, *37*, 778-786.
- 887 [52] Till, B.J.; Burtner, C.; Comai, L.; Henikoff, S. Mismatch cleavage by single-strand specific  
888 nucleases. *Nucleic Acids Res.* **2004**, *32*, 2632-2641.
- 889 [53] Triques, K.; Piednoir, E.; Dalmais, M.; Schmidt, J.; Le Signor, C.; Sharkey, M.; Caboche, M.;  
890 Sturbois, B.; Bendahmane, A. Mutation detection using ENDO1: Application to disease  
891 diagnostics in humans and TILLING and EcoTILLING in plants. *BMC Mol. Biol.* **2008**, *9*, 9.
- 892 [54] Barkley, N.A.; Wang, M.L. Application of TILLING and EcoTILLING as reverse genetic  
893 approaches to elucidate the function of genes in plants and animals. *Curr. Genomics*, **2008**, *9*, 212-  
894 226.
- 895 [55] Colbert, T.; Till, B.J.; Tompa, R.; Reynolds, S.; Steine, M.N.; Yeung, A.T.; McCallum, C.M.;  
896 Comai, L.; Henikoff, S. High-throughput screening for induced point mutations. *Plant Physiol.*,  
897 **2001**, *126*, 480-484.
- 898 [56] Perry, J.A.; Wang, T.L.; Welham, T.J.; Gardner, S.; Pike, J.M.; Yoshida, S.; Parniske, M. A  
899 TILLING reverse genetics tool and a web-accessible collection of mutants of the legume *Lotus*  
900 *japonicus*. *Plant Physiol.*, **2003**, *131*, 866-871.
- 901 [57] Caldwell, D.G.; McCallum, N.; Shaw, P.; Muehlbauer, G.J.; Marshall, D.F.; Waugh R. A  
902 structured mutant population for forward and reverse genetics in barley (*Hordeum vulgare* L).  
903 *Plant J.*, **2004**, *40*:143-150.
- 904 [58] Weil, C.F.; Monde, R.A. Getting to the point - mutations in maize. *Crop Sci.*, **2007**, *47*, S60-S67.
- 905 [59] Triques, K.; Sturbois, B.; Gallais, S.; Dalmais, M.; Chauvin, S.; Clepet, C.; Aubourg, S.; Rameau,  
906 C.; Caboche, M.; Bendahmane, A. Characterization of *Arabidopsis thaliana* mismatch specific



907 endonucleases: application to mutation discovery by TILLING in pea. *Plant J.*, **2007**, 51, 1116-  
908 1125.

909 [60] Dahmani-Mardas, F.; Troadec, C.; Boualem, A.; Lévêque, S.; Alsdon, A.A.; Aldoss, A.A.;  
910 Dogimont, C.; Bendahmane, A. Engineering melon plants with improved fruit shelf life using the  
911 TILLING approach. *PLoS One* **2010**, 5 (12), e15776.

912 [61] Kadaru, S.B.; Yadav, A.S.; Fjellstrom, R.G.; Oard J.H. Alternative EcoTILLING protocol for  
913 rapid, cost-effective single-nucleotide polymorphism discovery and genotyping in rice (*Oryza*  
914 *sativa* L.). *Plant Mol. Biol. Reporter* **2006**, 24, 3-22

915 [62] Mejlhede, N.; Kyjovska, Z.; Backes, G.; Burhenne, K.; Rasmussen, S.K.; Jahoor, A. EcoTILLING  
916 for the identification of allelic variation in the powdery mildew resistance genes *mlo* and *Mla* of  
917 barley. *Plant Breed.* **2006**, 125, 461-467.

918 [63] Wang, J.; Sun, J.Z.; Liu, D.C.; Yang, W.L.; Wang, D.W.; Tong, Y.P.; Zhang, A.M. Analysis of  
919 *Pina* and *Pinb* alleles in the micro-core collections of Chinese wheat germplasm by EcoTILLING  
920 and identification of a novel *Pinb* allele. *J. Cereal Sci.* **2008**, 48, 836-842.

921 [64] Ramos, M.L.; Huntley, J.J.; Maleki, S.J.; Ozias-Akins, P. Identification and characterization of a  
922 hypoallergenic ortholog of Ara h 2.01. *Plant Mol. Biol.* **2009**, 69, 325-335.

923 [65] Galeano, C.H.; Gomez, M.; Rodriguez, L.M.; Blair, M.W. CEL I Nuclease Digestion for SNP  
924 Discovery and Marker Development in Common Bean (*Phaseolus vulgaris* L.). *Crop Sci.* **2009**,  
925 49, 381-394.

926 [66] Wang, J.; Sun, J.Z.; Liu, D.C.; Yang, W.L.; Wang, D.W.; Tong, Y.P.; Zhang, A.M. Analysis of  
927 *Pina* and *Pinb* alleles in the micro-core collections of Chinese wheat germplasm by EcoTILLING  
928 and identification of a novel *Pinb* allele. *J. Cereal Sci.*, **2008**, 48, 836-842.

929 [67] Ramos, M.L.; Huntley, J.J.; Maleki, S.J.; Ozias-Akins, P. Identification and characterization of a  
930 hypoallergenic ortholog of Ara h 2.01. *Plant Mol. Biol.*, **2009**, 69:325-335.

931 [68] Piron, F.; Nicolai, M.; Minoia, S.; Piednoir, E.; Moretti, A.; Salgues, A.; Zamir, D.; Caranta, C.;  
932 Bendahmane, A. An induced mutation in tomato eIF4E leads to immunity to two potyviruses.  
933 *PLoS One* **2010**, 5 (6), e11313.

934 [69] Ibiza, V.P.; Cañizares, J.; Nuez, F. EcoTILLING in Capsicum species: searching for new virus  
935 resistances. *BMC Genomics* **2010**, 11, 631.

- 936 [70] Deschamps, S.; Campbell, M.A. Utilization of next-generation sequencing platforms in plant  
937 genomics and genetic variant discovery. *Mol. Breed.* **2010**, *25*, 553-570.
- 938 [71] Ganal, M.W.; Altmann, T.; Röder, M.S. SNP identification in crop plants. *Curr. Opin. Plant Biol.*,  
939 **2009**, *12*, 211-217.
- 940 [72] Weigel, D.; Mott, R. The 1001 genomes project for *Arabidopsis thaliana*. *Genome Biol.*, **2009**,  
941 *10*, 107.
- 942 [73] Gore, M.A.; Wright, M.H.; Ersoz, E.S.; Bouffard, P.; Szekeres, E.S.; Jarvie, T.P.; Hurwitz, B.L.;  
943 Narechania, A.; Harkins, T.T.; Grills, G.S.; Ware, D.H.; Buckler, E.S. Large-scale discovery of  
944 gene-enriched SNPs. *The Plant Genome* **2009**, *2*, 121-133.
- 945 [74] Myles, S.; Chia, J.M.; Hurwitz, B.; Simon, C.; Zhong, G.Y.; Buckler, E.; Ware, D. Rapid genomic  
946 characterization of the genus *Vitis*. *PLoS One* **2010**, *5* (1): e8219.
- 947 [75] McCouch, S.R.; Zhao, K.; Wright, M.; Tung, C.W.; Ebana, K.; Thomson, M.; Reynolds, A.;  
948 Wang, D.; DeClerck, G.; Ali, M.L.; McClung, A.; Eizenga, G.; Bustamante, C. Development of  
949 genome-wide SNP assays for rice. *Breed. Sci.* **2010**, *60*, 524-535.
- 950 [76] Wu, X.; Ren, C.; Joshi, T.; Vuong, T.; Xu, D.; Nguyen, H.T. SNP discovery by high-throughput  
951 sequencing in soybean. *BMC Genomics.* **2010**, *11*:469.
- 952 [77] Deschamps, S.; La Rota, M.; Ratashak, J.P.; Biddle, P.; Thureen, D.; Farmer, A.; Luck, S.; Beatty,  
953 M.; Nagasawa, N.; Michael, L.; Llaca, V.; Sakai, H.; May, G.; Lightner, J.; Campbell, M.A. Rapid  
954 genome-wide single nucleotide polymorphism discovery in soybean and rice via deep  
955 resequencing of reduced representation libraries with the Illumina genome analyzer. *The Plant*  
956 *Genome*, **2010**, *3*, 1.
- 957 [78] Bundock, P.C.; Elliott, F.G.; Ablett, G.; Benson, A.D.; Casu, R.E.; Aitken, K.S.; Henry, R.J.  
958 Targeted single nucleotide polymorphism (SNP) discovery in a highly polyploid plant species  
959 using 454 sequencing. *Plant Biotechnol. J.* **2009**, *7*, 347-354.
- 960 [79] Trick, M.; Long, Y.; Meng, J.; Bancroft, I. Single nucleotide polymorphism (SNP) discovery in  
961 the polyploid *Brassica napus* using Solexa transcriptome sequencing. *Plant Biotechnol. J.* **2009**, *7*,  
962 334-346.
- 963 [80] Oliver, R.E.; Lazo, G.R.; Lutz, J.D.; Rubenfield, M.J.; Tinker, N.A.; Anderson, J.M.; Morehead  
964 N.H.W.; Adhikary, D.; Jellen, E.N.; Maughan, P.J.; Brown Guedira, G.L.; Chao, S.; Beattie, A.D.;  
965 Carson, M.L.; Rines, H.W.; Obert, D.E.; Bonman, J.M.; Jackson, E.W. Model SNP development

966 for complex genomes based on hexaploid oat using high-throughput 454 sequencing technology.  
967 *BMC Genomics* **2011**, *12*, 77.

968 [81] Shen, Y.; Wan, Z.; Coarfa, C.; Drabek, R.; Chen, L.; Ostrowski, E.A.; Liu, Y.; Weinstock, G.M.;  
969 Wheeler, D.A.; Gibbs, R.A.; Yu, F. A SNP discovery method to assess variant allele probability  
970 from next-generation resequencing data. *Genome Res.*, **2010**, *20*, 273-280.

971 [82] Maughan, P.J.; Smith, S.M.; Fairbanks, D.J.; Jellen, E.N. Development, characterization, and  
972 linkage mapping of single nucleotide polymorphisms in the grain amaranths (*Amaranthus* sp.). *The*  
973 *Plant Genome* **2011**, *4*, 92-101.

974 [83] Hyten, D.L.; Cannon, S.B.; Song, Q.; Weeks, N.; Fickus, E.W.; Shoemaker, R.C.; Specht, J.E.;  
975 Farmer, A.D.; May, G.D.; Cregan, P.B. High-throughput SNP discovery through deep  
976 resequencing of a reduced representation library to anchor and orient scaffolds in the soybean  
977 whole genome sequence. *BMC Genomics* **2010**, *11*, 38.

978 [84] You, F.M.; Huo, N.; Deal, K.R.; Gu, Y.Q.; Luo, M.C.; McGuire, P.E.; Dvorak, J.; Anderson, O.D.  
979 Annotation-based genome-wide SNP discovery in the large and complex *Aegilops tauschii*  
980 genome using next-generation sequencing without a reference genome sequence. *BMC Genomics*  
981 **2011**, *12*, 59.

982 [85] Barbazuk, W.B.; Emrich, S.J.; Hsin, D.; Chen, L.; Li, P.; Schnable, P.S. SNP discovery via 454  
983 transcriptome sequencing. *Plant J.* **2007**, *51*, 910-918.

984 [86] Novaes, E.; Drost, D.R.; Farmerie, W.G.; Pappas, G.J.Jr.; Grattapaglia, D.; Sederoff, R.R.; Kirst,  
985 M. High-throughput gene and SNP discovery in *Eucalyptus grandis*, an uncharacterized genome.  
986 *BMC Genomics* **2008**, *9*, 312.

987 [87] Ueno, S.; Le Provost, G.; Léger, V.; Klopp, C.; Noirot, C.; Frigerio, JM.; Salin, F.; Salse, J.;  
988 Abrouk, M.; Murat, F.; Brendel, O.; Derory, J.; Abadie, P.; Léger, P.; Cabane, C.; Barré, A.; de  
989 Daruvar, A.; Couloux, A.; Wincker, P.; Reviron, M.P.; Kremer, A.; Plomion, C. Bioinformatic  
990 analysis of ESTs collected by Sanger and pyrosequencing methods for a keystone forest tree  
991 species: oak. *BMC Genomics* **2010**, *11*, 650.

992 [88] Vidal, R.O.; Mondego, J.M.; Pot, D.; Ambrósio, A.B.; Andrade, A.C.; Pereira, L.F.; Colombo,  
993 C.A.; Vieira, L.G.; Carazzolle, M.F.; Pereira, G.A. A high-throughput data mining of single  
994 nucleotide polymorphisms in *Coffea* species expressed sequence tags suggests differential

995 homeologous gene expression in the allotetraploid *Coffea arabica*. *Plant Physiol.* **2010**, *154*, 1053-  
996 1066.

997 [89] Schafleitner, R.; Tincopa, L.R.; Palomino, O.; Rossel, G.; Robles, R.F.; Alagon, R.; Rivera, C.;  
998 Quispe, C.; Rojas, L.; Pacheco, J.A.; Solis, J.; Cerna, D.; Kim, J.Y.; Hou, J.; Simon, R. A sweet  
999 potato gene index established by de novo assembly of pyrosequencing and Sanger sequences and  
1000 mining for gene-based microsatellite markers. *BMC Genomics* **2010**, *11*, 604.

1001 [90] Garg, R.; Patel, R.K.; Tyagi, A.K.; Jain, M. De novo assembly of chickpea transcriptome using  
1002 short reads for gene discovery and marker identification. *DNA Res.* **2011**, *18*, 53-63.

1003 [91] Shirasawa, K.; Isobe, S.; Hirakawa, H.; Asamizu, E.; Fukuoka, H.; Just, D.; Rothan, C.; Sasamoto,  
1004 S.; Fujishiro, T.; Kishida, Y.; Kohara, M.; Tsuruoka, H.; Wada, T.; Nakamura, Y.; Sato, S.;  
1005 Tabata, S. SNP discovery and linkage map construction in cultivated tomato *DNA Res.* **2010**, *17*,  
1006 381-391.

1007 [92] Guo, S.; Zheng, Y.; Joung, J.G.; Liu, S.; Zhang, Z.; Crasta, O.R.; Sobral, B.W.; Xu, Y.; Huang, S.;  
1008 Fei, Z. Transcriptome sequencing and comparative analysis of cucumber flowers with different sex  
1009 types. *BMC Genomics* **2010**, *11*, 384.

1010 [93] Blanca, J.; Cañizares, J.; Ziarsolo, P.; Esteras, C.; Mir, G.; Nuez, F.; García-Mas, J.; Picó, B.  
1011 Melon transcriptome characterization. SSRs and SNPs discovery for high throughput genotyping  
1012 across the species. *The Plant Genome*, **2011**, in press.

1013 [94] Narina, S.S.; Buyyarapu, R.; Kottapalli, K.R.; Sartie, A.M.; Ali, M.I.; Robert, A.; Hodeba, M.J.D.;  
1014 Sayre, B.L.; Scheffler, B.E. Generation and analysis of expressed sequence tags (ESTs) for marker  
1015 development in yam (*Dioscorea alata* L.). *BMC Genomics* **2011**, *12*, 100.

1016 [95] Yang, S.S.; Tu, Z.J.; Cheung, F.; Xu, W.W.; Lamb, J.F.S.; Jung, H.J.G.; Vance, C.P.; Gronwald,  
1017 J.W. Using RNA-Seq for gene identification, polymorphism detection and transcript profiling in  
1018 two alfalfa genotypes with divergent cell wall composition in stems. *BMC Genomics* **2011**, *12*,  
1019 199.

1020 [96] Huang, X.H.; Wei, X.H.; Sang, T.; Zhao, Q.A.; Feng, Q.; Zhao, Y.; Li, C.Y.; Zhu, C.R.; Lu, T.T.;  
1021 Zhang, Z.W.; Li, M.; Fan, D.L.; Guo, Y.L.; Wang, A.; Wang, L.; Deng, L.W.; Li, W.J.; Lu, Y.Q.;  
1022 Weng, Q.J.; Liu, K.Y.; Huang, T.; Zhou, T.Y.; Jing, Y.F.; Li, W.; Lin, Z.; Buckler, E.S.; Qian,  
1023 Q.A.; Zhang, Q.F.; Li, J.Y.; Han, B. Genome-wide association studies of 14 agronomic traits in  
1024 rice landraces. *Nature Genet.* **2010**, *42*: 961-976.

- 1025 [97] Massman, J.; Cooper, B.; Horsley, R.; Neate, S.; Dill-Macky, R.; Chao, S.; Dong, Y.; Schwarz, P.;  
1026 Muehlbauer, G.J.; Smith, K.P. Genome-wide association mapping of Fusarium head blight  
1027 resistance in contemporary barley breeding germplasm. *Mol. Breed.* **2011**, *27*: 439-454.
- 1028 [98] Close, T.J.; Bhat, P.R.; Lonardi, S.; Wu, Y.H.; Rostoks, N.; Ramsay, L.; Druka, A.; Stein, N.;  
1029 Svensson, J.T.; Wanamaker, S.; Bozdog, S.; Roose, M.L.; Moscou, M.J.; Chao, S.A.M.; Varshney,  
1030 R.K.; Szucs, P.; Sato, K.; Hayes, P.M.; Matthews, D.E.; Kleinhofs, A.; Muehlbauer, G.J.;  
1031 DeYoung, J.; Marshall, D.F.; Madishetty, K.; Fenton, R.D.; Condamine, P.; Graner, A.; Waugh, R.  
1032 Development and implementation of high-throughput SNP genotyping in barley. *BMC Genomics*  
1033 **2009**, *10*, 582.
- 1034 [99] Belo, A.; Zheng, P.Z.; Luck, S.; Shen, B.; Meyer, D.J.; Li, B.L.; Tingey, S.; Rafalski, A. Whole  
1035 genome scan detects an allelic variant of fad2 associated with increased oleic acid levels in maize.  
1036 *Mol. Genet. Genomics*, **2008**, *279*, 1-10.
- 1037 [100] Kump, K.L.; Bradbury, P.J.; Wissler, R.J.; Buckler, E.S.; Belcher, A.R.; Oropeza-Rosas, M.A.;  
1038 Zwonitzer, J.C.; Kresovich, S.; McMullen, M.D.; Ware, D.; Balint-Kurti, P.J.; Holland, J.B.  
1039 Genome-wide association study of quantitative resistance to southern leaf blight in the maize  
1040 nested association mapping population. *Nature Genet.* **2011**, *43*, 163-168.
- 1041 [101] Yan, J.B.; Yang, X.H.; Shah, T.; Sanchez-Villeda, H.; Li, J.S.; Warburton, M.; Zhou, Y.; Crouch,  
1042 J.H.; Xu, Y.B. High-throughput SNP genotyping with the GoldenGate assay in maize. *Mol.*  
1043 *Breed.* **2009**, *25*, 441-451.
- 1044 [102] Emanuelli, F.; Battilana, J.; Costantini, L.; Le Cunff, L.; Boursiquot, J.M.; This, P.; Grando, M.S.  
1045 A candidate gene association study on muscat flavor in grapevine (*Vitis vinifera* L.). *BMC Plant*  
1046 *Biol.* **2010**, *10*, 241.
- 1047 [103] Deulvot, C.; Charrel, H.; Marty, A.; Jacquin, F.; Donnadiou, C.; Lejeune-Hénaut, I.; Burstin, J.;  
1048 Aubert, G. Highly-multiplexed SNP genotyping for genetic mapping and germplasm diversity  
1049 studies in pea. *BMC Genomics.* **2011**, *11*, 468.
- 1050 [104] Neumann, K.; Kobiljski, B.; Dencic, S.; Varshney, R.K.; Borner, A. Genome-wide association  
1051 mapping: a case study in bread wheat (*Triticum aestivum* L.). *Mol. Breed.* **2011**, *27*, 37-58.
- 1052 [105] Chao, S.; Dubcovsky, J.; Luo, M.C.; Baenziger, S.P.; Matnyazov, R.; Clark, D.R.; Talbert, L.E.;  
1053 Anderson, J.A.; Dreisigacker, S.; Glover, K.; Chen, J.; Campbell, K.; Bruckner, P.L.; Rudd, J.C.;  
1054 Haley, S.; Carver, B.F.; Perry, S.; Sorrells, M.E.; Akhunov, E.D. Population- and genome-specific

1055 patterns of linkage disequilibrium and SNP variation in spring and winter wheat (*Triticum*  
1056 *aestivum* L.). *BMC Genomics*. **2011**, *11*, 727.

1057 [106] Beaulieu, J.; Doerksen, T.; Boyle, B.; Clément, S.; Deslauriers, M.; Beauseigle, S.; Poulin, P.;  
1058 Lenz, P.; Caron, S.; Rigault, P.; Bicho, P.; Bousquet, J.; Mackay, J. Association genetics of wood  
1059 physical traits in the conifer white spruce and relationships with gene expression. *Genetics*, **2011**,  
1060 *188*, 197-214.

1061 [107] Edenberg, H.J.; Yunlong, L. Laboratory methods for high-throughput genotyping. *Cold Spring*  
1062 *Harb Protoc*. **2009**, *16*, 183-193.

1063 [108] Appleby, N.; Edwards, D.; Batley, J. New technologies for ultra-high throughput genotyping in  
1064 plants. *Methods Mol Biol*. **2009**, *513*, 19-39.

1065 [109] Lin, C.H.; Yeakley, J.M.; McDaniel, T.K.; Shen, R. Medium- to high-throughput SNP genotyping  
1066 using VeraCode microbeads. *Methods Mol. Biol*. **2009**, *496*, 129-142.

1067 [110] Gabriel, S.; Ziaugra, L.; Tabbaa D. 2009. SNP genotyping using the Sequenom MassARRAY  
1068 iPLEX platform. *Curr. Protocols Human Genet.*, **Suppl. 60**, 2.12.1-2.12.18.

1069 [111] Eckert, A.J.; van Heerwaarden, J.; Wegrzyn, J.L.; Nelson, C.D.; Ross-Ibarra, J.; Gonzalez-  
1070 Martinez, S.C.; Neale, D.B. Patterns of population structure and environmental associations to  
1071 aridity across the range of loblolly pine (*Pinus taeda* L., Pinaceae). *Genetics* **2010**, *185*, 969-982

1072 [112] Eckert, A.J.; Pande, B.; Ersoz, E.S.; Wright, M.H.; Rashbrook, V.K.; Nicolet, C.M.; Neale, D.B.  
1073 High-throughput genotyping and mapping of single nucleotide polymorphisms in loblolly pine  
1074 (*Pinus taeda* L.). *Tree Genet. Genomes*, **2009**, *5*, 225-234.

1075 [113] Liu, S.; Chen, H.D.; Makarevitch, I.; Shirmer, R.; Emrich, S.J.; Dietrich, C.R.; Barbazuk, W.B.;  
1076 Springer, N.M.; Schnable, P.S. High-throughput genetic mapping of mutants via quantitative  
1077 single nucleotide polymorphism typing. *Genetics* **2010**, *184*, 19-26.

1078 [114] Neves, L.; Mamani, E.; Alfenas, A.; Kirst, M.; Grattapaglia, D. A high-density transcript linkage  
1079 map with 1,845 expressed genes positioned by microarray-based Single Feature Polymorphisms  
1080 (SFP) in *Eucalyptus*. *BMC Genomics* **2011**, *12*, 189.

1081 [115] Sato, K.; Nankaku, N.; Takeda, K. A high-density transcript linkage map of barley derived from a  
1082 single population. *Heredity* **2009**, *103*: 110-117.

- 1083 [116] Chutimanitsakun, Y.; Nipper, R.W.; Cuesta-Marcos, A.; Cistue, L.; Corey, A.; Filichkina, T.;  
1084 Johnson, E.A.; Hayes, P.M. Construction and application for QTL analysis of a Restriction Site  
1085 Associated DNA (RAD) linkage map in barley. *BMC Genomics* **2011**, *12*, 4.
- 1086 [117] Celton, J.M.; Christoffels, A.; Sargent, D.J.; Xu, X.M.; Rees, D.J.G. Genome-wide SNP  
1087 identification by high-throughput sequencing and selective mapping allows sequence assembly  
1088 positioning using a framework genetic linkage map. *BMC Biol.* **2010**, *8*, 155.
- 1089 [118] Huang, X.H, Feng, Q.; Qian, Q.; Zhao, Q.; Wang, L.; Wang, A.H.; Guan, J.P.; Fan, D.L.; Weng,  
1090 Q.J.; Huang, T.; Dong, G.J.; Sang, T.; Han, B. High-throughput genotyping by whole-genome  
1091 resequencing. *Genome Res.* **2009**, *19*: 1068-1076.
- 1092 [119] Xu, J.J.; Zhao, Q.A.; Du, P.N.; Xu, C.W.; Wang, B.H.; Feng, Q.; Liu, Q.Q.; Tang, S.Z.; Gu, M.H.;  
1093 Han, B.; Liang, G.H. Developing high throughput genotyped chromosome segment substitution  
1094 lines based on population whole-genome re-sequencing in rice (*Oryza sativa* L.). *BMC Genomics*  
1095 **2010**, *11*, 656.
- 1096 [120] Xie, W.B.; Feng, Q.; Yu, H.H.; Huang, X.H.; Zhao, Q.A.; Xing, Y.Z.; Yu, S.B.; Han, B.; Zhang,  
1097 Q.F. Parent-independent genotyping for constructing an ultrahigh-density linkage map based on  
1098 population sequencing. *Proc. Natl. Acad. Sci. USA* **2010**, *107*, 10578-10583.
- 1099 [121] Boerma, R.; Wilson, R.; Ready, E. Soybean genomics research program strategic plan  
1100 implementing research to meet 2012–2016 strategic milestones. *The Plant Genome* **2011**, *4*, 1-11.
- 1101 [122] Gore, M.A.; Chia, J.M.; Elshire, R.J.; Sun, Q.; Ersoz, E.S.; Hurwitz, B.L.; Peiffer, J.A.;  
1102 McMullen, M.D.; Grills, G.S.; Ross-Ibarra, J.; Ware, D.H.; Buckler, E.S. 2009. A first-generation  
1103 haplotype map of maize. *Science*, **2009**, *326*, 1115-1117.
- 1104 [123] Branca, A.; Paape, T.; Briskine, R.; Zhou, P.; Wang, S.; Denny, R.; Mudge, J.; Bharti, A.K.;  
1105 Farmer, A.; May, G.D.; Tiffin, P.L.; Young, N.D. The *Medicago truncatula* HapMap project: deep  
1106 coverage sequencing of 30 inbred lines using Illumina's Solexa technology. *Plant & Animal*  
1107 *Genomes XVIII Conference, San Diego, CA, 2010*, P417.
- 1108 [124] McNally, K.L.; Childs, K.L.; Bohnert, R.; Davidson, R.M.; Zhao, K.; Ulat, V.J.; Zeller, G.; Clark,  
1109 R.M.; Hoen, D.R.; Bureau, T.E.; Stokowski, R.; Ballinger, D.G.; Frazer, K.A.; Cox, D.R.;  
1110 Padhukasahasram, B.; Bustamante, C.D.; Weigel, D.; Mackill, D.J.; Bruskiewich, R.M.; Röttsch,  
1111 G.; Buell, C.R.; Leung, H.; Leach, J.E. Genome wide SNP variation reveals relationships among  
1112 landraces and modern varieties of rice. *Proc. Natl. Acad. Sci. USA* **2009**, *106*, 12273-12278.

- 1113 [125] Yamamoto, T.; Nagasaki, H.; Yonemaru, J.; Ebana, K.; Nakajima, M.; Shibaya, T.; Yano, M.  
1114 Definition of the pedigree haplotypes of closely related rice cultivars by means of genome-wide  
1115 discovery of single-nucleotide polymorphisms. *BMC Genomics* **2010**, *11*, 267.
- 1116 [126] Finkel E. Imaging. With 'phenomics,' plant scientists hope to shift breeding into overdrive.  
1117 *Science*, **2009**, *325*, 380-381.
- 1118 [127] Hyten, D.L.; Smith, J.R.; Frederick, R.D.; Tucker, M.L.; Song, Q.; Cregan, P.B. Bulk segregant  
1119 analysis using the GoldenGate assay to locate the *Rpp3* locus that confers resistance to soybean  
1120 rust in soybean. *Crop Sci.* **2009**, *49*, 265-271.
- 1121 [128] Kang, H.X.; Weng, Y.Q.; Yang, Y.H.; Zhang, Z.H.; Zhang, S.P.; Mao, Z.C.; Cheng, G.H.; Gu,  
1122 X.F.; Huang, S.W.; Xie, B.Y. Fine genetic mapping localizes cucumber scab resistance gene *Ccu*  
1123 into an R gene cluster. *Theor. Appl. Genet.* **2011**, *122*, 795-803.
- 1124 [129] Yu, H.H.; Xie, W.B.; Wang, J.; Xing, Y.Z.; Xu, C.G.; Li, X.H.; Xiao, J.H.; Zhang, Q.F. Gains in  
1125 QTL detection using an ultra-high density SNP map based on population sequencing relative to  
1126 traditional RFLP/SSR markers. *PLoS One* **2011**, *6* (3), e17595.
- 1127 [130] Rafalski JA. Novel genetic mapping tools in plants: SNPs and LD-based approaches. *Plant Sci.*  
1128 **2002**, *162*, 329-33.
- 1129 [131] Buckler, E.S.; Holland, J.B.; Bradbury, P.J.; Acharya, C.B.; Brown, P.J.; Browne, C.; Ersoz, E.;  
1130 Flint-Garcia, S.; Garcia, A.; Glaubitz, J.C.; Goodman, M.M.; Harjes, C.; Guill, K.; Koon, D.E.;  
1131 Larsson, S.; Lepak, N.K.; Li, H.H.; Mitchell, S.E.; Pressoir, G.; Peiffer, J.A.; Rosas, M.O.;  
1132 Rocheford, T.R.; Romay, M.C.; Romero, S.; Salvo, S.; Villeda, H.S.; da Silva, H.S.; Sun, Q.; Tian,  
1133 F.; Upadyayula, N.; Ware, D.; Yates, H.; Yu, J.M.; Zhang, Z.W.; Kresovich, S.; McMullen, M.D.  
1134 The genetic architecture of maize flowering time. *Science* **2009**, *325*, 714-718.
- 1135 [132] Andersen, J.R.; Zein, I.; Wenzel, G.; Darnhofer, B.; Eder, J.; Ouzunova, M.; Lubberstedt, T.  
1136 Characterization of phenylpropanoid pathway genes within European maize (*Zea mays* L.) inbreds.  
1137 *BMC Plant Biol.*, **2008**, *8*, 2.
- 1138 [133] Harjes, C.E.; Rocheford, T.R.; Bai, L.; Brutnell, T.P.; Kandianis, C.B.; Sowinski, S.G.; Stapleton,  
1139 A.E.; Vallabhaneni, R.; Williams M.; Wurtzel E.T.; Yan, J.B.; Buckler ES. Natural genetic  
1140 variation in lycopene epsilon cyclase tapped for maize biofortification. *Science* **2008**, *319*: 330-  
1141 333.



- 1142 [134] Zheng, P.; Allen, W.B.; Roesler, K.; Williams, M.E.; Zhang, S.; Li, J.; Glassman, K.; Ranch, J.;  
1143 Nubel, D.; Solawetz, W.; Bhatramakki, D.; Llaca, V.; Deschamps, S.; Zhong, G.Y.; Tarczynski,  
1144 M.C.; Shen, B. A phenylalanine in DGAT is a key determinant of oil content and composition in  
1145 maize. *Nature Genet.* **2008**, *40*: 367-372.
- 1146 [135] Manicacci, D.; Camus-Kulandaivelu, L.; Fourmann, M.; Arar, C.; Barrault, S.; Rousselet, A.;  
1147 Feminias, N.; Consoli, L.; Frances, L.; Mechin, V.; Murigneux, A.; Prioul, J.L.; Charcosset, A.;  
1148 Damerval, C. Epistatic interactions between Opaque2 transcriptional activator and its target gene  
1149 CyPPDK1 control kernel trait variation in maize. *Plant Physiol.* **2009**, *150*: 506-520.
- 1150 [136] Atwell, S.; Huang, Y.S.; Vilhjalmsson, B.J.; Willems, G.; Horton, M.; Li, Y.; Meng, D.Z.; Platt,  
1151 A.; Tarone, A.M.; Hu, T.T.; Jiang, R.; Mulyati, N.W.; Zhang, X.; Amer, M.A.; Baxter, I.; Brachi,  
1152 B.; Chory, J.; Dean, C.; Debieu, M.; de Meaux, J.; Ecker, J.R.; Faure, N.; Kniskern, J.M.; Jones,  
1153 J.D.G.; Michael, T.; Nemri, A.; Roux, F.; Salt, D.E.; Tang, C.L.; Todesco, M.; Traw, M.B.;  
1154 Weigel, D.; Marjoram, P.; Borevitz, J.O.; Bergelson, J.; Nordborg, M. Genome-wide association  
1155 study of 107 phenotypes in *Arabidopsis thaliana* inbred lines. *Nature* **2010**, *465*, 627-631.
- 1156 [137] Simko, I.; Pechenick, D.A.; McHale, L.K.; Truco, L.J.; Ochoa, O.E.; Michelmore, R.W.;  
1157 Scheffler, B.E. Association mapping and marker-assisted selection of the lettuce dieback  
1158 resistance gene *Tvr1*. *BMC Plant Biol.* **2009**, *9*: 135-151.
- 1159 [138] Andersen, J.R.; Lübberstedt, T. Functional markers in plants. *Trends Plant Sci.* **2003**, *8*, 554-560.
- 1160 [139] Hospital, F.; Chevalet, C.; Mulsant, P. Using markers in gene introgression breeding programs.  
1161 *Genetics*, **1992**, *132*, 1199-1210.
- 1162 [140] Gopalakrishnan, S.; Sharma, R.K.; Anand Rajkumar, K.A.; Joseph, M.; Singh, V.P.; Bhat, K.V.;  
1163 Singh, N.K.; Mohapatra, T. Integrating marker assisted background analysis with foreground  
1164 selection for identification of superior bacterial blight resistant recombinants in Basmati rice. *Plant*  
1165 *Breed.* **2008**, *127*, 131-139.
- 1166 [141] Cao, Z.; Tian, F.; Wang, N.; Jiang, C.; Lin, B.; Xia, W.; Shi, J.; Long, Y.; Zhang, C.; Meng, J.  
1167 Analysis of QTLs for erucic acid and oil content in seeds on A8 chromosome and the linkage drag  
1168 between the alleles for the two traits in *Brassica napus*. *J. Genet. Genomics* **2010**, *37*, 231-240.
- 1169 [142] Huang, N.; Angeles, E.R.; Domingo, J.; Magpantay, G.; Singh, S.; Zhang, G.; Kumaravadivel, N.;  
1170 Bennett, J.; Khush, G.S. Pyramiding of bacterial blight resistance genes in rice: marker-assisted  
1171 selection using RFLP and PCR. *Theor. Appl. Genet.* **1997**, *95*, 313-320.

- 1172 [143] Ishii, T.; Yonezawa, K. Optimization of the marker-based procedures for pyramiding genes from  
1173 multiple donor lines: II. strategies for selecting the objective homozygous plant. *Crop Sci.*, **2007**,  
1174 *47*, 1878-1886.
- 1175 [144] Wei, X.; Liu, L.; Xu, J.; Jiang, L.; Zhang, W.; Wang, J.; Zhai, H.; Wan, J. Breeding strategies for  
1176 optimum heading date using genotypic information in rice. *Mol Breeding* **2010**, *25*, 287-298.
- 1177 [145] Lü, H.Y.; Liu, X.F.; Wei, S.P.; Zhang, Y.M. Epistatic association mapping in homozygous crop  
1178 cultivars. *PLoS One* **2011**, *6* (3), e17773.
- 1179 [146] Jannink, J.L.; Lorenz, A.J. Iwata, H. Genomic selection in plant breeding: from theory to practice.  
1180 *Briefings Functional Genomics* **2010**, *9*, 166-177.
- 1181 [147] Meuwissen, T.H.E.; Hayes, B.J.; Goddard, M.E. Prediction of total genetic value using genome-  
1182 wide dense marker maps. *Genetics* **2001**, *157*, 1819-1829.
- 1183 [148] Heffner, E.L.; Sorrells, M.E.; Jannink, J.L. Genomic selection for crop improvement. *Crop Sci.*  
1184 **2009**, *49*: 1-12.
- 1185 [149] Bernardo, R. Genomewide selection for rapid introgression of exotic germplasm in maize. *Crop*  
1186 *Sci.* **2009**, *49*, 419-425.
- 1187 [150] Bernardo, R.; Yu, J. Prospects for genomewide selection for quantitative traits in maize. *Crop Sci.*,  
1188 **2007**, *47*, 1082-1090.
- 1189 [151] Lorenzana, R.E.; Bernardo, R. Accuracy of genotypic value predictions for marker-based selection  
1190 in biparental plant populations. *Theor. Appl. Genet.* **2009**, *120*: 151-161.
- 1191 [152] Bernardo, R. Molecular markers and selection for complex traits in plants: learning from the last  
1192 20 years. *Crop Sci.* **2008**, *48*, 1649-1664.
- 1193 [153] Araus, J.L.; Slafer, G.A.; Royo, C.; Serret, M.D. Breeding for yield potential and stress adaptation  
1194 in cereals. *Critical Rev. Plant Sci.* **2008**, *27*, 377-412.
- 1195

1196 Table 1. Comparison of the main characteristics of the conventional Sanger and some of the most  
1197 currently used next generation sequencing (NGS) technologies and approximate sequencing cost (in US \$  
1198 per Mbp).

| Technology          | Read length (bp) | Mbp per run | Cost (\$/Mbp) |
|---------------------|------------------|-------------|---------------|
| Sanger              | 1000             | 0.001       | 3000.00       |
| 454 Roche           | 450              | 450         | 66.00         |
| Illumina Hi-Seq2000 | 100              | 270000      | 0.07          |
| Solid 5500xl        | 50               | 270000      | 0.07          |

1199  
1200

1201 Table 2. Some important databases and repositories of genomic information of interest for breeders.

| Database           | Description                                  | URL   |
|--------------------|--|---|
| Genbank            | General public sequence repository           | <a href="http://www.ncbi.nlm.nih.gov/genbank/">http://www.ncbi.nlm.nih.gov/genbank/</a>                             |
| EMBL               | General public sequence repository           | <a href="http://www.ebi.ac.uk/embl/">http://www.ebi.ac.uk/embl/</a>   |
| DDBJ               | General public sequence repository           | <a href="http://www.ddbj.nig.ac.jp">http://www.ddbj.nig.ac.jp</a>   |
| UniProt            | Protein sequences and functional information | <a href="http://www.uniprot.org/">http://www.uniprot.org/</a>   |
| NCBI               | Biomedical and genomics information          | <a href="http://www.ncbi.nlm.nih.gov/">http://www.ncbi.nlm.nih.gov/</a>   |
| Gene Index Project | Transcriptome repository                     | <a href="http://compbio.dfci.harvard.edu/tgi/">http://compbio.dfci.harvard.edu/tgi/</a>                             |
| GOLD               | Repository of genomes databases              | <a href="http://genomesonline.org/cgi-bin/GOLD/bin/gold.cgi">http://genomesonline.org/cgi-bin/GOLD/bin/gold.cgi</a> |
| Phytozome          | Genomic plant database                       | <a href="http://www.phytozome.net/">http://www.phytozome.net/</a>   |
| Plantgdb           | Genomic plant database                       | <a href="http://www.plantgdb.org">http://www.plantgdb.org</a>   |
| CropNet            | Genomic plant database                       | <a href="http://ukcrop.net/">http://ukcrop.net/</a>   |
| SGN                | Solanaceae information resource              | <a href="http://solgenomics.net/">http://solgenomics.net/</a>   |
| Gramene            | Grass information resource                   | <a href="http://www.gramene.org/">http://www.gramene.org/</a>   |
| MaizeGDB           | Maize information resource                   | <a href="http://www.maizegdb.org/">http://www.maizegdb.org/</a>   |
| Tair               | Arabidopsis information resource             | <a href="http://www.arabidopsis.org/">http://www.arabidopsis.org/</a>   |
| CottonDB           | Cotton information resource                  | <a href="http://cottondb.org/">http://cottondb.org/</a>   |
| CPGR               | Phytopathogen genomic resource               | <a href="http://cpgr.plantbiology.msu.edu/">http://cpgr.plantbiology.msu.edu/</a>                                   |

1202

1203

1204

1205 Table 3. Some examples of the utility of molecular markers developed by means of high-throughput  
 1206 genomics techniques for the breeding of important crops.

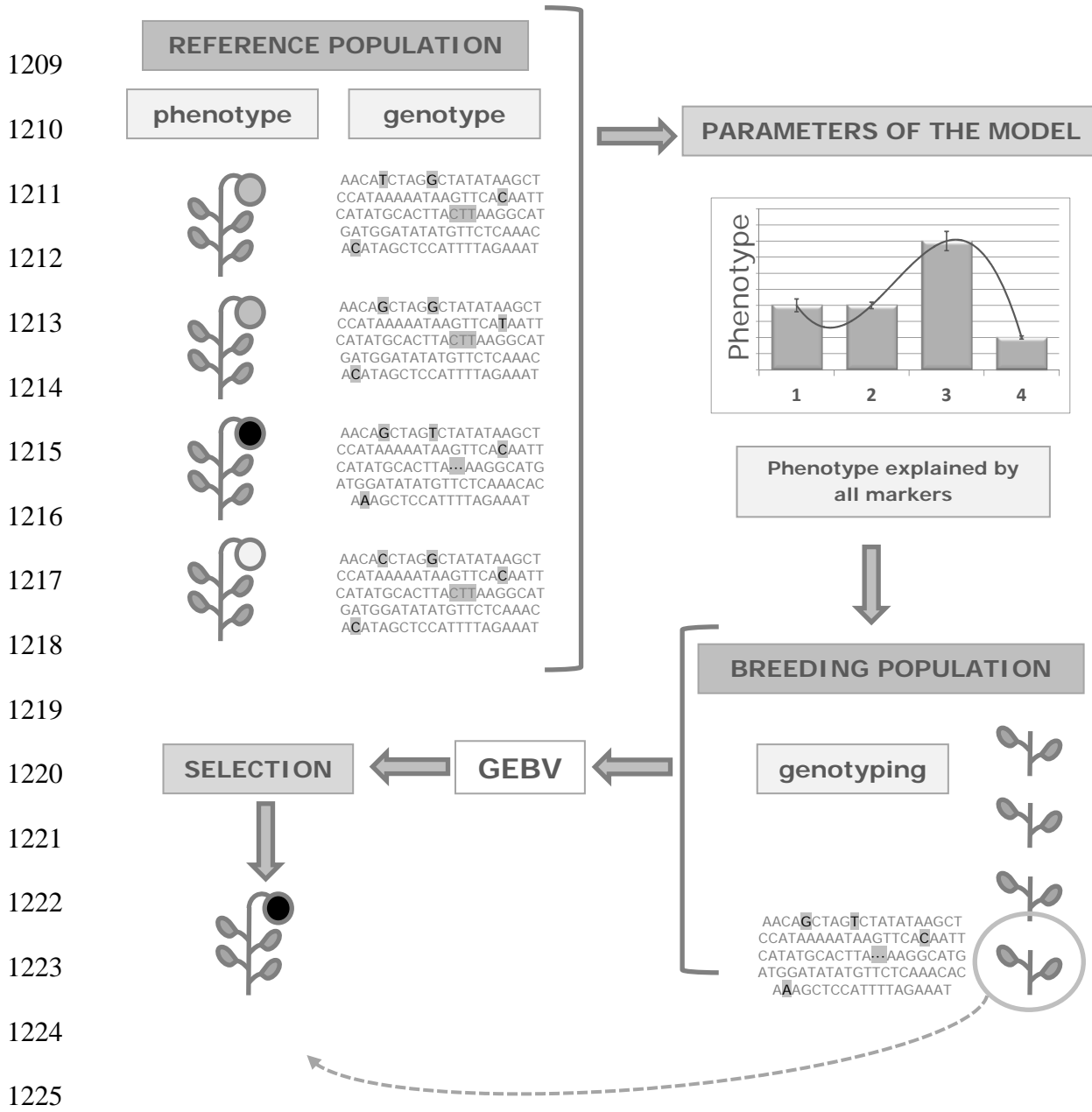
| Crop                                    | Markers                      | Plant material  | Use for breeding   | Reference |
|---|------------------------------|---|--|-----------|
| Rice ( <i>Oryza sativa</i> )            | ~3.6·10 <sup>6</sup><br>SNPs | 517 rice landraces  | Association studies for 14<br>agronomic traits                       | [96]      |
| Barley ( <i>Hordeum<br/>vulgare</i> )   | 1,536 SNPs                   | 768 breeding lines  | Association studies for<br><i>Fusarium</i> head<br>blight resistance | [97]      |
|   | 3,072 SNPs                   | 336 DH lines and<br>213 germplasm<br>selections                               | High-density genetic map<br>construction and MAF<br>estimation       | [98]      |
| Maize ( <i>Zea mays</i> )               | 8,590 SNPs                   | 553 elite maize<br>inbred lines   | Association studies for oleic<br>acid content                        | [99]      |
|   | 1,106 SNPs                   | 5,000 RILs  | Association studies for<br>resistance to southern<br>leaf blight     | [100]     |
|   | 1,536 SNPs                   | 154 maize inbred<br>lines   | Diversity studies  | [101]     |
| Grapevine ( <i>Vitis<br/>vinifera</i> ) | 94 SNPs and 7<br>indels      | 148 grape varieties   | Association studies for muscat<br>flavor candidate gene VvDXS        | [102]     |
|   | 9000 SNPs                    | 10 cultivated <i>Vitis</i><br><i>vinifera</i> and 7 wild<br><i>Vitis</i> spp. | Diversity and<br>population structure studies                        | [74]      |
| Pea ( <i>Pisum sativum</i> )            | 384 SNPs                     | 91 RIL mapping<br>population and 373<br><i>Pisum</i> accessions               | Linkage map construction and<br>diversity studies.                   | [103]     |
| Wheat ( <i>Triticum<br/>aestivum</i> )  | 874 DArT<br>markers          | winter<br>wheat core<br>collection of 96<br>accessions                        | Association studies for 20<br>agronomic traits                       | [104]     |

|  |            |   |   |       |
|--|------------|---|---|-------|
|  | 1,536 SNPs | 478 spring and<br>winter wheat<br>cultivars | Diversity studies   | [105] |
| White spruce ( <i>Picea<br/>glauca</i> ) | 944 SNPs   | 492 individuals                             | Association studies with<br>549 candidate genes and 25<br>wood quality traits | [106] |

---

1207

1208



1226 Figure 1. Genomic selection scheme. Information on phenotype and genotype for a reference population  
 1227 allows estimating parameters for the model. This model explains phenotype based on all markers  
 1228 analyzed. The model predicts the phenotype of plants in a breeding population on the basis of the  
 1229 genotyping results: this is the genomic estimated breeding value (GEBV), used to select the desired  
 1230 phenotypes.

1231