TECHNICAL ADVANCE

Generation of a 3D indexed *Petunia* insertion database for reverse genetics

Michiel Vandenbussche¹, Antoine Janssen², Jan Zethof¹, Nathalie van Orsouw², Janny Peters¹, Michiel J.T. van Eijk², Anneke S. Rijpkema¹, Harrie Schneiders², Parthasarathy Santhanam¹, Mark de Been³, Arjen van Tunen² and Tom Gerats^{1,2,*} ¹Radboud University, IWWR/Plant Genetics, Toernooiveld 1, 6525 ED Nijmegen, The Netherlands, ²Keygene NV, AgroBusiness Park 90, 6708 PW Wageningen, The Netherlands, and ³Centre for Molecular and Biomolecular Informatics (CMBI), Radboud University Nijmegen Medical Centre, PO Box 9101, 6500 HB Nijmegen, The Netherlands

Received 23 November 2007; revised 1 February 2008; accepted 18 February 2008. *For correspondence (fax +31 2436 52787; e-mail t.gerats@science.ru.nl).

Summary

BLAST searchable databases containing insertion flanking sequences have revolutionized reverse genetics in plant research. The development of such databases has so far been limited to a small number of model species and normally requires extensive labour input. Here we describe a highly efficient and widely applicable method that we adapted to identify unique transposon-flanking genomic sequences in Petunia. The procedure is based on a multi-dimensional pooling strategy for the collection of DNA samples; up to thousands of different templates are amplified from each of the DNA pools separately, and knowledge of their source is safeguarded by the use of pool-specific (sample) identification tags in one of the amplification primers. All products are combined into a single sample that is subsequently used as a template for unidirectional pyrosequencing. Computational analysis of the clustered sequence output allows automatic assignment of sequences to individual DNA sources. We have amplified and analysed transposon-flanking sequences from a Petunia transposon insertion library of 1000 individuals. Using 30 DNA isolations, 70 PCR reactions and two GS20 sequencing runs, we were able to allocate around 10 000 transposon flanking sequences to specific plants in the library. These sequences have been organized in a database that can be BLAST-searched for insertions into genes of interest. As a proof of concept, we have performed an in silico screen for insertions into members of the NAM/NAC transcription factor family. All in silico-predicted transposon insertions into members of this family could be confirmed in planta.

Keywords: insertion mutagenesis, Petunia, pyrosequencing, transposon, dTph1, NAC family.

Introduction

One of the classical challenges in biology is to understand how changes in gene function during evolution have contributed to the diversity of life as it exists today. To study the evolution of gene function, it is necessary to compare loss-of-function mutants for orthologous genes in different species. In Arabidopsis, the publicly available collections of insertion site-sequenced T-DNA lines (Alonso *et al.*, 2003) have revolutionized reverse-genetics approaches, as databases with insertion site flanking sequences can be searched for mutants by a simple sequence homology search instead of having to perform laborious PCR-based assays. Similar platforms are being developed for the monocot species rice and maize (Droc *et al.*, 2006; van Enckevort *et al.*, 2005; Kolesnik *et al.*, 2004; Settles *et al.*, 2007; Yazaki *et al.*, 2004). Plant comparative biology research would greatly benefit if such resources could be developed for a broader range of species. We wish to create such a platform for the model species *Petunia*, taking advantage of the endogenous, highly active *dTph1* transposable element system in line W138 (Gerats *et al.*, 1990). Petunia has several qualities that make

it a fairly amenable model system, amongst which are highly efficient forward and reverse transposon insertion screening methods (Gerats and Vandenbussche, 2005).

As a starting point for further improvements in these methods, we aimed to simultaneously sequence and order dTph1 flanking sequences amplified from a population of 1000 individuals, by using the powerful 454 liquid-phase massive parallel sequencing technique (Margulies *et al.*, 2005). For reverse-genetics purposes, the origin of every flanking sequence tag in the resulting database needs to be assigned to a specific individual within the population. As 454 sequencing in principle does not allow discrimination between individual DNA sources in a pooled sample, we have developed and applied a general methodology (Figure 1) that allows massive parallel amplification and



sequencing of thousands of different fragments that are automatically assigned to their specific origin. We have coupled the use of 5'-encoded PCR primers (see also Binladen *et al.*, 2007; van Orsouw *et al.*, 2007) with a threedimensional pooling system, and have adapted the original GS20 pyrosequencing protocol (Margulies *et al.*, 2005) to enable unidirectional sequencing of the amplified transposon tag library, followed by automated assignment of sequences to individuals or families within the population by a computational approach.

The three-dimensional (3D) pooling principle (Figure 1) allows identification of a specific individual within a group of up to several thousands of individuals in a single, smallscale PCR experiment (Zwaal et al., 1993) as commonly applied in PCR-based Petunia insertion mutagenesis screens (Koes et al., 1995; Vandenbussche et al., 2003). The DNA samples prepared from the 3D pooled plant material each serve as a separate template for PCR amplification of the desired targets. A different amplification primer is used for every PCR reaction, carrying a unique four-nucleotide sample identification tag at the 5' end (Keygene[™] SegTag; van Orsouw et al., 2007). The PCR products are subsequently pooled, and unidirectionally sequenced using a modified GS20 protocol. The resulting sequences can be clustered according to identity and reduced to a single database line by 'reading' the various sequence identification tags encountered within each cluster. Querying the resulting database with the predefined profiles of individuals and families allows automatic assignment of sequences to their source.

For the *Petunia* experiment presented here, a classic 3D pooling system was used (Koes *et al.*, 1995; Vandenbussche *et al.*, 2003). The size of the analysed population, 1000

Figure 1. Multi-dimensional indexed massive parallel amplification and sequencing principle.

DNA templates are obtained by organizing the source material in a virtual 3D cubic set-up and pooling material from each dimension before DNA extraction. For example, all individuals in the x1 plane are pooled together, and this is performed similarly for all planes in all dimensions. Each individual DNA source in such a cube is therefore represented in three pools, once in every dimension, where the pool numbers reflect the spatial positioning of the DNA source in the 3D grid. Such 3D coordinate sets can also be translated to binary code as a series of ones and zeroes in a database query line, indicating the presence or absence of individual DNA templates in every pool of the library. Small groups of related individuals, such as families, can be profiled in the same way (e.g. famB). Templates are PCR-amplified from every DNA pool separately, using amplification primers containing pool number information encrypted as a four-nucleotide 'barcode' at its 5' end. These amplification products are pooled together in one sample and sequenced by 454 sequencing. The resulting sequences are then clustered according to homology in the region between the amplification primers. Next, these clusters are then converted to single database lines containing a consensus sequence based on all sequences within each cluster, and the 5' four-nucleotide sample identification tags are back-translated/converted to a series of numbers indicating the pool numbers encountered within each cluster. These database lines are then cross-searched against the originally defined profiles of individuals and known families within the population, leading to automatic annotation of the resulting clusters.

individuals, organized in a $10 \times 10 \times 10$ grid, is moderate, but nevertheless represents a highly complex sample due to the high copy number of the endogenous Petunia dTph1 transposon system. Individual W138 plants, the standard line used in Petunia insertion mutagenesis screens, can contain up to an estimated 200 dTph1 transposable element copies in either a homozygous or heterozygous state (Gerats et al., 1990). Every generation, 5-10% of these transpose to a new location, creating a new insertion event, present in a single heterozygous individual (De Keukeleire et al., 2001; Van den Broeck et al., 1998). Upon selfing, such an insertion will segregate in the next generation. Thus, a population of 1000 plants could harbor 10 000-20 000 new insertions, hemizygously present in a single individual. In addition, as all W138 plants derive from a single plant (Doodeman et al., 1984), ancestral insertions shared between all individuals are also encountered. A *dTph1* insertion library is usually composed of a number of small families (e.g. 100 families of 10 plants each), with every family originating from a self of an individual W138 plant. As such, unique insertions that occurred in the parental lines used to set up a library now segregate in their respective progenies. These are called family insertions, and are of great interest for the coupling of forward phenotypic screens with candidate gene insertion sequences.

Due to the complex genealogy of W138 insertion libraries, resulting in unique, family and ancestral insertion sites, considerable differences in copy number per insertion site exist within the population: the most extreme difference (1:2000) is between a new unique insertion event, which is hemizygously present in a single plant, and the immobile ancestral transposon insertion sites homozygous in all 1000 individuals. Family insertion copy numbers vary between these extremes, depending on the family size and the relatedness of the various families within the library. The mass amplification and sequencing of a complete *dThp1* insertion library is thus comparable with sequencing a non-normalized EST library, in which highly abundant, moderately expressed and minimally expressed ESTs correspond to ancestral, family and unique insertion events. As a result, the unique insertion flanking sequences are amplified much less efficiently than high copy number transposon flanking sequences. For reverse-genetics purposes, the high copy number class is of limited value, as it represents a small (and relatively immobile) set of insertion loci amplified from many to all individuals of the library. A minimum of around 50 such ancestral immobile loci were estimated to be present in previous populations (De Keukeleire et al., 2001); therefore, this class of insertions unfortunately forms a large proportion of the total number of *dTph1* loci in a population $(50 \times 2 \times 1000 = 100\ 000)$. Therefore, the total number of sequences required to sufficiently resolve and annotate all low copy number flanking sequences (10 000-20 000) is

Petunia 3D indexed database for reverse genetics 3

very large, as omnipresent insertion sequences compete significantly for total sequencing capacity. Full coverage of all insertion flanking sequences in a limited number of sequence runs is therefore not expected. Moreover, a number of sequence clusters will remain unresolved, as one or more of the 3D coordinates will be missing. To at least partially reduce this problem, we implemented a normalization step for the amplified library before sequencing, in order to diminish the proportion of sequences derived from ancestral *dTph1* loci.

Results

3D indexed mass amplification and sequencing of Petunia dTph1 *transposon flanking sequences*

A graphical overview of the 3D indexed mass-amplification and sequence procedure, customized for analysis of Petunia dTph1 transposon flanking sequences, is shown in Figure S1, and technical details are described in Experimental procedures. Briefly, we have applied a modified transposon display protocol (Van den Broeck et al., 1998) to mass-amplify *dTph1* transposon flanking sequences from the 30 DNA pools harvested from the W138 library (Table S1). Four-nucleotide tags encrypting the 30 coordinates (Table S2) were incorporated in the transposon-derived primer used in the 2nd round of PCR amplification. The obtained PCR products were pooled for each dimension (x, y and z), resulting in three samples that were then subjected to one round of normalization using a hydroxyapatite column purification procedure (Bonaldo et al., 1996). The selected single-stranded fraction, enriched for low copy number sequences, was made double-stranded by primer extension, and subsequently digested to enable GS20/454 adapter ligation. Msel and Mun sites were added to the transposon primer (containing the four-nucleotide sample identification tag) and the extension primer, respectively. The GS20/454 adapters A and B (Table S2) were modified to have sticky Msel and Munl ends instead of blunt ends, allowing unidirectional ligation of the sequencing annealing site in adapter A, immediately adjacent to the four-nucleotide tag in the transposon-specific primer. As a consequence, all amplified fragments originating from the same template have the same orientation, facilitating subsequent computational clustering. Msel and Munl were initially used for genomic DNA digestion, guaranteeing that no Msel and Munl sites will be present within the amplified transposon flanking sequences during digestion. The adapter-ligated mixture was then used as a template for a final PCR amplification step using amplification primers A and B (Table S2). A 4 μ g aliquot of this mixture was then used as a template for two GS20 sequence runs, performed according to the manufacturer's instructions.

Clustering and automatic annotation of sequences to individual DNA sources

The combined output of the two GS20 runs yielded a total of 659 219 sequences, with an average read length of 101 bp. A BLAST search (Altschul et al., 1990) was used for detection of the 5' coded Petunia dTph1 primer (NNNN-IRoutw primer, Table S2) at the start of the read; this could be recognized in 94.8% of the sequences. This result indicates that the large majority of the sequenced amplicons are derived from transposon flanking loci. The coordinates of the start of the BLAST hit were subsequently checked for the four-base sample identification tag. Using a PERL script, this tag could be automatically translated in one of the 30 possible pool coordinates in 96.4% of the sequences, and incorporated into a database line. Sequences with unrecognizable sample identification tags were marked as 'unknown' and removed at a later stage. Finally, primer and sample identification tag sequences were trimmed and all resulting sequences shorter than 20 bp were removed from the dataset. Such fragments originate from transposon insertion sites with an *Mse*l restriction site in close proximity to the *dTph1* insertion site and are of limited value for in silico screenings due to their short read length.

Clustering of the remaining dataset (611 000 sequences in total) resulted in the generation of 82 570 unique clusters, each with a variable number of homologous sequences. Each individual cluster was subsequently converted into a single database line containing a consensus sequence and all numerical data required for further classification and annotation. For each pool coordinate, the number of fragments encountered within the cluster carrying that coordinate was summed, and the total number of fragments in each cluster was counted. A distribution analysis based on these numbers showed that extreme differences in copy number exist between clusters, ranging from one to over 6000 copies (Figure 2a), as expected due to simultaneous amplification of unique and highly redundant insertion sites in a large population. Around 25% of the sequences (611 000 in total) belong to one of 113 clusters, with a copy number varying from 500 up to 6800. These clusters represent the so-called ancestral insertions that are present in the majority of the individuals within the population. These redundant insertion loci make up a considerable part of the total sequence capacity, but represent only a limited number of insertion loci. The largest number of clusters can be found in the lower copy number classes, corresponding to unique (both somatic and germinal) and recent family insertions, with the latter segregating in a limited number of related individuals within the library. For further analysis, the number of detected pools per pool axis (x, y and z) was calculated. With the latter information, sequence clusters were further classified and annotated according to 1D, 2D, 3D and mD (multiD) coordinate classes (Figure 2b). At this

stage, 2759 of the clusters were discarded as they were composed solely of sequences previously marked as 'unknown' due to unreadable sequence tags. The 1D and 2D classes represent incomplete coordinate sets, missing a coordinate for two and one dimensions, respectively. The existence of incomplete coordinate sets is caused by either somatic insertion events or is due to under-sampling of the sequence library (see Discussion). On the other hand, 3D class clusters have exactly one x, one y and one z coordinate, and represent unique insertions present in one individual of the library. The mD class is composed of clusters having at least one x, y and z coordinate, excluding the 3D class clusters. These mD clusters therefore contain more than one pool coordinate in at least one of the dimensions, and correspond to insertion loci shared between at least two related individuals in the population.

Due to the 3D sampling strategy of the library, a sequence can only be assigned to an individual plant if it belongs to a cluster in which every dimension (x, y and z) is represented by at least one pool coordinate. All clusters fitting this criterion (3D and mD classes, 13 375 clusters in total) were therefore organized in a sub-database for final annotation.

In this set, a total of 5881 3D class clusters (Figure 2c), could be identified, which were automatically assigned a plant number between 1 and 1000 using the following formula:

$$((x-1) \times 10) + ((y-1) \times 1) + ((z-1) \times 100) + 1$$

with *x*, *y* and *z* all varying between 1 and 10 according to the library sampling used (Table S1).

The remaining set (mD class) represents so-called family insertions segregating in at least two related individuals within the library. As the relatedness of the lines in the library had been well documented during the previous two generations, a 'profile' for every family of plants in the library could be constructed as a translation of its spatial positioning within the library set-up (Table S1). Separate search profiles were created for small families and large families, for which the members display mainly first- and second-degree relatedness, respectively. These gueries (see Experimental procedures) resulted in unambiguous annotation of around 4000 insertions as segregating in either one of the defined small or large families (Figure 2c, Tables S3 and S4). In addition, a small number of clusters did not contain sufficient information to be assigned to either of the two related search profiles. For these cases, a double annotation was applied, combining both possibilities in the sequence header. The remaining set of unassigned clusters was not further analysed, and presumably contains family insertions reflecting ancestral and undocumented relationships between individuals and families within the library. In a final step, only the annotated insertion flanking sequences from lines for which sufficient seeds are available for distribution were selected, and grouped together in a final BLAST-



Petunia 3D indexed database for reverse genetics 5

Figure 2. Distribution analysis of 3D indexed mass-amplified and sequenced *dTph1* transposon flanking sequence library. (a) Cluster distribution analysis versus copy number classes.

(b) Number of clusters versus coordinate classes. Clusters have been classified according to the four possible dimensional classes. Question marks in the coordinates indicate the absence of pool coordinates for the respective dimension. The 'unknown' class indicates clusters containing only non-recognizable sequence tags.

(c) Number of assigned insertions versus cluster copy number classes. Total numbers of assigned insertions are indicated per category (total number analysed), together with the average copy number of clusters representing each type of insertion class.

searchable database, currently containing around 9500 non-redundant annotated transposon flanking sequences.

In silico screens for insertions into genes of interest

© 2008 The Authors

To evaluate whether the resulting annotated transposon flanking sequence database represents a reliable mutant identification source, we have performed an *in silico* screen for insertions into members of the *NAM/NAC* transcription factor family (Souer *et al.*, 1996). *NAC* genes have been identified as important regulators involved in various aspects of plant development (reviewed by Olsen *et al.*, 2005), and are easily recognized by their highly conserved NAC domain. The 22 known *Petunia NAC* family members were extracted from Genbank and subjected to a BLAST search against the insertion flanking sequence database. The search resulted in identification of nine putative insertions into eight *Petunia NAC* members (Figure 3a and Table S5). All these *dTph1* insertions reside in coding regions, which usually results in loss-of-function mutations. Five of these insertions correspond to known *Petunia NAC* members, while the remaining four represent insertions into

				de	non.	ert in	3 mot	ationatio		
(a)	Genelal	lele Acct	ession of	ert poe	atabasi atabasi	CR SCREEPIN	Proger	.,ч ^{сон} Х	Υ	Z
(u)	PhNH3-1	AF509866	323	3D	3D	<u>#912</u>	yes			
	PhNH3-2	AF509866	227	mD	mD	<u>F154/103</u>	yes			
I	PhNH10-1	AF509873	649	3D	3D	<u>#572</u>	yes			
1	PhNH17-1	AF510214	813	mD	mD	<u>F284</u>	yes			
I	PhNH22-1	AF510218	211	3D	3D	<u>#205</u>	yes			
1	PhNH23-1	EU249322	188	3D	3D	<u>#247</u>	yes	1 2 3 4 5 6 7 8 9 10	1 2 3 4 5 6 7 8 9 10	1 2 3 4 5 6 7 8 9 10
I	PhNH24-1	EU249325	1	mD	mD	<u>F167</u>	yes	1 2 3 4 5 6 7 8 9 10	1 2 3 4 5 6 7 8 9 10	1 2 3 4 5 6 7 8 9 10
1	PhNH25-1	EU249324	1	3D	3D	<u>#322</u>	yes			
1	PhNH26-1	EU249323	1	mD	mD	<u>F280/84</u>	yes		1 2 3 4 5 6 7 8 9 10	1 2 3 4 5 6 7 8 9 10
(1)									(2)	
(D)	PhNAM-1	X92204	490	1D	mD		yes		1 2 3 4 5 6 7 8 9 10	1 2 3 4 5 6 7 8 9 10
	PhNAM	X92204	705	1D	2D		no	1 2 3 4 5 6 7 8 9 10	1 2 3 4 5 6 7 8 9 10	1 2 3 4 5 6 7 8 9 10 (1)
	PhNH3	AF509866	236	1D	1D		no	1 2 3 4 5 6 7 8 9 10	1 2 3 4 5 6 7 8 9 10	1 2 3 4 5 6 7 8 9 10 (1)
I	PhNH16-1 /PhCUC3	AF510213	511	2D	3D		yes			
	PhNH17-2	AF510214	732	2D	3D		yes			
1	PhNH27-1	EU249326	1	2D	3D		yes	1 2 3 4 5 6 7 8 9 10	1 2 3 4 5 6 7 8 9 10	1 2 3 4 5 6 7 8 9 10

 $\mathbf{\Delta}$

-0

Figure 3. *In silico* and classic PCR screening for insertions into members of the *Petunia NAC* transcription factor family for complete gel images see Figure S2. (a) Overview of all *in silico*-identified and automatically assigned insertions into *NAC* genes.

(b) Selection of the tested *in silico*-identified *NAC* insertions with incomplete (non-assignable) coordinate sets. The numbers above the gel images represent the 30 pool samples derived from the 3D pooling of 1000 individuals. Numbers in parentheses above the pool numbering reflect *in silico* screening results and show the copy number of transposon flanking sequences with corresponding pool numbers as encountered in the sequence collection. The results of the automatic annotation are shown as plant numbers for unique insertions, or as family names (F code) for segregating insertions. The gel images show the results of a conventional PCR screening as described previously (Vandenbussche *et al.*, 2003) using gene-specific primers designed to flank the predicted insertion sites. Insertion positions are relative to the start of corresponding Genbank accessions. 1D, 2D, 3D and mD indicate the coordinate sets after database screening (*in silico*) and classic PCR screening. Assignable coordinate sets have been boxed. New NAC members isolated during this screening are underlined.

new *NAC* family members (named PhNH23–26), based on sequence conservation in their NAM domain (Figure S3), illustrating additional benefits of the database as a gene discovery source. The coordinates of five of the nine putative insertions identify a unique single plant (3D class insertion) within the population, while the remaining four insertions represent so-called family insertions (mD class), segregating in a small group of genetically related plants. In addition, three putative insertion events were identified that had one of their coordinates missing (2D hits), while 15 putative insertions were identified as 1D sets (not shown). Incomplete coordinate sets may represent somatic insertion events, present in only one or two of the three samples taken from each individual, or might represent true germinal insertion events for which one or more of the coordinates is missing for technical reasons. To investigate this further, we included three 2D and three 1D sets in the analysis presented below.

As an independent confirmation of the automatic assignment of insertions to individuals and small families, we screened the 30 DNA pools of the library using a conventional gene-specific screening method (Vandenbussche et al., 2003) for each of the nine automatically assigned insertions (Figure 3a) and six selected incomplete sets (Figure 3b) (see also Figure S2a-e). All in silico-assigned insertion events could be confirmed by PCR screening, as shown by the appearance of specific amplification products in all of the identified pool coordinates. Further, presumably due to the higher sensitivity of the gene-specific screening method, the missing coordinate for all three 2D sets could be determined (identifying a new gene: PhNH27-1; Figure 3b), but this was not possible for two of the three 1D hits (Figure 3b). These results confirm both the occurrence of somatic insertion events, and the incompleteness of clusters due to under-sampling of the sequence library. Similarly, for three family insertions, additional positive pools were identified that were not found in the sequence database; these pools fitted within the previously identified family profile. These results clearly indicate that not all templates are amplified or not all amplified templates are represented in our current sequence collection. Finally, we were able to confirm the presence of the respective insertion events in progenies of all positive individuals or families identified in this screen (results not shown); a number of these lines are currently being analysed for phenotypes.

Discussion

The use of 5' coded primers has been proposed previously to distinguish a limited number of DNA sources within a 454 sequence template (Binladen *et al.*, 2007; van Orsouw *et al.*, 2007). Here we show that the combined use of 5' coded primers with a multidimensional pooling and annotation strategy results in a extremely powerful method, enabling highly complex templates to be analysed by 454 sequencing while the source is still identifiable, and with a minimal labour investment in terms of template preparation and processing. This approach can be applied to any system for which multiple templates from multiple sources have to be analysed on a truly large scale. The generation of hundreds of thousands of sequences per 454 run provides sufficient over-sampling capacity to analyse and annotate tens of thousands of templates simultaneously.

We developed and applied this procedure to generate a database of annotated transposon flanking sequences, mass-amplified from a *Petunia* W138 library of 1000 individuals, which can then be used for *in silico* reverse-genetics screens. However, the nature of the endogenous

dTph1 transposon system in *Petunia* created some extra complicating factors.

The extreme differences in copy number between ancestral and new unique insertion sites within such a library made it necessary to normalize the samples before sequencing. Even so, 25% of the final sequence reads appeared to be derived from insertions present in many to all individuals of the library. Other normalization procedures might turn out to be more efficient (see, for example, Cheung et al., 2006; Shagin et al., 2002). Other undesired sequence tags are those that are derived from somatic insertions and thus will not be transmitted to the next generation. The presence of these sequences in the final database would lead to false identification of putative insertion lines. In Petunia, it has been demonstrated that a 3D pooling strategy allows discrimination between somatic and germinal insertion events (Koes et al., 1995). Unfortunately, a 3D indexed pooling strategy implies that a sequence can only be annotated if it has been amplified and sequenced at least three times independently in the total sequence collection, at least once from every dimensional pool.

It is obvious that, with the same total sequence capacity, a 2D pooling strategy would have resulted in a larger number of annotated sequences, as two coordinates each time would have sufficed to appoint a sequence to a source. We nevertheless chose a 3D indexed mass amplification and sequencing strategy, so as to be able to exclude all somatically derived transposon flanking sequences from the database. This problem is specific to *Petunia*; for most other applications, a 2D pooling strategy might be a better choice.

Despite the complications inherent in *dTph1* transposon biology, we identified almost 10 000 different insertion flanking sequences that could be automatically assigned to specific individuals of the library using only 30 DNA preparations, 70 PCR reactions and two GS20 sequencing runs. Almost 6000 (5881) of the sequence tags represent unique insertions present in only one individual of the library. Using defined family history profiles based on three generations of known family relationships between the various individuals, we could automatically annotate a further 4000 insertions as segregating in small families in the library.

The successful assignment of these family insertions in this study also indicates that the methodology is well suited to large-scale parallel mapping of family-specific genotypes in heterogeneous populations. Combined with the fast progress in 454 sequencing technology in terms of both capacity and read length, and with further improvement and application of the proposed strategy, we expect that creation of a saturated insertion sequence library is realistic for *Petunia* as well as for other model systems. As a starting point, we have organized the approximately 10 000 identified insertion sequence tags in a database. BLAST searches for insertions into genes of interest will be performed on request, and seeds of positively identified insertion lines will

be provided free for academic purposes. We expect this to be a valuable resource for functional genomics research in *Petunia* and other Solanaceae.

Experimental procedures

Library set-up

Seeds representing 1000 progenies from self-pollinations (mainly W138 but also W138/Mitchell hybrids) were sown, and usually one individual of each progeny was grown for the transposon library. These 1000 plants were grown under standard greenhouse conditions in trays, and after harvesting for DNA sampling were transferred to individual pots and maintained until sufficient seeds were produced, resulting from manual pollinations. A schematic representation of the library set-up is shown in Table S1. Equal colouring within blocks (*z* dimension) between individuals indicates a family relationship in the 2nd degree. Identical parental accession codes at different positions indicate direct sibling relationships.

DNA preparation

The 1000 selected library plants were sampled in a 3D fashion, resulting in 30 DNA samples representing 100 plants each, as described previously (Koes *et al.*, 1995; Vandenbussche *et al.*, 2003). The DNA samples were further treated using a modified version of the transposon display method originally developed to selectively amplify all dTph1 flanking sequences from various samples of a one single plant (Van den Broeck *et al.*, 1998).

A 5 μ g DNA aliquot for each pool sample was incubated with *Mun*l (Biolabs; http://www.neb.com/nebecomm) and *Msel* (Biolabs; http://www.neb.com/nebecomm) for 1.5 h at 37°C in 70 μ l restriction mix (5 units *Mun*l, 5 units *Msel*, 1× NEB4 and 1× BSA). Adapters were ligated by incubation for 3 h at 37°C after addition of 30 μ l ligation mix (15 Weiss units of T4 DNA ligase (Fermentas; http:// www.fermentas.com), 40 pmol biotinylated *Mun*l adapter, 400 pmol *Msel* adapter, 1 mM ATP, 1× NEB4 and 1× BSA).

Selection of biotinylated DNA fragments to enrich for transposon flanking fragments

To remove the excess of biotinylated adapters, DNA samples were purified using the QIAquick PCR purification kit (Qiagen, http:// www.qiagen.com/) and DNA was eluted using 50 μ I EB buffer (supplied elution buffer; Qiagen PCR purification kit). DNA samples were incubated with approximately 0.1 mg MyOne streptavidine beads C1 (Dynal; http://www.invitrogen.com/dynal) for 1 h on a rotator at room temperature, in 500 μ I binding buffer (10 mM Tris–CI pH 8, 2 μ NaCl, 1 mM EDTA and 0.1% Triton X-100). Beads were collected with a magnet and the supernatant removed. Beads were washed once using 200 μ I STEX buffer (10 mM Tris–CI pH 8, 1 μ NaCl, 1 mM EDTA and 0.1% Triton X-100), and transferred to another tube prior to three additional washes with 200 μ I STEX buffer each. Finally the beads were resuspended in 50 μ I T₀₁E buffer (10 mM Tris–CI pH 8, 0.1 mM EDTA) and transferred to another tube.

Selective pre-amplification

To enrich for transposon-flanking sequences, DNA-bead samples were subjected to two rounds of PCR amplification. For the first PCR

amplification, 2 µl of DNA-bead sample was mixed with 6 pmol *Mun*l-ACAC- and 6 pmol *Mse*l primer (Table S2) in a PCR buffer with 0.6 units of Red Hot DNA polymerase (Abgene; http://www.abgene. com), 7.5 mM Tris–Cl pH 8.8, 20 mM (NH₄)₂SO₄, 0.01% Tween-20, 2.5 mM MgCl₂ and 0.2 mM dNTPs in a final volume of 20 µl. The samples were incubated in a PE9600 (Perkin Elmer; http://www. appliedbiosystems.com) according to the following PCR profile: one cycle comprising a 15 sec step at 94°C, a 30 sec step at 65°C and a 60 sec step at 72°C; the annealing temperature was then lowered each cycle by 0.7°C for 13 cycles, and then kept at 56°C for another 22 cycles. The resulting PCR products were checked by electrophoresis of 5 µl of reaction mix on a 1.5% agarose gel (a weak low-molecular-weight smear should be visible). The remaining 15 µl was diluted 10 times with H₂O.

Specific PCR amplification of dTph1 flanking sequences

To amplify the *dTph1* flanking sequences and to incorporate the pool-specific four-nucleotide sample identification tag, 5 μ l of the 10-fold diluted pre-amplified material was mixed with 15 pmol NNNN-IR_{outw} primer and 15 pmol Mse primer (Table S2) in a PCR buffer with 1 unit Red Hot DNA polymerase (Abgene), 7.5 mM Tris-Cl pH 8.8, 20 mM (NH₄)₂SO₄, 0.01% Tween-20, 2.5 mM MgCl₂ and 0.2 mM dNTPs in a final volume of 50 μ l. The samples were incubated in a PE9600 (Perkin Elmer) according to the following PCR profile: one cycle comprising a 15 sec step at 94°C, a 30 sec step at 65°C and a 60 sec step at 72°C; the annealing temperature was then lowered each cycle by 0.7°C for 13 cycles, and then kept at 56°C for another 22 cycles.

The resulting PCR products from the ten samples from each dimension were pooled to create three samples: an x, y and z sample.

Normalization

In order to optimize the amount of unique fragments against a backdrop of fragments shared by many to all individuals, we normalized the three pooled samples. The procedure involves hybridization and purification steps to obtain single-stranded DNA. The approach relies on the differential re-annealing kinetics between highly abundant and rare DNA templates, by which the single-stranded fraction will become enriched for low copy number templates. The single-stranded fraction was isolated by hydroxyapatite chromatography (Bonaldo *et al.*, 1996).

The pooled PCR products (x, y and z samples, containing)around 10 μg of DNA each) were precipitated and dissolved in 30 μ l formamide. The DNA was melted at 80°C for 3 min (under mineral oil) and hybridized for 16 h at 30°C, as described in Bonaldo et al. (1996), in a final volume of 50 µl. Thereafter, the samples were diluted by addition of 500 µl pre-warmed (60°C) loading buffer (0.04 M sodium phosphate buffer, pH 6.8), and single-stranded DNA was purified by hydroxyapatite chromatography at 60°C, using a 0.1 ml column pre-equilibrated with the loading buffer. After a 2 ml wash with loading buffer, the singlestranded DNA was eluted from the column using eight repeats of 100 µl pre-warmed 0.12 м sodium phosphate buffer (pH 6.8), and elution fractions 2-8 were pooled. The phosphate buffer was removed using nucleospin columns (Macherey Nagel; http:// www.mn-net.com). The DNA was eluted four times using 200 μ l NE buffer, precipitated with iso-propanol, and finally dissolved in 50 μ l H₂O for each elution. The effectiveness of the normalization was verified by a PCR competition experiment between a high copy number flanking sequence and a low copy sequence, and a transposon display experiment on normalized and non-normalized samples (results not shown). In view of the results, further optimization of the normalization would be useful, and could be achieved by varying parameters or using other procedures (see, for example, Cheung *et al.*, 2006; Shagin *et al.*, 2002).

GS20 template preparation and sequencing

Primer extension. For conversion to double-stranded DNA, 50 µl of single-stranded DNA was mixed with 50 pmol extension primer (Table S2) in a PCR buffer with 1 unit Platinum Taq DNA polymerase (Invitrogen, http://www.invitrogen.com/), 20 mm Tris–Cl pH 8.4, 50 mm KCl, 1.5 mm MgCl₂ and 0.2 mm dNTPs in a final volume of 75 µl. The three samples were then incubated in a PE9600 (Perkin Elmer) for 2 min at 94°C, 1 min at 56°C and 10 min at 72°C to complete the primer extension.

Digestion and GS20 adapter ligation. The three normalized (double-stranded) DNA pools were digested with *Mun*l and *Mse*l to enable directional GS20 adapter ligation. The total amount of normalized double-stranded DNA was incubated for 1.5 h at 37°C in 90 μ l restriction mix (20 units *Mun*l, 20 units *Mse*l, 1× NEB4 and 1× BSA). The GS20 adapters were ligated by incubation for 4 h at 37°C after addition of 20 μ l ligation mix (15 Weiss units of T4 DNA ligase, 200 pmol GS20-*Mun*l adapter B, 200 pmol GS20-*Mse*l adapter A, 1 mm ATP, 1× NEB4 and 1× BSA).

If there is enough DNA (a few micrograms; when normalization isn't required), it is possible to ligate the BioTEG-GS20-*Mun*l adapter B (Table S2) instead of the non-biotinylated adapter. In this case, the samples can be pooled after adapter ligation to create one 'superpool', and are ready for the 'library immobilization step' of the GS20 sequencing protocol.

PCR re-amplification before sequencing. After normalization, we re-amplified the ligation mixture to obtain sufficient and biotinylated DNA for GS20 sequencing. Five microlitres of template were mixed with 15 pmol of adapter primer A and 15 pmol BioTEG-adapter B primer (Table S2; http://www.biolegio.com) in a PCR buffer with 1 unit of Red Hot DNA polymerase (Abgene), 7.5 mM Tris–Cl pH 8.8, 20 mm (NH₄)₂SO₄, 0.01% Tween-20, 2 mM MgCl₂ and 0.2 mm dNTPs in a final volume of 50 µl. The samples were incubated in a PE9600 (Perkin Elmer) according to the following PCR profile: one cycle with a 15 sec step at 94°C, a 30 sec step at 65°C and a 60 sec step at 72°C; the annealing temperature was then lowered each cycle by 0.7°C for 13 cycles, and then kept at 56°C for another 22 cycles.

The samples were pooled to create one superpool. Approximately 4 μ g of the superpool fragments were used as input for GS20 library construction. Several adaptations to the standard library preparation protocol were made. First, no shearing, end-polishing, phosphorylation or ligation of the A and B adapters was carried out; the samples entered at the 'library immobilization step' of the protocol. Second, the Bst DNA polymerase fill-in step of the protocol was omitted.

Sequencing. After library construction, a titration run was performed to estimate the optimal amount of single-stranded fragments containing the A and the B adapter for the actual sequencing run. After the titration run, emulsion PCR and bead enrichment were carried out according to the standard GS20 protocol (Roche Applied Science; http://www.roche-applied-science.com). Two full picotitre plates (70×75 mm) with two regions each were used. Sequencing was performed according to the manufacturer's instructions (Roche Applied Science).

Clustering of the GS20 sequence output

The cleaned and trimmed dataset was clustered using a combined approach. The first step was a fast clustering using CDHIT (Weizhong *et al.*, 2001), using a percentage identity of 93%. From each resulting cluster, a consensus fragment was generated as a representative. All CDHIT cluster representatives were further iteratively clustered using BLASTCLUST (Altschul *et al.*, 1990) (percentage identity of 93%), until the number of clusters did not further decrease.

Annotation of segregating insertions in the library

For each of the 30 pool coordinates, these small and large family profiles (Table S3 and S4) contain the value *x* (where the number of fragments within the cluster assigned to the given pool can be 0 or more), 0 (where the number of fragments within the cluster assigned to the given pool is 0), >1 (which is the number of fragments within the cluster assigned to the given pool) or *A* (where the sum of the number of fragments encountered within the pools with this value *A* is 1).

A number of search profiles, especially in the large families, display closely related patterns, due to inter-mixed positional arrangement of the corresponding families within the library. For a more unambiguous distinction between these patterns, we included the operators 'A' and '>1' in the search profiles of some families.

The search profiles were translated into SQL queries using a PERL script. The script executed the queries, and all resulting entries were annotated as a family insertion (small or large class). This step was consecutively performed for small and large families, and only the remaining non-annotated insertions were searched for in the 2nd round. This strict order of extraction is necessary as the majority of small family profiles reside within larger family profiles.

Sequence deposition

Sequence data for *PhNH23*, *PhNH24*, *PhNH25*, *PhNH26* and *PhNH27* have been deposited with the EMBL/Genbank libraries under accession numbers EU249322–EU249326.

Patents and trademarks

The method for high-throughput screening of transposon tagging populations using massively parallel sequence identification of insertion sites and KeyGene[™] SeqTag technology are covered by patents and patent applications owned by Keygene NV (http:// www.keygene.com). Application for trademark registration for Seq Tag technology has been filed by Keygene NV.

Acknowledgements

We thank Hans Sommer and Zsuzsanna Schwarz-Sommer (both Max Planck Institute, Cologne, Germany) for sharing their normalization protocol. M.V. was funded by a HORIZON grant (050-71-036) from the Netherlands Genomics Initiative. A.S.R. was funded by the Netherlands Organization for Scientific Research (grant 814.02.009).

Supplementary Material

The following supplementary material is available for this article online:

Figure S1. Overview of 3D indexed mass amplification and directional sequencing of *Petunia* right border *dTph1* flanking fragments. **Figure S2.** PCR confirmation of *in silico*-predicted *NAC* insertion alleles for 3D pooled DNA samples from the library.

Figure S3. Protein alignment of translated *dTph1* flanking sequences representing new *Petunia NAC* members, together with a selection of known *Petunia NAC* genes.

Table S1. Library set-up.

Table S2. Oligo and adapter sequences.

 Table S3. Small family search profiles.

Table S4. Large family search profiles.

 Table S5. Database representation of the nine in silico-identified

 NAC insertions.

This material is available as part of the online article from http:// www.blackwell-synergy.com

Please note: Blackwell Publishing are not responsible for the content or functionality of any supplementary materials supplied by the authors. Any queries (other than missing material) should be directed to the corresponding author for the article.

References

- Alonso, J.M., Stepanova, A.N., Leisse, T.J. et al. (2003) Genomewide insertional mutagenesis of Arabidopsis thaliana. Science, 301, 653–657.
- Altschul, S.F., Gish, W., Miller, W., Myers, E.W. and Lipman, D.J. (1990) Basic local alignment search tool. *J. Mol. Biol.* **215**, 403– 410.
- Binladen, J., Gilbert, M.T., Bollback, J.P., Panitz, F., Bendixen, C., Nielsen, R. and Willerslev, E. (2007) The use of coded PCR primers enables high-throughput sequencing of multiple homolog amplification products by 454 parallel sequencing. *PLoS*, 14, e197.
- Bonaldo, M., Lennon, G. and Bento Soares, M. (1996) Normalization and subtraction: two approaches to facilitate gene discovery. *Genome Res.* 6, 791–806.
- Cheung, F., Haas, B.J., Goldberg, S.M.D., May, G.D., Xiao, Y. and Town, C.D. (2006) Sequencing *Medicago truncatula* expressed sequenced tags using 454 life sciences technology. *BMC Genomics*, 7, 272–282.
- De Keukeleire, P., Maes, T., Sauer, M., Zethof, J., Van Montagu, M. and Gerats, T. (2001) Analysis by transposon display of the behavior of the dTph1 element family during ontogeny and inbreeding of *Petunia hybrida*. *Mol. Genet. Genomics*, 265, 72–81.
- Doodeman, M., Boersma, E.A., Koomen, W. and Bianchi, F. (1984) Genetic analysis of instability in *Petunia hybrida*: a highly unstable mutation induced by a transposable element in the An1 locus for flower colour. *Theor. Appl. Genet.* **67**, 345–355.
- Droc, G., Ruiz, M., Larmande, P., Pereira, A., Piffanelli, P., Morel, J.B., Dievart, A., Courtois, B., Guiderdoni, E. and Perin, C. (2006) OryGenesDB: a database for rice reverse genetics. *Nucleic Acids Res.* 34, D736–D740.
- van Enckevort, L.J., Droc, G., Piffanelli, P. et al. (2005) EU-OSTID: a collection of transposon insertional mutants for functional genomics in rice. Plant Mol. Biol. 59, 99–110.

- Gerats, T. and Vandenbussche, M. (2005) A model for comparative research: Petunia. *Trends Plant Sci.* **10**, 251–256.
- Gerats, A.G., Huits, H., Vrijlandt, E., Marana, C., Souer, E. and Beld, M. (1990) Molecular characterization of a nonautonomous transposable element (dTph1) of petunia. *Plant Cell*, 2, 1121–1128.
- Koes, R., Souer, E., Houwelingen, A. et al. (1995) Targeted gene inactivation in Petunia by PCR-based selection of transposon insertion mutants. Proc. Natl Acad. Sci. USA, 92, 8149– 8153.
- Kolesnik, T., Szeverenyi, I., Bachmann, D., Kumar, C.S., Jiang, S., Ramamoorthy, R., Cai, M., Ma, Z.G., Sundaresan, V. and Ramachandran, S. (2004) Establishing an efficient Ac/Ds tagging system in rice: large-scale analysis of Ds flanking sequences. *Plant J.* 37, 301–314.
- Margulies, M., Egholm, M., Altman, W.E. et al. (2005) Genome sequencing in microfabricated high-density picolitre reactors. *Nature*, **437**, 376–380.
- Olsen, A.N., Ernst, H.A., Leggio, L.L. and Skriver, K. (2005) NAC transcription factors: structurally distinct, functionally diverse. *Trends Plant Sci.* **10**, 79–87.
- van Orsouw, N.J., Hogers, R.C.J., Janssen, A. *et al.* (2007) Complexity reduction of polymorphic sequences (CRoPS[™]): a novel approach for large-scale polymorphism discovery in complex genomes. *PLoS ONE*, **2**(11), e1172.
- Settles, A.M., Holding, D.R., Tan, B.C. et al. (2007) Sequenceindexed mutations in maize using the UniformMu transposontagging population. *BMC Genomics*, 8, 116.
- Shagin, D.A., Rebrikov, D.V., Kozhemyako, V.B., Altshuler, I.M., Shcheglov, A.S., Zhulidov, P.A., Bogdanova, E.A., Staroverov, D.B., Rasskazov, V.A. and Lukyanov, S. (2002) A novel method for SNP detection using a new duplex-specific nuclease from crab hepatopancreas. *Genome Res.* 12, 1935–1942.
- Souer, E., van Houwelingen, A., Kloos, D., Mol, J. and Koes, R. (1996) The no apical meristem gene of Petunia is required for pattern formation in embryos and flowers and is expressed at meristem and primordia boundaries. *Cell*, 85, 159–170.
- Van den Broeck, D., Maes, T., Sauer, M., Zethof, J., De Keukeleire, P., D'Hauw, M., Van Montagu, M. and Gerats, T. (1998) Transposon display identifies individual transposable elements in high copy number lines. *Plant J.* **13**, 121–129.
- Vandenbussche, M., Zethof, J., Souer, E., Koes, R., Tornielli, G.B., Pezzotti, M., Ferrario, S., Angenent, G.C. and Gerats, T. (2003) Toward the analysis of the Petunia MADS box gene family by reverse and forward transposon insertion mutagenesis approaches: B, C, and D floral organ identity functions require SEPALLATA-like MADS box genes in Petunia. *Plant Cell*, **15**, 2680– 2693.
- Weizhong, L., Jaroszewski, L. and Godzik, A. (2001) Clustering of highly homologous sequences to reduce the size of large protein database. *Bioinformatics*, 17, 282–283.
- Yazaki, J., Kojima, K., Suzuki, K., Kishimoto, N. and Kikuchi, S. (2004) The rice PIPELINE: a unification tool for plant functional genomics. *Nucleic Acids Res.* **32**, D383–D387.
- Zwaal, R.R., Broeks, A., van Meurs, J., Groenen, J.T. and Plasterk, R.H. (1993) Target-selected gene inactivation in *Caenorhabditis elegans* by using a frozen transposon insertion mutant bank. *Proc. Natl Acad. Sci. USA*, **90**, 7431–7435.

The EMBL/Genbank accession numbers for the PhNH23-27 sequences are EU249322-EU249326, respectively.



Supplementary Figure1: Overview of 3D-indexed mass amplification and directional sequencing of petunia right border *dtph1* flanking fragments.

Supplementary Figures 2A-E: PCR confirmation of nam insertion alleles on 3D pooled DNA samples of the library, complete gel images.

M = Size marker

Arrows indicate fragments representing NAC insertion alleles. Sequences of the gene specific primers used in these screenings are provided in Supplementary Table 2. Corresponding insertion alleles are indicated.

Fig. 2A Χ У Ζ Ζ Х У Μ Μ Μ Μ Μ Μ Μ Μ PhNAM-1 PhNH26-1 PhNAM

Fig. 2B



Fig. 2C



Fig. 2D







Supplementary Figure 3: Protein alignment of translated *dTph1* flanking sequences representing new *Petunia* NAC members together with a selection of known Petunia NAC genes.

Supplementary Table 1: Library set-up



Supplementary Table 2: Oligo and adapter sequences

Oligo Name	Sequence (5'-3')
Bio-MunI adapter:	
Bio-Mun top	BIOTIN-CTCGTACACTACC
Mun-hottom	
	AATTCGTACGCAGTC
Msel adapter:	
Mse-top	GACGATGAGTCCTGAG
Mse-bottom	TACTCAGGACTCAT
MunI-ACAC primer	AGACTGTGTACGAATTGACAC
MseI primer	GACGATGAGTCCTGAGTAA
Inverted repeat primers:	
NNNN-Iroutw primer	CATATATTAANNNNGTAGCTCCGCCCCTG
NNNN-IRoutw-pool X1 primer	CATATATTAAACACGTAGCTCCGCCCCTG
NNNN-IRoutw-pool X2 primer	CATATATTAA ACAG GTAGCTCCGCCCCTG
NNNN-IRoutw-pool X3 primer	CATATATTAA ACGA GTAGCTCCGCCCCTG
NNNN-IRoutw-pool X4 primer	CATATATTAA ACGT GTAGCTCCGCCCCTG
NNNN-IRoutw-pool X5 primer	CATATATTAA ACTC GTAGCTCCGCCCCTG
NNNN-IRoutw-pool X6 primer	CATATATTAAACTGGTAGCTCCGCCCCTG
NNNN-IRoutw-pool X7 primer	
NNNN-IRoutw-pool X0 primer	
NNNN-IRoutw-pool X10 primer	CATATATAAAGCTGTAGCTCCGCCCCTG
NNNN-IRoutw-pool Y1 primer	CATATATTAAAGTCGTAGCTCCGCCCCTG
NNNN-IRoutw-pool Y2 primer	CATATATTAAAGTGGTAGCTCCGCCCCTG
NNNN-IRoutw-pool Y3 primer	CATATATTAA ATCG GTAGCTCCGCCCCTG
NNNN-IRoutw-pool Y4 primer	CATATATTAA ATGC GTAGCTCCGCCCCTG
NNNN-IRoutw-pool Y5 primer	CATATATTAACACAGTAGCTCCGCCCCTG
NNNN-IRoutw-pool Y6 primer	CATATATTAA CACT GTAGCTCCGCCCCTG
NNNN-IRoutw-pool Y7 primer	CATATATTAA CAGA GTAGCTCCGCCCCTG
NNNN-IRoutw-pool Y8 primer	CATATATTAA CAGT GTAGCTCCGCCCCTG
NNNN-IRoutw-pool Y9 primer	CATATATTAACATCGTAGCTCCGCCCCTG
NNNN-IRoutw-pool Y10 primer	CATATATTAACATGGTAGCTCCGCCCCTG
NNNN-IRoutw-pool Z1 primer	
NNNN-IROULW-pool 22 primer	
NNNN-IROUTW-pool Z3 primer	
NNNN IRoutw-pool Z5 primer	CATATATTAACTCAGTAGCTCCGCCCCTG
NNNN-IRoutw-pool Z6 primer	CATATATTAA CTCT GTAGCTCCGCCCCTG
NNNN-IRoutw-pool Z7 primer	CATATATTAA CTGA GTAGCTCCGCCCCTG
NNNN-IRoutw-pool Z8 primer	CATATATTAA CTGT GTAGCTCCGCCCCTG
NNNN-IRoutw-pool Z9 primer	CATATATTAA GACA GTAGCTCCGCCCCTG
NNNN-IRoutw-pool Z10 primer	CATATATTAA GACT GTAGCTCCGCCCCTG
extension primer	CATATACAATTGGACGATGAGTCCTGAGTAA
GS20 MunI adapter B:	
GS20 - Mun - top	CCTATCCCCTGTGTGCCTTGCCTTATCCCCCTGTTGCGTGTCTCAG
GS20 - Mun - bottom	A A TTCTCA CA CA CCCA A CA CCCA A A CA CA CA CA
	AATTCTGAGACACGCAACAGGGGATAGGCAAGGCACACAGGGGA
GS20 MSel adapter A:	
GS20-Mse-top	CCATCTCATCCCTGCGTGTCCCCATCTGTTCCCTCCCTGTCTCAG
GS20-Mse-bottom	TACTGAGACAGGGAGGGAACAGATGGGACACGCAGGGATGAG
Amplification adapter primer A & B:	
adapter primer B	
adapter primer A	
sequence primer A	CCATCTGTTCCCTCCCTGTC
NAM screening:	
IRoutw	GAATTCGCTCCGCCCCTG
PhNAM_REV	TTTGTTGAGGTCAACTTCAGCAATG
PhNH26_REV	
PNNH3_REV	ATACTCGTGCATTATCCAGTTAGTC
	AILIGAIAACTGUTUTGAAAUTGTG CCCATTTTCCATACTTTCTAAAUTGTG
PhNH15 FW	
PhNH16/PhCUC3 FW	CAAGCAGTTGGACAAAATTGTGG
PhNH17 FW	ACTACAAGCAAGATCAGCGGTTAC
PhNH22 REV	ACACAAGTGTTTTCTTCATGCCAAC
PhNH24 REV	TGGCCATGGGGAGCTCTTCC
PhNH27_REV/PhNH25_REV	CCTGTTGCCTTCCAGAATCCT

Total indiv.	FamilyName	X1 X2	X3	X4 🛛	X5 X	6 X7	X8	X9 (X10	Y1	Y2	Y3	Y4	Y5	Y6	Y7	Y8	Y9 Y1	Z1	Z2	Z3	Z4	Z5	Z6	Z7	Z8 2	Z9 Z10
2	fam229-7	x 0	0	0	0 ×	0	0	0	0	x	0	0	0	0	0	0	0	0 0	x	0	0	0	0	0	0	0	0 0
2	fam229-8/fam229-10	x 0	0	0	0 x	0	0	0	0	0	x	0	0	0	0	0	0	0 0	x	0	0	0	0	0	0	0	0 0
3	fam229-10	x A	0	0	0 x	0	0	0	0	0	x	0	A	0	0	0	0	0 0	x	0	0	0	0	0	0	0	0 0
2	fam220.0 w		0	0	0	0	ő	0	0	0	0	v	0	õ	0	0	0	0 0	Û	ñ	0	0	õ	0	0	0	0 0
2	fom220_11_W		0	0		0	0	0	0	0	0	_	0	0	0	0	0	0 0		0	0	0	0	0	0	0	0 0
3	1a11229-11-W		0	0		0	0	0	0	0	0	0	0	~	×	0	0	0 0	×	0	0	0	0	0	0	0	0 0
3	fam229-16	x 0	0	0	0 ×	0	0	0	0	0	0	0	0	0	Х	X	0	0 0	х	0	0	0	0	0	0	0	0 0
2	fam229-17	x 0	0	0	0 ×	0	0	0	0	0	0	0	0	0	0	0	Х	0 0	Х	0	0	0	0	0	0	0	0 0
4	famK218-6	хх	0	Х	0 0	0	0	X	0	0	0	0	0	0	0	0	0	х х	х	0	0	0	0	0	0	0	0 0
2	fam229-23-W	0 x	0	0	0 0	X	0	0	0	0	0	х	0	0	0	0	0	0 0	X	0	0	0	0	0	0	0	0 0
2	fam229-24-W	0 x	0	0	0 0	x	0	0	0	0	0	0	х	0	0	0	0	0 0	x	0	0	0	0	0	0	0	0 0
2	fam113-7-W/fam113-4-W	0 0	х	0	0 0	0	х	0	0	0	0	0	0	0	Х	0	0	0 0	x	0	0	0	0	0	0	0	0 0
3	fam113-4-W	0 A	x	0	0 0	0	x	0	0	0	0	0	0	0	x	А	0	0 0	x	0	0	0	0	0	0	0	0 0
3	fam113_3_W	0 0	Ŷ	0 0	0 0	Y	Ŷ	Ő	0	Ő	0 0	õ	ñ	0	0	0	Y	0 0	Ŷ	ñ	Ő	0	õ	Ő	ñ	ñ	0 0
5	fam113.2 W	V O		0	0 0		Û	v	0	0	0	0	0	0	ň	0	0		Û	0	0	0	0	0	0	0	0 0
5	fami(210.0			0	0 0		^		0	0	0	0	0	0	0	0	0			0	0	0	0	0	0	0	0 0
3	Iamk218-8		0	x	0 0	0	0	X	0	0	0	0	0	0	0	0	x	0 0	x	0	0	0	0	0	0	0	0 0
2	TamK218-7	0 0	0	X	X U	0	0	0	0	0	0	0	0	0	0	0	X	X 0	X	0	0	0	0	0	0	0	0 0
3	famK262-9-W	0 x	0	0	x 0	0	0	0	х	Х	0	0	0	0	0	0	0	0 0	х	0	0	0	0	0	0	0	0 0
2	famK262-3-W	0 0	0	0	x 0	0	0	0	х	0	Х	0	0	0	0	0	0	0 0	х	0	0	0	0	0	0	0	0 0
2	famK262-2	0 0	0	0	x 0	0	0	0	х	0	0	х	0	0	0	0	0	0 0	Х	0	0	0	0	0	0	0	0 0
2	famK262-1	0 0	0	0	x 0	0	0	0	x	0	0	0	х	0	0	0	0	0 0	x	0	0	0	0	0	0	0	0 0
4	fam212-10-W	0 x	0	0	0 0	X	х	0	0	0	0	0	0	х	Х	0	0	0 0	x	0	0	0	0	0	0	0	0 0
2	fam212-18	0 x	0	0	0 0	x	0	0	0	0	0	0	0	0	0	X	0	0 0	x	0	0	0	0	0	0	0	0 0
2	fam216-6-12D	0 0	0	X	0 0	0	0	X	0	0	0	0	0 I	x	0	0	0	0 0	x	0	0	0	0	0	0	0	0 0
3	famK234	0 0	× ×	0	0 0	n n	n	0	õ	ñ	n I	γ	y y	Y	ñ	õ	0	0 0	¥	ñ	õ	0	ñ	õ	õ	0	0 0
2	famK111_6_\//		0	×.	0 0		n	y l	0	v	0	0	0	0	0	õ	ñ	0 0	Û	0	0	ñ	0	0	õ	ñ	0 0
2	fomk111 = W/		0	<u> </u>	0 0		0	<u>,</u>	0	ĉ	0	0	0	0	0	0	0	0 0	Û.	0	0	0	0	0	0	0	0 0
2			0	~	0 0		U	×	0		~	U	0	0	0	0	0	0 0	×	0	0	0	0	0	0	0	
3	famK111-3-W	0 0	0	x	0 0	0	X	x	0	0	0	х	0	0	0	0	0	0 0	X	0	0	0	0	0	0	0	0 0
3	famK111-1-W	0 0	0	x	0 0	0	Х	X	0	0	0	0	Х	0	0	0	0	0 0	х	0	0	0	0	0	0	0	0 0
4	famK111-7-W	0 0	Х	Х	0 ×	X	Х	0	0	0	0	0	0	0	0	0	0	0 x	х	0	0	0	0	0	0	0	0 0
4	fam198-15/212-5	0 0	0	х	x 0	0	0	X	0	0	0	0	0	0	Х	x	0	x 0	x	0	0	0	0	0	0	0	0 0
2	famE210-9	0 0	0	0	x 0	0	0	0	x	0	0	0	0	x	0	0	0	0 0	x	0	0	0	0	0	0	0	0 0
2	famE210-5	0 0	0	0	x 0	0	0	0	x	0	0	0	0	0	0	х	0	0 0	x	0	0	0	0	0	0	0	0 0
2	famK270-6-W	0 0	х	0	0 0	0	х	0	0	х	0	0	0	0	0	0	0	0 0	x	0	0	0	0	0	0	0	0 0
3	famK272?-7-W	0 0	x	0	0 0	x	x	0	0	0	x	0	0	0	0	0	0	0 0	x	0	0	0	0	0	0	0	0 0
2	fam255-25	0 0	0	v	0 0		0	v	0	0	0	ñ	ň	v	l ñ	ñ	ñ	0 0	0	v	n n	ñ	ñ	ñ	ñ	ñ	0 0
2	fom218 21v270 E	0 0	0	Û -			0	0	0	0	0	0	0	0	0	0	0	0 0	0	Û	0	0	0	0	0	0	0 0
3	lalli318-21x270-3	0 0	0	×			0	0	<u>`</u>	0	0	0	×	0	0	0	0	0 0	0	×	0	0	0	0	0	0	0 0
2	fam270-2	0 0	0	0	X U	0	0	0	х	0	0	0	0	0	0	Х	0	0 0	0	х	0	0	0	0	0	0	0 0
2	fam270-7	0 0	0	0	x 0	0	0	0	х	0	0	0	0	0	0	0	0	x 0	0	х	0	0	0	0	0	0	0 0
2	fam270-8/270-3	0 0	0	0	x 0	0	0	0	х	0	0	0	0	0	0	0	0	0 x	0	х	0	0	0	0	0	0	0 0
3	fam270-3	0 0	Α	0	x 0	0	0	0	x	0	0	0	0	0	0	0	А	0 x	0	х	0	0	0	0	0	0	0 0
3	fam318-22	0 0	x	0	x 0	0	0	0	x	0	0	0	0	x	0	0	x	0 0	0	х	0	0	0	0	0	0	0 0
3	fam116-14x120-8	x 0	x	0	0 0	0	х	0	0	х	0	0	0	0	0	0	0	0 0	0	x	0	0	0	0	0	0	0 0
2	fam116-15x120-8	0 0	x	0	0 0	0	x	0	0	0	х	0	0	0	0	0	0	0 0	0	x	0	0	0	0	0	0	0 0
2	fam116-13x120-8	0 x	0	0	0 0	x	0	0	0	0	0	0	0	0	0	0	0	0 x	0	x	0	0	0	0	0	0	0 0
3	fam116-23x100-8	V O	v	0	0 0		v	l n	0	n i	v	0	ñ	0	v	0	0		0	Ŷ	0	0	0	0	0	0	0 0
5	fam110-20x100-0		<u> </u>	0	0 0		<u>.</u>	0	0		0	0	0	~	_	0	0	0 0	0	<u></u>	0	0	0	0	0	0	0 0
3	1211116-182199-1	X U	X	0	0 0		X	0	0	0	0	x	0	0	0	0	0	0 0	0	x	0	0	0	0	0	0	0 0
5	fam 199-1	X X	X	0	0 0	0	X	0	0	0	0	х	0	0	0	0	0	0 0	0	x	0	0	0	0	0	0	0 0
3	fam116-20x199-10	x 0	X	0	0 0	0	Х	0	0	0	0	0	Х	0	0	0	0	0 0	0	х	0	0	0	0	0	0	0 0
3	fam116-8x120-6	хх	0	0	0 0	X	0	0	0	0	0	0	0	Х	х	0	0	0 0	0	х	0	0	0	0	0	0	0 0
2	fam120-9	x 0	0	0	0 ×	0	0	0	0	0	0	0	0	0	Х	0	0	0 0	0	х	0	0	0	0	0	0	0 0
4	fam120-12	x 0	0	0	0 ×	0	0	0	0	0	0	0	0	0	0	х	x	0 x	0	х	0	0	0	0	0	0	0 0
2	fam120-23	0 x	0	0	0 0	x	0	0	0	x	0	0	0	0	0	0	0	0 0	0	Х	0	0	0	0	0	0	0 0
3	fam116-1x199-9	0 x	0	0	0	x	0	0	0	0	Х	Х	0	0	0	0	0	0 0	0	х	0	0	0	0	0	0	0 0
2	fam116-3x199-1	0 x	0	0	0 0	x	0	0	0	0	0	x	0	0	0	0	0	0 0	0	x	0	0	0	0	0	0	0 0
3	fam116-4x199-17	0	0	0	0	x	0	0	0	0	0	0	X	0	0	0	0	0 0	0	X	0	0	0	0	0	0	0 0
3	fam116-5x199-10		0 0	0 0	0	Ŷ	ő	ñ	0	Ő	0 0	0	0	Ŷ	l ñ	0	ñ	0 0	0 0	Ŷ	0 0	0	õ	Ő	ñ	ñ	0 0
3	fam116 10x100 2		0	0	0	Ĵ	ő	0	0	0	0	0	0	0	0	v	0	0 0	0	Û	0	0	0	0	0	0	0 0
3	fom116 11 x100 10		0	0		÷.	0	0	0	0	0	0	0	0	0	^	0	0 0	0	<u> </u>	0	0	0	0	0	0	0 0
2	1a11110-11-X199-10		0	0	0 0		0	0	0	0	0	0	0	0	0	0	~	0 0	0	X	0	0	0	0	0	0	0 0
2	fam116-12	U X	0	0	0 0	X	0	0	0	0	0	0	0	0	0	0	0	X U	0	x	0	0	0	0	0	0	0 0
2	fam116-13x120-8	0 x	0	0	0 0	X	0	0	0	0	0	0	0	0	0	0	0	0 x	0	х	0	0	0	0	0	0	0 0
2	fam116-21-W	0 0	X	0	0 0	0	X	0	0	0	0	0	0	Х	0	0	0	0 0	0	Х	0	0	0	0	0	0	0 0
2	fam116-24x199-16	0 0	X	0	0 0	0	Х	0	0	0	0	0	0	0	0	X	0	0 0	0	Х	0	0	0	0	0	0	0 0
2	fam295-8	0 0	x	0	0 0	0	X	0	0	0	0	0	0	0	0	0	0	X 0	0	Х	0	0	0	0	0	0	0 0
3	fam318-16x198-7	0 0	0	х	X 0	0	0	0	х	х	0	0	0	0	0	0	0	0 0	0	х	0	0	0	0	0	0	0 0
3	fam318-20	0 0	0	x	X C	0	0	0	x	0	X	х	0	0	0	0	0	0 0	0	х	0	0	0	0	0	0	0 0
3	fam295-18-W	0 0	0	x	0 0	0	0	X.	0	0	X	х	0	0	0	0	0	0 0	0	X	0	0	0	0	0	0	0 0
2	fam318-9-W	0 0	0	x	0 0	n n	ñ	×.	0	0	0	0	0	0	x.	0	0	0 0	ñ	x	0	0	0	0	0	0	0 0
2	fam318-10-W/		0	X	0 0		ñ	Ŷ.	ñ	ň	ñ	ñ	ñ	0	0	×	õ	0 0	ñ	Y.	ñ	õ	ñ	ñ	õ	õ	0 0
2	fam210 11 M/		0	Ŷ	0 0		0	Û	0	0	0	0	0	0	0	0	v	0 0	0	Ŷ	0	0	0	0	0	0	0 0
2	fom240_42		0	~	0 0		0		0		0	0	0	0	0	0	^			~	0	0	0	0	0	0	
2	iam318-13	0 0	U	x	0 0		0	X	U	0	U	U	U	U	U	U	U	× 0	U	X	U	U	U	U	U	U	0 0
2	tam318-15	0 0	0	X	υΟ	0	0	X	0	0	Ú	0	0	0	0	0	0	0 x	0	х	0	0	0	0	0	0	υ 0
2	fam318-17-W	0 0	0	0	X C	0	0	0	Х	0	X	0	0	0	0	0	0	0 0	0	Х	0	0	0	0	0	0	0 0
2	fam318-24	0 0	0	0	X 0	0	0	0	X	0	0	0	0	0	Х	0	0	0 0	0	Х	0	0	0	0	0	0	0 0
2	fam304-5-W	x x	0	0	0 0	0	0	0	0	Х	0	0	0	0	0	0	0	0 0	0	0	Х	0	0	0	0	0	0 0

Supplementary Table 3: Small family search profiles

<u> </u>	famili 0			
2	lambiz			0 0
2	fam304-7-W	x 0 0 0 0 0	0 x x 0 0 0 0 0 0 0 0 x 0 0 0 0 0	0 0
2	fam169	x x 0 0 0 0	0 <u>0 0 x</u> 0 0 0 x 0 0 0 0 x 0 0 0 0 0	0 0
-	fam204.40			0 0
2	fam304-19	0 0 0 0 0 0		0 0
2	fam129-9	0 0 0 0 0 0	0 0 0 0 <mark>x x</mark> 0 0 0 0 0 0 x 0 0 0 0 0	0 0
2	fam172-1/fam172	0 0 0 0 0 0	0 0 0 0 0 x x 0 0 0 0 x 0 0 0 0 0	0 0
4	fom172			0 0
4	Idili172	A 0 0 0 0 0		0 0
2	fam288-19	0 x 0 0 0 0	0 0 0 x x 0 0 0 0 0 0 0 x 0 0 0 0 0	0 0
2	fam312-11	0 0 x 0 0 0	0 0 0 0 <mark>x x</mark> 0 0 0 0 0 0 x 0 0 0 0 0	0 0
2	fam290-5			0 0
2	fam200_10			0 0
2	tam290-18			0 0
2	fam310-3-W	0 0 0 0 × 0	0 0 0 0 0 <mark>x x</mark> 0 0 0 0 0 x 0 0 0 0 0	0 0
2	fam324-12	0 0 0 0 X	0 0 0 0 0 0 0 x x 0 0 x 0 0 0 0 0	0 0
2	fam126-11			0 0
2	100.10	0 0 0 0 0		0 0
3	fam126-18	0 0 0 0 0 0		0 0
2	fam114-6	0 0 0 0 0 0	0 0 x x 0 0 0 0 0 0 0 0 x 0 0 0 0	0 0
2	fam114-14	0 0 0 0 0 0	0 0 <u>0 0 0 0 0 x x</u> 0 0 0 0 x 0 0 0 0	0 0
2	fam307-3-W	x 0 0 0 0 0		0 0
-	6			0 0
2	lam307-21	XXUUUU		0 0
2	fam240-18	0 x x 0 0 0	0 0 x 0 0 0 x 0 0 0 0 0 0 x 0 0 0 0	0 0
3	fam216-6	0 x x 0 0 0	0 0 0 0 x 0 0 0 0 x x 0 0 x 0 0 0 0	0 0
2	fam216-7			0 0
-	fom194.21			0 0
2	laiii 104-21			0 0
4	fam184-14	0 0 0 0 x 0	0 0 0 0 0 x x x x 0 0 0 0 x 0 0 0 0	0 0
4	fam306	× 0 0 0 0 0	0 0 0 0 0 <mark>x x x x</mark> 0 0 0 <mark>x</mark> 0 0 0 0	0 0
2	fam112-11	0 0 0 0 0	0 0 0 0 x 0 0 0 0 0 x x 0 0 0 0	0 0
-	fam252.24 \\/			0 0
3	1a111233-24-VV			0 0
2	fam253-1-W	000000	0 0 0 <mark>x x 0</mark> 0 0 0 0 0 0 0 x 0 0 0	0 0
2	fam131-12	x 0 0 0 0 0	0 0 0 <mark>x 0 x</mark> 0 0 0 0 0 0 0 0 x 0 0 0	0 0
2	fam161-19	0 × 0 0 0 0	0 0 0 0 0 0 x x 0 0 0 0 x 0 0 0	0 0
-	fam269_10			0 0
2	18/11208-19	0 0 0 x 0 0		0 0
3	fam164-20	0 0 0 0 × 0	x 0 0 0 0 0 0 x 0 0 0 0 0 x 0 0 0	0 0
2	fam164-18x164-10	0 0 0 <mark>x</mark> 0 0	0 0 0 0 0 0 x 0 x 0 0 0 0 x 0 0 0	0 0
3	fam133-19	0 0 0 0 X	x 0 0 0 0 0 0 x 0 0 0 0 x 0 0 0	0 0
õ	fam140.04			0 0
2	Idiii 140-21	0 0 0 0 0 0		0 0
2	fam148-20	0 0 0 0 0 ×	0 0 0 0 0 0 0 0 x x 0 0 0 0 x 0 0 0	0 0
2	fam148-13	0 0 0 0 0 x	0 0 x 0 0 0 0 x 0 0 0 0 0 0 x 0 0 0	0 0
2	fam134-10	0 0 0 0 0	0 0 x 0 0 0 0 0 0 0 0 0 0 x 0 0 0	0 0
2	fam134.4			0 0
2	Iall1134-4	0 0 0 0 0 0		0 0
2	fam132-20	0 0 0 0 0 0	0 0 0 0 0 0 0 x x 0 0 0 0 0 x 0 0	0 0
2	fam127-1	x 0 0 0 0 0	0 0 x 0 0 x 0 0 0 0 0 0 0 0 x 0 0	0 0
_				
2	fam144-19	0 0 × 0 0 0	0 0 0 0 0 <mark>x x 0</mark> 0 0 0 0 0 0 x 0 0	0 0
2	fam147-13	0 0 0 <mark>x x</mark> 0	0 0 0 0 0 0 0 x 0 0 0 0 0 0 0 x 0 0	0 0
4	fam271-7			0 0
- -	fam147.47			0 0
2	lam147-17			0 0
2	fam273-6	0 0 0 0 0 x	0 0 0 <u>0 0 0 0 0 0 x x 0 0 0 0 x 0 0</u>	0 0
2	fam273-18	0 0 0 0 0 0	0 0 0 <mark>x x</mark> 0 0 0 0 0 0 0 0 0 x 0 0	0 0
3	fam273-19	0 0 0 0 0 0	0 0 0 0 x x 0 0 0 0 0 0 0 0 x 0 0	0 0
2	fom246P			0 0
2	Talli240B	0 0 0 0 0 0		0 0
2	fam173-13	0 0 0 x x 0	0 0 x x 0 0 0 0 0 0 0 0 0 0 0 0 x 0	0 0
2	fam138-21	0 0 0 0 0 0	0 0 0 0 0 0 0 0 x x 0 0 0 0 0 x 0	0 0
2	fam149-14	0 0 0 0 0 0	0 0 x 0 0 x 0 0 0 0 0 0 0 0 0 x 0	0 0
2	fam149-15	0 0 0 0 0 0		0 0
2	fam207 10			0 0
2	lam287-18			0 0
2	fam243-12	0 0 0 × 0 0	x x 0 0 0 0 0 0 0 0 0 0 0 0 0 x 0	0 0
2	fam243-21	0 0 0 x x 0	0 0 0 0 0 x 0 0 0 0 0 0 0 0 0 x 0	0 0
3	fam106-20	0 0 x 0 0 0	0 0 0 0 0 0 0 0 x x x 0 0 0 0 x 0	0 0
2	fam106_11			0 0
2	laii100-11			0 0
2	tam106-19			0 0
2	fam106-8	0 0 0 <mark>x x</mark> 0	0 <u>0 0 0 0 0 0 0 x 0 0 0 0 x 0</u>	0 0
2	fam155-24	× 0 0 0 0 0	0 x 0 x 0 0 0 0 0 0 0 0 0 0 0 x	0 0
2	fam272-17	0 0 0 0 0 0	x x 0 0 0 0 0 0 0 0 0 0 0 0 0 x	0 0
-	fam272 17			
4	lam272-15	x 0 0 0 0 0		0 0
2	fam272-6	0 0 0 0 0 0	0 0 0 x 0 0 0 0 0 0 0 0 0 0 0 x	0 0
2	fam272-8x272-7	0 0 0 0 0 0	0 0 0 0 x 0 0 0 0 0 0 0 0 0 0 x	0 0
2	fam272_0	0 0 0 0 0 0		
2	1011272-9	0 0 0 0 0 0		0 0
2	tam2/2-12			0 0
2	fam272-13	0 0 0 0 0	0 <u>0 0 0 0 0 x 0 0 0 0 0 0 0 x 0 0 0 0 x 0 0 0 x 0 0 0 x 0 0 0 x 0 0 0 0 x 0 0 0 0 x 0 0 0 0 x 0 0 0 0 x 0 0 0 0 x 0 0 0 0 0 x 0 0 0 0 0 x 0 0 0 0 0 0 0 x 0</u>	0 0
2	fam163-11-W	× 0 0 0 0 0	0 x x 0 0 0 0 0 0 0 0 0 0 0 0 x	0 0
2	fam163_6	0 x 0 0 0 0		
-	fam:005 0414			
2	tam285-6x5/285-21W	x x U U U O	<u> </u>	0 0
4	fam285-21-W	x x A 0 0 0	X A O O A A O O O O O O O O O X	0 0
2	fam285-4	x 0 0 0 0 0	0 0 0 0 0 0 0 x x 0 0 0 0 0 x	0 0
2	fam154_6v285_21			
~	10111104-0X200-21			
2	tam154-12	υυυχχΟ	<u> </u>	0 0
2	fam154-21x103-4	0 0 0 0 <mark>x x</mark>	x 0 0 0 x 0 <u>0</u> 0 <u>0</u> 0 0 0 0 0 0 0 x	0 0
2	fam154-5	x 0 x 0 0 0	0 0 0 0 0 x 0 x 0 0 0 0 0 0 0 x	0 0
2	fam103-6			
~				
2	tam103-13	U U U U U X	x x v v v v v v v v v v v v v v v 0 0 x	υυ

· ·	ferm102.11		•	0	0	0				0	•		0	0	0			0	~	~	~	•	0	•	•	0	0	0			~
3	ram103-14	0	0	U	U	U	X	X	X	υ	U	х	0	0	0	X	X	U	0	U	υ	U	U	0	U	U	U	U	X	U	υ
3	fam160-11	0	0	0	0	0	0	0	х	х	0	0	Х	Х	0	0	0	0	0	0	0	0	0	0	0	0	0	0	Х	0	0
3	fam160-21	0	0	0	0	0	0	0	х	х	Х	х	0	0	х	х	0	0	0	0	0	0	0	0	0	0	0	0	x	0	0
2	fam160-22	0	0	0	0	0	0	0	0	х	х	0	х	0	x	0	0	0	0	0	0	0	0	0	0	0	0	0	x	0	0
2	199-2x160-19	0	0	0	0	0	0	0	0	х	0	0	0	0	0	0	0	0	Х	х	0	0	0	0	0	0	0	0	x	0	0
3	fam160-5	0	0	0	0	0	0	0	х	0	0	0	0	0	0	0	х	х	x	0	0	0	0	0	0	0	0	0	x	0	0
3	fam118-5	0	0	0	0	0	0	0	0	х	Х	0	0	0	0	0	х	х	0	0	0	0	0	0	0	0	0	0	x	0	0
2	fam118-10	x	х	0	0	0	0	0	0	0	0	0	х	х	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	x	0
2	fam118-20	x	0	0	0	0	0	0	0	0	0	х	0	0	0	0	0	0	0	0	Х	0	0	0	0	0	0	0	0	x	0
4	fam152-W	0	0	0	0	х	0	0	0	0	0	х	x	0	Х	х	0	0	0	0	0	0	0	0	0	0	0	0	0	x	0
2	fam182-18	0	0	0	0	Y	0	0	0	0	0	0	0	0	0	0	0	0	0	x	Y	0	0	0	0	0	0	0	0	Y	0
2	fam281.8	0	0	õ	0	0	v	0 V		ñ	0	v	v	0	0	0	0	0	0	0	0	0	0	0	õ	0	0	0	0	Û	õ
2	fam100_10_W	0	0	0	0	0	<u>^</u>	Ê.	0	0	0	^	· .	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		0
2	fam109-10-vv	0	0	0	0	0	0	X	0	0	0	0	X	Х	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	X	0
2	fam279-1	0	0	0	0	0	0	0	Х	0	0	0	Х	Х	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	X	0
3	fam217-17-W	0	0	0	0	0	0	0	0	х	Х	х	0	0	Х	0	0	0	0	0	0	0	0	0	0	0	0	0	0	X	0
2	fam217-3	0	0	0	0	0	0	0	х	0	0	0	0	0	0	0	0	0	Х	Х	0	0	0	0	0	0	0	0	0	x	0
4	fam194	0	0	0	0	0	0	0	0	0	Х	0	0	0	0	Х	Х	0	х	x	0	0	0	0	0	0	0	0	0	x	0
2	fam270-9	x	0	0	0	0	0	0	0	0	0	х	х	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	x
2	fam245-8-w	0	х	0	0	0	0	0	0	0	0	0	0	0	х	Х	0	0	0	0	0	0	0	0	0	0	0	0	0	0	x
2	fam282-8-w	0	0	х	х	0	0	0	0	0	0	х	х	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	x
2	fam282-14	0	0	0	x	0	0	0	0	0	0	0	0	0	0	0	х	х	0	0	0	0	0	0	0	0	0	0	0	0	x
3	fam278-24	0	0	0	0	0	X	X	X	0	0	х	0	х	0	X	0	0	0	0	0	0	0	0	0	0	0	0	0	0	x
2	fam278-22	0	0	0	0	0	0	x	x	0	0	0	х	0	х	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	x
2	fam286-16-W	0	0	0	0	0	X	0	0	X	0	0	0	0	0	x	0	0	0	X	0	0	0	0	0	0	0	0	0	0	x



Supplementary Table 4: Large family search profiles

Supplementary Table 5: Database lines representing assigned NAC insertions

							Х	(-di	me	nsi	on			Y-dimension								Z-dimension										
alleles	Annotation	id	Х	Ϋ́	Z	12	3	4	56	37	8	91	0	12	3	4 !	56	7	8 9	ə 10) 1	12	3	4	5	67	8	9	10	-	CI. Size	seq
phnh3-1	#912	15756	1	1	1	03	0	0	0 0	0 (0	0 () (3	0	0 0	0 (0	0 0	0 ((0 0	0	0	0	0 0	0 (0	2		8	CCCAGGACCTTAGGGATCAAGAAGGCACTAGT
phnh3-2	F154/103	19127	2	1	1	0 0	0	0	0 5	51	0	0 (0 0	0 0	0	0 0	0 (0	0 0	3	(0 0	0	0	0	0 0	5	0	0		14	CCAAGAGATAGAAAATACCCAAATGGTTCACGO
phnh10-1	#572	16952	1	1	1	0 0	0	0	0 0	0 (5	0 (0 0	3 3	0	0 0	0 (0	0 0	0 ((0 0	0	0	0	20	0	0	0		10	ACCGAAAGCCAGAGATTCTACCACCGCTACCA
phnh17-1	F284	36289	2	3	1	02	8	0	0 0	0 (0	0 (° C	10	1	2 (0 (0	0 0	0 0	(0 0	0	0	0	0 0	0	6	0	1	20	ATAACTACTGGAACTGGCAAGGTGGTTGTTGTG
phnh22-1	#205	13798	1	1	1	10	0	0	0 0	0 (0	0 (0 0	0 0	0	0 4	10	0	0 0	0 (0	0 0	1	0	0	0 0	0	0	0		6	GAAGGCAACAGGACTTGACAAGCAAATAGTGA
phnh23-1	#247	30744	1	1	1	0 0	0	0	6 (0 (0	0 (0 0	0 0	0	0 0	0 (2	0 0	0 (0	0 0	2	0	0	0 0	0	0	0		10	GCAGCTAAGGCGACGTTTGGAGAACAAGGAAT
phnh24-1	F167	7367	2	3	1	0 0	0	0	0 0) 3	3	0 (D C	0 0	0	0 0	0 0	0	4 3	32	0	0 0	0	0	0	0 0	0	5	0		20	GGGAGGGGACAAATGCATAAGGGAACAGCTTC
phnh25-1	#322	45951	1	1	1	0 0	2	0	0 0	0 (0	0 (0 0	3 3	0	0 0	0 (0	0 0	0 (0	0 0	0	1	0	0 0	0	0	0		6	AGCAGCAGGATTTTGGAAGGCAACAGGAAGAG
phnh26-1	F280/84	25950	2	2	1	02	1	0	0 0	0 (0	0 (0 0	0 0	0	0 0) 1	1	0 0	0 0	(0 0	0	0	0	0 0	0	1	0		6	TACTCAACAAGTGAATCATTCTGGCATGGTGAG