

# *Molecular Characterization of Mutant Germplasm*

*A Manual*

*Prepared by the  
Joint FAO/IAEA Programme of Nuclear Techniques in Food and Agriculture*



**Joint FAO/IAEA Division**  
of Nuclear Techniques in Food and Agriculture

*Plant Breeding and Genetics Laboratory, Seibersdorf, 2015*



## FOREWORD

Plant biotechnology applications must not only respond to the challenges of improving food security and fostering socio-economic development, but in doing so, promote the conservation, diversification and sustainable use of plant genetic resources for food and agriculture. Today the biotechnology toolbox available to plant breeders offers many new possibilities for accelerating the breeding process, and increasing productivity, crop diversification and production, while developing a more sustainable agriculture. The early versions of this manual provided a companion to training courses on plant mutant germplasm characterization. As such, the content was tailored to the curricula of the course. It has now developed to include new technologies as they emerge in providing a contemporary tool kit for genotypic analysis and selection in plant breeding and genetics.

The first print of this manual on selected molecular marker techniques was prepared using the hand-outs and other materials distributed to participants of the FAO/IAEA Interregional Training Course on "Mutant germplasm characterisation using molecular markers". The course was hosted by the Joint FAO/IAEA Programme of Nuclear Techniques in Food and Agriculture at the Plant Breeding and Genetics Laboratory (PBGL, formerly the Plant Breeding Unit) of the Agriculture and Biotechnology Laboratory at the IAEA Laboratories in Seibersdorf, Austria, in 2001. Messrs J. Bennetzen (USA), K. Devos (UK), G. Kahl (Germany), U. Lavi (Israel), M. Mohan (ICGEB) and S. Nielen (FAO/IAEA) contributed protocols to the first print version. These contributions and others were formally compiled into the first early editions of the manual by Messrs P. Gustafson (USA), B. Forster (UK), M. Gale (UK), R. Adlam (UK), M. Maluszynski and S. Nielen of the Joint Programme. In later editions, J Fernandez-Manjarres (Colombia) provided the section on population genetics, and Plant Breeding and Genetics Section Head Pierre Lagoda provided the protocol on multivariate analysis. While this series of courses ended in 2007, there has been a continual demand from trainees for a codified set of standard protocols, and so the Plant Breeding and Genetics Laboratory (PBGL) has continued adapting this book by incorporating new protocols with the aim of assisting Member States in the appropriate application of molecular tools with minimal costs. These include protocols for TILLING/Scotilling, DNA quantification, low-cost and low toxicity DNA extraction, alternative enzymology for enzymatic mismatch cleavage (new in 2013), methods for rapid bench-top purification of single-strand-specific nucleases used in mutation discovery assays (new in 2014). Of note in this the successful implementation of low-cost and non-toxic DNA extraction methods developed by the PBGL and first delivered to the Member States in the 2013 edition of this manual. These methods have been successfully adapted for 20 different crops. For 2015 we've added a PBGL protocol for molecular validation of production of doubled haploid plants. This been validated in barley, Tef, and sugarbeet. Particular thanks for work on recent editions (since 2010) go to PBGL staff Owen Huynh, Joanna Jankowicz-Cieslak, and Bradley Till.

We strive to improve the manual with each edition. We very much appreciate feedback, suggestions and comments, which could further improve and enrich the contents of this manual. Correspondence should be addressed directly to Mr. P.J.L Lagoda, Head of Plant

Breeding and Genetics Section, Joint FAO/IAEA Division of Nuclear Techniques in Food and Agriculture, P.O. Box 100, Vienna, Austria, Telephone: +43 1 2600 21626; email [P.Lagoda@iaea.org](mailto:P.Lagoda@iaea.org).

A hard copy with attached CD-ROM will be distributed, free of charge, to interested scientists from FAO and IAEA Member States. Requests for the manual should be sent to Ms. K. Allaf, Plant Breeding and Genetics Section, Joint FAO/IAEA Division of Nuclear Application in Agriculture, P.O. Box 100, Vienna, Austria, Telephone: +43 1 2600 21621 or by email: [K.Allaf@iaea.org](mailto:K.Allaf@iaea.org).

## **LIST OF ACRONYMS**

<b>AFLP</b>	<b>A</b> mplified <b>F</b> ragment <b>L</b> ength <b>P</b> olymorphism
<b>CAPS</b>	<b>C</b> leaved <b>A</b> mplified <b>P</b> olymorphic <b>S</b> equences
<b>CJE</b>	<b>C</b> elery <b>J</b> uice <b>E</b> xtract
<b>DH</b>	<b>D</b> oubled Haploid
<b>EST</b>	<b>E</b> xpressed <b>S</b> equences <b>T</b> ag
<b>IPCR</b>	<b>I</b> nverse <b>P</b> olymerase <b>C</b> hain <b>R</b> eaction
<b>IRAP</b>	<b>I</b> nter- <b>R</b> etrotransposon <b>A</b> mplified <b>P</b> olymorphism
<b>ISSR</b>	<b>I</b> nter- <b>S</b> imple <b>S</b> equences <b>R</b> epeat amplification
<b>PCR</b>	<b>P</b> olymerase <b>C</b> hain <b>R</b> eaction
<b>RAPD</b>	<b>R</b> andom <b>A</b> mplified <b>P</b> olymorphic <b>D</b> N
<b>REMAP</b>	<b>R</b> etrotransposon- <b>M</b> icrosatellite <b>A</b> mplified <b>P</b> olymorphism
<b>RFLP</b>	<b>R</b> estriction <b>F</b> ragment <b>L</b> ength <b>P</b> olymorphism
<b>SCAR</b>	<b>S</b> equences <b>C</b> haracterized <b>A</b> mplified <b>R</b> egion
<b>SNP</b>	<b>S</b> ingle <b>N</b> ucleotide <b>P</b> olymorphism
<b>SSCP</b>	<b>S</b> ingle <b>S</b> tranded <b>C</b> onformation <b>P</b> olymorphism
<b>SSR</b>	<b>S</b> imple <b>S</b> equences <b>R</b> epeat
<b>STS</b>	<b>S</b> equences <b>T</b> agged <b>S</b> ite
<b>TILLING</b>	<b>T</b> argeting <b>I</b> nduced <b>L</b> ocal <b>L</b> esions <b>I</b> N <b>G</b> enomes
<b>NGS</b>	<b>N</b> ext <b>G</b> eneration <b>S</b> equencing

# TABLE OF CONTENTS

FOREWORD.....	I
LIST OF ACRONYMS .....	III
TABLE OF CONTENTS .....	IV
1. INTRODUCTION TO MOLECULAR MARKERS .....	1-1
1.1. Use of molecular markers: A cautionary tale .....	1-2
1.1.1. An example of how not to use molecular markers.....	1-2
1.1.2. An example of efficient application of markers.....	1-3
1.2. A Summary of Marker Techniques .....	1-4
1.3. Ideal genetic markers.....	1-4
1.4. Marker application suitability .....	1-5
1.5. Implementation.....	1-8
1.6. Requirements .....	1-8
1.7. Comparison of different marker systems.....	1-9
2. LOW COST DNA EXTRACTION WITHOUT TOXIC ORGANIC PHASE SEPARATION 2-1	
2.1. Materials.....	2-1
2.2. Solutions to Prepare .....	2-3
2.3. Methods (for centrifuge tubes).....	2-3
2.4. Example Data.....	2-6
2.5. Conclusions .....	2-8
3. DNA QUANTIFICATION .....	3-1
3.1. Protocol for gel electrophoresis .....	3-1
3.1.1. Preparation of DNA concentration standards.....	3-1
3.1.2. Preparing agarose gels.....	3-2
3.1.3. Preparing samples for loading into gels.....	3-2
3.1.4. Running the gel .....	3-2

3.1.5. Photographing the gel .....	3-3
3.2. Quantification of DNA using image analysis software.....	3-3
<b>4. RESTRICTION ENZYME DIGEST.....</b>	<b>4-1</b>
<b>5. FINDING CANDIDATE GENES AND PRIMER DESIGN FOR MOLECULAR TESTING: AN EXAMPLE FROM THE ANNOTATED <i>SORGHUM BICOLOR</i> GENOME.....</b>	<b>5-1</b>
5.1. Overview .....	5-1
<b>6. SSR.....</b>	<b>6-1</b>
6.1. Protocol.....	6-2
6.1.1. PCR reaction mix .....	6-2
6.1.2. PCR amplification.....	6-3
6.1.3. Separation of the amplification products in agarose gel .....	6-3
6.1.4. Denaturing gel electrophoresis .....	6-4
6.1.5. Assembling the glass plate sandwich.....	6-4
6.1.6. Casting gel .....	6-5
6.2. Setting up the operation .....	6-5
6.3. Polyacrylamide gel running conditions.....	6-6
6.4. Silver-staining .....	6-6
6.5. References .....	6-7
6.6. Reagents needed .....	6-8
<b>7. ISSR .....</b>	<b>7-1</b>
7.1. Protocol.....	7-1
7.1.1. Prepare 20µl reaction mix.....	7-2
7.1.2. PCR amplification.....	7-2
7.1.3. Separation and visualization of the amplification products.....	7-2
7.1.4. Gel running conditions.....	7-3
7.1.5. Silver-staining.....	7-3
7.2. Primers available at Plant Breeding & Genetics Laboratory (FAO/IAEA) .....	7-3
7.3. References .....	7-4
7.4. Reagents needed .....	7-4

8. AFLP .....	8-1
8.1. Protocol.....	8-2
8.1.1. Restriction of genomic DNA and ligation of adapters to the DNA fragments.....	8-2
8.1.2. Pre-amplification.....	8-3
8.1.3. PCR pre-amplification .....	8-3
8.1.4. Check-step.....	8-3
8.1.5. Selective pre-amplification.....	8-4
8.1.6. PCR mix for selective amplification, products to be visualized on PAGE.....	8-5
8.1.7. PCR profile for Selective amplification, products to be visualised on PAGE .....	8-5
8.1.8. Polyacrylamide Gel Electrophoresis (PAGE) .....	8-5
8.1.9. Silver staining of PAG.....	8-6
8.1.10. PCR mix for selective amplification, products to be visualized on an automated DNA analyser .....	8-6
8.1.11. PCR profile for selective amplification, products to be visualized on an automated DNA analyser.....	8-6
8.1.12. Electrophoresis using an automated DNA analyser.....	8-7
8.1.13. Production of single primer, linear PCR products .....	8-7
8.1.14. PCR amplification to produce single stranded DNA.....	8-7
8.2. Required enzymes and primer sequences for AFLP assays.....	8-8
8.2.1. Restriction enzymes.....	8-8
8.3. Preparation of adapters .....	8-8
8.4. Reagents needed .....	8-8
8.5. Sequence information of adapters and primers used for AFLP .....	8-9
8.6. References .....	8-10
9. REMAP & IRAP .....	9-1
9.1. Protocol.....	9-1
9.1.1. Prepare a 50µl reaction mix.....	9-3
9.1.2. PCR amplification .....	9-3
9.1.3. Separation and visualization of the amplification products.....	9-3
9.2. References .....	9-4
9.3. Reagents needed .....	9-4



10.	SINGLE NUCLEOTIDE POLYMORPHISMS (SNPs).....	10-1
10.1.	References.....	10-2
11.	TILLING.....	11-1
11.1.	Protocol.....	11-1
11.1.1.	PCR reaction with IRDye-labeled primers .....	11-1
11.1.2.	Heteroduplex digestion, preparation of Sephadex spin plates.....	11-3
11.1.3.	Agarose gel analysis of enzymatic mismatch cleavage, and sample purification ....	11-4
11.1.4.	Sample purification and volume reduction .....	11-5
11.1.5.	Preparing, loading, and running LI-COR gels .....	11-6
11.1.6.	Data Analysis .....	11-7
11.2.	Computation tools .....	11-8
11.2.1.	Selecting the best region to screen and designing primers .....	11-8
11.3.	Data analysis .....	11-10
11.4.	Additional info.....	11-13
11.4.1.	List of consumables and equipment .....	11-13
11.5.	Frequently asked questions.....	11-15
11.6.	Additional protocols .....	11-16
11.6.1.	Sequencing .....	11-16
11.7.	EMS mutagenesis of Arabidopsis seed.....	11-17
11.7.1.	Materials.....	11-18
11.7.2.	Standard size batch.....	11-18
11.7.3.	A note on technique.....	11-18
11.7.4.	DNA extraction.....	11-19
11.8.	References.....	11-20
12.	ALTERNATIVE ENZYMOLOGY FOR MISTMATCH CLEAVAGE FOR TILLING AND ECOTILLING: EXTRACTION OF ENZYMES FROM WEEDY PLANTS.....	12-1
12.1.	Objective .....	12-1
12.2.	Materials .....	12-1
12.3.	Methods.....	12-2
12.3.1.	Enzyme extraction.....	12-2

12.3.2. Concentration of enzyme extractions .....	12-3
12.3.3. Test of Mismatch Cleavage Activity.....	12-4
12.4. Example results.....	12-5
12.5. Conclusions.....	12-6
13. LOW-VOLUME, NON-TOXIC AND RAPID EXTRACTION OF SINGLE-STRAND-SPECIFIC NUCLEASES FROM CELERY .....	13-1
13.1. Objective .....	13-1
13.2. Materials .....	13-1
13.3. Methods.....	13-1
13.3.1. CEL I preparation.....	13-1
13.3.2. Activity tests .....	13-5
13.4. Conclusions.....	13-7
13.5. Contributors .....	13-7
14. A PROTOCOL FOR VALIDATION OF DOUBLED HAPLOID PLANTS BY ENZYMATIC MISMATCH CLEAVAGE .....	14-1
14.1. Abstract .....	14-1
14.2. Introduction .....	14-1
14.3. Materials .....	14-4
14.3.1. PCR amplification .....	14-4
14.3.2. Enzymatic mismatch cleavage .....	14-4
14.3.3. Agarose gel electrophoresis .....	14-4
14.4. Methods.....	14-4
14.4.1. PCR amplification .....	14-4
14.4.2. Enzymatic mismatch .....	14-5
14.4.3. Agarose gel electrophoresis and data analysis.....	14-5
14.5. Notes.....	14-7
14.6. Acknowledgments .....	14-9
14.7. Contributors .....	14-9
14.8. References.....	14-9

15. MULTIVARIATE ANALYSIS – PHYLOGENETICS AND PRINCIPAL COMPONENT ANALYSIS .....	15-1
15.1. Phylogenetics.....	15-1
15.2. Inferring phylogeny from pairwise distances: construction of a distance tree using clustering with the unweighted pair group method with arithmetic mean (UPGMA).....	15-2
15.3. Distance measures.....	15-2
15.4. Some reflexions on the comparison between genetic distances.....	15-8
15.5. What genetic distance estimator to choose for essential derivation?.....	15-8
15.6. Genetic distances between populations .....	15-9
15.7. Protocol: tree reconstruction.....	15-10
15.8. UPGMA exercise.....	15-16
15.9. Principal Component Analysis (PCA) .....	15-20
15.9.1. Considerations and references .....	15-22
15.10. References.....	15-24
16. POPULATION GENETICS.....	16-1
16.1. Reading and coding genetic data.....	16-1
16.1.1. Presence/absence coding of dominant data .....	16-1
16.1.2. Allele size coding for microsatellites.....	16-2
16.1.3. Categorical coding.....	16-4
16.1.4. Presence/absence coding of co-dominant data.....	16-4
16.1.5. Formatting dominant data as co-dominant.....	16-5
16.1.6. Notes of formatting diploid data with spread sheets .....	16-6
16.1.7. Transforming data types using software.....	16-6
16.1.8. The FSTAT data file.....	16-7
16.2. Genetic diversity.....	16-8
16.3. Genetic structure.....	16-11
16.3.1. Nei’s population genetics parameters: $G_{st}$ family.....	16-11
16.3.2. Sewall Wright’s F-statistics .....	16-11
16.4. Population and individual divergence and phylogenetic trees.....	16-12
16.5. Web resources and software – non-exhaustive.....	16-13
16.6. References.....	16-17

16.7. Some key concepts.....	16-19
16.8. Equations.....	16-20
<b>17. APPENDICES .....</b>	<b>17-1</b>
17.1. General DNA extraction techniques .....	17-1
17.1.1. Phenol/chloroform extraction .....	17-1
17.1.2. Ethanol precipitation .....	17-1
17.1.3. Solutions.....	17-2
17.2. Polymerase chain reaction protocol .....	17-2
17.2.1. References .....	17-6
17.3. Plant genome database contact information .....	17-7
17.4. Acronyms of chemicals and buffers.....	17-8
<b>NOTES .....</b>	<b>1</b>

## 1. INTRODUCTION TO MOLECULAR MARKERS

Traditionally, molecular markers have played a major role in the genetic characterization and improvement of many crop species. They have also contributed to, and greatly expanded, our abilities to assess biodiversity, reconstruct accurate phylogenetic relationships, and understand the structure, evolution and interaction of plant and microbial populations. Molecular markers systems reveal variation in genomic DNA sequence and allow the tracking of this variation, ideally linked to phenotypic trait variation, in crossing programmes. The first generation of molecular markers, RFLPs, were based on DNA-DNA hybridisation and were slow and expensive. The invention of the polymerase chain reaction (PCR) to amplify short segments of DNA gave rise to a second generation of faster and less expensive PCR-based markers, which became popular in genotyping of many species. Today, next generation sequencing technologies have become the dominant tool for marker assisted breeding in developed countries and biotechnology companies. While incredibly powerful, these techniques are still cost-limiting and carry a heavy bioinformatics load, making use difficult in developing countries. This will likely change in the future as sequencing technologies and analysis tools increase in power and decrease in cost. Until then, we provide in this manual a series of low cost marker systems that are applicable in many laboratories with infrastructure for basic molecular biology.

Molecular markers are being used extensively to investigate the genetic basis of agronomic traits and to facilitate the transfer and accumulation of desirable traits between breeding lines. They are used both to tag target genes and to monitor the genetic background. A number of techniques have been particularly useful for genetic analysis. For example, collections of RFLP probes have been very versatile and important for the generation of genetic maps, construction of physical maps, the establishment of syntenic relationships between genomes, and marker assisted breeding. Numerous examples of specific genes that have been identified as tightly linked to RFLP markers are available for the improvement of specific agronomic traits in almost all major crops. Specific examples include viral, fungal and bacterial resistance genes in maize, wheat, barley, rice, tomatoes and potatoes. Additional examples include insect resistance genes in maize, wheat and rice as well as drought and salt tolerance in sorghum. These markers often used in conjunction with bulked segregant analysis and detailed genetic maps, provide a very efficient method of characterizing and locating natural and induced mutated alleles at genes controlling interesting agricultural traits. Markers have also been used to identify the genes underlying quantitative variation for height, maturity, disease resistance and yield in virtually all major crops. In particular, the PCR-based techniques have been useful in the assessment of biodiversity, the study of plant and pathogen populations and their interactions; and identification of plant varieties and cultivars. Amplified DNA techniques have produced sequence-tagged sites that serve as landmarks for genetic and physical mapping. It is envisioned that emerging oligonucleotide-based technologies derived from the use of hybridization arrays, the so-called DNA chips and oligonucleotide arrays, will become important in future genomic studies. However, many of these are still under development, are proprietary, or require the use of expensive equipment, and are therefore not yet suitable or cost-effective for adequate transfer to developing countries. Clearly, the initial transfer of technology has only involved a selected group of techniques that are well established and/or seem to have a broad application (*e.g.*, RFLP,

SSR, ISSR, AFLP, RAPD, IRAP and REMAP and SNPs). However, techniques are continuously changing and evolving, so technology transfer needs to keep pace with current developments in genomics. Capacity for handling molecular marker data has been identified as a bottleneck to the integration of molecular techniques in germplasm management. A module on population genetics, dealing specifically with the analysis of molecular marker data is included in this edition of the manual.

## 1.1. Use of molecular markers: A cautionary tale

Molecular biology is an exciting discipline with new techniques constantly being developed and high impact publications coming from the work. As such, it is tempting for the junior scientists to think of molecular tools as a starting point for their breeding objectives. The downside, however, is that these tools are often challenging to master, expensive and easy to mis-apply. It is important that experiments are carefully designed with proper controls and that the researcher understands the strengths and limitations of the chosen application. In this section we focus on the use of molecular markers. These tools can provide rapid, valuable information on the nucleotide diversity of collections allowing deductions of evolutionary relationships and gene flow. However, this manual is focused on *mutant* germplasm characterization, and when applying these tools for evaluation of induced mutant populations, an understanding of the genetics of the species and heritability of variation is required for proper application. To highlight this, we offer two different examples of application of markers; one correct, the other incorrect. If you are uncertain if molecular markers are right for you, please feel free to contact the Plant Breeding and Genetics Laboratory for further advice.

### 1.1.1. An example of how not to use molecular markers.

A research group is starting a new project to use induced mutations to breed for improved disease resistance in barley. They have never used induced mutations before and would like to use molecular markers to track disease resistance because it is very time consuming and expensive for them to test their material phenotypically at every generation.

The group produces a large M1 population that was treated with gamma rays. They self-fertilize the barley and grow the M2 in the next generation. They apply pathogen to the plants and score resistance. Of 10,000 plants, they find 50 with some increase in resistance to the pathogen. These 50 plants come from 20 different M1 parents. They collect tissue from these 50 plants, along with 10 mutagenized plants that are susceptible and 10 plants that were not mutagenized. They extract DNA, and perform an AFLP marker analysis. They hope to find bands that are common in the resistant plants but not in the control. Their data is not conclusive, so they decide to look at even more plants.

### WHY IS THIS A BAD IDEA?

Current data suggests that most mutagenesis is random. In other words, different plants will have different changes in the DNA. Therefore, you don't expect the same mutations to be found in progeny from different M1 plants. Applying statistical probability, you might see this once or twice in a large population, but never 20 times. Therefore you don't expect to find bands in the mutants that arise due to common mutations.

### BUT, HOW COME THEY ARE ALL DISEASE RESISTANT?

If a trait is polygenic, there may be many genes involved in a trait. Different plants in the example population may have mutations in different genes that give a similar phenotypic response. So, you don't need to mutate the same gene to get a similar phenotype. Additionally, there may be many possible mutations within the same gene that could give you a phenotype. The different alleles may not give the same signal in a marker assay.

### 1.1.2. An example of efficient application of markers

The researchers working with the barley population above have produced one line that is highly disease resistant after backcrossing to the parental line and applying selective pressure through five generations.

The issue with the parental line and the mutant line is that they are low yielding. The researchers would like to introgress the disease resistance into a high yielding cultivar that farmers are growing. To aid in this, the researchers apply a set of SSR markers to 300 plants from the disease resistance line, 300 parents and 300 of the elite variety. They identify one new band with a set of SSR primers that is present in all mutants but not in either the parent or the elite variety. They set out a crossing plan where they cross the mutant line with the elite variety. They self the F1s and then select only plants with the mutant SSR band. Starting in the F2, they select plants for disease resistance. They also apply AFLP and choose disease resistant plants that share the majority of markers with the elite variety.

### WHY IS THIS A GOOD APPROACH?

The researchers have developed a marker by evaluating plants that are genetically related and harbouring the same mutation. Evaluation of a large number of plants allows the establishment that the marker is genetically linked to the mutation causing the phenotype. The lack of such bands in the control material reduces the risk that the marker is from some source of natural genetic variation. In the end, using AFLP allows for a high density of information on the genetic background of the selected individuals. It should be fairly straightforward to determine which plants have mostly elite variety background. This is what the breeder wants,

the elite variety with only that small amount of DNA conferring disease resistance introgressed, and not a lot of other DNA from the less suitable parent.

## 1.2. A Summary of Marker Techniques

Table 1.2–1. List of marker techniques

Marker/technique	PCR-based	Polymorphism (abundance)	Dominance
RFLP	No	Low-Medium	Co-dominant
RAPD	Yes	Medium-High	Dominant
SSR	Yes	High	Co-dominant
ISSR	Yes	High	Dominant
AFLP	Yes	High	Dominant
IRAP/REMAP	Yes	High	Co-dominant
Additional marker systems			
Morphological	No	Low	Dominant/Recessive/Co-dominant
Protein/isozyme	No	Low	Co-dominant
STS/EST	Yes	High	Co-dominant/Dominant
SNP	Yes	Extremely High	Co-dominant
SCARS/CAPS	Yes	High	Co-dominant
Microarray	Yes	High	

## 1.3. Ideal genetic markers

(highly dependent on application and species involved)

- No detrimental effect on phenotype
- Co-dominant in expression
- Single copy
- Economic to use
- Highly polymorphic
- Easily assayed
- Multi-functional
- Highly available (un-restricted use)
- Genome-specific in nature (especially when working with polyploids)
- Can be multiplexed
- Ability to be automated
- A perfect marker for the gene of interest, though for practical plant breeding a tightly linked marker is usually good enough.



## 1.4. Marker application suitability

RFLP	<p>Comparative maps                      Framework maps, bin mapping                      Genetic maps                      Breeding                      Varietal/line identification                      (multiplexing of probes necessary)                      Marker-assisted selection                      F<sub>1</sub> identification                      Diversity studies                      Novel allele detections                      Gene tagging                      Bulk segregant analysis                      Map-based gene cloning</p>
SSR	<p>This marker system is not suggested due to major issues in the lack of reproducibility.                      Fingerprinting                      Varietal/line identification (multiplexing of primers necessary)                      Framework/region specific mapping                      Genetic maps                      F<sub>1</sub> identification                      Comparative mapping                      Breeding                      Bulk segregant analysis                      Diversity studies                      Novel allele detections                      Marker-assisted selection                      High-resolution mapping                      Seed testing                      Map-based gene cloning</p>
ISSR	<p>Fingerprinting                      Varietal/line identification                      Genetic maps                      F<sub>1</sub> identification                      Gene tagging                      Breeding                      Bulk segregant analysis                      Diversity studies                      Marker-assisted selection                      High-resolution mapping                      Seed testing</p>

AFLP	<ul style="list-style-type: none"> <li>Fingerprinting</li> <li>Very fast mapping</li> <li>Region-specific marker saturation</li> <li>Varietal identification</li> <li>Genetic maps</li> <li>F<sub>1</sub> identification</li> <li>Gene tagging</li> <li>Breeding</li> <li>Bulk segregant analysis</li> <li>Diversity studies</li> <li>Marker-assisted selection</li> <li>High-resolution mapping</li> <li>Map-based gene cloning</li> </ul>
IRAP/REMAP	<ul style="list-style-type: none"> <li>Fingerprinting</li> <li>Varietal identification</li> <li>F<sub>1</sub> identification</li> <li>Gene tagging</li> <li>Bulk segregant analysis</li> <li>Diversity studies</li> <li>Marker-assisted selection</li> <li>High-resolution mapping</li> <li>Seed testing</li> </ul>
Morphological	<ul style="list-style-type: none"> <li>Genetic maps</li> <li>Alien gene introduction</li> <li>Varietal/line identification</li> <li>F<sub>1</sub> identification</li> <li>Novel phenotypes</li> <li>Breeding</li> </ul>
Protein and Isozyme	<ul style="list-style-type: none"> <li>Genetic maps</li> <li>Quality trait mapping</li> <li>Varietal/line identification (multiplexing of proteins or isozymes necessary)</li> <li>F<sub>1</sub> identification</li> <li>Breeding</li> <li>Seed testing</li> </ul>
STS/EST	<ul style="list-style-type: none"> <li>Fingerprinting</li> <li>Varietal identification</li> <li>Genetic maps</li> <li>F<sub>1</sub> identification</li> <li>Gene tagging and identification</li> <li>Bulk segregant analysis</li> <li>Diversity studies</li> </ul>

	<ul style="list-style-type: none"> <li>Marker-assisted selection</li> <li>Novel allele detection</li> <li>High-resolution mapping</li> <li>Map-based cloning</li> </ul>
SNP	<ul style="list-style-type: none"> <li>Genetic maps</li> <li>F<sub>1</sub> identification</li> <li>Breeding</li> <li>Gene tagging</li> <li>Alien gene introduction</li> <li>Bulk segregant analysis</li> <li>Diversity studies</li> <li>Novel allele detections</li> <li>Marker-assisted selection</li> <li>High resolution mapping</li> </ul>
SCARS/CAPS	<ul style="list-style-type: none"> <li>Framework mapping</li> <li>Can be converted to allele-specific probes</li> <li>F<sub>1</sub> identification</li> <li>Gene tagging</li> <li>Bulk segregant analysis</li> <li>Diversity studies</li> <li>Marker-assisted selection</li> <li>Map-based cloning</li> </ul>
Microarray	<ul style="list-style-type: none"> <li>Fingerprinting</li> <li>Sequencing</li> <li>Transcription</li> <li>Varietal identification</li> <li>Genetic maps</li> <li>F<sub>1</sub> identification</li> <li>Gene tagging and identification</li> <li>Bulk segregant analysis</li> <li>Diversity studies</li> <li>Marker-assisted selection</li> <li>High-resolution mapping</li> </ul>

## 1.5. Implementation

Table 1.5–1. Relative costs of marker techniques.

Marker/techniques	Development costs	Running costs per data point	Portability (Lab/Crops)
RFLP	Medium	High	High/High
RAPD	Low	Low	Low/Low
SSR	High	Medium	High/Low
ISSR	Low	Low	High/Low
AFLP	Medium-High	Low	High/Low
IRAP/REMAP	High	Medium	High/Low

*Additional marker systems not covered in the course*

Morphological	Depends	Depends	Limited to breeding aims
Protein and isozyme	High	Medium	High/High
SCARS/CAPS	High	Medium	High/Low
STS/EST	High	Medium	Medium/High
SNP	High	Medium-Low	Unknown
Microarray	Medium	Low	Unknown

## 1.6. Requirements

Table 1.6–1. Requirements for marker techniques.

Marker/technique	Amount/ quality of DNA	DNA Sequence Required	Radioactive detection	Gel system
RFLP	High/High	No	Yes/No	Agarose
RAPD	Low/Low	No	No	Agarose
SSR	Low/Medium	Yes	No	Acrylamide/ Agarose
ISSR	Low/Medium	Yes/No	No	Acrylamide/ Agarose
AFLP	Low/High	No	Yes/No	Acrylamide
IRAP/REMAP	Low/Medium	Yes	No	Acrylamide/ Agarose

*Additional marker systems not covered in the course*

Morphological	No	No	No	None
---------------	----	----	----	------

Protein/isozyme	No	No	No	Agarose/ Acrylamide
STS/EST	Low/High	Yes	Yes/No	Acrylamide/ Agarose
SNP	Low/High	Yes	No	Sequencing required
Microarray	Low/High	Yes	No	None
SCARS/CAPS	Low/High	Yes	Yes/No	Agarose

### 1.7. Comparison of different marker systems

Table 1.7–1. Advantages and disadvantages of various marker techniques.

Marker	Advantages	Disadvantages
<b>RFLP</b>	<ul style="list-style-type: none"> <li>• Unlimited number of loci</li> <li>• Codominant</li> <li>• Many detection systems</li> <li>• Can be converted to SCARs</li> <li>• Robust in usage</li> <li>• Good use of probes from other species</li> <li>• Detects in related genomes</li> <li>• No sequence information required</li> </ul>	<ul style="list-style-type: none"> <li>• Labour intensive</li> <li>• Fairly expensive</li> <li>• Large quantity of DNA needed</li> <li>• Often very low levels of polymorphism</li> <li>• Can be slow (often long exposure times)</li> <li>• Needs considerable degree of skill</li> </ul>
<b>RAPD</b>	<ul style="list-style-type: none"> <li>• Results obtained quickly</li> <li>• Fairly cheap</li> <li>• No sequence information required</li> <li>• Relatively small DNA quantities required</li> <li>• High genomic abundance</li> <li>• Good polymorphism</li> <li>• Can be automated</li> </ul>	<ul style="list-style-type: none"> <li>• Highly sensitive to laboratory changes</li> <li>• Low reproducibility within and between laboratories</li> <li>• Cannot be used across populations nor across species</li> <li>• Often see multiple loci</li> <li>• Dominant</li> </ul>
<b>SSR</b>	<ul style="list-style-type: none"> <li>• Fast</li> <li>• Highly polymorphic</li> <li>• Robust</li> </ul>	<ul style="list-style-type: none"> <li>• High developmental and start-up costs</li> <li>• Species-specific</li> <li>• Sometimes difficult interpretation because of stuttering</li> </ul>

Marker	Advantages	Disadvantages
<b>ISSR</b>	<ul style="list-style-type: none"> <li>• Can be automated</li> <li>• Only very small DNA</li> <li>• Codominant</li> <li>• Multiallelic</li> <li>• Multiplexing possible</li> <li>• Does not require radioactivity</li> <li>• Highly polymorphic</li> <li>• Robust in usage</li> <li>• Can be automated</li> </ul>	<ul style="list-style-type: none"> <li>• Usually single loci even in polyploids</li> <li>• Usually dominant</li> <li>• Species-specific</li> </ul>
<b>AFLP</b>	<ul style="list-style-type: none"> <li>• Small DNA quantities required</li> <li>• No sequence information required</li> <li>• Can be automated</li> <li>• Can be adapted for different uses, e.g. cDNA-AFLP</li> </ul>	<ul style="list-style-type: none"> <li>• Evaluation of up to 100 loci</li> <li>• Marker clustering</li> <li>• Dominant</li> <li>• Technique is patented</li> <li>• Can be technically challenging</li> </ul>
<b>IRAP/ REMAP</b>	<ul style="list-style-type: none"> <li>• Highly polymorphic depends on the transposon</li> <li>• Robust in usage</li> <li>• Can be automated</li> <li>• Species-specific</li> </ul>	<ul style="list-style-type: none"> <li>• Alleles cannot be detected</li> <li>• Can be technically challenging</li> </ul>
<b>Additional marker systems</b>		
<b>Morphological</b>	<ul style="list-style-type: none"> <li>• Usually fast</li> <li>• Usually cheap</li> </ul>	<ul style="list-style-type: none"> <li>• Few in number</li> <li>• Often not compatible with breeding aims</li> <li>• Need to know the genetics</li> </ul>
<b>Protein and Isozyme</b>	<ul style="list-style-type: none"> <li>• Fairly cheap</li> <li>• Fairly fast analysis</li> <li>• Protocol for any species</li> <li>• Codominant</li> <li>• No sequence information required</li> </ul>	<ul style="list-style-type: none"> <li>• Often rare</li> <li>• Often different protocol for each locus</li> <li>• Labour intensive</li> <li>• Sometimes difficult to interpret</li> </ul>

<b>Marker</b>	<b>Advantages</b>	<b>Disadvantages</b>
<b>STS/EST</b>	<ul style="list-style-type: none"> <li>• Fast</li> <li>• cDNA sequences</li> <li>• Non-radioactive</li> <li>• Small DNA quantities required</li> <li>• Highly reliable</li> <li>• Usually single-specific</li> <li>• Can be automated</li> </ul>	<ul style="list-style-type: none"> <li>• Sequence information required</li> <li>• Substantially decreased levels of polymorphism</li> </ul>
<b>SNP</b>	<ul style="list-style-type: none"> <li>• Robust in usage</li> <li>• Polymorphism are identifiable</li> <li>• Different detection methods available</li> <li>• Suitable for high throughput</li> <li>• Can be automated</li> </ul>	<ul style="list-style-type: none"> <li>• Very high development costs</li> <li>• Requires sequence information</li> <li>• Can be technically challenging</li> </ul>
<b>SCARS/CAPS</b>	<ul style="list-style-type: none"> <li>• Codominant</li> <li>• Small DNA quantities required</li> <li>• Highly reliable</li> <li>• Usually single locus</li> <li>• Species-specific</li> </ul>	<ul style="list-style-type: none"> <li>• Very labour intensive</li> </ul>
<b>Microarray</b>	<ul style="list-style-type: none"> <li>• Single base changes</li> <li>• Highly abundant</li> <li>• Highly polymorphic</li> <li>• Codominant</li> <li>• Small DNA quantities required</li> <li>• Highly reliable</li> <li>• Usually single locus</li> <li>• Species-specific</li> <li>• Suitable for high throughput</li> <li>• No gel system</li> <li>• Can be automated</li> </ul>	<ul style="list-style-type: none"> <li>• Very high development and start-up costs</li> <li>• Portability unknown</li> </ul>

## 2. LOW COST DNA EXTRACTION WITHOUT TOXIC ORGANIC PHASE SEPARATION

One of the most common activities of molecular biology is the extraction of genomic DNA from cells. Traditional methods utilized lysis followed by organic phase separation to remove unwanted molecules such as proteins. Commercialized kits from companies such as Qiagen have circumvented unwanted toxic organic phase separation by using methods that employ DNA binding to silica with the use of chaotropic salts. This approach has proven superior in terms of speed and quality of product and has become the industry standard. The main issue with these commercial kits is that costs can become prohibitively expensive for large scale applications. The protocol below describes a home-made silica DNA binding protocol that costs about 1/10th that of a commercial kit and produces DNA quality suitable for TILLING and other high-throughput molecular applications.

### 2.1. Materials

	Company
<b>MATERIALS FOR LOW-COST DNA EXTRACTIONS</b>	
Celite 545 silica powder (Celite 545-AW reagent grade)	Supelco 20199-U
SDS (Sodium dodecyl sulfate) for mol biol approx 99%	Sigma L-4390-250G
Sodium acetate anhydrous	Sigma S-2889 (MW=82.03g/mol)
NaCl (Sodium chloride)	Sigma S-1314-1KG (MW=58.44g/mol)
RNase A	10 microgram per ml.
Ethanol	Ethanol absolute for analysis (Merck 1.00983.2500)
Nuclease-free H <sub>2</sub> O	Gibco ultrapure distilled water (DNase, RNase-free)
Guanidine thiocyanate	Sigma G9277 (MW=118.2g/mol)
Microcentrifuge tubes (1.5mL, 2.0mL)	Any general laboratory supplier
Micropipettes (1000µL, 200µL, 20µL)	Any general laboratory supplier
Microcentrifuge	Eppendorf Centrifuge 5415D
Optional: Shaker for tubes	Eppendorf Thermomixer comfort for 1.5mL tubes
<b>MATERIALS FOR GRINDING OF LEAF MATERIAL (depending on grinding method)</b>	
Liquid nitrogen	
Mortar and pestle or, TissueLyser, ...	e.g. Qiagen TissueLyser II



Metal beads (tungsten carbide beads, 3mm)	Qiagen Cat.No. 69997 (for TissueLyser)
<b>EVALUATION OF DNA YIELD AND QUALITY</b>	
DNA concentration	ND-NanoDrop Spectrophotometer (optional) 1000
Agarose gel equipment	Any supplier providing horizontal mini-gels
<b>TILLING-PCR</b>	
Thermocycler	Biorad C1000 Thermal cycler, or equivalent
PCR tubes	Life Science No 781340
TaKaRa Ex Taq™ Polymerase (5U/ul)	TaKaRa
10X Ex Taq™ Reaction Buffer	TaKaRa
dNTP Mixture (2.5mM of each dNTP)	TaKaRa
Agarose gel equipment	Any supplier providing horizontal mini-gels

## 2.2. Solutions to Prepare

Buffer	Receipt	Comments
<b>Stock solutions</b>		
5M NaCl stock solution	MW=58.44g/mol 29.22g / 100mL	If keeping stocks for a long period, check to make sure high molarity stocks stay in solution. If precipitate forms, warm solution until back in solution, or discard and make fresh.
3M Sodium acetate (pH = 5.2)	MW=82.03g/mol 24.61g / 100mL	Adjust pH value with glacial acetic acid
95% (v/v) Ethanol	95 mL ethanol abs + 5 mL H <sub>2</sub> O	
Tris-EDTA (TE) buffer (10x)	<u>Composition:</u> For 100mL 100mM Tris-Cl, pH8.0 10 mM EDTA Tris-Cl stock 2mL of 0.5M EDTA stock	Tris and EDTA can be prepared from powder. Note that the pH of tris changes with temperature.
LYSIS BUFFER (standard)	0.5% SDS (w/v) in 10x TE 0.5g SDS /100mL	PBGL has developed a range of lysis buffers for different crops. If performance is poor, contact PBGL for modified buffers.
DNA BINDING BUFFER	6M Guanidine thiocyanate MW = 118.2 g/mol 70.92 g / 100mL (6M)	!!! it takes several hours until dissolved (leave it approx. 4-5 hours)
WASH BUFFER	1mL of 5M NaCl + 99mL of 95% EtOH	!!! PREPARE FRESH, because the salt precipitates during storage
DNA ELUTION BUFFER	depending on application (e.g. TE-buffer; Tris-HCl buffer)	

## 2.3. Methods (for centrifuge tubes)

### PREPARATION OF SILICA POWDER-DNA BINDING-SOLUTION

- Fill silica powder (Celite 545 silica) into 50 mL-Falcon-tube (to about 2.5-mL = approx. 800mg)
- Add 30 mL dH<sub>2</sub>O
- Shake vigorously (vortex and invert)
- Let slurry settle for approx. 15 min
- Remove (pipette off) the liquid
- Repeat 2 times (a total of 3 washes)
- After last washing step: resuspend the silica powder in about the same amount of water
- (up to about 5 mL)

- STORE the silica solution at RT until further use (silica : H<sub>2</sub>O = 1 : 1)

**Before use:**

- suspend stored silica solution (silica : H<sub>2</sub>O = 1 : 1) by vortexing
- Transfer ~50 µL of silica solution to 2mL-tubes (prepare 1 tube per sample)
- **NB** try to keep the silica suspended during pipetting to ensure an equal distribution
- Add 1mL H<sub>2</sub>O (a final wash step)
- Mix by vortexing
- Centrifuge: full speed (13.200) for 10-20 sec
- Pipette off liquid
- Add 700 µL DNA binding buffer (6M Guanidine thiocyanate)
- Suspend the silica powder in DNA binding buffer
- The silica binding solution is now ready for further use in the protocol (see Methods)

**PREPARATIONS**

- For TissueLyser: Prepare 2 mL-tubes (1 per sample): add 3 metal beads (tungsten carbide beads, 3mm) per tube
- Harvest leaf material (starting amount of material: about 100 mg fresh weight)

**GRINDING**

Use appropriate / available grinding protocol (mortar & pestle, Qiagen TissueLyser,)

**For the TissueLyser:**

- Freeze 2-mL tubes containing leaf material and 3 metal beads in liquid nitrogen
- Grind in TissueLyser by shaking (10 sec at 1/30 speed)
- Re-freeze in liquid nitrogen (>30 sec)
- Grind again in TissueLyser by shaking (10 sec at 1/30 speed)
- Re-freeze in liquid nitrogen (>30 sec)
- Store in liquid nitrogen until lysis buffer is added

**LYSIS**

- Add 800µ Lysis buffer
- Add 4 µL RNaseA (10 µg/ml)
- Vortex (~2 min until the powder is dissolved in the buffer)
- Incubate: 10min at room temperature
- Add 200 µL 3M Sodium Acetate (pH 5.2)
- Mix by inversion of tubes
- Incubate on ice for 5 min
- Centrifuge 13,200 rpm / 5 min / RT (pellet the leaf material)

**DNA BINDING**

- prepare 700 µL silica binding solution (see above)
- transfer 800 µL of the supernatant to the tubes containing silica binding solution)

- !! Do not transfer leaf material!
- Completely resuspend the silica powder by vortexing and inversion of tubes (approx. 20 sec)
- incubate 15 min at RT (on a shaker at 400 rpm and/or invert tubes from time to time)
- Centrifuge 13,200 rpm / 3 min / RT (pellet the silica)
- Remove the supernatant (with pipette)

#### **WASHING (2 times washing)**

- Add 500 mL **wash buffer**
- !! Prepared fresh (see above)!
- Completely resuspend the silica powder by vortexing and inversion of tubes (approx. 20 sec)
- Centrifuge 13,200 rpm / 3 min / RT (pellet the silica)
- Repeat the washing step (optional: a third washing step)
- Remove the supernatant with pipette (as complete as possible)
- optional: short spin and remove residual liquid
- After last washing step: dry the silica in the hood up to 1 hour at RT (make sure there is no wash buffer left)

#### **RESUSPENSION**

- Add 200uL TE buffer or 10mM Tris buffer
- Completely resuspend the silica powder by vortexing and inversion of tubes (approx. 20 sec)
- Incubate: 20 min / RT / with gentle agitation (on a shaker at 400 rpm and/or invert tubes from time to time)
- Centrifuge (for tubes): 13,200 rpm / 5 min / RT (pellet the silica)
- transfer 180 µL supernatant to new tube (avoid transferring silica powder!)
- optional: if there is still silica powder in the preps – repeat the centrifugation
- check for concentration and integrity of DNA
- store the genomic DNA at -20°C for long-term storage or 4°C for short-term storage

#### **VALIDATION OF LOW-COST DNA PREPARATIONS FOR TILLING APPROACHES**

Follow the protocol contained in “Positive control for mutation discovery using agarose gels, version 2.4” available at <http://mvgs.iaea.org/LaboratoryProtocols.aspx> , to test that your DNA is suitable for TILLING and Ecotilling applications.

## 2.4. Example Data

Table 1. Different combinations of self-made (low-cost) buffers and buffers from Qiagen DNeasy Plant Mini kit tested with barley tissue

Sample	1		2		3		4		5		6		7		8	
	+	-	+	-	+	-	+	-	A	B	A	B	A	B	A	B
Lysis	Dneasy kit* +Shredder columns -Shredder columns		Dneasy kit* +Shredder columns -Shredder columns		Dneasy kit* +Shredder columns -Shredder columns		Dneasy kit* +Shredder columns -Shredder columns		<b>Lysis buffer (PBGL)</b>		<b>Lysis buffer (PBGL)</b>		<b>Lysis buffer (PBGL)</b>		<b>Lysis buffer (PBGL)</b>	
DNA binding buffer	Buffer AP3/E*		Buffer AP3/E*		<b>6M Guanidine thiocyanate</b>		<b>6M Guanidine thiocyanate</b>		Buffer AP3/E*		Buffer AP3/E*		<b>6M Guanidine thiocyanate</b>		<b>6M Guanidine thiocyanate</b>	
DNA wash buffer	Buffer AW*		<b>Wash buffer – PBGL</b>		Buffer AW*		<b>Wash buffer-PBGL</b>		Buffer AW*		<b>Wash buffer-PBGL</b>		Buffer AW*		<b>Wash buffer-PBGL</b>	
DNA concentration (ng/μL)	14	13	34	41	7	8	4	10	12	11	12	20	10	16	13	17
Total yield (μg)	2.6	2.4	6.2	7.3	1.3	1.5	0.7	1.9	2.2	2.0	2.2	3.5	1.8	2.8	2.4	3.0
260/280 value	1.95	1.83	1.81	1.91	1.37	1.52	1.41	1.73	1.66	1.63	1.64	1.83	1.75	1.55	1.76	1.71

\*components of Qiagen DNeasy Plant Mini kit

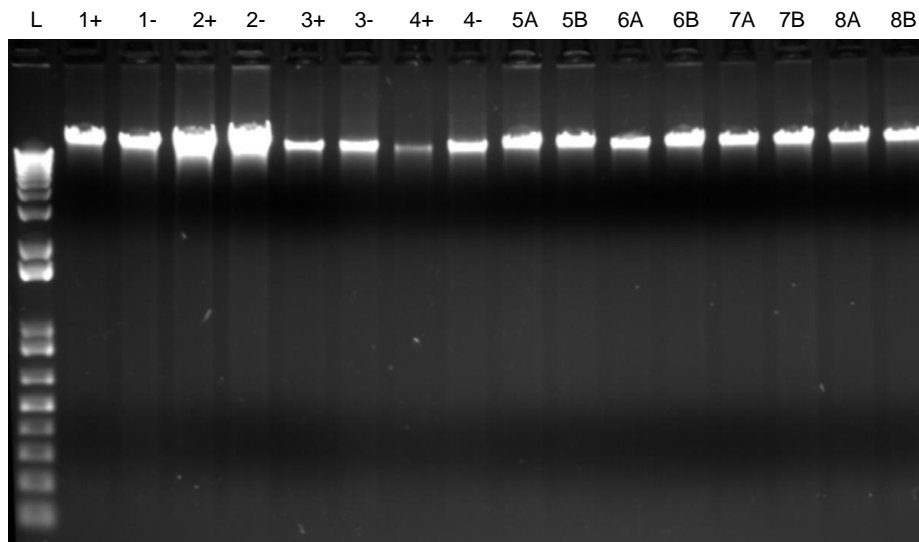


Figure 1. Quality of barley genomic DNA extractions using silica powder and different combinations of self-made (low-cost) buffers and buffers provided by Qiagen DNeasy kit. 8  $\mu$ L of each genomic DNA extraction were separated on a 0.7% agarose gel.

- 1-8: Barley genomic DNA preparation
- +: using QIAshredder columns for the preparation of barley leaf lysates (lysis procedure following the kit instructions)
- : preparation of leaf lysates using the kit instruction (but without using QIAshredder columns)
- A, B: technical replicates
- L: size standard (1 kB Plus DNA ladder - Invitrogen)

All of the genomic DNA preparations show similar DNA concentrations (Table 1) and a good quality of the genomic DNA on the agarose gel (Figure 1). Only the DNA preparations “2+” and “2-” (buffer components from the kit in combination with our wash buffer) show clearly higher concentrations and yields (about 2-3 times higher) than all other DNA preparations. These results indicate that by modifications of the protocol (i.e. modifications of buffers) some improvements of the DNA yields are possible.

The DNA preparations of samples 8A and 8B were extracted exclusively with self-made (low-cost) buffers and show a comparable concentration and yield as the other extractions.

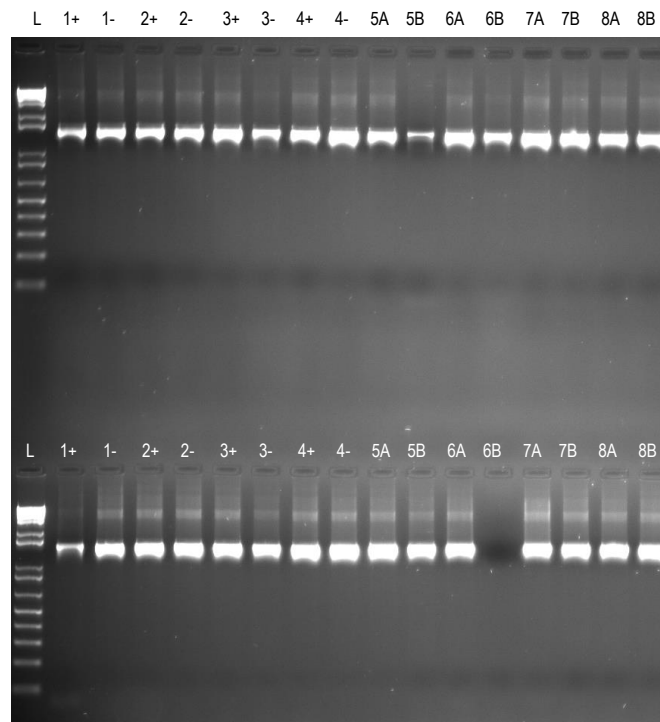


Figure 2. TILLING-PCR products amplified from genomic DNA extractions of barley (obtained by silica-based, low-cost DNA isolation method using different combinations of self-made buffers and buffers provided by Qiagen DNeasy kit). An aliquot of 5uL of each PCR reaction was separated on a 1.5% agarose gel.

top half – Target gene: nb2-rdg2a (1500bp-PCR product);

bottom half – Target gene: nbs3-rdg2a (1491bp-PCR product)

1-8: Barley genomic DNA preparation (see Table 1)

+: using QIAshredder columns for the preparation of barley leaf lysates – Lysis procedure following the kit instructions;

-: preparation of leaf lysates using the kit instruction (but without using QIAshredder columns

A, B: technical replicates

L: size standard (1 kb Plus DNA ladder - Invitrogen)

## 2.5. Conclusions

The DNA extractions from barley using the silica-based, low-cost method provided high-quality genomic DNA and sufficient yield suitable for standard PCR application such as molecular markers and TILLING.

### 3. DNA QUANTIFICATION

This protocol is designed to provide a standardized method for evaluating the quality and quantity of genomic DNA samples extracted from different plant species. Proper quantification and normalization of DNA samples to a common concentration is necessary prior to pooling samples for TILLING or Eco-tilling. A failure to combine genomes at an equal concentration can increase the false positive error rate because some polymorphisms will be represented at a concentration below the limits of detection.

#### 3.1. Protocol for gel electrophoresis

##### 3.1.1. Preparation of DNA concentration standards.

Lambda DNA (Invitrogen cat. # 25250-010) is used as a concentration standard.

- A. Estimate how much concentration standard will be needed for a project (same organism, DNA prepared using the same methods, see 1.B.). Take this volume of DNA and vortex using the same settings as the genomic DNA extraction protocol used. This should shear the DNA to the approximate same size fragments as the genomic DNA. It is important to get the standard near to the same size as the genomic DNA because the intensity of ethidium bromide staining is a product of the size of DNA fragments.
- B. Using the sheared DNA from 1.A, prepare DNA concentration standards at 115 ng/ $\mu$ l, 76.9 ng/ $\mu$ l, 51.3 ng/ $\mu$ l, 34.2 ng/ $\mu$ l, 22.8 ng/ $\mu$ l, 15.2 ng/ $\mu$ l, 10.1ng/ $\mu$ l, 6.8ng/ $\mu$ l, 4.5 ng/ $\mu$ l, and 3 ng/ $\mu$ l. These are derived from the formula:  $3 \times 1.5^i$ ,  $i =$  integers from 0 through 7. This is intended to provide the most accurate binning of DNA concentration estimates when performing visual analysis. Prepare the standards as independent dilutions from the stock of shaken Lambda to avoid cumulative error in low concentration DNA references. Prepare enough of each standard so that you have at least 3  $\mu$ l for every 14 samples. Note that the concentration of lambda DNA may vary from batch to batch. Make sure to calculate dilutions based on the information printed on the stock tube.



### 3.1.2. Preparing agarose gels.

Prepare a 1.5% Agarose gel in 0.5x TBE buffer with 0.15 µg/ml ethidium bromide. Use at least a 24 tooth comb when preparing the gel. Place the solid gel into a rig containing 0.5x TBE buffer with 0.15 µg/ml ethidium bromide.

**CAUTION:** Ethidium bromide is mutagenic. Wear gloves, lab coat and goggles. Dispose of gloves in toxic trash when through. Avoid contaminating other lab items (equipment, phones, door handles, light switches) with ethidium bromide.

### 3.1.3. Preparing samples for loading into gels.

**NOTE:** When you have many samples to quantify, it is best to first test ~28 to determine the range of DNA concentrations from your extraction method. Samples above 62 ng/µl will be diluted to ~ 20 ng/µl for accurate quantification. If the majority of the small test subset have concentrations > 62 ng/ µl, you may want to dilute the rest of the samples prior to the agarose gel assay. This will save a gel run and the time required to estimate DNA concentrations.

- A. Add 3 µl of DNA sample plus 2 µl DNA load dye (30% glycerol plus bromophenol blue – Do not add xylene cyanol as it migrates near the genomic fragment and can interfere with quantification). Use the same volumes for the DNA concentrations standards.
- B. Load the gel. When using a 28 tooth comb, lanes 1-14 should contain genomic DNA samples and lanes 15-28 the concentration standard. Lane 15 should contain the 3 ng/ µl standard, lane 16 the 4.5 ng/ µl standard and so on with lane 28 containing the 115 ng/ µl standard.

### 3.1.4. Running the gel

Run gel at 5-6 V/cm (160V on a large Owl A2 rig, should be about the same for our rigs) for 30-60 min. The DNA sample should be completely out of the well and into the gel about 0.2 cm. Do not run the gel too long as the genomic DNA band will become diffuse and hard to quantify.

**NOTE:** Degraded samples (those producing smeary bands with standard agarose gels) should be run on a 3% MetaPhor agarose gel (~10.5g MetaPhor (Cambrex) in 350ml 0.5x TBE). The preparation of the MetaPhor gel is more specific in that it must be allowed to hydrate in the 0.5x TBE for ~15 min prior to melting. After melting and pouring, allow to set at room temperature, then put in the cold room (4°C) for 15-30 min. This final step is critical for proper setting of the gel.

### 3.1.5. Photographing the gel

It is important to get a proper exposure of the gel that shows a difference in ethidium staining in the concentration ranges you are assaying. For example, if all of your samples are at 20 ng/  $\mu\text{l}$ , you should be able to observe a noticeable difference in the 34.2 ng/  $\mu\text{l}$ , 22.8 ng/  $\mu\text{l}$ , and 15.2 ng/  $\mu\text{l}$  concentration standards. Make sure this is clear on the gel printout.

- A. Adjust the image so as to take the longest possible exposure that does not saturate the image of any of the samples being assayed. It is all right to saturate the image of a reference sample that has higher [DNA] than any of the samples being assayed. Save this image in TIFF format. Print this image.
- B. It may not be possible to set the exposure such that all bands can be visualized without saturating the higher concentration samples. In such a case, a second exposure is required for the notebook, but not for the scoring protocol on the gel documentation system as the computer can score samples that may be difficult to see by eye. Adjust the exposure of the gel so as to allow for the visualization of the lowest [DNA] samples. This will cause the saturation of the images of the highest [DNA] samples. Save this image as a TIFF file. Print this image.

### 3.2. Quantification of DNA using image analysis software

DNA concentrations can be estimated manually by comparing band intensity to the intensity of DNA standards of known concentration. A computer programme that capable of measuring pixel density can provide a more accurate and objective estimation of DNA concentration. In this method a standard curve is created with the DNA concentration standards and sample concentrations are estimated using the standard curve. Many GelDoc systems provide software for automated or semi-automated determination of DNA concentration based on pixel density. We provide here an alternative that will work on any digital tiff image using free image analysis software and Microsoft excel. The method can thus be applied to most labs.

1. The free programme ImageJ (<http://rsbweb.nih.gov/ij/>), is a public domain program developed by Wayne Rasband of the National Institutes of Health, USA Download this onto your computer. Full documentation can be obtained from the website.
2. Open ImageJ
3. Open the tiff image to be analysed (File>Open). A demonstration image titled "Cassava\_DNA\_test2c.tif" can be found on (URL) for practice.

**CAUTION:** Do not use compressed file formats such as jpeg.

4. Straighten the image so that the lanes are parallel with the image dialog box (Image>Rotate>Arbitrarily). In the rotate dialog box, select preview, set the

Grid Lines number to 30, and adjust the angle in degrees until the bands are in line with the grid lines. The Interpolate feature should be selected. Note that you can set negative degrees by placing a minus (-) sign before the angle degree number. You may have to use a decimal setting to get the lanes to line up. When finished, click OK.

5. Subtract background noise (Process>Subtract Background). Deselect “light background”. It is important that you don’t set the rolling ball radius too small. It should be no more than half the width of the box you draw for the band (see step 7).
6. Select the rectangle tool in the ImageJ toolkit dialog box.
7. Find the highest intensity band on the gel to be analysed and draw a box around it. Make sure that the box surrounds the entire signal but does not overlap on the signal from another band. Check the height (h) and width (w) values and make sure that the larger of the two values is not more than 2x the size of the rolling ball radius chosen in step 5.

**TIP:** Select the magnifying tool and make the gel image as large as reasonable.

8. Left click and hold the mouse over the box and move it so that it is positioned around lane 1.

**CAUTION:** The box should contain only signal from the lane to be measured. Failure to do so will lead to an inaccurate reading.

9. Measure the box by hitting the m key. A full screen table should appear with columns for sample #, Area, Mean, min and max values. Minimize the table so that you can again view the gel image.
10. Move the box to lane 2 and hit the m key.

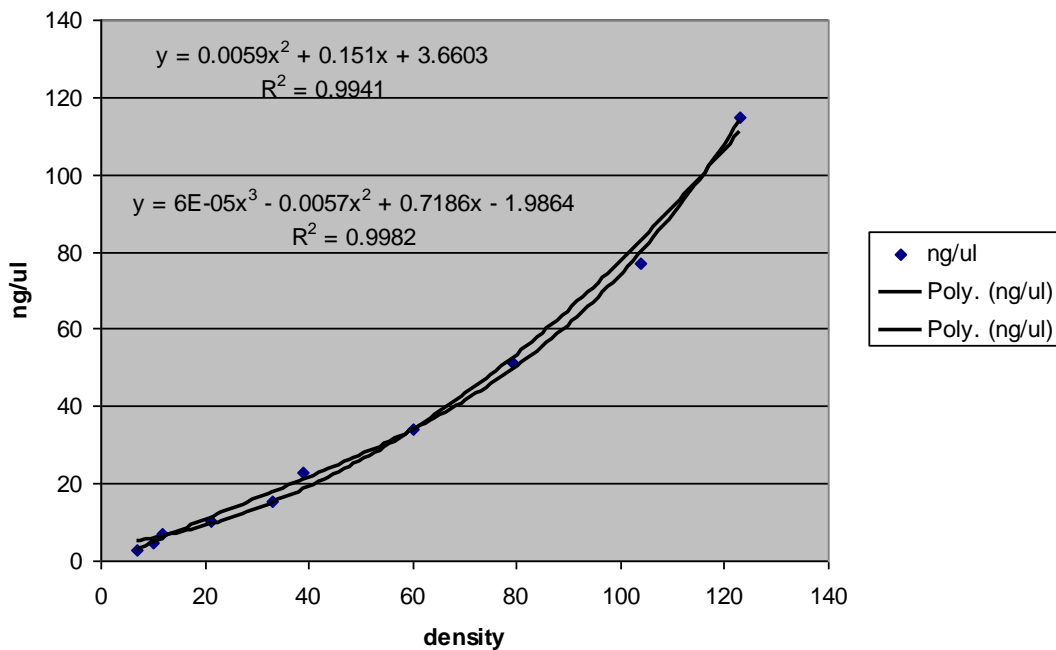
**CAUTION:** Do not change the size of the box. You must measure the same volume of box for each lane. If you accidentally change the size of the box while measuring lanes, start over.

11. Continue to move the box and hit the m key until all the lanes in a gel tier are measured, including the standards.
12. Evaluate the table. Does every sample have the same area value? If not, you have changed the size of the box and you need to start over. Does the number of samples equal the number of lanes on the gel? If not, you either missed a lane or counted a lane more than once. If so, you need to start over.
13. When you are satisfied that the table is correct. Select the entire contents of the table (control A), copy and paste into the raw data section of the excel worksheet.

14. Copy the density (area) from the last 6 samples representing the standards of known concentration in the test image. Paste these data into the density column just below the raw data. The excel table for the test gel image is found on (URL)

**CAUTION:** If you used less than the normal complement of standards, or put the standards in a different order than is represented in the “ng/μl” column, you will need to modify this section appropriately.

15. Select the density and ng/μl columns including the title cells (A, B 41-47 in excel). Click the “Chart Wizard” button. Select XY (Scatter) as chart type and scatter with no point connection as sub type. Click next
16. Select the series in columns. Click next and fill out the title (Gel #), X axis (density) and Y axis (ng/μl). Click next and save the graph as an object in the workbook. Click finish. Move this graph to the graph section of the worksheet.
17. Inspect the graph. Are there any points that are clearly off of the trend? If so, consider removing this data point and re-drawing the graph. This may become more evident once you have drawn the trendline (Step 18).
18. Add a trendline (Chart>Add Trendline). Under type, click polynomial and select 2nd order. Click Options and select “Display equation on chart, and display r-squared value on chart. Click ok. OPTIONAL: You may try a higher order polynomial to evaluate how differences in curve fitting can affect your concentration estimation (see figure below showing second and third order polynomial).



19. Fill in the sample # next to the lane number in the DNA concentration table to the right of the raw data section.
20. Copy and paste the density from the raw data into the density column of the concentration table to the right of the raw data.
21. Insert the formula for the second order polynomial into the first cell of the second order polynomial column. Copy the formula from the graph, then click on the cell, type the equal (=) symbol in the formula box and paste the formula. Replace x<sup>2</sup> with the density data from the first sample. This sample should be in cell J7, so you would replace x<sup>2</sup> with \*j7\*j7. Replace x with \*j7. When finished, press the enter key. The value should appear in the cell.
22. Click on the cell. Pull the right corner so that the box extends over the entire column. You should see all the cells in that column fill with the appropriate values.

Optional: Repeat Step 21 and 22 for the third order polynomial. For x<sup>3</sup>, use \*j7\*j7\*j7. For many cases the second order polynomial will be sufficient. The main differences will be in estimating high (>50 ng/μl) concentration samples.

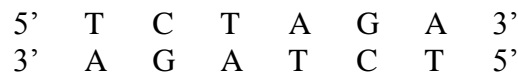
23. Save the gel image in ImageJ as a tif image in a new folder labelled with the gel image name.
24. In the excel workbook, import the tif gel image and place it near the Gel Image field.
25. Compare the band intensities on the image with the concentrations estimated from the standard curve. Do you agree with the estimations? If not, consider repeating the measurement.
26. Compare your data with the data provided in the sample data tab of the excel sheet. Did you get the same results?

## 4. RESTRICTION ENZYME DIGEST

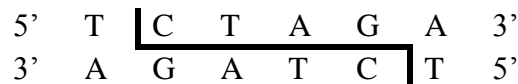
Restriction enzymes are produced by various bacterial strains. In these bacterial strains they are responsible for limiting attack from certain bacteriophages. They act by cutting (“restricting”) the phage DNA at a sequence-specific point, thereby destroying phage activity. Sequence-specific cutting is a fundamental tool in molecular biology. DNA fragments can be ligated back together (“recombined”) by T4 DNA ligase. In addition to cloning and molecular marker applications, restriction digestion is being used for new techniques such as for creation of restriction phased libraries for Next Generation Sequencing (NGS). Many restriction enzymes have been cloned and are available in a commercially pure form. They are named after their bacterial origin: e.g. *EcoRI* from *E. coli*.

The known restriction enzymes recognize four or six bases (eight in the case of “very rare cutters” like *NotI* and *SfiI*). Recognition sequences are almost always “palindromic” where the first half of the sequence is reverse-complementary to the second:

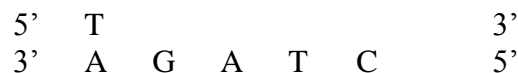
e.g. the *XbaI* site is



The position of the actual cut is enzyme dependent and symmetrical on the opposite strand:



leaving cohesive termini (sticky ends) at the 5' end:



The commercially available restriction enzymes are supplied with the appropriate restriction buffers (10 x concentrated). The enzymes are adjusted to a specific activity per  $\mu\text{l}$ , usually 10 U/ $\mu\text{l}$ . (1 Unit is the amount of enzyme needed to cut 1  $\mu\text{g}$  of lambda DNA in one hour at 37°C).

A typical restriction digestion is performed using between 20 $\mu\text{l}$  and 100 $\mu\text{l}$  reaction volume per 5  $\mu\text{g}$  and more of plant DNA. For purified plasmid DNA 2 U per  $\mu\text{g}$  DNA is sufficient, for plant DNA 4 U per  $\mu\text{g}$  should be used.

For example: digestion of 5 µg DNA in 40 µl reaction volume:

Restriction buffer (10x)	4 µl
DNA 1 µg/µl	5 µl
Doubled distilled H <sub>2</sub> O	29 µl
Enzyme (10 U/µl)	2 µl

Incubate for at least 1 hour at 37°C. The restriction enzyme can be inactivated by heating to 65°C for 10 minutes or by adding 1.0µl 0.5 M EDTA.

Note however, that protein engineering and advanced biochemistry have allowed major improvements from the canonical restriction digestions above. For example, ThermoScientific have developed a suite of fast enzymes that can digest complete genomes in 15 minutes, versus the traditional overnight digestion. Such digestions can be accomplished with no star activity.

## 5. FINDING CANDIDATE GENES AND PRIMER DESIGN FOR MOLECULAR TESTING: AN EXAMPLE FROM THE ANNOTATED *SORGHUM BICOLOR* GENOME.

### 5.1. Overview

There are several levels of genome annotation. The goal of this method is to quickly identify annotated genes and recover gene and transcript/protein sequences from the Sorghum genome that have potentially interesting biological function, without extensive bioinformatics expertise or tools. The same methods can be applied to many other annotated genomes. Genome project websites typically have text files of genome annotations. Many genome projects use the same generic genome browser architecture, and so retrieval of sequences described here will work for different species. For example, there are many genomes available on Phytozome.

Retrieve a list of annotated genes in the Sorghum genome. This file:

[ftp://ftp.jgi-psf.org/pub/compngen/phytozome/v8.0/Sbicolor/annotation/Sbicolor\\_79\\_annotation\\_info.txt](ftp://ftp.jgi-psf.org/pub/compngen/phytozome/v8.0/Sbicolor/annotation/Sbicolor_79_annotation_info.txt)

while not the most verbose annotation it is easily opened and searchable.

- Open this file up and hit control F, you can do a quick text match search for keywords like disease. If you search for disease, you get >100 hits.
- The first hit for a text search of disease is Sb0019s003010.1.
- Recover sequences for your favourite gene
- There are (at least) two ways to retrieve the sequence for primer design.
  1. First, you can search NCBI (<http://www.ncbi.nlm.nih.gov>). You need to remove the “.1” at the end because this delineation is not in NCBI. What you’ll get is an 800,000 bp scaffold that contains the gene sequence. Unfortunately, it contains many predicted proteins, but the annotation isn’t there. Which means that it is very hard to find the protein you’re looking for unless you blast all the hypothetical peptides. This isn’t very convenient.
  2. To retrieve genomic, cDNA and protein sequences, goto the genome website <http://www.phytozome.net/sorghum>. Click “Browse Genome” and then enter Sb0019s003010.1 into the landmark or region window and click search. You’ll get the gene model back with blast hits to other plant proteins. Move the mouse over this pile up and you’ll get individual annotations from



the different species (this is good to do to double check you have the correct gene).

- Download sequences for downstream analysis and primer design.
- In many cases (such as TILLING/Scotilling) it is best to be searching for potentially functional variation. So, it will be more efficient to screen exonic regions. In this example notice the exonic regions are mostly on the left side.
- It is not very intuitive how to get both the genomic and transcript sequence from this graphical output. Put your mouse over the transcript and right click. A new window will appear from phytozome and you can get the sequences you need from the sequencing tab.

**FOR TILLING and Scotilling applications design primers following protocol in chapter section 13.2.1.**

## 6. SSR

**SSR (Microsatellite) definition:** Any one of a series of very short (2-10 bp), middle repetitive, tandemly arranged, highly variable (hypervariable) DNA sequences dispersed throughout fungal, plant, animal and human genomes (Kahl, 2001).

Simple sequence repeats (SSR) or microsatellites are a class of repetitive DNA elements (Tautz and Rentz, 1984; Tautz, 1989). The di-, tri- or tetra-nucleotide repeats are arranged in tandem arrays consisting of 5 – 50 copies, such as (AT)<sub>29</sub>, (CAC)<sub>16</sub> or (GACA)<sub>32</sub>. SSRs are abundant in plants, occurring on average every 6-7 kb (Cardle et al., 2000). These repeat motifs are flanked by conserved nucleotide sequences from which forward and reverse primers can be designed to PCR-amplify the DNA section containing the SSR. SSR alleles, amplified products of variable length, can be separated by gel electrophoresis and visualised by silver-staining, autoradiography (if primers are radioactively labelled) or via automation (if primers are fluorescently labelled) (Figures 6.1 and 6.2). SSR analysis is amenable to automation and multiplexing (Figure 6.2), and allows genotyping to be performed on large numbers of lines, and multiple loci to be analysed simultaneously. SSRs can be identified by searching among DNA databases (e.g. EMBL and Genebank), or alternatively small insert (200-600bp) genomic DNA libraries can be produced and enriched for particular repeats (Powell et al., 1996). From the sequence data, primer pairs (of about 20 bp each) can be designed (software programmes are available for this).

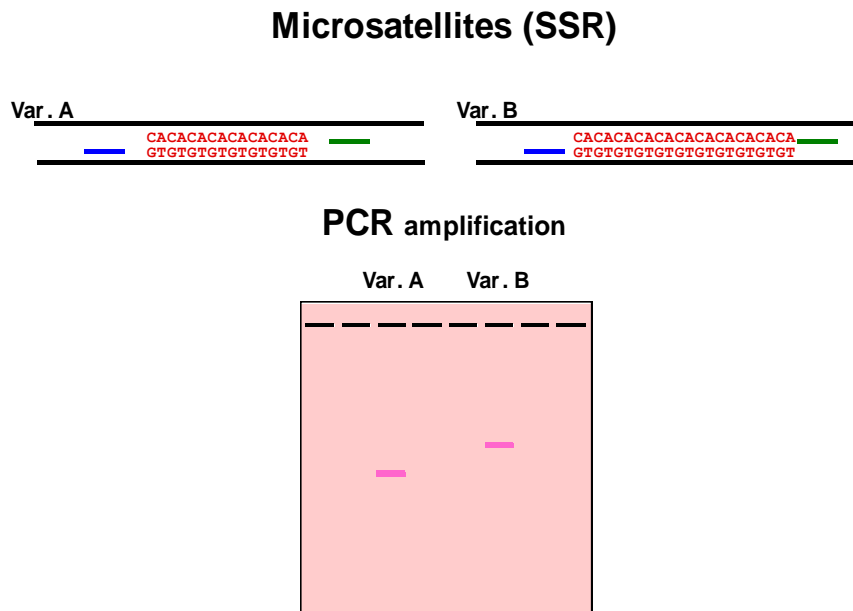


Figure 6-1. The schematic above shows how SSR variation (short A and long B) can be detected using gel electrophoresis after PCR with forward (blue) and reverse primers (green) (with permission, K. Devos).

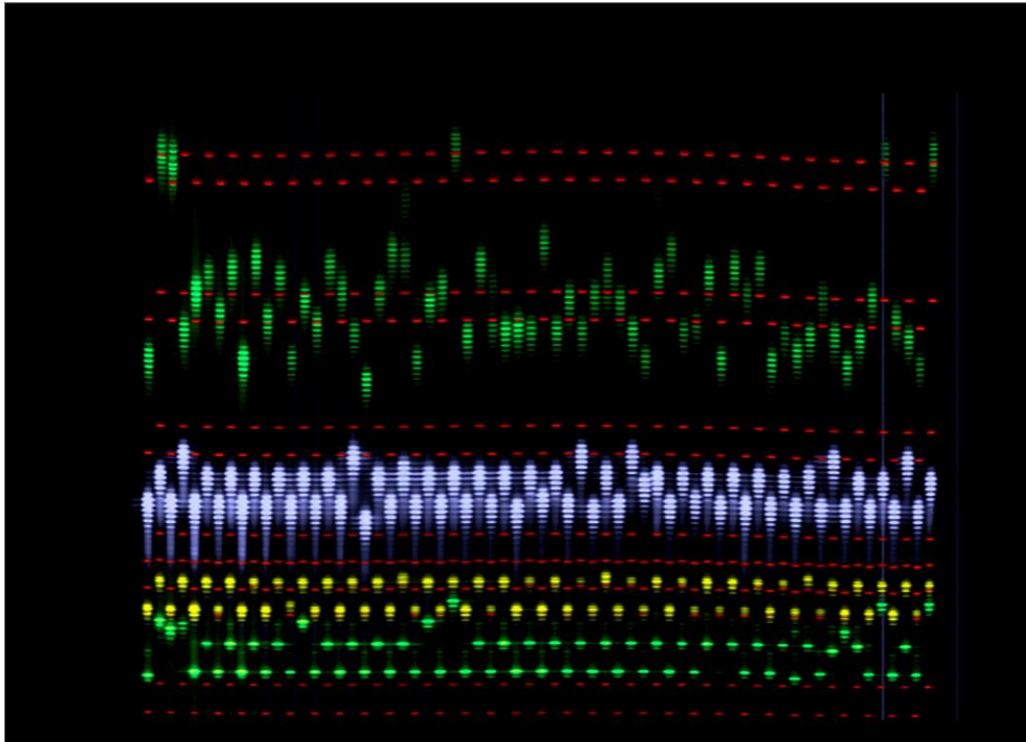


Figure 6-2. A computer image showing an example of SSR multiplexing with different colours (with permission, J. Kirby and P. Stephenson).

## 6.1. Protocol

### 6.1.1. PCR reaction mix

Microsatellite primers are specific for each individual genome or species. It is essential to know that the primer pairs chosen will work for your given species.

NOTE: Wear gloves and lab coat at all times for safety and to prevent contamination.

#### Prepare 25 µl Reaction Mix

1. Take four sterile PCR tubes and to each add the following:

10 x Taq buffer	2.5 µl
MgCl <sub>2</sub> (25mM)	1.5µl
dNTPs (10 mM)	1.0 µl
Forward primer (10 µM)	0.8 µl
Reverse primer (10 µM)	0.8 µl
Taq DNA polymerase (5U/µl)	0.25 µl
DNA (20ng/µl)	1.0 µl
*Add sterile distilled water up to	<b>25µl</b>

2. Mix by gently tapping against the tube.
3. Centrifuge briefly (~14,000 rpm for 5 seconds).

NOTE: Keep all reagents and reaction mix on ice until used.

### 6.1.2. PCR amplification

Place tubes in a PCR machine and amplify using a programme designed for the primers being used; an example is given below:

Step 1	Initial denaturing	94°C	5 minutes
Step 2	Denaturing	94°C	1 minute
Step 3	Annealing*	55°C	1 minute
Step 4	Extension	72°C	2 minute
Step 5	Cycling	repeat steps 2-5 for 34 cycles	
Step 8	Final extension	72°C	5 minutes
Step 9	Hold	4°C	forever

\*NOTE: The annealing temperature (*Step 3*), in particular, can and does vary with primers used. Please note this when changing primers.

### 6.1.3. Separation of the amplification products in agarose gel

NOTE: Where SSR polymorphism is large, bands can be separated in agarose gels, however small base-pair differences among alleles require separation in polyacrylamide gels.

1. Take 5µl of the PCR product into a fresh tube.
2. Add 2 µl 5X loading buffer containing dye.
3. Centrifuge briefly (14,000 rpm for 5 seconds).
4. Load all 7µl of the mixture into a 1.5 % agarose gel (which is made up of 25% fine agarose and 75% normal agarose with 2µl/100ml ethidium bromide for staining DNA).
5. Run gel until dark blue colour marker has run two thirds of the gel.

NOTE: Do not run the dye off the gel or you will also lose your DNA samples.

NOTE: See Section of RFLP Protocol (Agarose gel electrophoresis) for details of gel preparation and running.

6. Stain gel with ethidium bromide (*Caution: ethidium bromide is toxic: wear gloves and lab coat and avoid inhalation*).
7. Visualise under UV light (*Caution: wear gloves, and UV protective glasses or a shield over your face when you are exposed to the UV light of the transilluminator*).

### 6.1.4. Denaturing gel electrophoresis

NOTE: Denaturing the samples produces single-stranded DNA, which is used for detection in polyacrylamide gels (see below). Single-stranded detection is preferred as it results in a greater clarity in band separation for detection. Setting up and casting a polyacrylamide gel using sequencing apparatus involves the followings.

### 6.1.5. Assembling the glass plate sandwich

1. Wear gloves and lab coat, and place the Integral Plate Chamber (IPC), i.e. the big plate on the bench, horizontally, glass side up. Clean the upper surface of the glass plate using Alconox and warm water. Rinse and dry the plate.
2. Clean the upper surface with 95% ethanol. Apply a thin film of Sigmacote (2ml) to the upper surface of the plate and spread evenly using blue roll and dry. Repel silane or Repelcote are other brand names of the same product.

NOTE: Change gloves between working with the bigger and smaller plates as you will be using 2 different chemicals, bind silane and repel silane that must not contaminate the unintended glass plate. One is a 'binder' while the other repels and when properly applied ensure that the gel sticks only one surface and not the other. Contamination can be brought about by not changing gloves and this will lead to breakage of the gel between the 2 plates!

3. Clean the smaller plate using Alconox and water (you may also need to use a razor blade to remove old bits of gel that have stuck). Rinse and dry the plate, clean the upper surface only with 95% ethanol.
4. Prepare fresh bind silane solution by adding 3 $\mu$ l of binding solution to 1ml of 95% ethanol mixed with 5 $\mu$ l of glacial acetic acid.
5. Apply prepared bind silane solution to the upper surface of the plate and spread evenly using blue roll.

NOTE: Clean everything following use, and dispose of materials carefully according to the regulations of your organization.

NOTE: The glass plates must be meticulously clean. Detergent microfilm left on the glass plate may result in a high (brown coloured) background for the stained gel.

6. Place clean, dry spacer on the long edge of the IPC plate. Make sure that there is no untrimmed adhesive underneath the spacer.
7. Place the outer glass plate on the top of the spacers. The raised plastic edges on the IPC plates will help position the spacer and plate. Align the outer plate and spacer with bottom edge. Precise alignment is necessary.
8. Slide clamps over the gel plate assembly, one clamp at a time. This can be done while holding the IPC vertically. Start each clamp (there is right and left clamp) near the bottom end first, then slide the clamp on to the IPC assembly until it snaps into a place along the entire length.

NOTE: The clamps must fit reasonably tightly to prevent the spacer from leaking. Make sure the clamps are all the way on, with the spacer and outer glass plate flush at the bottom.

### 6.1.6. Casting gel

1. Prepare 100 ml of gel solution per plate by adding together:

*Acrylamide/bis solution 19:1 (40 %)	15 ml
TBE (10X)	10 ml
Urea 8 M	50 g
Make up to 100 ml with distilled water.	

\* *Caution: acrylamide is toxic*

NOTE: An alternative option is to use a pre-mixed solution, SequaGel®XR (National Diagnostics, Inc.), which gives sharper bands

2. Filter the solution and keep at 4°C and take as required when ready to cast a gel.
3. Add 28 µl TEMED (Caution: TEMED is corrosive) and 800 µl 10% fresh ammonium persulphate solution (Caution: ammonium persulphate, APS, is harmful) to 100 ml of the gel mix,
4. Gently draw up acrylamide solution into a 100ml syringe, avoiding air bubbles.
5. Adjust angle of plates so gel solution flows slowly down one side. Keep the acrylamide solution flow consistent by varying the flow rate by tilting the gel assembly. This reduces the formation of bubbles during the filling. Perfect clean plates will not allow bubbles to form. If bubbles do form, tap the glass plate gently to dislodge them.

NOTE: Gel will start to polymerize after adding APS, be prepared to move quickly.

6. Insert the flat side of a 0.4mm shark's tooth comb between plates before the gel polymerizes. Place the binder clamps over the glass plates to insure that the plates are held firmly against the comb
7. Leave to polymerise for approximately 1 hour.

NOTE: Make up the developer for silver-staining while the gel is polymerising, see section 5.1.6 below.

### 6.2. Setting up the operation

1. Place the IPC assembly into the universal base, against the back of the wall. Stick a gel temperature indicator on to the outer plate, somewhere near the centre of the gel, to monitor the temperature during electrophoresis.
2. Fill the upper buffer chamber with 1X TBE buffer. The level of the buffer should be about 1cm from the top all the time during the run.
3. Fill the lower buffer chamber and adjust the levelling screws. Do not fill the lower chamber with more than 500ml of buffer

4. Remove the comb from the gel and clean the well space using distilled water. Replace comb carefully, teeth first this time.

NOTE: You can only replace the comb once, so be very careful!

5. Pull the plastic hood over the gel tank and insert the electrodes. Switch on the power pack and adjust the reading roughly to 900-1500 V and 70W.
6. Pre-run the gel at 125 watts. The gel temperature will stabilize near 55°C. Pre-running the gel at 45°C for an hour or two may result in better resolution, particularly if you use high catalyst concentration

### 6.3. Polyacrylamide gel running conditions

1. Prepare samples by adding 2 µl of formamide dye mix to 8 µl of your PCR reaction (second half). Denature the samples for 5 minutes and place on ice (Caution: formamide is harmful).
2. Load 1 kb marker ladder (10 µl 1 kb ladder (50 ng/µl) add 6 µl formamide loading buffer); load 5 µl into first lane (and at convenient intervals across the gel).
3. Load 8 µl of each sample containing the formamide dye mix into individual wells of the gel.
4. Run gel for approximately 1 hour and 20 minutes at 75 watts or until just before the dark blue runs off the bottom of the gel. You will need to quantify the best time for your particular PCR products.

NOTE: Do not run the dye off the gel or you will also run your sample off the gel and lose it.

### 6.4. Silver-staining

1. While the gel is polymerising, prepare the developer solution: Dissolve 60 g sodium carbonate in 2 litres of distilled water then add 400 µl of sodium thiosulphate solution (10 mg/ml) and 3 ml formaldehyde (37% solution) and store at 4°C (*Caution: Both sodium carbonate and formaldehyde are toxic, avoid inhalation and wear gloves and lab coat*). For best results, the developer must be chilled.
2. While the gel is running, prepare the fixative (10 % acetic acid): Add 200 ml glacial acetic acid to 1.8 litres distilled water (*Caution: acetic acid is corrosive, gloves should be worn*).
3. Prepare the silver-stain (*toxic, wear gloves*): Add 2g silver nitrate (AgNO<sub>3</sub>) solution in 2 litres of distilled water (*Caution: silver nitrate is corrosive, gloves should be worn*). Then add 3 ml formaldehyde (37% solution) and mix (*Caution: formaldehyde solution is toxic, Wear gloves and lab coat, and avoid inhalation*). Silver nitrate is light sensitive so store in an opaque bottle or wrap aluminium foil around the bottle.
4. Remove the gel from the rig and separate the plates. Place the gel in a tray with the fixative and leave shaking in a fume hood for 20 minutes.

NOTE: Do not pour solutions directly onto the gel as it may come off the plate! When running

5. Remove the gel and stand on a rack. Pour off fixative and save it as it can be used for up to 10 times. Wash the gel three times (2 min) in water. Remove the gel and stand. Pour out the water and replace with silver-stain, introduce the gel again and leave shaking for 30 minutes. For best results, cover the tray as light affects the AgNO<sub>3</sub> solution

NOTE: Silver stain (AgNO<sub>3</sub> and formaldehyde solution) can be re-used up to 10 times

NOTE: The next few procedures have to be followed quickly and carefully so make sure you have everything set up and ready.

6. Remove gel from the silver-stain solution and rest it on a tray containing water (do not put it in the water yet). Dispose of spent stain according to the regulations of your organization. Rinse the box that contained the silver-stain with water.
7. Set a timer for 10 seconds. Start the timer and quickly lower the gel into the water. Agitate several times to remove all excess silver-stain. When 10 seconds is up quickly drain the gel and place it in the developing solution.
8. Agitate the gel in developer solution and, use a piece of white paper placed behind the gel to check progress of the band development. Keep an eye on the gel as it develops. Stop the reaction when bands start to appear near the bottom of the gel (*i.e.*: 70 bp marker on the 1 kb ladder) by taking the gel out of the developer solution.
9. Put the gel into a tray containing 2 litres of stop solution (10% glacial acetic acid) for 5 minutes.

NOTE: The stop solution could be what was saved from earlier (first step fixative) if there is no need for re-use. If re-use is desired, it is best to have separate fixative and stop solutions as the latter contains AgNO<sub>3</sub> and therefore not suitable for use again as fixative.

10. Rinse gel in water for 5 minutes and leave it to dry standing vertically.
11. Gels can be recorded or documented using Kodak duplicating film.
  - 11.1. Place the glass plate upside down on the film.
  - 11.2. Expose to room light for 15-17 seconds (depending on the room light intensity).

NOTE: The longer the light exposure, the brighter the film gets following development.

Gels can be scanned or photocopied.

## 6.5. References

- Cardle, L., L Ramsay, D. Milbourne, M. Macaulay, D. Marshall, and R. Waugh, 2000. Computational and experimental characterisation of physically clustered simple sequence repeats in plants. *Genetics*. **156**: 847-854.
- Kahl, G., 2001. *The Dictionary of Gene Technology*. Wiley-VCH, Weinheim.
- Powell, W., G. C. Machray, and J. Provan, 1996. Polymorphism revealed by simple sequence repeats. *Trends in Plant Sci.* **1**(7): 215-222.



Tautz, D., 1989. Hypervariability of simple sequences as a general source for polymorphic DNA markers. *Nucleic Acids Res.* **17**: 6463-6471

Tautz, D. and M. Rentz, 1984. Simple sequences are ubiquitous repetitive components of eukaryotic genomes. *Nature.* **322**: 652-656.

## 6.6. Reagents needed

- Use only sterile distilled water for all solutions.

- *Taq* buffer

- dNTPs

- Alconox

- Repel silane (Repelcote, Sigmacote)

- Bind silane

- Sterile distilled water

- Primers

- *Taq* DNA polymerase (5U/ $\mu$ l)

- DNA (10-20ng/ $\mu$ l)

- 10 x loading buffer

Glycerol (80%)	600 $\mu$ l
----------------	-------------

Xylene cyanol	2.5 mg
---------------	--------

Bromophenol blue	2.5 mg
------------------	--------

Distilled water	400 $\mu$ l
-----------------	-------------

- 5 x loading buffer

Glycerol (80%)	300 $\mu$ l
----------------	-------------

Xylene cyanol	1.3 mg
---------------	--------

Bromophenol blue	1.3 mg
------------------	--------

Distilled water	400 $\mu$ l
-----------------	-------------

- Ethidium bromide

- Agarose

- Acrylamide

- Bis-acrylamide

- TEMED

- Ammonium persulphate

- Sodium thiosulphate

- TBE

H <sub>2</sub> O	~800 ml
------------------	---------

Tris base	108 g
-----------	-------

Boric acid	55 g
------------	------

EDTA	9.3 g
------	-------

ddH <sub>2</sub> O	Adjust volume to 1 litre
--------------------	--------------------------

- 100% ethanol

- Bind silane

- Sodium carbonate

- Glacial acetic acid

- Formamide dye mix (for 1 ml)

Formamide (deionized)	950µl
dd H <sub>2</sub> O	30µl
EDTA (0.5 M)	20µl
Bromophenol blue	1 mg
Xylene cyanol	1 mg

Mix and store at -20°C

## 7. ISSR

ISSR amplification definition: A variant of the polymerase chain reaction that uses simple sequence repeat primers (e.g.  $[AC]_n$ ) to amplify regions between their target sequences (Kahl, 2001).

Inter-SSR (ISSR) amplification is an example (one of many) of a PCR-based fingerprinting technique. The technique exploits the abundant and random distribution of SSRs in plant genomes by amplifying DNA sequences between closely linked SSRs (Figure 6.1). The method used in the FAO/IAEA course used 3'-anchored primers to amplify regions between two SSRs with compatible priming sites (Yang *et al.*, 1996). More complex banding patterns can be achieved using 5'-anchored primers that incorporate the SSR regions in their amplification products, and by combining 3'- and 5'- primers (Zietkiewicz *et al.*, 1994).

Other methods of fingerprinting using primers complementary to SSR motifs involve using SSR specific primers in combination with an arbitrary primer (Davila *et al.*, 1999), or in combination with primers that target other abundant DNA sequences such as retrotransposons (Provan *et al.*, 1999).

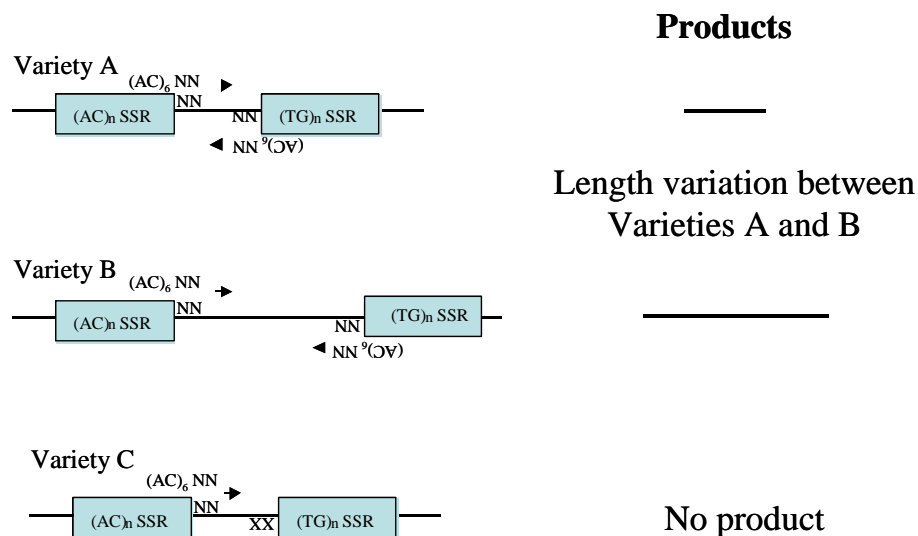


Figure 7-1. The above scheme shows how sequence variation between two SSRs results in variation in PCR products in varieties A, B and C. The figure shows variation at only one ISSR locus, amplification of all compatible ISSR loci among the genomes of a range of varieties will result in complex, fingerprinting, banding patterns.

### 7.1. Protocol

In the example below, one of three primers given in the ISSR protocol of Yang *et al.*, (1996) is used; this produces a relatively simple fingerprint (small number of bands). In more recent applications two or more primers have been used to produce more complex banding profiles (similar to AFLP profiles).

NOTE: Wear gloves and lab coat at all times for safety and to prevent contamination.

### 7.1.1. Prepare 20µl reaction mix

1. Take one PCR tube and add:

10x PCR buffer	2.5 µl
MgCl <sub>2</sub> (25mM)	1.5µl
Primer (10 mM)	2.5 µl
dNTPs (10mM)	0.8 µl
DNA (20ng/ µl)	1.25 µl
<i>Taq</i> DNA polymerase (5 U/ µl)	0.2 µl
Add sterile distilled water to bring volume to 20 µl	

2. Mix by tapping bottom of tube.
3. Centrifuge briefly (14,000 rpm for 5 seconds)

NOTE: Keep all reagents and reaction mix on ice.

### 7.1.2. PCR amplification

Place tube in a PCR machine and amplify using a programme designed for the primer(s). In this example the following programme can be used:

<i>Step 1</i>	Initial denaturing	94°C	7 minutes
<i>Step 2</i>	Denaturing	94°C	30 seconds
<i>Step 3</i>	Annealing*	54°C	45 seconds
<i>Step 4</i>	Extension	72°C	2 minute
<i>Step 5</i>	Cycling	repeat steps 2-4 for 30 cycles	
<i>Step 6</i>	Final extension	72°C	7 minutes
<i>Step 7</i>	Hold	4°C	forever

### 7.1.3. Separation and visualization of the amplification products

1. Add 2 µl of 5x loading buffer to 8 µl of PCR sample.
2. Vortex briefly.
3. Centrifuge briefly (14,000 rpm for 5 seconds)
4. Load samples into a non-denaturing 6% polyacrylamide gel/3M urea gel (see Section 5.1.4. of SSR protocol for preparation of 6% acrylamide gel. [Step 4: Use 180 g urea (3M) instead of 480 g (8M)!])

NOTE: Where the running of polyacrylamide gels is not feasible, 1.5% agarose gel may be used for fragment separation. For this, load sample into 1.5% agarose gel. A mixture of 25%

fine agarose and 75% routine agarose works very well (see Section 6.1.3. of SSR protocol for preparation of agarose gel [Step 4]).

#### 7.1.4. Gel running conditions

1. Run gel under non-denaturing condition at 12 V/cm for 10-13 hours.

NOTE: This is normally done overnight.

NOTE: Non-denaturing gels are run at low voltages and 1 x TBE to prevent denaturation of small fragments of DNA by the heat generated in the gel during electrophoresis.

2. Run agarose gel at 120V for at least 2 hours

NOTE: Do not run the bands off of the bottom of the gel.

#### 7.1.5. Silver-staining

Follow Section 6.1.6 of SRR Protocol (silver-staining).

### 7.2. Primers available at Plant Breeding & Genetics Laboratory (FAO/IAEA)

Primers ID	Sequence information	Primers ID	Sequence information
ISSR-1	(CAC) <sub>7</sub> T	ISSR-27	(GT) <sub>8</sub> G
ISSR-2	(GA) <sub>9</sub> C	ISSR-28	(AC) <sub>8</sub> T
ISSR-3	GT) <sub>9</sub> G	ISSR-29	(AC) <sub>8</sub> C
ISSR-4	(CAC) <sub>7</sub> G	ISSR-30	(AC) <sub>8</sub> G
ISSR-5	GT(CAC) <sub>7</sub>	ISSR-31	(TG) <sub>8</sub> A
ISSR-6	GTG) <sub>7</sub> C	ISSR-32	(TG) <sub>8</sub> G
ISSR-7	(CA) <sub>10</sub> G	ISSR-33	AG) <sub>8</sub> Y T
ISSR-8	(CT) <sub>9</sub> G	ISSR-34	(GA) <sub>8</sub> Y T
ISSR-9	(GA) <sub>9</sub> A Y	ISSR-35	(CT) <sub>8</sub> R A
ISSR-10	BDB(TCC) <sub>5</sub>	ISSR-36	(CT) <sub>8</sub> R C
ISSR-11	HVH(TCC) <sub>5</sub>	ISSR-37	(CA) <sub>8</sub> R T
ISSR-12	(AG) <sub>8</sub> T	ISSR-38	(CA) <sub>8</sub> R C
ISSR-13	(AG) <sub>8</sub> G	ISSR-39	(GT) <sub>8</sub> Y A
ISSR-14	(GA) <sub>8</sub> T	ISSR-40	(GT) <sub>8</sub> Y G
ISSR-15	(GA) <sub>8</sub> C	ISSR-41	(TC) <sub>8</sub> R T
ISSR-16	(GA) <sub>8</sub> A	ISSR-42	(AC) <sub>8</sub> Y G
ISSR-17	(CT) <sub>8</sub> A	ISSR-43	(AC) <sub>8</sub> Y A
ISSR-18	(CT) <sub>8</sub> G	ISSR-44	(AC) <sub>8</sub> Y T
ISSR-19	(CT) <sub>8</sub> T	ISSR-45	(TG) <sub>8</sub> R T

ISSR-20	(CA) <sub>8</sub> A	ISSR-46	(TG) <sub>8</sub> RC
ISSR-21	(CA) <sub>8</sub> G	ISSR-47	(ACC) <sub>6</sub>
ISSR-22	(GT) <sub>8</sub> A	ISSR-48	(ATG) <sub>8</sub>
ISSR-23	(GT) <sub>8</sub> C	ISSR-49	(CTC) <sub>6</sub>
ISSR-24	(GT) <sub>8</sub> T	ISSR-50	(GAA) <sub>6</sub>
ISSR-25	(TC) <sub>8</sub> A	ISSR-51	(GACA) <sub>6</sub>
ISSR-26	(GT) <sub>8</sub> C	ISSR-52	(TCC) <sub>5</sub> RY

Y=C/T

R=A/G

### 7.3. References

- Davila, J. A., Y. Loarce, and E. Ferrer, 1999. Molecular characterization and genetic mapping of random amplified microsatellite polymorphism in barley. *Theor. Appl. Genet.* 98: 265-273
- Provan, J., W. T. B. Thomas, B. P. Forster, and W. Powell, 1999. Copia-SSR: a simple marker technique which can be used on total genomic DNA. *Genome.* 42: 363-366
- Yang, W., A. C. De Olivera, I. Godwin, K Schertz, and J. L. Bennetzen, 1996. Comparison of DNA marker technologies in characterizing plant genome diversity: variability in Chinese sorghums. *Crop Sci.* 36: 1669-1676
- Zietkiewicz, E., A. Rafalski, and D. Labuda, 1994. Genome fingerprinting by simple sequence repeat (SSR)-anchored Polymerase Chain Reaction Amplification. *Genomics.* 20: 176-183

### 7.4. Reagents needed

Use only sterile distilled water for all solutions:

- *Taq* buffer
- dNTPs
- Sterile distilled water
- Primer(s)
- *Taq* DNA polymerase (5U/ µl)
- DNA (10-20 ng/ µl)
- 10 x loading buffer
- Glycerol (80%)                      600 µl
- Xylene cyanol                        2.5 mg
- Bromophenol blue                   2.5 mg
- Water                                    400 µl
- 5 x loading buffer
- Glycerol (80%)                      300 µl
- Xylene cyanol                        2.5 mg
- Bromophenol blue                   2.5 mg
- Water                                    400 µl
- Ethidium bromide

- Agarose
- Acrylamide
- Bis-acrylamide
- TEMED
- Ammonium Persulphate
- Alconox
- TBE (see 5.3)
- Ethanol(95%)
- Repelcote (Symacote)
- Bind silane
- Sodium carbonate
- Glacial acetic acid
- Sodium thiosulphate

## 8. AFLP

**Amplified Fragment Length Polymorphism (AFLP)** is basically a fingerprinting technique. It is a method by which selection of restricted fragments of a total genomic DNA digest is detected by PCR amplification. It is a combination of hybridisation and amplification-based strategies.

The AFLP technique combines components of RFLP analysis with PCR technology (Vos *et al.*, 1995). Total genomic DNA is digested with a pair of restriction enzymes, normally a frequent and a rare cutter. Adaptors of known sequence are then ligated to the DNA fragments. Primers complementary to the adaptors are used to amplify the restriction fragments. The PCR-amplified fragments can then be separated by gel electrophoresis and banding patterns visualized (Figure 8-1). A range of enzymes and primers are available to manipulate the complexity of AFLP fingerprints to suit application. Care is needed in selection of primers with selective bases.

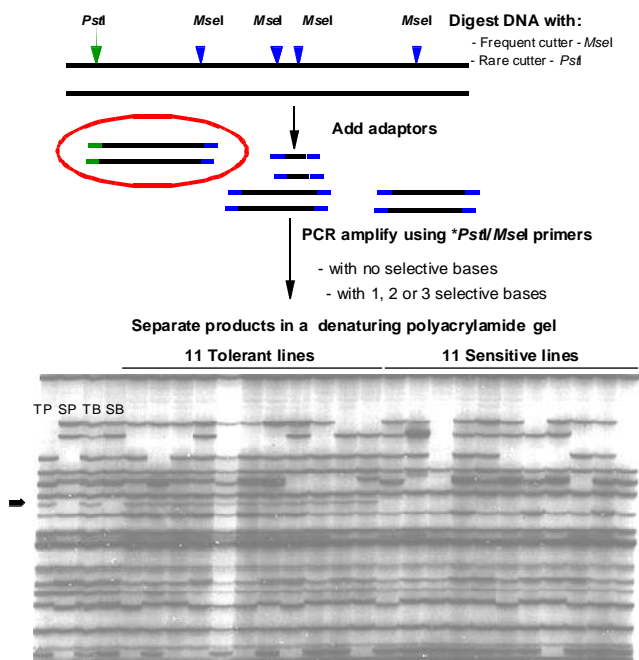


Figure 8-1. In the figure above AFLP profiles have been used in bulk segregant analysis to detect a band associated with tolerance to aluminium in rye, the arrow shows the presence or absence of a band in the tolerant (TP) and susceptible (SP) parents, tolerant (TB) and susceptible (SB) bulks, and 11 tolerant and 11 susceptible individuals (scheme and data with permission, K. Devos and Miftahudin, respectively).



## 8.1. Protocol

AFLP involves four major steps:

- I\* Cutting genomic DNA with restriction enzymes
- II\* Ligating double-strand adaptors to the restriction fragments
- III Amplifying (pre- and selective amplification) restriction fragments using primers
- IV Gel analysis of the amplified products

\*OPTIONAL: these two steps can be performed in one reaction

**NOTE:** Wear gloves and lab coat at all times for safety and to prevent contamination.

### 8.1.1. Restriction of genomic DNA and ligation of adapters to the DNA fragments

Two pairs of restriction enzymes, *MseI/Tru9I* and *PstI/EcoRI*, were used to digest the genomic DNA. *MseI/Tru9I* is a frequent cutter with a T↓TAA cutting site, whereas *PstI* and *EcoRI* are 6-base rare cutters with a CTGCA↓G (*PstI* is methylation sensitive) and G↓AATTC (*EcoRI*)

1. Put on gloves (to protect yourself and the reaction mix) and add the following to a 0.5 ml Eppendorf tube:

Restriction-ligation reaction mixture

Genomic DNA(20ng/μl)	~150ng
5x RL buffer	2μl
Rare cutting enzyme EcoRI (10U/μl)	0.10 μl
Frequent cutting enzyme Tru9I (10U/μl)	0.10 μl
EcoRI adaptor mix (50 pmole/μl)	0.5 μl
Tru9I adapter mix (50 pmole/μl)	0.5 μl
rATP (10 mM)	0.2 μl
T4 DNA ligase (5U/μl)	0.13μl
Sterile distilled water	Up to 10μl

2. Mix by tapping the bottom of the tube.
3. Centrifuge briefly (14,000 rpm for 5 seconds).
4. Incubate the resulting reaction mixture for a minimum of 3 hours at 37°C.
5. Inactivate the restriction endonuclease by incubating the mixture at 70°C for 15 min.
6. Place tubes on ice and do brief centrifugation to collect contents.

### 8.1.2. Pre-amplification

Pre-amplification is performed with primers having one selective nucleotide. The aim of pre-amplification is to generate enough template DNA for selective amplification step.

1. Set up the PCR reaction (on ice):

10 x PCR buffer	5 µl
Restriction-ligation reaction (from 7.1.1)	5 µl
EcoRI primer (10µM/µl)	1.5 µl
MseI/Tru91 primer (10µM/µl)	1.5 µl
dNTPs (10 mM)	1 µl
Taq DNA polymerase (5U/µl)	0.5 µl
Sterile distilled water	Up to 50µl

2. Mix by tapping the bottom of the tube.
3. Centrifuge briefly (14,000 rpm for 5 seconds).

**NOTE:** The *EcoRI* and *Tru91* primers used in pre-amplification are non- selective in that they recognise all *EcoRI* and *Tru91* priming sites.

### 8.1.3. PCR pre-amplification

This step amplifies all of the DNA fragments carrying PstI and TruI terminal adaptors, and provides sufficient template for subsequent selective amplification.

Place the tube in the PCR machine and amplify using the following programme:

<i>Step 1</i>	Denaturing	94°C	30 seconds
<i>Step 2</i>	Annealing	65°C (-0.7 °C/cycle)	30 seconds
<i>Step 3</i>	Extension	72°C	1 minute
<i>Step 4</i>	Cycling	repeat steps 1-3 for 11 cycles	
<i>Step 5</i>	Denaturing	94°C	30 seconds
<i>Step 6</i>	Annealing	56°C	30 seconds
<i>Step 7</i>	Extension	72°C	1 minute
<i>Step 8</i>	Cycling	repeat steps 5-7 for 22 cycles	
<i>Step 9</i>	Hold	4°C	forever

### 8.1.4. Check-step

It is important to check that everything has worked in the previous steps before proceeding.

1. Take a 5 µl aliquot of the PCR-amplified product from 7.1.3 above and place in a fresh 0.5 ml tube, and add 2 µl 5x loading buffer.
2. Vortex briefly.

3. Centrifuge briefly (14,000 rpm for 5 seconds).
4. Load the sample into a 1.2 % agarose gel.
5. Run gel at 50V for 30 minutes.
6. Visualise DNA by UV illumination (Figure 8-2).

(Caution: wear gloves, and UV protective glasses and shields over your face when you are exposed to the UV light of the transilluminator)

NOTE: If previous steps have worked you should see a clear DNA band (Figure 8-2).

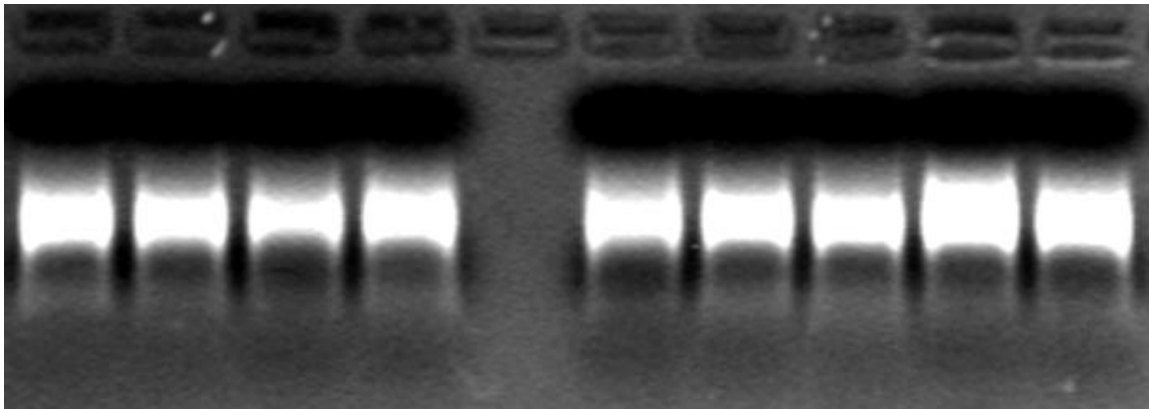


Figure 8-2.

7. Dilution of pre-amplified DNA:
  - For silver staining, dilute 5µl of pre-amplified DNA sample 1:50 with water (50 µl sample + 245 µl water).
  - For fluorescent labelling, dilute pre-amplified DNA to 1:10 with TE (10 µl sample + 90 µl water).
  - Store this dilution and the remaining pre-amplification product at -20°C (long term).

NOTE: The dilution of sample depend of amplified products (S.7.) that is used in selective amplification (8.1.3) PCRs, and now termed 'Test DNA'.

### 8.1.5. Selective pre-amplification

In this section, specific subsets in the test DNA are amplified using EcoRI and Tru91 primers that are extended with one to three selective nucleotides. Silver staining of the amplified fragments that have been electrophoresed on PAGE is commonly used for detection of DNA banding patterns. Alternatively, fluorescence-labelled primers can be used in the selective amplification PCR step and the products visualised on an automated DNA analyser. These two options are described below.

### 8.1.6. PCR mix for selective amplification, products to be visualized on PAGE

1. Put on gloves and in a PCR tube add:

Test DNA (diluted pre-amplified DNA from 8.1.4:Step 7)	5.0 µl
10 x PCR buffer	2.5 µl
<i>Eco</i> RI selective primer (10 µmol)	0.25 µl
<i>Tru</i> 91 selective primer (10 µmol)	0.75 µl
dNTPs (10 mM)	0.5 µl
<i>Taq</i> DNA polymerase (5U/µl)	0.2 µl
Sterile distilled water	Up to 25.0µl

2. Mix by gently tapping against the tube.
3. Centrifuge briefly (14,000 rpm for 5 seconds).

### 8.1.7. PCR profile for Selective amplification, products to be visualised on PAGE

Place tube in the PCR machine and amplify using the following programme:

<i>Step 1</i>	Denaturing	94°C	30 seconds
<i>Step 2</i>	Annealing	65°C (-0.7 °C/cycle)	30 seconds
<i>Step 3</i>	Extension	72°C	1 minute
<i>Step 4</i>	Cycling	repeat steps 1-3 for 13 cycles	
<i>Step 5</i>	Denaturing	94°C	30 seconds
<i>Step 6</i>	Annealing	56°C	30 seconds
<i>Step 7</i>	Extension	72°C	1 minute
<i>Step 8</i>	Cycling	repeat steps 5-7 for 23 cycles	
<i>Step 9</i>	Hold	4°C	forever

### 8.1.8. Polyacrylamide Gel Electrophoresis (PAGE)

The single-stranded AFLPs are separated in long, denaturing polyacrylamide gels (often referred to as sequencing gels).

1. Take a 5 µl aliquot of the PCR-amplified product from 10.1.3 above and place in a fresh 0.5 ml tube, and add 2 µl formamide loading buffer. The number of samples will be determined by the number of wells you have in your polyacrylamide gel.
2. Denature for 5 minutes at 95°C - 100°C, and snap-cool on ice.
3. Centrifuge briefly (14,000 rpm for 5 seconds).

- Run 5µl samples in denaturing 6% polyacrylamide gels. SequaGel®XR (<http://www.nationaldiagnostics.com/electroproducts/ec842.html>)

### 8.1.9. Silver staining of PAG

Follow the procedure given in the SSR Protocol (6.1.6. Silver-staining).

### 8.1.10. PCR mix for selective amplification, products to be visualized on an automated DNA analyser

- Put on gloves and in a PCR tube add:

Test DNA (diluted DNA from .7.1.2.2:S7)	5.0 µl
10 x PCR buffer (with Mg <sup>2+</sup> )	2.0 µl
Fluorescent EcoRI Primer (1µmol)	1.0µl
<i>Tru91</i> selective primer (5µmol)	1.0µl
dNTPs (10 mM)	0.40µl
<i>Taq</i> DNA polymerase (5U/µl)	0.20 µl
Sterile distilled water up to	<b>20.0µl</b>

- Mix by gently tapping against the tube.
- Centrifuge briefly (14,000 rpm for 5 seconds).

### 8.1.11. PCR profile for selective amplification, products to be visualized on an automated DNA analyser

Place tube in the PCR machine and amplify using the following programme:

<i>Step 1</i>	Denaturing	94°C	30 seconds
<i>Step 2</i>	Annealing	65°C (-0.7 °C/cycle)	30 seconds
<i>Step 3</i>	Extension	72°C	1 minute
<i>Step 4</i>	Cycling	repeat steps 1-3 for 11 cycles	
<i>Step 5</i>	Denaturing	94°C	30 seconds
<i>Step 6</i>	Annealing	56°C	30 seconds
<i>Step 7</i>	Extension	72°C	1 minute
<i>Step 8</i>	Cycling	repeat steps 5-7 for 29 cycles	
<i>Step 9</i>	Hold	4°C	forever

### 8.1.12. Electrophoresis using an automated DNA analyser

The single-stranded AFLPs are separated through electrophoresis on a capillary type automated DNA analyser (ABI Prism 3100 is used in the Plant Breeding and Genetics Laboratory).

1. Put on gloves and in a “sequencer” plate, add for each sample:

PCR-amplified product from 7.1.3.1	1.0µl
Formamide	13.0 µl
ROX standard	0.25 µl

2. Denature for 5 minutes at 95°C - 100°C, and snap-cool on ice.
3. Centrifuge briefly (14,000 rpm for 5 seconds) and check for air bubbles.
4. Load plate on the DNA analyser according to User’s manual and select the option for AFLP fragment separation.

### 8.1.13. Production of single primer, linear PCR products

NOTE: This procedure is used to avoid doubled stranded DNA fragments and results in a greater clarity of band separation.

1. Put on gloves and add in a PCR tube:

10X PCR buffer	2 µl
Selective amplification DNA (produced in Step 6)	2 µl
<i>Pst</i> I selective primer (50 ng/µl)	1.5 µl
dNTPs (2 mM)	2.5 µl
<i>Taq</i> DNA polymerase (5U/µl)	0.1 µl
Add sterile distilled water to make up to	<b>20 µl</b>

2. Mix gently by tapping the tube.
3. Centrifuge briefly (14,000 rpm for 5 seconds).

### 8.1.14. PCR amplification to produce single stranded DNA

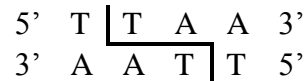
Put on gloves and place tube from 10.2..3. into a PCR machine and amplify using the following programme:

<i>Step 1</i>	Denaturing	94°C	30 seconds
<i>Step 2</i>	Annealing	56°C	30 seconds
<i>Step 3</i>	Extension	72°C	1 minute
<i>Step 4</i>	Cycling	repeat steps 1-3 for 22 cycles	
<i>Step 5</i>	Denaturing	94°C	30 seconds
<i>Step 6</i>	Hold	4°C	hold

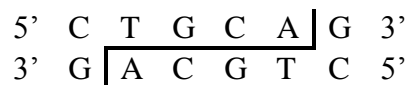
## 8.2. Required enzymes and primer sequences for AFLP assays

### 8.2.1. Restriction enzymes

MseI/Tr91



PstI



## 8.3. Preparation of adapters

**Tru9I adapter-oligos** have 16 and 14 nucleotides



Take 15µl of each to get the final concentration of 50pmol/µl in 30µl water.

**Pst1 adapter-oligos** have 21 and 14 nucleotides



Take 15µl of each to get the final concentration of 50pmol/µl in 30µl water.

## 8.4. Reagents needed

- Use only sterile distilled water for all solutions
- 5x RL buffer
  - 50 mM TrisAc pH7.5
  - 50 mM MgAc
  - 250 mM KAc
  - 25 mM DTT
  - 250 ng/µl BSA
- Rare cutting enzyme, *PstI* (5U/µl)
- Frequent cutting enzyme, *Tru9I* (5U/µl)
- *PstI* adaptor (5 pmole/µl) or *EcoRI* adaptor (5 pmole/µl)
- *Tru9I* adaptor (50 pmole/µl)
- rATP (10 mM)
- T4 DNA ligase
- 10 x PCR buffer
- *PstI* or *EcoRI* non-selective primer (50 ng/µl)

- *Tru9I* non-selective primer (50 ng/μl)
- *Taq* DNA polymerase (5U/μl)
- Agarose
- T<sub>0.1</sub>E buffer
- *Pst*I or *Eco*RI selective primer
- ***Tru9I*** selective primer
- dNTPs (10 mM)
- Formamide
- ROX Standard

### 8.5. Sequence information of adapters and primers used for AFLP

<i>Tru9I</i> -Adapter sequence 1	5'-GACGATGAGTCCTGAG-3'
<i>Tru9I</i> -Adapter sequence 1.	3'-TACTCAGGACTCAT-5'
<i>Eco</i> RI:	5'- CTCGTAGACTGCGTACC -3'
<i>Eco</i> RI	5'- AATTGGTACGCAGTCTAC -3'
Primers for pre-amplification	
<i>Tru9I</i> -primer	5'-GACGATGAGTCCTGAGTAA-3'
Eco-P0:	5'- GACTGCGTACCAATTC -3'
<i>Tru9I</i> -P0:	5'- GATGAGTCCTGAGTAA -3'
<i>Tru9I</i> -PC:	5'- GATGAGTCCTGAGTAAC -3'
<i>Tru9I</i> Selective primers**	
<i>Tru9I</i> -CAC	5-GATGAGTCCTGAGTAACAC-3'
<i>Tru9I</i> -ACC	5'-GATGAGTCCTGAGTAAACC-3'
<i>Tru9I</i> -CCA	5'-GATGAGTCCTGAGTAACCA-3'
<i>Tru9I</i> -CAA	5'-GATGAGTCCTGAGTAACAA-3'
<i>Tru9I</i> -ACG	5-GATGAGTCCTGAGTAAACG-3'
<i>Tru9I</i> -CAG	5'-GATGAGTCCTGAGTAACAG-3'
<i>Tru9I</i> -CAT	5'-GATGAGTCCTGAGTAACAT-3'
<i>Tru9I</i> -CGA	5'-GATGAGTCCTGAGTAAACGA-3'
<i>Tru9I</i> -CGT	5'-GATGAGTCCTGAGTAAACGT-3'
<i>Tru9I</i> -CCT	5'-GATGAGTCCTGAGTAACCT-3'
<i>Tru9I</i> -CTA-	5'- GATGAGTCCTGAGTAACTA 3'
<i>Tru9I</i> -CTC	5'- GATGAGTCCTGAGTAACTC -3'
<i>Tru9I</i> -CTG:	5'- GATGAGTCCTGAGTAACTG -3'
<i>Tru9I</i> -CTT:	5'- GATGAGTCCTGAGTAACTT -3'
<i>Tru9I</i> -GAA	5'- GATGAGTCCTGAGTAAGAA -3'
<i>Tru9I</i> -GAC:	5'- GATGAGTCCTGAGTAAGAC -3'
<i>Tru9I</i> -GAG	5'- GATGAGTCCTGAGTAAGAG -3'
<i>Tru9I</i> -GAT	5'- GATGAGTCCTGAGTAAGAT -3'
<i>Tru9I</i> -GTA:	5'- GATGAGTCCTGAGTAAGTA -3'
<i>Tru9I</i> -GTC:	5'- GATGAGTCCTGAGTAAGTC -3'
<i>Tru9I</i> -GTG	5'- GATGAGTCCTGAGTAAGTG -3'
<i>Tru9I</i> -GTT:	5'- GATGAGTCCTGAGTAAGTT -3'



<i>EcoRI</i> Selective primers**	
<i>EcoRI</i> AA	5'- GACTGCGTACCAATTCAA -3'
<i>EcoRI</i> AT	5'- GACTGCGTACCAATTCAT -3'
<i>EcoRI</i> TA	5'- GACTGCGTACCAATTCTA -3'
<i>EcoRI</i> TT	5'- GACTGCGTACCAATTCTT -3'
<i>EcoRI</i> AC	5'- GACTGCGTACCAATTCAC -3'
<i>EcoRI</i> AG	5'- GACTGCGTACCAATTCAG -3'
<i>EcoRI</i> TG:	5'- GACTGCGTACCAATTCTG -3'
<i>EcoRI</i> TC	5'- GACTGCGTACCAATTCTC -3'
<i>EcoRI</i> CTG	5'- GACTGCGTACCAATTCCTG -3'
<i>EcRI</i> GAC	5'- GACTGCGTACCAATTCGAC -3'
<i>EcoRI</i> GAA	5'- GACTGCGTACCAATTCGAA -3'
<i>EcoRI</i> CTA	5'- GACTGCGTACCAATTCCTA -3'
<i>EcoRI</i> AAC	5'- GACTGCGTACCAATTCAAC -3'
<i>EcoRI</i> AAG	5'- GACTGCGTACCAATTCAAG -3'
<i>EcRI</i> ACA	5'- GACTGCGTACCAATTCACA -3'
<i>EcoRI</i> ACC	5'- GACTGCGTACCAATTCACC -3'
<i>EcoRI</i> ACG	5'- GACTGCGTACCAATTCACG -3'
<i>EcoRI</i> ACT	5'- GACTGCGTACCAATTCACT -3'
<i>EcRI</i> AGC	5'- GACTGCGTACCAATTCAGC -3'
<i>EcoRI</i> AGG	5'- GACTGCGTACCAATTCAGG -3'
<i>EcoRI</i> GAT	5'- GACTGCGTACCAATTCGAT -3'
<i>EcRI</i> GAG	5'- GACTGCGTACCAATTCGAG -3'
<i>EcRI</i> CTT	5'- GACTGCGTACCAATTCCTT -3'
<i>EcoRI</i> CTC	5'- GACTGCGTACCAATTCCTC -3'

\*\*The same PCR primers are used for both the silver stained PAGE and automated DNA analyser options except that for the latter, primers labelled with either HEX or FAM fluorescent dye are used.

## 8.6. References

Vos, P., R. Hogers, M. Bleeker, M. Reijans, T. van de Lee, M. Hornes, A. Frijters, J. Pot, J. Peleman, M. Kuiper, and M. Zabeau, 1995. AFLP: a new technique for DNA fingerprinting. *Nucleic Acids Res.* 23(21): 4407-4414.

## 9. REMAP & IRAP

**REMAP definition:** Any difference in DNA sequence between two genomes, detected by polymerase chain reaction-mediated amplification of the region between a long terminal repeat of a retrotransposon and a nearby microsatellite (Kahl, 2001).

The dispersion, ubiquity and prevalence of retrotransposon-like elements in plant genomes can be exploited for DNA-fingerprinting. Two DNA techniques based on retrotransposon-like elements are introduced here: IRAP and REMAP (Kalendar *et al.*, 1999). The **IRAP** (**I**nter-**R**etrotransposon **A**mplified **P**olymorphism) markers are generated by the proximity of two retrotransposons using outward facing primers annealing to their long terminal repeats (LTRs). In REMAP (**R**etrotransposon-**M**icrosatellite **A**mplified **P**olymorphism) the DNA sequences between the LTRs and adjacent microsatellites (SSRs) are amplified using appropriate primers.

The principle of IRAP und REMAP is shown in Figure 9-1 below:

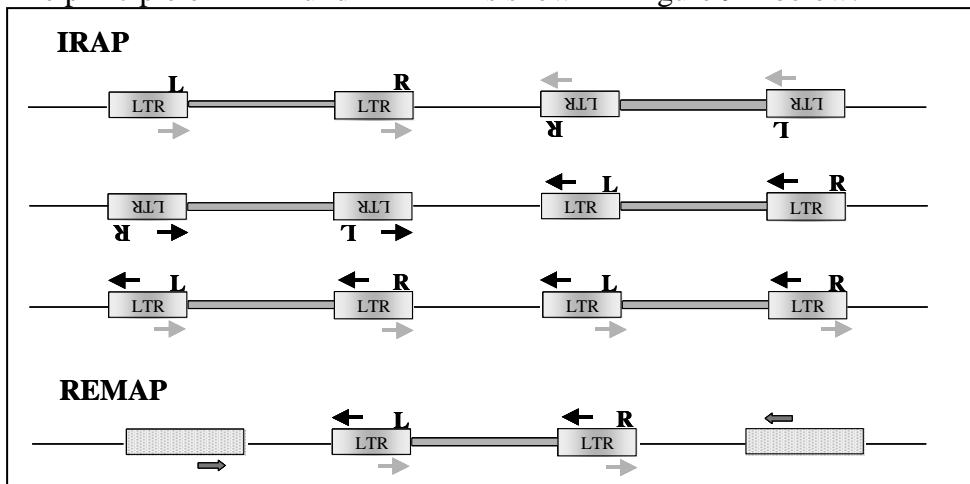


Figure 9-1. Principle of the IRAP und REMAP strategy. **IRAP:** PCR primers facing outward from the 5' (black arrows) and 3' (grey arrows) ends of LTRs will amplify intervening DNA from the retrotransposon in any of the three possible orientations (tail-to-tail, head-to-head, head-to-tail). **REMAP:** LTR primers are used together with a primer consisting of simple sequence repeats (blank boxes) (Kalendar *et al.*, 1999)

### 9.1. Protocol

REMAP and IRAP markers are species specific. In the FAO/IAEA course the following primers for rice and barley were available and used in conjunction with rice and barley DNA.

**Table 9.1.** LTR primers from the rice retrotransposon Tos17 (Hirochika *et al.*, 1996), sequence and PCR annealing temperatures ( $T_a$ ).

Primer	Sequence	$T_a$
TOS17LTR-1 (outward 3' end of LTR)	TTGGATCTTGTATCTTGTATATAC	56°C
TOS17LTR-2 (outward 3' end of LTR)	GCTAATACTATTGTTAGGTTGCAA	56°C
TOS17LTR-3 (outward 5' end of LTR)	CCAATGGACTGGACATCCGATGGG	56°C
TOS17LTR-4 (outward 5' end of LTR)	CTGGACATGGGCCAACTATACAGT	56°C

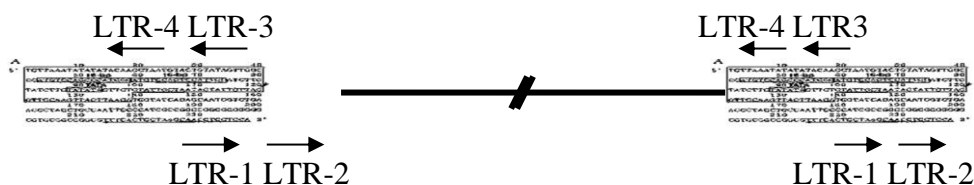
**Table 9.2.** LTR primers from the barley BARE-1 (Kalendar *et al.*, 1999), sequence and PCR annealing temperatures ( $T_a$ ).

Primer	Sequence	$T_a$
BARLTR-2(LTR forward) - IRAP	CTCGCTCGCCCACTACATCAACCGCGTTT ATT	60°C
BARLTR-3(LTR reverse) - IRAP/REMAP	GGAATTCATAGCATGGATAATAAACGAT TATC	60°C

**Table 9.3.** Microsatellite (SSR) primers and PCR annealing temperatures ( $T_a$ ).

Sequence	$T_a$
(GA) <sub>9</sub> C; (CT) <sub>9</sub> G; (CA) <sub>10</sub> G	54°C
(CAC) <sub>7</sub> G; (GTG) <sub>7</sub> C; (CAC) <sub>7</sub> T; GT(CAC) <sub>7</sub>	58°C

NOTE: It is very important to try different combinations of LTR- and microsatellite (SSR) primers for REMAP and LTR-primers for IRAP. Choose primers that have been derived from the species you are working with. The figure below shows you the orientation of only the TOS17-LTR-primers:



NOTE: Gloves and lab coat should be worn throughout.

### 9.1.1. Prepare a 50µl reaction mix

1. Take a sterile PCR tube and add:

10 x <i>Taq</i> buffer	5.0 µl
dNTPs (10 mM)	1.0 µl
Primer 1 (100 pmol/µl)	0.5 µl
Primer 2 (100 pmol/µl)	0.5 µl
DNA (100 ng/µl)	1.0 µl
<i>Taq</i> DNA polymerase (5 U/µl)	0.5 µl
Add ddH <sub>2</sub> O to bring volume to	<b>50 µl</b>

2. Mix by tapping against the tube.
3. Centrifuge briefly (14,000 rpm for 5 seconds).

### 9.1.2. PCR amplification

The PCR amplification programme used for the Tos17 sequence was:

<i>Step 1</i>	Initial denaturation	94°C	2 minutes
<i>Step 2</i>	Denaturation	94°C	30 seconds
<i>Step 3</i>	Primer annealing*	T <sub>a</sub>	30 seconds
<i>Step 4</i>	Ramp	0.5°C per second to 72°C	
<i>Step 5</i>	Primer extension	72°C	2 minutes
<i>Step 6</i>	Cycling	repeat steps 2-5 for 29 cycles	
<i>Step 7</i>	Final extension	72°C	8 minutes
<i>Step 8</i>	Hold	4°C	forever

\* See tables above for appropriate annealing temperatures (T<sub>a</sub>).

### 9.1.3. Separation and visualization of the amplification products

1. Place 15 µl of PCR into a fresh Eppendorf tube.
2. Add 3 µl of 5 X loading buffer containing dye.
3. Vortex briefly.
4. Centrifuge briefly (14,000 rpm for 5 seconds).
5. Load sample into a 2% NuSieve® agarose gel.

NOTE: NuSieve® agarose provides a good separation gel.

6. Run gel for approximately 80 minutes at 80 W (power limiting) or until dark blue front has run 2/3 down the gel.

NOTE: See Section 1 of RFLP Protocol (Agarose gel electrophoresis) for details of gel preparation and running.

7. Stain gel with ethidium bromide (*Caution: ethidium bromide is toxic wear gloves and avoid inhalation*).
8. Visualise bands under UV light (*Caution: wear UV protective glasses and shield your face when you are exposed to the UV light of the transilluminator*).

## 9.2. References

- Hirochika, H., K. Sugimoto, Y. Otsuki, H. Tsugawa, and M. Kanda, 1996. Retrotransposons of rice involved in mutations induced by tissue culture. *Proc.Natl.Acad.Sci.USA*. **93**: 7783-7788
- Kalendar, R., T. Grob, A. Regina, A. Suoniemi, and A. Schulman, 1999. IRAP and REMAP: two new retrotransposon-based DNA fingerprinting techniques. *Theor.Appl.Genet*. **98**: 704-711.

## 9.3. Reagents needed

Use only sterile distilled water for all solutions.

- *Taq* buffer
- dNTPs
- Primers
- *Taq* DNA polymerase (5U/μl)
- DNA (10-20 ng/μl)
- 10 x loading buffer:
 

Glycerol (80%)	600 μl
Xylene cyanol	2.5 mg
Bromophenol blue	2.5 mg
Water	400 μl
- 5 x loading buffer
 

Glycerol (80%)	300 μl
Xylene cyanol	1.3 mg
Bromophenol blue	1.3 mg
Water	400 μl
- Ethidium bromide
- Agarose
- Acrylamide
- Bis-acrylamide
- TBE

## 10. SINGLE NUCLEOTIDE POLYMORPHISMS (SNPS)

**SNP definition:** Any polymorphism between two genomes that is based on a single nucleotide exchange, small deletion or insertion. (Kahl, 2001).

Small nucleotide polymorphism (SNP) is a relatively new marker technology originally developed in human. SNPs are the most abundant polymorphic marker with 2 – 3 polymorphic sites every kilobase (Cooper et al., 1985). Originally discovered in humans, SNPs have now been developed for genotyping in plants. SNP technology is heavily dependent upon sequence data. Several methods are available for SNP detection including automated fluorescent sequencing denaturing high-performance liquid chromatography (DHPLC, Underhill *et al.*, 1996), DNA microarrays (Hacia and Collins, 1999), single-strand conformational polymorphism-capillary electrophoresis (SSCP-CE, Ren, 2001; Figure 1), microplate-array diagonal-gel electrophoresis (MADGE, Day *et al.*, 1998) and matrix-assisted laser desorption/ionisation time-of-flight (MALDI-TOF, Griffin and Smith, 2000).

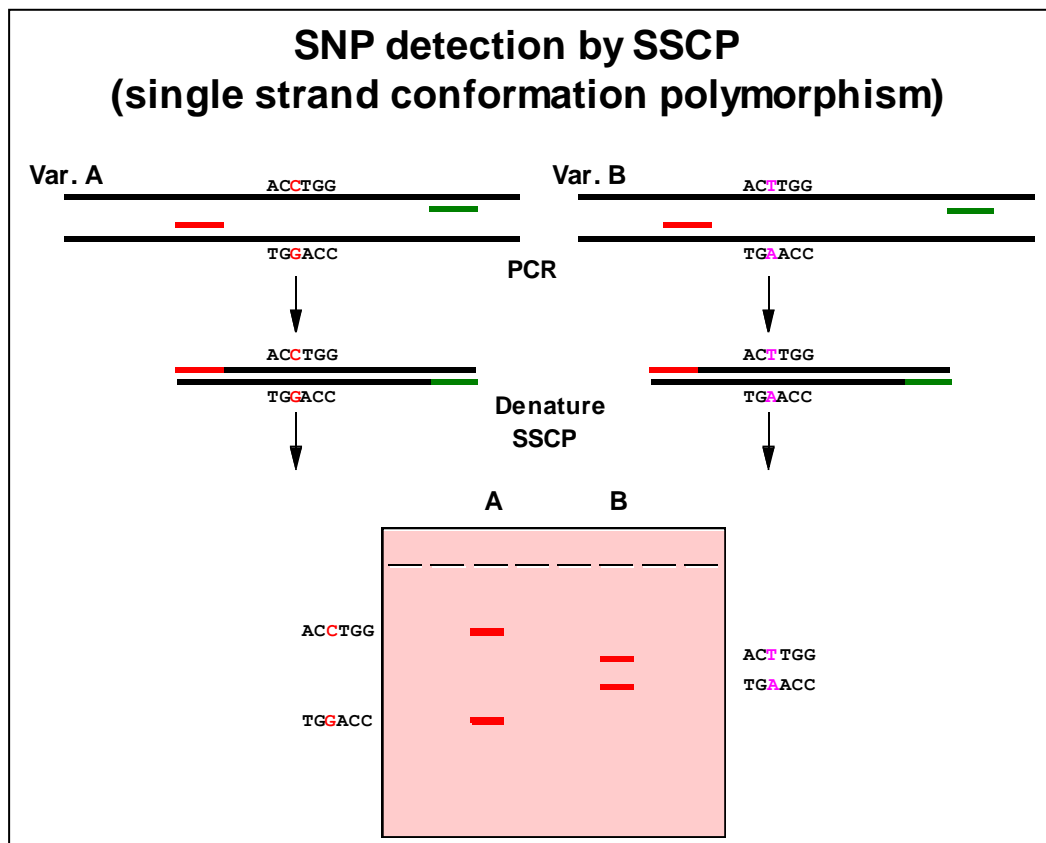


Figure 10-1. The scheme above shows how SNP variation can be detected between varieties A and B (with permission K. Devos).

## 10.1. References

- Cooper, D. N., B. A. Smith, H. J. Cooke, S. Niemann, and J. Schmidtke, 1985. An estimate of unique DNA sequence heterozygosity in the human genome. *Hum.Genet.* **69**(3): 201-205
- Day, I. N., E. Spanakis, D Palamand, G. P. Weavind, and S. D. O'Dell, 1998. Microplate-arrays diagonal-gel electrophoresis (DADGE) and melt-MADGE: tool for molecular genetic epidemiology. *Trends in Biotech.* **16**: 287-290
- Griffin, T. J. and L. M. Smith, 2000. Single-nucleotide polymorphism analysis by MALDI-TOF mass spectrometry. *Trends in Biotech.* **18**: 77-84
- Hacia, J. G. and F. S. Collins, 1999. Mutational analysis using oligonucleotide microarrays. *J.Med.Genet.* **36**: 730-736
- Kahl, G., 2001. *The Dictionary of Gene Technology.* Wiley-VCH, Weinheim.
- Ren, J., 2001. High-throughput single-strand conformation polymorphism analysis by capillary electrophoresis. *J.Chromatography B.Biomed.Science Appl.* **741**: 115-128
- Underhill, P. A., L. Jin, R Zemans, P. J. Oefner, and L. L. Cavalli-Sforza, 1996. A pre-Columbian Y chromosome-specific transition and its implications for human evolutionary history. *Proc.Natl.Acad.Sci.USA.* **93**: 196-200.

## 11. TILLING

TILLING (Targeting Induced Local Lesions IN Genomes) is a general strategy for the discovery of induced point mutations (MCCALLUM *et al.* 2000; COLBERT *et al.* 2001). The procedure consists of: setting up and running PCR using gene specific primers, denaturing and annealing PCR products to create heteroduplexes between mutant and wild-type strands, digesting heteroduplexes with a single-strand specific nuclease, purifying the products and reducing sample volume, loading sample onto a membrane comb, running the samples on a gel and processing and examining the gel images to identify mutations. The same methods can be used to identify naturally occurring polymorphisms in populations, called Ecotilling, (COMAI *et al.* 2004).

For this training course, we will be using primers for the Arabidopsis OX11 gene and eight genomic DNA samples, each containing a unique single nucleotide point mutation. The protocol has been scaled down from the standard high throughput TILLING protocol for the discovery of mutations in a large number of pooled samples (TILL *et al.* 2003; TILL *et al.* 2006b). Primers and genomic DNA samples are described in a publication on the use of single-strand specific nucleases for mismatch cleavage (TILL *et al.* 2004a). The standard high-throughput TILLING protocol will be followed using fluorescently labelled primers and a LI-COR DNA analyser. Additionally, students will analyse mutations using lower cost and lower throughput agarose gels (for examples see (SATO *et al.* 2006; GARVIN and GHARRETT 2007; GALEANO *et al.* 2009)). The goal of this section of the training course is to familiarize you with the bench and computational techniques that have been developed for TILLING. The hope is that students will leave with a firm understanding of TILLING and the ability to critically evaluate the usefulness of TILLING in his or her research program.

### 11.1. Protocol

Each group will receive a box containing samples, buffers and solutions for this section of the course. All materials are provided in the box except Ex-Taq polymerase. This will be distributed by the instructor.

#### 11.1.1. PCR reaction with IRDye-labeled primers

Make the following PCR master mix on ice:

Water	72 $\mu$ l
10x PCR buffer	11.4 $\mu$ l
25 mM MgCl <sub>2</sub>	13.6 $\mu$ l
2.5 mM each dNTP	18.4 $\mu$ l
primer cocktail *	8.0 $\mu$ l
Ex-Taq hot start version	1.2 $\mu$ l



Add 10 µl of PCR mix to each DNA sample (10 µl). Mix sample by pipetting up and down three times.

Place your set of 8 samples in the thermal cycler. Once all teams have deposited their samples, run the PCR cycling program (titled PCRTM70.cyc):

<i>Step 1</i>	Initial denaturation	95°C	2 minutes
<i>Step 2</i>	Denaturation	94°C	20 seconds
<i>Step 3</i>	Primer annealing	73°C (-1°C/cycle)	30 seconds
<i>Step 4</i>	Ramp	0.5°C per second to 72°C	
<i>Step 5</i>	Primer extension	72°C	1 minute
<i>Step 6</i>	Cycling	repeat steps 2-5 for 7 cycles	
<i>Step 7</i>	Denaturation	94°C	20 seconds
<i>Step 8</i>	Primer annealing	65°C	30 seconds
<i>Step 9</i>	Ramp	0.5°C per second to 72°C	
<i>Step 10</i>	Primer extension	72°C	1 minute
<i>Step 11</i>	Cycling	repeat steps 7-10 for 44 cycles	
<i>Step 12</i>	Final extension	72°C	5 minutes
<i>Step 13</i>	Denaturation	99°C	10 minutes
<i>Step 14</i>	Cooling	72°C	20 seconds
<i>Step 15</i>	Cycling	repeat step 14 for 70 cycles (-0.3°C/ cycle)	
<i>Step 16</i>	Hold	4°C	forever

NOTES: For purposes of training, we increase the volume of the master mix so that you have more than is needed. Normally this is not done, but the excess volume controls for pipetting errors and if one group makes a mistake, excess from the other groups can be provided to them.

\* The primer cocktail was made in advance as follows:

3 µl forward primer labeled with IRD700 dye (100µM)  
 2 µl unlabeled forward primer (100µM)  
 4 µl reverse primer labeled with IRD800 dye (100µM)  
 1 µl unlabeled reverse primer (100µM)

This mix was stored at -80°C. Prior to use, the mix is thawed on ice, diluted 1:10 with TE (10 mM Tris-HCl, 1 mM ethylene diamine tetraacetic acid (EDTA), pH 7.4) and distributed to each team.

Remove 4µl of samples #7 and #8 and put into new tubes for analysis of PCR product on agarose gel (Step 11.1.3).

### 11.1.2. Heteroduplex digestion, preparation of Sephadex spin plates

#### Heteroduplex digestion

Add 4 µl of water to samples #7 and #8 to bring the volume back to 10 µl. Because DNA has been removed for the agarose gel test, these samples should appear weaker on the LI-COR gel.

Prepare the following mix on ice:

Water	326µl
10X CEL I TILLING buffer	60µl
CJE nuclease <sup>#</sup>	14µl

#### NOTES:

\*10X CEL I buffer is:

5 ml 1M MgSO<sub>4</sub>

100 µl 10% Triton X-100

5 ml 1M Hepes pH 7.5

5 µl 20 mg/ml bovine serum albumen

2.5 ml 2M KCl

37.5 ml water

# The amount of enzyme required will vary depending on nuclease source or possibly from batch to batch of the same enzyme from the same source.

Mix components on ice. Add 40µl of mix to the PCR product and mix by pipetting 2-3 times. Incubate at 45°C for 15 min (in thermal cycler). Cool to 8°C and stop reaction by adding 10 µl of 0.25M EDTA to each sample.

Label a new 8-strip of PCR tubes a set 2 and transfer 35 µl of samples to these tubes. Divide samples by transferring into a new set of 8-tube strip. Set one will be used in Step 11.1.3 onwards.

#### Preparation of Sephadex spin plates

Prior to loading nuclease digested samples onto the denaturing polyacrylamide gel, salts must be separated from the DNA and sample volume reduced to 1.5 µl. There are several methods that can be used to accomplish this. The one you might be most familiar with is alcohol precipitation. For TILLING, we use a different method: size exclusion chromatography using Sephadex G50 medium beads. This is much faster than alcohol precipitation and provides consistent and high recovery of DNA. 96-well plates containing hydrated Sephadex can be prepared up to one week in advance.

Each team will practice preparing a Sephadex plate during the 90°C incubation in Step 12.1.3.1. Pour dry G50 (medium) powder into a 96-hole metal plate and distribute evenly using plastic scraper. Fit a 96-well membrane plate on top, then invert and tap to fill wells with powder. Use a multichannel pipette to add 300 µl water to the top of each well to hydrate, then cover and let sit at least 1 hr at room temperature. Plates are usually made in advance and stored at 4°C in a moist environment for up to one week.

### 11.1.3. Agarose gel analysis of enzymatic mismatch cleavage, and sample purification

#### Agarose gel analysis

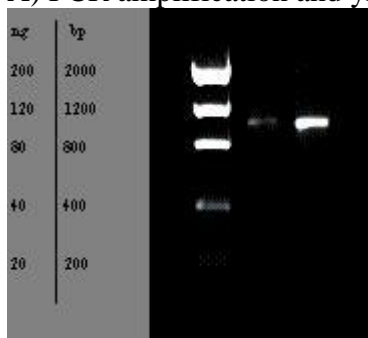
DNA samples are electrophoresed through an agarose gel to verify that (a) PCR was successful in Step 11.1.1 and (b) digestion of mutant DNA by CELI has occurred in Step 11.1.2.

Load samples in the following order:

Lane	1	2	3	4	5	6	7	8	9	10	11
Sample	Low DNA mass ladder	#7 from section 3.1	#8 from section 3.1	#1 from strip 2, section 3.2	#2	#3	#4	#5	#6	#7	#8
Volume (µl)	4	4	4	10	10	10	10	10	10	10	10

#### Data analysis

##### A) PCR amplification and yield



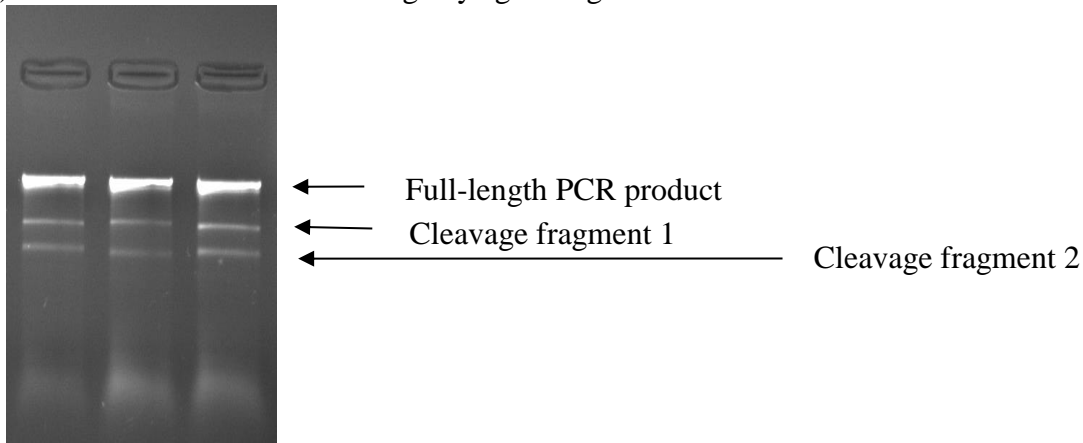
The figure above shows example data of what your first three gel lanes should look like. You should see a single band of the correct size (992 bp). The yield should be at least 7-10 ng/µl of PCR product. The Invitrogen low DNA mass ladder is quantitative and yields are determined by estimating the intensity of amplified PCR products. For example the intensity of the band in the first PCR sample is between 40 and 20 ng, so the concentration is 30 ng/4 µl or 7.5

ng/μl. The second sample is around 25 ng/μl. Both samples indicate that PCR yield is sufficiently robust for TILLING.

NOTES:

Primer yields are typically not assayed before CEL I digestion of samples. This is done here to evaluate your work. The PBGL typically performs PCR amplification tests on all gene-specific primers prior to purchasing expensive fluorescently labelled primers. Primers passing standardized quality control tests almost always perform well in TILLING experiments.

B) Evaluation of mutation cleavage by agarose gel



DNA used for PCR amplification of samples 1-8 each contains a single point mutation. Cleavage of the mutation creates two fragments of lower molecular weight that migrate faster than the full-length PCR product on the agarose gel. The size of these two fragments equals the size of the full-length PCR product. The eight samples have mutations at different positions on the PCR fragment and so will produce different sized fragments. Take some time to determine where you think mutations are based on the size of your bands.

**11.1.4. Sample purification and volume reduction**

All of the workshop samples will be loaded onto a single Sephadex plate. Visually check the Sephadex plate for moistness, and also check underneath for loose Sephadex. If there is any, lightly wipe the bottom with a wet paper towel and gently rinse the bottom holding the plate on its side. Assemble Sephadex plate, blue plate adaptor, and 96-well skirted 0.2 ml plate (this plate is the “waste” plate).

Spin 2 min at 440g.

Replace the waste plate with a sample catch plate containing 1.5 µl formamide load dye\* and 2 µl 200bp marker† in row D. Transfer the entire CEL I reaction sample to each spin plate well. Use a 20-200 µl 8-channel multi-pipettor. Caution: Be sure to dispense liquid to the middle of each well in the Sephadex spin plate, and do not touch the surface of the Sephadex.

Spin 2 min at 440g.

**NOTES:**

\* Formamide load dye is:

250 ml deionized formamide  
5 ml 0.5 M EDTA pH 8  
60 mg bromophenol blue

† 200 bp marker is made by PCR using gene specific IRD labeled primers that amplify a 200 bp target region. Perform PCR and Sephadex purification as outlined in this protocol. Dilute product to 0.5ng/µl in TE.

The instructor will re-array samples so that all eight samples from a group are adjacent on the LI-COR gel.

Incubate samples at 90°C for approx. 45 min until volume reduced to 1.5 µl.

### 11.1.5. Preparing, loading, and running LI-COR gels

All student samples will be run on a single gel. The instructor will demonstrate gel preparation.

Clean and assemble glass plates. Prepare the following mixture:

20 ml acrylamide gel mix (6.5%)  
15 µl TEMED  
150 µl fresh 10% ammonium persulfate

Fill a 20 ml syringe with acrylamide solution. Dispense along the top, avoiding bubbles by rapping just above the liquid edge whenever it appears one might get trapped. If any bubbles appear, remove them quickly after the gel is poured with a thin wire tool. Leaving a little excess at the well, insert the top spacer all the way and centered. Insert the Plexiglas pressure plate between the glass plate and casting rails. Tighten the top screws as soon the spacer is inserted, compressing the rubber pads on the pressure plate a little. Add acrylamide to the top glass edge where the comb is inserted and on the edges to assure that polymerization is not inhibited within the gel. Let the gel set at least 30 min before putting it into the gel box. Gels can be poured in advance and stored wrapped in a damp paper towel at 4°C for several days.

### Loading samples onto membrane combs

All samples will be loaded onto a single loading tray. Each team will load 0.25 µl of sample into the membrane comb loading tray. The instructor will dip the comb into the tray to absorb the sample. The sample should run 1/2 to 2/3 up the length of the comb.

#### NOTES:

Membrane combs are expensive. To reduce the costs, combs can be reused many times. After the comb has been used, rinse thoroughly with deionized water, soak in water for at least 30 minutes, and allow to dry completely before reuse.

### Running LI-COR gels

Pre-run gel 20 min. Gel settings: 1500 V, 40 mA, 40 W, Temp = 50°C, Width = 1028, Speed =2, Channels= 700 & 800

Make sure the back plate is clean and clear of any scratches in the data collection window. Check that the machine is properly focused before loading samples.

Clean the gel slot out with a syringe and drain the top buffer reservoir until the level is below the glass edge. Wick out the remaining buffer, first with a paper towel and then with a 6 inch wide strip of Whatman 1 paper, sliding it into the slot left by the spacer. Using a Pipetteman P1000, fill the slot with 1 ml of 1% Ficoll leaving just a thin bead, ~1 mm above the slot. Hold the comb at a 45°C vertical angle with lane 1 on the left, aim for the slot and insert rapidly by pushing gently until it just touches the gel surface along its length. Gently fill the reservoir to the fill line, insert the electrode/cover, close the top and then click on “Collect image”. From the time the comb touches the slot until the time the current is applied should be no more than about 20 min or so to prevent diffusion. After 10 min, open the LI-COR (be sure that you hear the ‘pling’ signal and the high voltage light goes off), remove the comb and gently rinse the slot with buffer. Replace the upper electrode, close the door and resume the run for 3hrs 45min.

#### 11.1.6. Data Analysis

This component of the TILLING exercise is intended to be performed by students on computers with internet access. Programs and training files along with the protocol below can be downloaded here: [http://tilling.fhcrc.org/tillingdemo/computational\\_tools.shtml](http://tilling.fhcrc.org/tillingdemo/computational_tools.shtml).

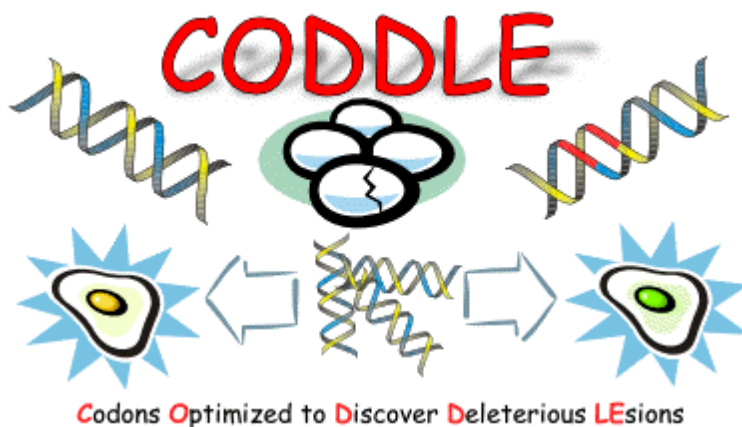
By following the instructions on the webpage, you can easily access all the links described in the protocol below.

## 11.2. Computation tools

### 11.2.1. Selecting the best region to screen and designing primers

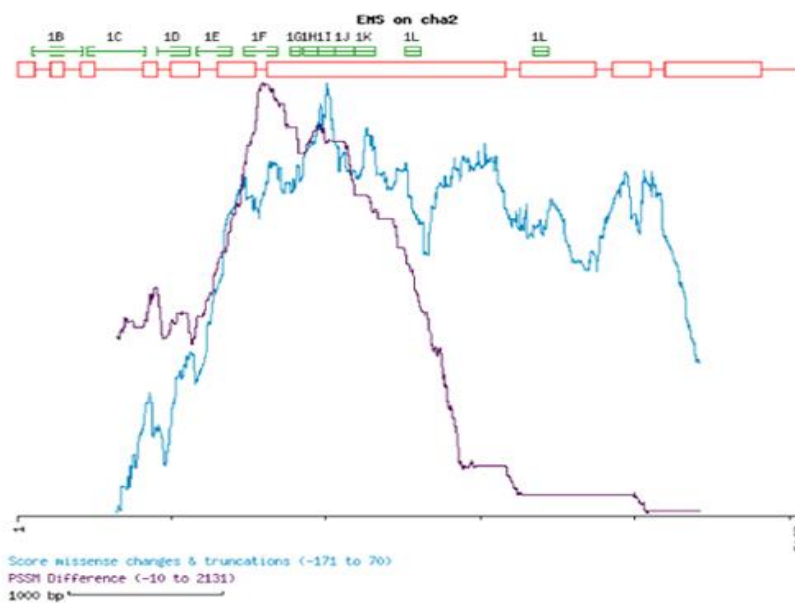
The current PCR target size for TILLING is between 725 and 1600 bp, with the optimum being around 1.5 kbp. The average gene size in Arabidopsis is 3-4 kb and thus a single PCR amplicon will not cover a whole gene. For genes larger than 1.6 kb, one can either screen the entire gene with overlapping primer pairs (TILLING by tiling), or one can choose the region of a gene with the highest number of possible deleterious changes. For projects where there are a large number of targets, or where the cost of screening could become prohibitive, choosing a “best” screening region is a good approach. This is the approach that STP takes for its public services. For this section of the course, students will use computational tools to choose a target region for TILLING, design primers, and place an order with STP. There are three important components necessary for the optimal TILLING order: 1) a good gene model (intron/exon positions), 2) a good protein sequence homology model, and 3) a good PCR primer pair.

These choices are facilitated by the CODDLE Input Utility, (<http://www.proweb.org/input/>) which accepts genomic, cDNA and/or protein sequences from your own files or via links from public databases.



1. Open the **Test Genes page** (<http://tilling.fhcrc.org/tillingdemo/CODDLEtestgenes/>) in a new browser window. Select a gene by clicking on the gene name. Here you will find both genomic and protein sequence information.
2. Select and copy the genomic DNA sequence.
3. Open **CODDLE Input Utility** (<http://www.proweb.org/input/>) in a new browser window.
4. In the CODDLE input page, enter the gene name and paste in the genomic sequence information.
5. Go back to the gene page and copy the protein sequence.
6. Paste the protein sequence in the appropriate window.

7. Click the “Begin Processing” button. The CODDLE input utility is now creating a gene model and searching for homology information that will help identify regions that are likely to be important for protein function.
8. A new window should appear with a summary of the Blocks family protein homology, an intron/exon join statement and the amino acid sequence. Click the “Proceed with CODDLE” button.
9. In the CODDLE page, select “TILLING w/EMS (plants)” as the mutation method, then click “CODDLE your gene”. CODDLE will now evaluate every possible mutation and provide a high scoring window where the highest number of deleterious changes are likely to be found. A new window will open with the CODDLE output. The graphical output shows the gene model (red boxes and lines), protein homology (green boxes) and the score of the gene (purple and blue lines). The purple line indicates the score for predicted deleterious missense changes, and the blue line is the score for the total number of non-silent changes. In the example below, the highest scoring window for missense and truncation changes is centred at position 2008.



Below the graph is information on the Blocks protein homology and an additional options box where you can examine a region of the gene that was not selected as the high scoring region. Below this, the changes and predicted effect of the changes can be seen at the sequence level. For a complete description of the symbols used, and more detailed information on CODDLE, please visit the **CODDLE glossary**.

10. When you are satisfied with the CODDLE output, click “Create primers for this window”.
11. Evaluate the information in the Primer3 window. Note that the optimum T<sub>m</sub> for primers is 70°C. Click “pick primers”.
12. In the output page, click “display this pair of primers” for your favorite set of primers.
13. You will now be directed to a page summarizing your primer choices. Note that the percentage of each type of change is listed.
14. When satisfied, click “order TILLING of this region”.



15. You are now directed to an STP order page. Enter the following email address: [test@fhcrc.org](mailto:test@fhcrc.org) and select Arabidopsis as the organism.
16. Click “place order”. Your order will now be searched in the STP database. If the target has been previously screened, you will be provided with information on found mutations. If it is a new target, it will be blasted against the Arabidopsis genome to ensure that the primers are designed to the correct organism. Once ready, click “store” to store the order in the database.
17. NOTES: The CODDLE input utility, CODDLE and Primer3 are all general tools that are available on the World Wide Web. You may find them useful for non-TILLING applications. Steps 14-16 have been included to illustrate that placing, verifying and confirming orders are tasks that have been automated by STP.
18. Additional Exercises: Once you have familiarized yourself with CODDLE and primer design, try inputting other information in the CODDLE input utility such as the Genbank URL of your favorite sequence (step 4). Also, try making additional Blocks with the SIFT programme (step 8). Finally, use the additional options window of the CODDLE output (step 9) to design primers to a different region of the gene.

### 11.3. Data analysis

The programme GelBuddy has been created to assist the discovery of mutations and polymorphisms ((ZERR and HENIKOFF 2005). It is available as a free download (<http://www.proweb.org/gelbuddy/>). This program should already be loaded onto the training course computers. For this exercise, download sample images from here (<http://tilling.fhrc.org/tillingdemo/ImagesforFAOgelBud/>). Be sure to download both the IRD700 and IRD800 images. The protocol below uses the basic Gelbuddy features for analysis of a standard TILLING gel. Tools are provided for the analysis of EcoTILLING or two dimensionally pooled gels that are not described. More information can be found at the GelBuddy page.

1. Download IRD700 and 800 jpeg or tiff images to your desktop. For example, download both 43ugfp115a\_bt.7 and 43ugfp115a\_bt.8.
2. Open Gel Buddy.
3. Import images. Under file, choose “Open 700 and 800 channel images”.
4. Select the first image to load. While holding down the shift key, select the second image. Click “open”.
5. Adjust the 700 channel image to the desired intensity using the slider bars located on the upper region of the GelBuddy window.



6. Adjust the 800 channel image. Click the 700-800 box at the top of the window to switch to the 800 channel. With the 800 channel selected, adjust the image as in step 5.



7. Call lanes. Click the “find lanes” box located in the tool bar at the top of the window.



8. Set the number of sample lanes in the “find lanes” pop up window (the default is 96 for a standard TILLING run). Select segmented lane tracks. Unless one of the channels is very bad, use the both channels for detecting lanes. Click “ok”.
9. Editing lanes. The blue lane markers should run through the lanes with the 200 bp marker. If they do not, or one or more lanes are called wrong, click the “edit lanes mode” in the toolbar.



10. Select the lane you wish to edit or the lane adjacent to the area where you wish to add a lane. Under the edit menu, select insert or delete lanes as required. If a lane merely needs to be “straightened”, select the boxed regions and drag to the desired location.
11. Click the “show lanes box” to remove lines.
12. Set the molecular weight migration. Click the “show calibration information” box. Vertical lines will appear.



13. Place the mouse over one of the numbers in blue and drag that number to the desired location on the gel. The 700 should align with the highest band in the ladder lanes. The 200 should align with the 200 bp marker.
14. Now set the 0% and 100% migration by dragging the red numbers to the bottom of the signal on the gel image (100%) and to the top of the full length product (0%). When complete, click the “calibration information” box again to make lines disappear.
15. Select mutations. Select the “record signals mode” box. Using the 700-800 box, switch between channels to find mutations. You will be prompted to enter the size of the full length product (0% migration). Enter the number at 0% and click “ok”. Enter your initials in the “created by” box. The signal grouping should be set to “all lanes”. Click the mouse over the mutation to select the mutation. When selecting mutations, note that mutations in the 700 channel are marked red and those in the 800 are marked with a blue box. If you are unsure of a mutation, note that the size of the band is given at the bottom of the window when your mouse is over the mutation. For any one lane, the sizes of bands in the blue and red boxes should equal the full length product. Do not be alarmed if the sizes are up to 100 base pairs off. To delete a box, hold down the option key and click the box.



16. Once you have selected all of the mutations, select the “show signals” box to remove the boxes. Look at the gel again to be sure you have selected real mutations. Select the box again to make the boxes reappear.



17. To zoom in to a region of the gel, select the “zoom in mode” box and click on the region you wish to enlarge. To zoom out, select the “zoom out mode” box. To fit the image back to the original window, select the “zoom to window” box.

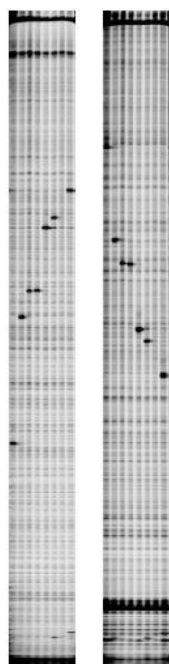


18. When you have finished analysing the gel, click the log box to see a report. Inspect the signals sorted by lane table. True mutations should have paired signals in the 700 and 800 channel that add up to the full-length product size.



19. Compare your data with what was found by STP. At STP, data from GelBuddy is directly posted to the program Squint in the STP database using the GelBuddy autopost function. You can view squint files for this exercise here (<http://tilling.fhcr.org/cgi-bin/displayWorkshop.pl?form=newSquint>). Under “squinting”, click “new/modify/view”. In the LI-COR run name field, enter the run name. The run name does not include .7.jpg or .8.jpg. For instance, for the first set of images on the images page, you would type 42600mla\_eb as the run name. Select “list current squint file” in the select a squint action box. Click the submit button to view the squint file. Did you find all the mutations? Did you find more than were reported? Note that mutations are given a confidence score based on quality. Confidence level A: the bands in both channels are clear and add up to the full-length size; level B: there are two corresponding bands but one of the bands is questionable; level C: data is available for one of the two channels but the band is most likely a mutation; level D: data is available for one of the two channels and the band is weak.
20. Try some other features in GelBuddy. For weak bands, try the “show inverted image” box to view the inverted image. Click the calibration box to show the horizontal calibration lines. Under the options pull down menu, try changing some of the calibration settings and see what happens to the lines. Notice that GelBuddy is compensating for lane to lane variation such as gel smiling. Want to see what the samples you processed should look like? Below is a test gel of these samples run in Seattle.





Mutations in the Arabidopsis OXII gene.  
IRD700 left, IRD800 right. Lanes 1-8  
equal samples 1-8 in the training course.

Lane	Fragment size in bp (IRD700)
1	288
2	441
3	476
4	477
5	566
6	580
7	603 *
8	624

\*The mutation in lane 7 is homozygous and thus is not detected when screened alone. For the course, an equal amount of wild-type DNA has been added to sample 7 so that heteroduplexes will be formed.

## 11.4. Additional info

### 11.4.1. List of consumables and equipment

Note that not all equipment is necessary for a successful TILLING operation, and not all equipment may be available. For instance, the comb-loading robot is no longer being sold by MWG, and neither are the thermal cyclers. Manual comb-loading is relatively easy, and most thermal cyclers should work for TILLING, so lower cost options are available.

#### Lab Supplies

Product	Company	Catalog Number
LI-COR 4300 S DNA analyzer	LI-COR	4300-02
Apricot pipettor	Perkin Elmer	PP-550
Combloder	MWG	Comblod
Centrifuge 5804 (Cel I)	Brinkman	2262250-1
Thermocycler Primus 96	MWG	4000-000005
Centrifuge 5810 (Genomic)	Brinkman	2262500-4
Nanopure (Water Treatment)	VWR (Barnstead)	13500-866
Centrifuge 5417C (PCR bench)	Brinkman	2262170-0
Equalizer (electric pipettes)	Matrix	2139
Heat plate sealer	Marsh (AB Gene)	AB-0384
pH meter	Fisher	13-636-AR10

Heat blocks	VWR	52434-232
Pipettes LTS multi-channel 20-1000µl	Rainin	L8-20, L8-200...
APC surge protector	CDWG	323633
Multi heat block	Fisher	NC9800611
Pipettes LTS single channel 20-1000µl	Rainin	L-20, L-200...
Stir plate	Fisher	11-500-49SH

Consumables

Product	Company	Catalog Number
Membrane combs	Gel Company	CAJ96
MWG 96-well plates	MWG	4050-000003
QT tip, 250ul clear, sterile filter tip	Molecular Bio Products	1043-60-5
QT tip, 500ul clear, sterile non-filter tip	Molecular Bio Products	1043-61-7
Acrylamide	Li Cor	82705607
Buffer reservoirs	Apogent Discoveries	8094
Sephadex G-50	A.Pharmacia	17-0043-02
EDTA	Research Organics	3002E
Ficol	Fisher	BP525-25
Tris	Research Organics	30960T
Boric acid	Research Organics	1748B
Milipore plates	Fisher	MAHVN 4550
Formamide	Sigma	F5786
Sealing tape PCR	Island Scientific	IS-609
Sealing tape non-PCR	Island Scientific	IS-SEAL
IRD 700	LI-COR	4200-60
IRD 800	LI-COR	4000-45
Taq, dNTP, PCR buffer	Pan Vera	TAK RR001C
Seq direct clean-up kit	Qbiogene	9904-200
EZPeel clear heat seal	Marsh Bio Prod	AB-0812
EZPeel aluminium heat seal	Marsh Bio Prod	AB-0745
LTS tips 10F	Rainin	GP-L10F
LTS tips 10S	Rainin	GP-L10S
LTS tips 200F	Rainin	GP-L200F
LTS tips 250S	Rainin	GP-L250S
LTS tips 1000F	Rainin	GP-L1000F
LTS tips 1000S	Rainin	GP-L1000S
20uL LTS tips spacesaver	Rainin	GPS-L10
200uL LTS tips spacesaver	Rainin	GPS-L250S
1000uLLTS tips spacesaver	Rainin	GPS-L1000S
Sephadex column loader 45ul	Fisher	MACL09645
Sephadex scraper replacement	Fisher	MACL0SC03

## 11.5. Frequently asked questions

*Will TILLING work in my favourite organism?*

TILLING is a general method and should work for most organisms. Requirements include the ability to induce mutations, propagate and/or store mutant organisms and PCR amplify gene specific targets.

*What about polyploids or duplicated gene targets?*

STP has successfully screened polyploid species. Additionally, Slade, *et al.*, have published TILLING data for polyploid wheat (SLADE *et al.* 2005). For polyploids and duplicated gene targets, a good approach is to pre-test unlabeled primers before purchasing IRD labeled primers. This is the approach taken for the Maize TILLING Project (<http://genome.purdue.edu/maizetilling/>). Following PCR and agarose gel analysis, products are sequenced. Primer pairs are selected for TILLING if they produce at least 7 ng/μl of product and sequence analysis indicates the amplification of a single target.

*What if there is no genomic sequence available for my organism?*

Short of cloning genes, you can design primers to EST data (or whatever is available) and pre-screen the primers. Sequencing the PCR products will provide genomic sequence information. It is important to select primers that yield products within the appropriate size range for your assay. Also, you may wish to avoid TILLING large amounts of intron as mutations in introns are likely to be non-functional. You may be able to use genomic sequence from a related organism to guess at the position of introns in your organism.

*I do not have access to a LI-COR, can I still TILL?*

The choice of read out platform (the machine used), can affect the level of allowable pooling, rate of false positives and negatives, robustness of the assay, as well as other factors. Thus, the choice of read out platform can have a large impact on the cost and throughput of your operation. STP has exclusively used LI-CORs and therefore it is difficult to comment directly on other platforms. Perry *et al.* published TILLING work using an ABI 377 (PERRY *et al.* 2003). Other end labeling strategies, such as using radioactivity, should work. Again, the throughput, efficiency and screening cost associated with the platform should be considered.

An alternative to end labeling is body labeling. Body labeling DNA may not be as efficient as end labeling either the DNA or a probe. That said, one can use single-strand specific nucleases to induce double strand breaks in DNA, allowing visualization on native agarose gels (BURDON and LEES 1985; CHAUDHRY and WEINFELD 1995; HOWARD *et al.* 1999; SOKURENKO *et al.* 2001) Most likely, this will prove to be a lower throughput option.

*I am more interested in EcoTILLING. How is it different?*

EcoTILLING is a method for the discovery and genotyping of natural polymorphisms (COMAI *et al.* 2004). The starting material for EcoTILLING is DNA from “natural” populations rather than mutagenized ones. Depending on the population, one might expect a substantially higher frequency of polymorphisms than the rare induced mutations found in a chemically mutagenized population. The wet bench protocols used for TILLING and EcoTILLING are the same. GelBuddy has been designed to work with EcoTILLING data and some EcoTILLING-specific features are available in GelBuddy.

Will a chemical mutagen be effective on all genes? What about background mutations in the lines? Do I need a license to TILL?

For answers to these questions, please see the STP FAQ page (<http://tilling.fhcr.org/files/FAQ.html>).

## 11.6. Additional protocols

### 11.6.1. Sequencing

This protocol is a scaled down version of the standard high-throughput sequencing protocol.

H <sub>2</sub> O	54.8µl
Ex Taq buffer	10µl
dNTP	8µl
forward primer (10 µM)	1µl
reverse primer (10 µM)	1µl
HS Ex Taq	0.25µl

Add 15 µl mix to 5 µl DNA and mix well.

Run the following programme:

<i>Step 1</i>	Initial denaturation	95°C	2 minutes
<i>Step 2</i>	Denaturation	94°C	20 seconds
<i>Step 3</i>	Primer annealing	73°C (-1°C/cycle)	30 seconds
<i>Step 4</i>	Ramp	0.5°C per second to 72°C	
<i>Step 5</i>	Primer extension	72°C	1 minute
<i>Step 6</i>	Cycling	repeat steps 2-5 for 7 cycles	
<i>Step 7</i>	Denaturation	94°C	20 seconds
<i>Step 8</i>	Primer annealing	65°C	30 seconds
<i>Step 9</i>	Ramp	0.5°C per second to 72°C	
<i>Step 10</i>	Primer extension	72°C	1 minute

<i>Step 11</i>	Cycling	repeat steps 7-10 for 44 cycles	
<i>Step 12</i>	Final extension	72°C	5 minutes
<i>Step 13</i>	Hold	4°C	forever

Quantify yield on an agarose gel (this is normally done only on 1 row of a 96 well plate).

#### Pre-sequencing clean-up:

To 10 µl PCR product add and mix well:

\*4 µl Shrimp alkaline phosphatase

\*1 µl Endonuclease I (keep enzymes on ice at all times)

\*Check company protocol for units/µl

Incubate 37°C for 15 min., 80°C for 15 min. (Follow manufacturer's suggestion).

The pre-sequencing amplification is performed with the unlabeled primers used in the TILLING screen. Following the manufacturer's protocol, HS Ex-Taq (Takara) is used in a 20 µl final reaction volume with 0.005 ng genomic DNA (for Arabidopsis).

Sequencing RXN (Big Dye version 3.0 or higher/ ABI 3100 or higher)

Add 5 µl of 5% DMSO to PCR product and mix

To new set of tubes add:

4 µl diluted Big Dye (version 3.0 or higher) (1:1 dilution with PCR H<sub>2</sub>O)

1 µl forward primer (3 µM)

5 µl PCR product (diluted with DMSO)

Mix well and spin down.

<i>Step 1</i>	Initial denaturation	95°C	5 minutes
<i>Step 2</i>	Denaturation	95°C (ramp at 1°C/sec)	10 seconds
<i>Step 3</i>	Primer annealing	50°C (ramp at 1°C/sec)	5 seconds
<i>Step 4</i>	Primer extension	60°C (ramp at 1°C/sec)	4 minutes
<i>Step 5</i>	Cycling	repeat steps 2-4 for 24 cycles	
<i>Step 6</i>	Hold	8°C (ramp at 1°C/sec)	forever

Big dye removal and running the ABI is performed by a core facility.

Sequence trace analysis is performed using Sequencher™ 4.5 software (Gene Codes). Both heterozygous and homozygous mutations can be confirmed utilizing the mapping information gathered in the TILLING screens.

## 11.7. EMS mutagenesis of Arabidopsis seed

EMS mutagenesis of maize pollen for the population used in the Maize TILLING Project has been described (TILL *et al.* 2004b).



### 11.7.1. Materials

Orbital shaker: Aros 160 with a 1.25 cm radius of gyration  
 10-15 L tub  
 Microfuge tubes with 50 mg of seed each  
 Stir plate and stir bar  
 1000 ml beaker  
 1 L 2 N NaOH  
 Squeeze bottle of di H<sub>2</sub>O  
 10% Tween 20  
 P-1,000 pipetter with barrier tips. Some of these ought to have notches cut in Tip as per  
 “A Note on Technique.”  
 P-20 pipetter with barrier tips  
 EMS (methanesulphonic acid ethyl ester), Sigma  
 Glass scintillation vials I.D. = 2.5 cm  
 Box for dry hazardous materials disposal  
 Plastic bag for hazardous materials disposal  
 Box of nitrile (not latex) gloves and a lab coat

### 11.7.2. Standard size batch

In order to avoid variation in mutation rate that could arise from scaling properties, the first 10 mutagenesis procedures for this project except the 6th were done in standard batches of 50 mg seed in 4ml of EMS solution. Only flat-bottomed glass scintillation vials of 2.5 cm ID were used so as to avoid subtle variations in the agitation of the seeds. This standard procedure did not make the concentration of EMS a good predictor of the EL count. Because of this, and to allow reducing the number of people needed to care for a batch of M<sub>1</sub> plants, quantities of seeds less than 50 mg are now allowed.

### 11.7.3. A note on technique

Before beginning this procedure, cut a couple of notches in the tip of several of the P1000 tips. If the notch is too small to allow seeds to pass through, the tip can be pressed against the bottom of the scintillation vial and the supernatant can be efficiently aspirated without loss of seeds.

#### Day 1:

1. Preparation of Fume Hood for procedure.
  - 1.1. Label each scintillation vial with the concentration of EMS that is to be used in it.
  - 1.2. Warn all personnel that a dangerous procedure is about to be performed in the hood.
  - 1.3. Place all materials in hood.
  - 1.4. Put 125 ml of 2 N NaOH and 375 ml of H<sub>2</sub>O in beaker with stir bar slowly rotating.  
Place remaining 875 ml of 2 N NaOH in tub with 2.6 L H<sub>2</sub>O.

2. Add 4 ml of H<sub>2</sub>O to each vial and mark level with a fine tip marker then empty vial of H<sub>2</sub>O.
3. Rinse seed into each vial with 4 ml of diH<sub>2</sub>O. Add 40 ml of 10% Tween 20 to each vial and agitate at 180 RPM for 15 sec.
4. Pipette off Tween/ H<sub>2</sub>O and add 4ml DI-H<sub>2</sub>O to each vial. Agitate for 5 min at 180 RPM. Repeat for 4 total washes.
5. Add DI-H<sub>2</sub>O to each vial to 4 ml line made in 2) in order to achieve a total volume of 4 ml.
6. Use gloves, lab jacket, and fear for following steps.
7. Add .425 X (ml) EMS to each vial with barrier tip P-20s. X is desired [EMS] (mM). Dispose of tips in beaker of 0.5 N NaOH.
8. Agitate for 17 hr at 180 RPM at room temperature.

### Day 2:

1. Pipette off EMS solution from each vial and dispose in flask of 0.5 N NaOH.
2. Fill each vial to shoulder with di H<sub>2</sub>O from squeeze bottle, swirl by hand, then pipette off supernatant and dispose as in 1). Repeat 5 times.
3. Add diH<sub>2</sub>O to vial to achieve 4 ml and agitate 15 sec.
4. Pipette off as in 2) and repeat.
5. Store at 4°C until sown.
6. Allow NaOH that has been used for EMS disposal to stir for 30 min., then gently pour contents of beaker into tub of 0.5 N NaOH, placing beaker in tub as well, then pour down drain and flush with cold running water for 15 min.
7. Wipe off pipettes and inside of hood with dil NaOH, and call Hazardous Materials Disposal to remove solid waste.

### 11.7.4. DNA extraction

DNA isolation is done per FastDNA a kit protocol (revision #6540-999-1D04, <http://www.qbiogene.com/fastprep/protocols.shtml>), with the following variations and warnings:

1. Use only one ceramic bead per shaker tube.
2. Run shaker for 45 min at 4.5 m/s.
3. The first centrifuge spin should be at 14,000 × g for up to 30 min. Draw off as much as 800 - 900 ml supernatant from shaker tube in Step 3 of the FastDNA protocol.
4. After DNA is bound in a pellet to the Binding Matrix, take care not to disturb the pellet when discarding supernatants in Step 4 of the FastDNA protocol.
5. To make re-suspension easier, all spins before a re-suspension (both 1-minute spins in Step 4 of the Fast DNA protocol) should be at 9,000-10,000 × g for 3 min.
6. To re-suspend a pellet (Steps 4 and 5 of the Fast DNA protocol), use the vortex or noisily rake the tube across a tube rack, a practice known as “ducking” for the quack-like sound made. When ducking, take care to hold down the cap of the tube to prevent it from popping open.

In Step 5 of the Fast DNA protocol, elute binding matrix with 200 ml DES. Spin at  $14000 \times g$  for ~5 min. Then pipette off 180 ml of supernatant, taking extreme care not to draw up particles of Binding Matrix, and transfer supernatant to a sterile screw-top tube. Add 20 ml of 10x TE @ 3.2 mg/ml RNase A.

## 11.8. References

- BURDON, M. G., and J. H. LEES, 1985 Double-strand cleavage at a two-base deletion mismatch in a DNA heteroduplex by nuclease S1. *Biosci Rep* **5**: 627-632.
- CHAUDHRY, M. A., and M. WEINFELD, 1995 Induction of double-strand breaks by S1 nuclease, mung bean nuclease and nuclease P1 in DNA containing abasic sites and nicks. *Nucleic Acids Res* **23**: 3805-3809.
- COLBERT, T., B. J. TILL, R. TOMPA, S. REYNOLDS, M. N. STEINE *et al.*, 2001 High-throughput screening for induced point mutations. *Plant Physiol* **126**: 480-484.
- COMAI, L., K. YOUNG, B. J. TILL, S. H. REYNOLDS, E. A. GREENE *et al.*, 2004 Efficient discovery of DNA polymorphisms in natural populations by Ecotilling. *Plant J* **37**: 778-786.
- GALEANO, C. H., M. GOMEZ, L. M. RODRIGUEZ and M. W. BLAIR, 2009 CEL I Nuclease Digestion for SNP Discovery and Marker Development in Common Bean (*Phaseolus vulgaris* L.). *Crop Science* **49**: 381-394.
- GARVIN, M. R., and A. J. GHARRETT, 2007 DEco-TILLING: an inexpensive method for single nucleotide polymorphism discovery that reduces ascertainment bias. *Molecular Ecology Notes* **7**: 735-746.
- HOWARD, J. T., J. WARD, J. N. WATSON and K. H. ROUX, 1999 Heteroduplex cleavage analysis using S1 nuclease. *Biotechniques* **27**: 18-19.
- MCCALLUM, C. M., L. COMAI, E. A. GREENE and S. HENIKOFF, 2000 Targeted screening for induced mutations. *Nat Biotechnol* **18**: 455-457.
- PERRY, J. A., T. L. WANG, T. J. WELHAM, S. GARDNER, J. M. PIKE *et al.*, 2003 A TILLING reverse genetics tool and a web-accessible collection of mutants of the legume *Lotus japonicus*. *Plant Physiol* **131**: 866-871.
- SATO, Y., K. SHIRASAWA, Y. TAKAHASHI, M. NISHIMURA and T. NISHIO, 2006 Mutant Selection from Progeny of Gamma-ray-irradiated Rice by DNA Heteroduplex Cleavage using Brassica Petiole Extract. *Breeding Science* **56**: 179-183.
- SLADE, A. J., S. I. FUERSTENBERG, D. LOEFFLER, M. N. STEINE and D. FACCIOTTI, 2005 A reverse genetic, nontransgenic approach to wheat crop improvement by TILLING. *Nat Biotechnol* **23**: 75-81.
- SOKURENKO, E. V., V. TCHESNOKOVA, A. T. YEUNG, C. A. OLEYKOWSKI, E. TRINTCHINA *et al.*, 2001 Detection of simple mutations and polymorphisms in large genomic regions. *Nucleic Acids Res* **29**: E111.
- TILL, B. J., C. BURTNER, L. COMAI and S. HENIKOFF, 2004a Mismatch cleavage by single-strand specific nucleases. *Nucleic Acids Res* **32**: 2632-2641.

- TILL, B. J., S. H. REYNOLDS, E. A. GREENE, C. A. CODOMO, L. C. ENNS *et al.*, 2003 Large-scale discovery of induced point mutations with high-throughput TILLING. *Genome Res* **13**: 524-530.
- TILL, B. J., S. H. REYNOLDS, C. WEIL, N. SPRINGER, C. BURTNER *et al.*, 2004b Discovery of induced point mutations in maize genes by TILLING. *BMC Plant Biol* **4**: 12.
- TILL, B. J., T. ZERR, L. COMAI and S. HENIKOFF, 2006 A protocol for TILLING and Ecotilling in plants and animals. *Nat Protoc* **1**: 2465-2477.
- ZERR, T., and S. HENIKOFF, 2005 Automated band mapping in electrophoretic gel images using background information. *Nucleic Acids Res* **33**: 2806-2812.

## 12. ALTERNATIVE ENZYMOLOGY FOR MISMATCH CLEAVAGE FOR TILLING AND ECOTILLING: EXTRACTION OF ENZYMES FROM WEEDY PLANTS

### 12.1. Objective

A crude celery extract containing the single-strand-specific nuclease CEL1 has been widely used in TILLING and Ecotilling projects around the world. Yet, celery is hard to come by in some Member States. Based on previous studies and bioinformatic analysis suggestion homologies exist to CEL1 in all plants. Therefore, we developed a protocol for extraction of active enzyme from plants common across the world: weeds. We isolated weed plants from the grassland around the Seibersdorf laboratories and isolated a crude enzyme extract (in parallel to the enzyme extracts from celery). Since, there was no or only very low mismatch digestion activity in the crude extract, we applied a centrifuge-based filter method to concentrate the enzyme extract.

### 12.2. Materials

MATERIALS / BUFFERS FOR ENZYME EXTRACTIONS	Notes
hand-held mixer (or juicer)	From any supplier
STOCK: 100mM PMSF (stock in isopropanol)	To prepare an aqueous solution of 100µM PMSF (for buffers A and B), add 1 ml 0.1M PMSF per liter of solution immediately before use.
STOCK: 1M Tris-HCl, pH 7.7.	
Buffer A: 0.1 M Tris-HCl, pH 7.7, 100 µM PMSF.	
Buffer B: 0.1 M Tris-HCl, pH7.7, 0.5 M KCl, 100 µM PMSF.	
Dialysis tubing with a 10,000 Dalton molecular weight cut off (MWCO)	e.g. Spectra/PorR Membrane MWCO: 10,000, Spectrum Laboratories, Inc.
(NH <sub>4</sub> ) <sub>2</sub> SO <sub>4</sub> (Ammonium sulphate)	
Sorvall Centrifuge	Or equivalent centrifuge/rotor combination to achieve needed gravitational force

MATERIALS FOR CONCENTRATION OF ENZYME EXTRACTS	
Amicon Ultra Centrifugal Filters (0.5mL, 10K)	Millipore Amicon Ref.No. UFC501024 24Pk
Refrigerated (4°C) Microcentrifuge	e.g. Eppendorf 5415R

<b>TILLING-PCR</b>	
Thermocycler	e.g. Biorad C1000 Thermal cycler
PCR tubes	Life Science No 781340
TaKaRa Ex Taq™ Polymerase (5U/ul)	TaKaRa
10X Ex Taq™ Reaction Buffer	TaKaRa
dNTP Mixture (2.5mM of each dNTP)	TaKaRa
Agarose gel equipment	

## 12.3. Methods

### 12.3.1. Enzyme extraction

1. Collect approximately 200 grams of mixed monocot and dicot weedy plants were collected that were growing on the periphery of our sorghum field.
2. Wash material 3x in water and then ground using a hand-held mixer and by adding about 300 mls of water to facilitate tissue disruption (or optional in a juicer)
3. Add 1M Tris-HCl (pH7.7) and 100mM PMSF to a final concentration of Buffer A (0.1M Tris-HCl and 100µM PMSF) (NOTE: Stocks and water should be kept at 4°C, perform subsequent steps at 4°C)
4. Centrifuge for 20 min at 2600 x g in Sorvall GSA rotor to pellet debris. Save supernatant.
5. Bring the supernatant to 25% ammonium sulphate (add 144 g per liter of solution). Mix gently at 4°C (cold room) for 30 min.
6. Centrifuge for 40 min at 4°C at ~14,000 x g in sorvall GSA rotor (~9000 rpm). Discard the pellet.
7. Bring the supernatant to 80% ammonium sulphate (add 390 g per liter of solution). Mix gently at 4°C for 30 min.
8. Centrifuge for 1.5 hours at 4°C at ~14,000 x g in sorvall GSA rotor. SAVE the pellet. Discard the supernatant (be careful in decanting the supernatant!) The pellet can be stored at -80°C for at least two weeks.
9. OPTIONAL: Pellets can be frozen at -80°C for months.
10. Resuspend the pellets in ~ 1/10 the starting volume with Buffer B (Frozen pellets of the weed juice extract were suspended in 15mL Buffer B and pellets of the celery juice extract in 10 mL Buffer B). Ensure the pellet is thoroughly resuspended.

11. Dialyze against Buffer B at 4°C (2 Liters per 10mls of resuspended solution).. Use e.g. Spectra por 7 MWCO 10000 tubing. (NOTE: Soak the dialysis tubing in nanopure water for 30 min. before use.)
12. Dialyze for 1 hour against Buffer B at 4°C
13. Repeat for a total of 4 dialysis steps with a minimum of 4 hours dialysis. (NOTE: Longer dialysis is better, it is often convenient to perform the third dialysis overnight).
14. Remove liquid from dialysis tubing. It is convenient to store ~75% of the liquid in a single tube at -80°C and the remainder in small aliquot for testing. This protein mixture does not require storage in glycerol and remains stable through multiple freeze-thaw cycles, however, limiting freeze thaw cycles to 5 limits the chance of reduced enzyme activity
15. Perform activity test (step 3.3, or proceed immediately to enzyme concentration, step 3.2)



Figure 1. Mixture of different plant species (weedy plants) from the grassland around the Seibersdorf laboratories used for the isolation of an enzyme extract for mismatch cleavage.

### 12.3.2. Concentration of enzyme extractions

Concentration of weed and celery enzyme extracts is done using Amicon Ultra 10K centrifugal filter devices (for 0.5mL starting volume; in 1.5-mL tubes).

1. Perform with 600µL of protein extract after dialysis
2. Clear extract by centrifugation at 30 min / 10,000 x g / 4°C (to pellet plant material) in refrigerated microcentrifuge
3. Transfer 500 µL of the (cleared) supernatant to a filter device (keep the rest of the supernatant as control “before concentration”).

4. Centrifuge the filter device with a collection tube inserted per manufacturer's instructions for 30 min / 14,000 x g / 4°C
5. Remove filter device, invert and place in new collection tube.
6. Centrifuge for 2 min / 1,000 x g / 4°C
7. Measure the recovered volume. This is your concentrated protein. Calculate the concentration factor with the following formula: Starting volume/Final volume = concentration factor

### 12.3.3. Test of Mismatch Cleavage Activity

1. Produce TILLING-PCR products for mismatch cleavage tests with the concentrated enzyme extracts. The example below is for barley.

GENES/PRIMER: nb2-rdg2a (1500bp-PCR product)  
 nb2-rdg2a\_F2           TCCACTACCCGAAAGGCACTCAGCTAC  
 nb2-rdg2a\_R2           GCAATGCAATGCTCTTACTGACGCAA

TILLING PCR REACTIONS (TaKaRa ExTaq enzyme): total volume: 25uL

10x ExTaq buffer (TaKaRa)	2.5 µL
dNTP mix (2.5 mM)	2.0 µL
Primer forward (10 µM)	0.3 µL
Primer reverse (10 µM)	0.3 µL
TaKaRa Taq (5U /µl)	0.1 µL
Barley genomic DNA (5 ng/µL)	5.0 µL
H2O (to 25 µL)	14.8 µL

TILLING PCR cycling program for TILLING ("PCRTM70"):

95°C for 2 min;

loop 1 for 8 cycles (94°C for 20 s, 73°C for 30 s, reduce temperature 1°C per cycle, ramp to 72°C at 0.5°C/s, 72°C for 1 min);

loop 2 for 45 cycles (94°C for 20 s, 65°C for 30 s, ramp to 72°C at 0.5°C/s, 72°C for 1 min);

72°C for 5 min;

99°C for 10 min;

loop 3 for 70 cycles (70°C for 20 s, reduce temperature 0.3°C per cycle); hold at 8°C.

2. Mix 10µL of PCR product with 10uL weed digestion mix to a volume of 20µL
3. Incubate at 45°C for 15 min
4. Add 2.5µL of 0.5M EDTA (pH 8.0) – to stop reaction
5. Load a 10µL aliquot on an agarose gel



## 12.4. Example results

Concentrations of protein extracts:

Table 1. Calculations of concentration factors after centrifugation with Amicon Ultra 10K – Starting volume: 500 µL (“Before” centrifugation = considered as 1x concentrated)

	<b>Recovered volume</b>	<b>Concentration factor (calculated from 500 µL starting volume)</b>
Weed	~42 µL	11.9x
Cell	~33 µL	15.2x

Mismatch digestions using celery and weed enzyme extracts:

Table 2. Set-up of mismatch digestions using celery and weed enzyme before and after centrifugation with Amicon Ultra 10K. The enzyme concentration in the extracts were calculated using the calculated concentration factors from Table 1.

	<b>1 - BEFORE</b>	<b>2 - after</b>	<b>3 - after</b>	<b>4 - after</b>
Enzyme	3.5 uL (1x)	0.5 uL	3 uL	6 uL
Cell buffer	1.5 uL	1.5 uL	1.5 uL	1.5 uL
H2O	5 uL	8.0 uL	5.5 uL	2.5 uL
Tot. Volume	10 uL	10 uL	10 uL	10 uL
<b>Celery</b> enzyme concentration in relation to extract before centrifugation (3.5uL – before = 1x)	<b>1x</b>	7.6 uL <b>2.2x</b>	45.6 uL <b>13.0x</b>	91.2 uL <b>26.1x</b>
<b>Weed</b> enzyme concentration in relation to extract before centrifugation (3.5uL - before = 1x)	<b>1x</b>	5.95 uL <b>1.7x</b>	35.7 uL <b>10.2x</b>	71.4 uL <b>20.4x</b>

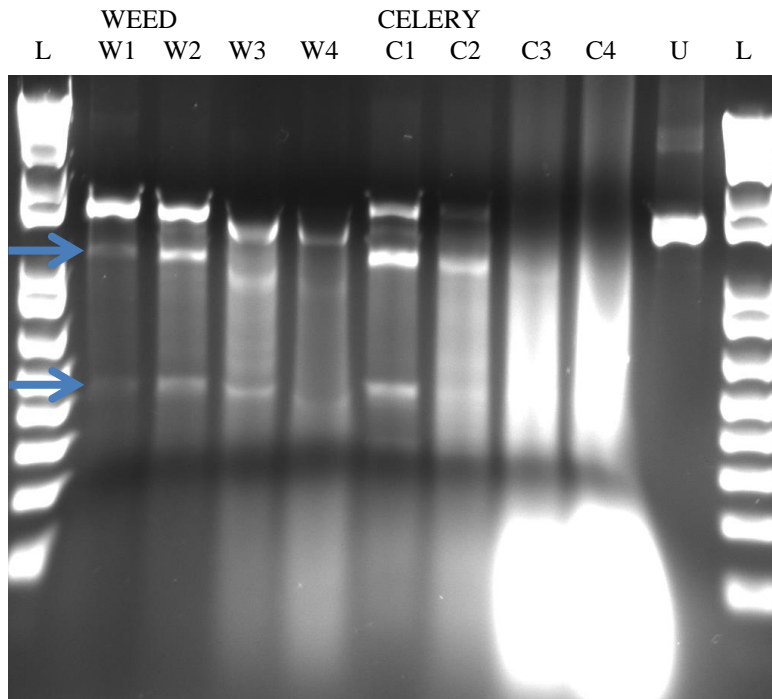


Figure 2. Mismatch cleavage with celery and weed enzyme extracts. TILLING-PCR products of the target gene *nb2-rdg2a* (1500bp-PCR product) were produced from genomic DNA of barley. The PCR products were digested with weed and celery enzyme extracts before and after concentration by centrifugation with Amicon Ultra 10K. 10  $\mu$ L of the digested PCR products were separated on a 1.5% agarose gel. Position of SNPs are marked with blue arrows. Concentrations of Weed (W) and Celery (c) extracts are listed above the lanes. A 1kb ladder is loaded on either side of the samples.

## 12.5. Conclusions

Crude enzyme extracts of weeds show a similar activity to that of celery extract for the cleavage of single nucleotide polymorphisms. The per unit activity, however, was lower than than for CEL I, likely owing to the co-precipitation of other plant proteins in weeds, presumably including a larger amount of RUBISCO. This limitation can be overcome through the use of a simple centrifugation based protein concentration step. 150 ml of weed extract produces enough enzyme for approximately 2000 reactions with this protocol.

## 13. LOW-VOLUME, NON-TOXIC AND RAPID EXTRACTION OF SINGLE-STRAND-SPECIFIC NUCLEASES FROM CELERY

### 13.1. Objective

The aim of this protocol is to provide a quick method for crude celery juice (CJE) extraction for 5000 reactions or more that removes the toxicity and use of specialized equipment and methods (preperatory centrifuge and dialysis), so that it can be performed in a standard molecular biology laboratory. This enzyme is used for SNP and small indel discovery and genotyping applications such as TILLING and Ecotilling.

### 13.2. Materials

1. Juicer (*e. g.*, Le Quipe).
2. Celery.
3. 1M Tris-HCl, pH 7.7.
4. Buffer A: 0.1 M Tris-HCl, pH 7.7,
5. Buffer B: 0.1 M Tris-HCl, pH7.7, 0.5 M KCl,
6. Amicon Ultra 0.5ml 10K Centrifugal filters (Millipore Amicon Ref.No. UFC501024 24Pk).

### 13.3. Methods

#### 13.3.1. CEL I preparation

1. Perform all steps at 4°C when possible. Most steps can be performed at room temperature.
2. Rinse desired amount of celery with water. One bunch (~1 lb) yields approximately enough CEL I for 500,000 standard TILLING reactions. Remove any leaves and cut off tough tissue at base of stalk. For this protocol we aim for the production of about 15mls of juice with 0.5kg of material typically giving 200-400mls.
3. Juice the desired amount of material.
4. Add 1M Tris to a final concentration of Buffer A.
5. 18 mL celery juice + 2 mL 1M Tris-HCl buffer (pH=7.7)
6. Distribute liquid into 10 2.0 ml microcentrifuge tubes.

7. Spin the juice for 20min at 2600 g to pellet debris at 4C if possible.
8. Pour supernatant into a beaker.
9. Bring the supernatant to 25%  $(\text{NH}_4)_2\text{SO}_4$  by adding 144 g per liter of solution. Mix gently at 4°C for 30 min. Using stir plate and magnetic stir bar.  
*Total volume (from 10 tubes): 18.5 mL – 2.66g  $(\text{NH}_4)_2\text{SO}_4$  added*



Figure 1. Protein precipitation with  $(\text{NH}_4)_2\text{SO}_4$

10. Distribute liquid into 10 2.0 ml microcentrifuge tubes. Spin at 15000 g at 4°C for 40 min.

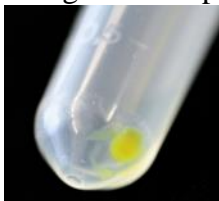


Figure 2. Protein pellet after 25%  $(\text{NH}_4)_2\text{SO}_4$  precipitation (discard)

11. Pour supernatant into clean beaker. Discard pellet.
12. Bring the supernatant from 25% to 80%  $(\text{NH}_4)_2\text{SO}_4$  by adding 390 gram per liter of solution. Mix gently at 4°C for 30 min.  
*Total volume (from 10 tubes): 18 mL – 7.02g  $(\text{NH}_4)_2\text{SO}_4$  added*
13. Distribute liquid into 10 11 2.0 ml microcentrifuge tubes. Spin 15000 x g for 1.5 hr. Save the pellet and discard the supernatant, being careful in decanting the supernatant. The pellet can be stored at -80°C for months.



Figure 3. Protein pellet after 80%  $(\text{NH}_4)_2\text{SO}_4$  precipitation (keep and resuspend)

14. Resuspend the pellets in  $\sim 1/10$  the starting volume with Buffer B, ensuring the pellet is thoroughly resuspended. Target final volume for all 10 pellets is 1.5mls. Add 1.5mls buffer B to tube #1, resuspend by pipetting up and down or vortexing. Then transfer this liquid to tube #2 and repeat, continue until the last tube.

*In total it were 11 tubes with pellets: 5 pellets were resuspended in 750uL and 6 pellets in 750uL – then combined to 1.5 mL (total volume of liquid + pellets  $\sim 2$  mL)*



Figure 4.  $\sim 2$  mL liquid after re-suspension and combination of 11 pellets (derived from 80%  $(\text{NH}_4)_2\text{SO}_4$  precipitation)

15. Desalting: Use Amicon ultra filters. Distribute liquid into four filter devices, making sure not to exceed 500ul in any filter. Attach collection tube and spin at 14000g for 30minutes. When complete, remove liquid from collection tube and add 500ul buffer B and repeat. Repeat this step a total of 4 times.

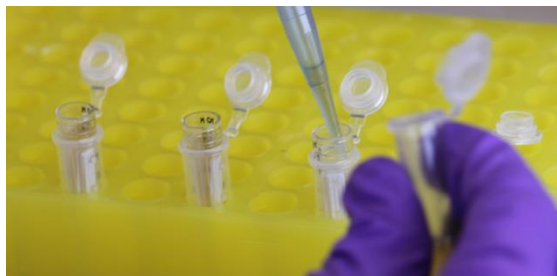


Figure 5. Transfer of liquid after resuspension to 4 Amicon Ultra filter devices

**Table 1.** Volume of retained liquid in the Amicon Ultra filter devices after the 5 centrifugation steps and a resulting (calculated) desalting factor.

Centrifugation (30 min at 15000 g)	Starting volume	End volume (+ buffer added)	Desalting	Desalting (calculated)	factor
1	500	100 + 400	5x	5x	
2	500	$\sim 50 + 450$	10x	50x	
3	500	$\sim 35 + 470$	14x	700x	
4	500	$\sim 30 + 475$	16.6x	<b>11620x</b>	
5	500	For elution	no	-	

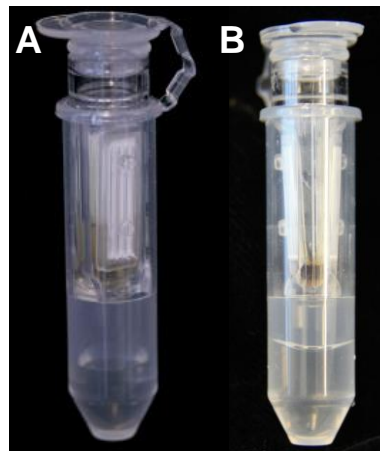


Figure 6. Amicon Ultra filter device (A) after 1<sup>st</sup> centrifugation (~100uL liquid retained in the filter device and ~400uL flow-through) and (B) after 5<sup>th</sup> centrifugation (~35uL liquid retained in the filter and ~465uL flow-through)

16. Elute sample: To elute sample, invert the filter and place inverted into a new collection tube. Centrifuge at 1000g for 2min.

**Table 2.** Volumes of recovered liquid from each Amicon Ultra filter device after inverted centrifugation and (calculated) concentration factor of the enzyme.

Eluate	Starting vol	Elution volume	
1		35	
2		45	
3		40	
4		35	
<b>Total</b>	<b>2000 uL</b>	<b>155 uL</b>	<b>12.9 x</b>

17. Combine all eluates

*Final volume of 4 tubes: 155 uL*



Figure 7. Recovered eluates after centrifugation of the inverted Amicon filter devices

18. Centrifuge (remove solid material) - 4°C for 30 min at 10.000g  
*Centrifuged an aliquot of 70uL (other 70 uL frozen without centrifugation)*

19. Use supernatant for activity test

### 13.3.2. Activity tests

For standard TILLING applications, test a range of amounts of CEL I for activity with known mutations/polymorphisms following the high throughput TILLING protocol. Target amount per reaction  $X = 7.5 \times 10^{-8}$  x total amount juice in  $\mu\text{l}$ . For example, the target range for a bunch of celery giving 400mls juice is  $400000 \times 7.5 \times 10^{-8}$  or 0.03  $\mu\text{l}$  per reaction. To assay activity, perform a standard titration curve with the outermost points flanking the target range on either side by a factor of 100. With excess enzyme, full length PCR product will disappear; as the amount of enzyme falls below the target range, PCR product and background bands will become increasingly dark to the point where the image becomes difficult to interpret.

Synthesis of TILLING-PCR products for mismatch cleavage tests

(barley TILLING primer #13): GENES/PRIMER: Mlo9 (1476 bp-PCR product)

#13 HV\_Mlo9-F2 CATTGTGTCGCAAAACAGCAAGTTCGAC  
HV\_Mlo9-R2 TTGTCTCATCCCTGGCTGAAGGAAAA

TEMPLATE: 1:1-mixture of Golden Promise and HOR-1606 gDNA – mixture gives mismatch cleavage

TILLING PCR REACTIONS (TaKaRa ExTaq enzyme): total volume: 25uL

10x ExTaq buffer (TaKaRa)	2.5 uL
dNTP mix (2.5 mM)	2.0 uL
Primer forward (10 uM)	0.3 uL
Primer reverse (10 uM)	0.3 uL
TaKaRa Taq (5U /ul)	0.1 uL
Barley genomic DNA (5 ng/uL)	5.0 uL
H2O (to 25 uL)	14.8 uL

TILLING PCR cycling program for TILLING (“PCRTM70”)

95°C for 2 min;

loop 1 for 8 cycles (94°C for 20 s, 73°C for 30 s, reduce temperature 1°C per cycle, ramp to 72°C at 0.5°C/s, 72°C for 1 min);

loop 2 for 45 cycles (94°C for 20 s, 65°C for 30 s, ramp to 72°C at 0.5°C/s, 72°C for 1 min);

72°C for 5 min;

99°C for 10 min;

loop 3 for 70 cycles (70°C for 20 s, reduce temperature 0.3°C per cycle);

hold at 8°C

- CELI-digestions: mix 10uL of PCR product with 10uL digestion mix to a volume of 20uL (see Table 3 for set-up of digestion mixes)
- Incubate at 45<sup>0</sup>C for 15 min
- Add 2.5uL of 0.5M EDTA (pH 8.0) – to stop reaction
- Load a 10uL aliquot on an 1.5% agarose gel

SERIAL DILUTIONS OF CELI enzyme

**Table 3.** Serial dilutions of isolated Cell enzyme and set-up of Cel digestion mixes.

Dilution factor		Enzyme (uL)	Cell buffer	H2O (uL)
<b>(A) CELI from dialysis</b>	<b>1x</b>			
0x		0	1.5 uL	8.5
0.1x		0.35 (1:10)	1.5 uL	8.15
<b>1x</b>	<b>= 0.35 uL</b>	<b>3.5 (1:10)</b>	1.5 uL	5
5x		1.75	1.5 uL	6.75
10x		3.5	1.5 uL	5
24x		8.5	1.5 uL	0
<b>(B) CELI from Amicon filters*</b>	Dilution factor – <b>12.9x</b>			
0x		0	1.5 uL	8.5
0.1x		0.25 (1:100)	1.5 uL	8.25
<b>1x</b>	<b>= 0.027 uL</b>	<b>0.25 (1:10)</b>	1.5 uL	8.25
5x		1.25 (1:10)	1.5 uL	7.25
10x		2.5 (1: 10)	1.5 uL	6.0
20x		0.5	1.5 uL	8.0
50x		1.25	1.5 uL	7.25
100x		2.5	1.5 uL	6.0

\*1x concentration of the CELI purified with Amicon filter devices were calculated using the concentration factor 12.9x



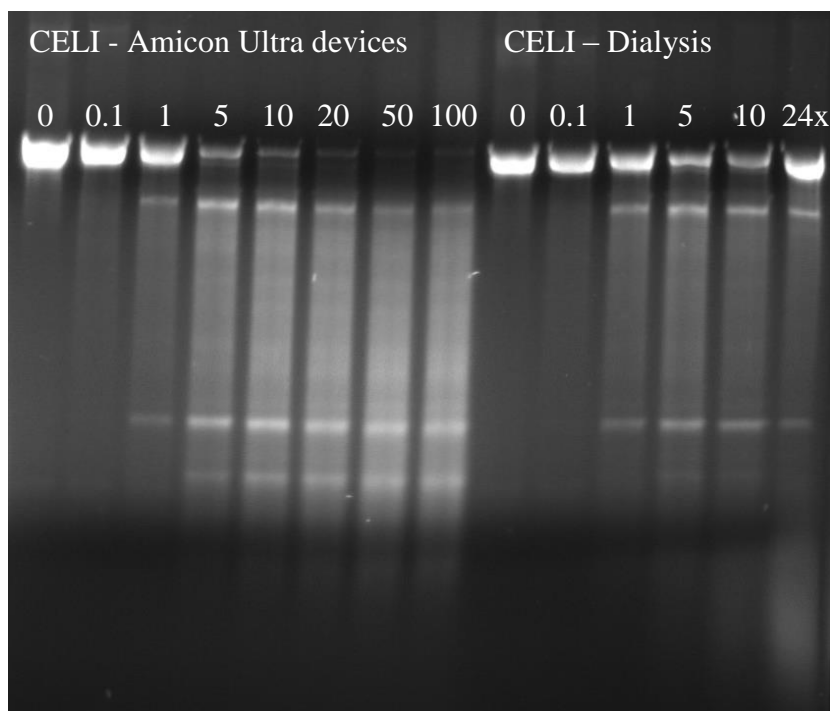


Figure 8. Activity tests of CELI enzyme isolated with Amicon Ultra filter devices (left) and with dialysis method (right). A serial dilution of enzyme activity is shown. There are no cleavage bands present at the control without CELI enzyme (0x) and at the lowest dilution (0.1x) in both extracts. Both enzyme extracts show activity from 1x diluted. However, the background in the CELI extracts purified with Amicon filter devices seems to be stronger than in the activity assays of CELI enzyme purified with dialysis.

### 13.4. Conclusions

The activity tests showed that mismatch cleavage activity could be detected in celery extracts purified with Amicon Ultra (0.5mL, 10k) centrifugal filter devices (and omitting the dialysis step). The whole isolation procedure could be carried out within 1 day using standard laboratory equipment (i.e. a cooled microcentrifuge). However, a stronger background on the agarose gels (possibly originating from salt remnants retained in the filter devices) is an issue. Number of reactions (obtained from 18 mL celery juice and using 4 Amicon Ultra filter devices): we have recovered a total volume of 155 uL from the 4 filter devices. The activity assay shows a clear cleavage pattern with the 5x concentrated enzyme (0.125 uL per reaction). This would allow a total of at least 1240 reactions. However, a lower amount of enzyme (between 1x and 5x) seems to work either and would increase the number of reactions accordingly.

### 13.5. Contributors

Experimental design, data interpretation: Bernhard Hofinger and Bradley Till  
 Experimental execution: Bernhard Hofinger  
 Celery enzyme extraction: Bernhard Hofinger, Owen Huynh, Biguang Huang, Bradley Till  
 Manuscript preparation and editing: Bernhard Hofinger and Bradley Till  
 Development of TILLING protocol: Owen Huynh, Bradley Till

## 14. A PROTOCOL FOR VALIDATION OF DOUBLED HAPLOID PLANTS BY ENZYMATIC MISMATCH CLEAVAGE

### 14.1. Abstract

Doubled haploidy is an important tool for plant breeders. It provides a rapid means of developing recombinant populations that are homozygous and therefore genetically fixed. Homozygosity is also important in plant mutation breeding where many induced mutations are predicted to be recessive and mutant alleles need to be in a homozygous state before new traits are expressed. While production of doubled haploids has been described for many plant species, efficient means to validate that produced materials are indeed homozygous are needed. Polymorphism discovery methods utilizing enzymatic mismatch cleavage are ideally suited for validation of doubled haploid plants. We describe here a low-cost protocol that utilizes self-extracted single-strand specific nucleases, standard PCR and agarose gels that can be applied to most plant species.

### 14.2. Introduction

First reported in the early 1920s, methods for the production of haploid plants have now been described for more than 250 species (BLAKESLEE *et al.* 1922; MALUSZYNSKI *et al.* 2003). In many cases, haploid plants either spontaneously become diploid or this state can be induced with the treatment of chemicals such as colchicine. It remains a popular and powerful tool in plant sciences and breeding because once plants are doubled haploid, they are homozygous, genetic variants are fixed, and plants are true-breeding. A wide range of methods have been described for producing doubled haploid plants, and efficiencies can vary dramatically, with less than 10% haploid embryo formation in some cases (MALUSZYNSKI *et al.* 2003; FORSTER *et al.* 2007; GEIGER and GORDILLO 2009; RAVI and CHAN 2010). For successful and efficient research and breeding it is therefore necessary to validate that materials are truly doubled haploid and homozygous before plants are selected for further evaluation and use.

Enzymatic mismatch cleavage for discovery and genotyping of single nucleotide polymorphisms (SNPs) and small insertion/deletions (indels) is ideally suited to detect loss of heterozygosity in doubled haploid plants. Widely used in Targeting Induced Local Lesions IN Genomes (TILLING) reverse-genetic studies, the procedure begins with the selection of gene-specific primers for PCR amplification of ~1 to 1.5 kb regions. This is followed by denaturation and annealing to create heteroduplexed DNA in samples where heterozygous variation exists. Samples are then treated with a self-prepared enzyme mixture containing single-strand specific nucleases that cleaves DNA at mismatched regions. The resulting products are electrophoresed and the presence of cleaved DNA fragments of lower molecular

weight than the original PCR product indicates the presence and approximate location of heterozygous polymorphisms (TILL *et al.* 2006c). The process can be made low-cost by self-extraction of nuclease and the use of standard agarose gel electrophoresis as a gel readout platform (Springer book and Huynh *et al.*). The use of enzymatic mismatch cleavage has also been widely used in EcoTILLING studies for discovery of natural nucleotide variation in populations (TILL 2014). This has been shown to be highly accurate with less than 5% false discovery and false negative error rates, even in highly heterozygous polyploids (TILL *et al.* 2006a; TILL *et al.* 2010).

The same methods used for TILLING and EcoTILLING can be easily adapted to evaluate homozygosity in putatively doubled haploid plants. As a proof of principle, Hofinger and colleagues adapted this method for barley and compared it to another approach for validation of doubled haploid (DH) plants: Simple Sequence Repeat (SSR) markers (HOFINGER *et al.* 2013). In this work the authors showed that 11/26 primer pairs were suitable for DH validation by enzymatic mismatch cleavage, while only 3/32 previously characterized SSR primer sets were suitable. Thus the enzymatic mismatch cleavage approach was ~4.5x more efficient. Furthermore, because enzymatic mismatch cleavage utilizes gene or target-specific primers, the approach allows for accurate DH production error estimations per plant and per experiment taking into account genetic linkage. In addition, using gene-specific primers allows selection of specific alleles and for diversification of allele types in applications such as those involving F1 hybrids. The entire process can be completed in one day (Figure 1). Enzymatic mismatch cleavage has been used in over 20 plant and animal species for TILLING and EcoTILLING and thus this adaptation for DH validation is expected to be broadly applicable (KUROWSKA *et al.* 2011; TILL 2014). Indeed, in addition to barley we show data for screening of putative DH *Eragrostis tef* plants (Figure 2, Gugsa *et al.* 2006). While the methods described in this protocol are straight-forward and low-cost, a key component of successful application is proper experimental design. Example data and a discussion of experimental design are provided along with a detailed protocol suitable for many crop species.

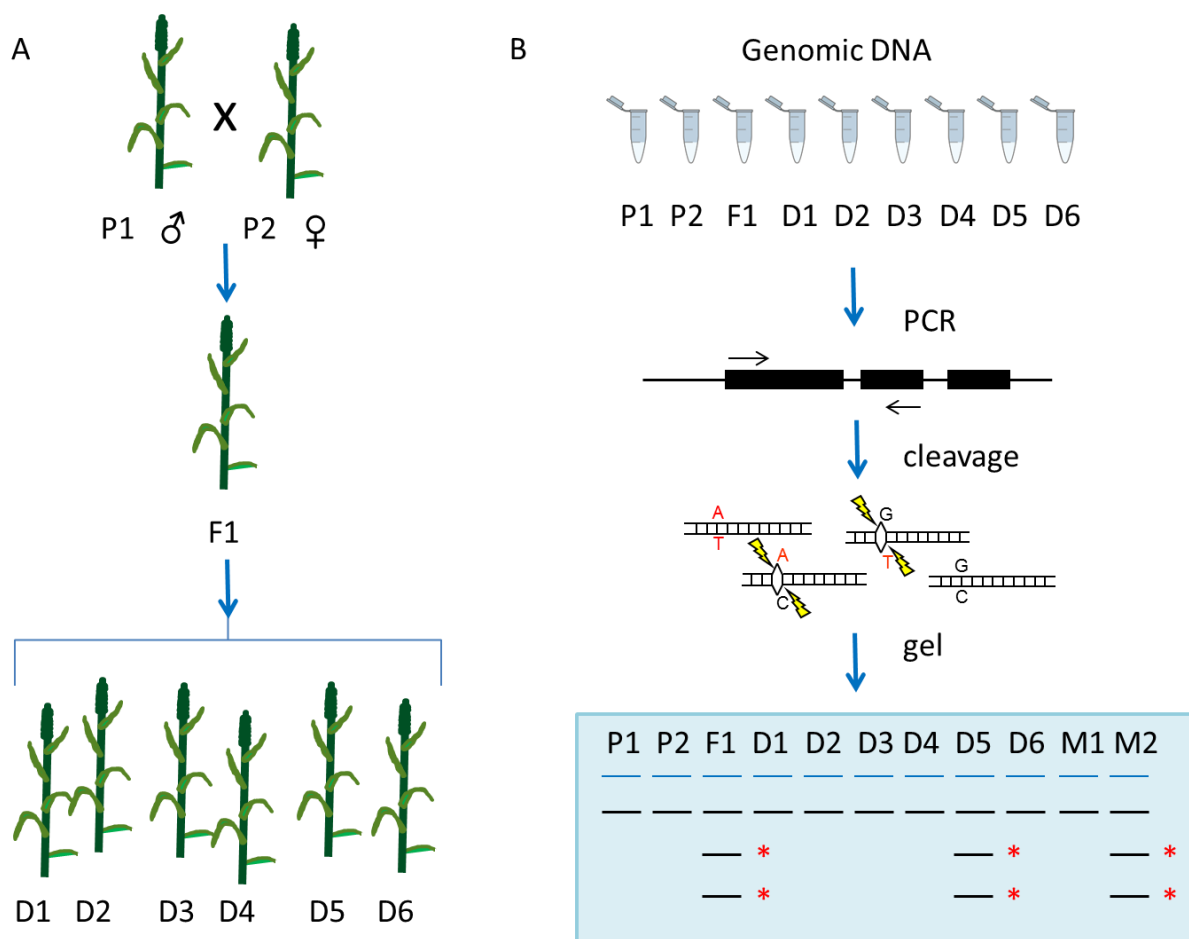


Figure 1. The use of enzymatic mismatch cleavage for validation of doubled haploid plants in F1 hybrid studies. A. In a typical F1 hybrid approach, two genetically diverse parental plants (P1 and P2) are crossed to create an F1 hybrid (F1). Putatively doubled haploids are produced from the F1 plant (D1 through D6). B. To validate the production of DH plants, genomic DNA is prepared from all materials. PCR is performed to amplify a specific target. PCR products are denatured and annealed to form mismatches where polymorphisms exist and samples are incubated with a nuclease that cleaves DNA that is not base-paired. DNA fragments are evaluated by agarose gel electrophoresis. Suitable primer pairs show cleavage products indicating heterozygosity in the F1 sample (marked by asterisks). In the example drawn, all putatively DH plants are homozygous except sample D5. Lanes marked M1 and M2 represent mixtures of sample D1 and P1 and D1 and P2 (respectively). This allows assignment of the parental allele in the DH line. In this example D1 contains the allele from P1 because there was no heterozygosity observed in the mixture of the two samples. More primer pairs should be screened to increase confidence that the remaining plants are truly DH.

## 14.3. Materials

### 14.3.1. PCR amplification

1. Taq polymerase & Taq polymerase buffer (See note 1)
2. dNTP mix
3. MgCl<sub>2</sub> (25 mM stock), if not supplied in Taq polymerase buffer
4. Forward and reverse primers (See Note 2)
5. Laboratory grade water (distilled or deionized and autoclaved)
6. Thermalcycler with heated lid

### 14.3.2. Enzymatic mismatch cleavage

1. MgSO<sub>4</sub>
2. Triton X-100
3. HEPES (4-(2-hydroxyethyl)-1-piperazineethanesulfonic acid )
4. BSA (bovine serum albumen)
5. KCl
6. 10x cleavage buffer (prepare a stock solution of: 5 ml 1 M MgSO<sub>4</sub> ,100 µl 10% Triton X-100, 5 ml 1M HEPES (pH 7.4), 5 µl 20 mg/mL bovine serum albumen, 2.5 ml 2M KCl, 37.5 ml water
7. Single-strand-specific nuclease (See Note 3).
8. EDTA (Ethylenediaminetetraacetic acid), 0.25 M.

### 14.3.3. Agarose gel electrophoresis

1. Standard agarose (e.g. Sigma A9539)
2. TBE buffer
3. Ethidium bromide (10 mg/ml) (See note 4)
4. Horizontal gel electrophoresis system with power supply
5. Gel photography system
6. Gel loading dye (30% glycerol plus xylene cyanol) (See note 5)
7. DNA molecular weight ladder (e.g. 1 kb ladder, Invitrogen 10787-018)

## 14.4. Methods

### 14.4.1. PCR amplification

1. Select samples for analysis based on guidelines for experimental design (see Note 6).

2. Prepare a master mix according to manufacturer's guidelines for the Taq polymerase used. For example for 10 samples with Takara Ex Taq: 109.5  $\mu\text{L}$  water, 20  $\mu\text{L}$  10 $\times$  Ex Taq buffer, 16  $\mu\text{L}$  2.5 mM dNTP mix, 2  $\mu\text{L}$  Forward primer (10  $\mu\text{M}$ ), 2  $\mu\text{L}$  reverse primer (10  $\mu\text{M}$ ), 0.5  $\mu\text{L}$  TaKaRa HS taq (5 U/ $\mu\text{L}$ ).
3. Add 5  $\mu\text{L}$  of DNA to PCR tubes (see note 7).
4. Add 15  $\mu\text{L}$  of PCR mastermix to PCR tubes and mix by pipetting.
5. Centrifuge tubes for 1 minute at 5000 x g.
6. Incubate samples in thermal cycler with following conditions: Step 1, Initial denaturation, 95°C, 2 minutes, Step 2, Denaturation, 94°C, 20 seconds, Step 3, Primer annealing, 73°C (-1°C/cycle), 30 seconds, Step 4, Ramp, 0.5°C per second to 72°C, Step 5, Primer extension, 72°C, 1 minute, Step 6, Cycling, repeat steps 2-5 for 7 cycles (8 cycles in total), Step 7, Denaturation, 94°C, 20 seconds, Step 8, Primer annealing, 65°C, 30 seconds, Step 9, Ramp, 0.5°C per second to 72°C, Step 10, Primer extension, 72°C, 1 minute, Step 11, Cycling, repeat steps 7-10 for 44 cycles (45 cycles in total), Step 12, Final extension, 72°C, 5 minutes, Step 13, Denaturation, 99°C, 10 minutes, Step 14, Cooling, 72°C, 20 seconds, Step 15, Cycling & Touchdown, repeat step 14 for 69 cycles (-0.3°C/ cycle), Step 16, Hold, 8°C, forever (see note 8).
7. Samples can be stored at -20°C for months.

#### 14.4.2. Enzymatic mismatch

1. Prepare the following enzyme digestion master mix: 115  $\mu\text{L}$  water, 35  $\mu\text{L}$  10x cleavage buffer, 25  $\mu\text{L}$  single-strand-specific nuclease (see note 3)
2. Centrifuge PCR products for 1 minute at 5000 x g.
3. Add 20  $\mu\text{L}$  of enzyme digestion master mix to PCR reactions.
4. Incubate reactions in a thermal cycler at 45°C for 15 minutes (no heated lid) and then hold at 8°C.
5. Add 6  $\mu\text{L}$  of 0.25 M EDTA to each tube to stop the reaction.

#### 14.4.3. Agarose gel electrophoresis and data analysis

1. Prepare a 1.5% agarose gel in TBE buffer. Heat mixture in microwave until all agarose is dissolved. Take care to avoid boiling agarose.
2. Cool agarose to approximately 55°C (until you can hold the flask in your hand for 5 seconds without burning your hand).
3. Add ethidium bromide to the warm gel solution to obtain a final concentration of 0.5  $\mu\text{g}/\text{mL}$  in the gel. Mix by swirling.
4. Pour gel solution into gel apparatus and insert comb according to the instructions of the manufacturer.
5. Add 2  $\mu\text{L}$  6 $\times$  loading dye to 10  $\mu\text{L}$  PCR product.
6. Run gel at 100 V for 1 hour or until suitable separation of bands is achieved.
7. Photograph the gel under UV light.
8. Evaluate data (Figures 2, 3 and Note 9).

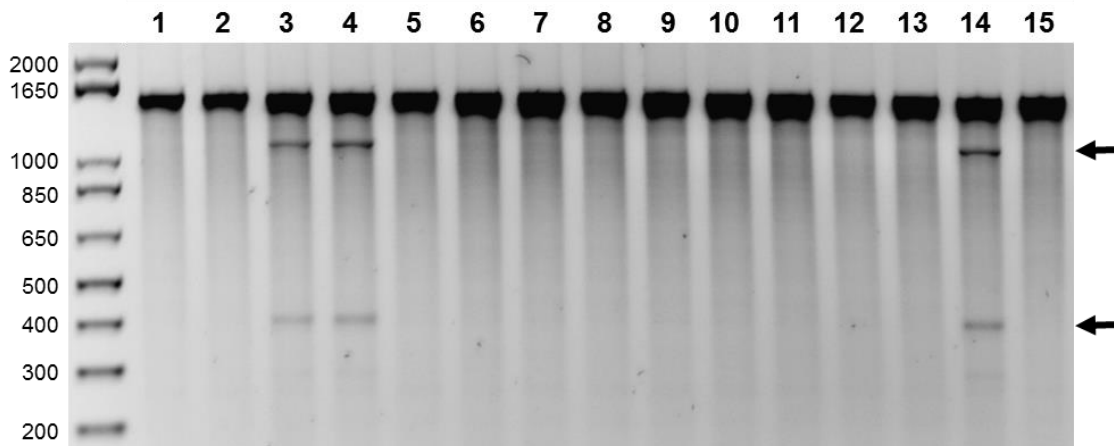


Figure 2: Example data of validation of doubled haploid production in barley. A 1476 bp fragment of the barley *Mlo9* gene was PCR amplified and digested with a crude celery juice extract (CJE) containing single-strand-specific nuclease activity followed by agarose gel analysis. The top band in lanes 1-15 represents undigested PCR product. The cleavage products of DNA-heteroduplexes exclusively present in heterozygous samples are marked with arrows. Parental lines Golden Promise (GP) and HOR1606 are homozygous for this gene region (lanes 1 and 2 respectively). A synthetic mixture of parental DNA and also the F1 sample from crossing of the two parents show cleavage fragments resulting from a heterozygous SNP (lanes 3 & 4). Doubled haploid plants (lanes 5-13) are homozygous. Mixtures of genomic DNA from a DH plant and GP show cleavage products while mixture of the same material with HOR1606 does not, indicating the DH harbors the GP allele (lanes 14 & 15). This figure and legend reproduced from (Hofinger et al. 2013).

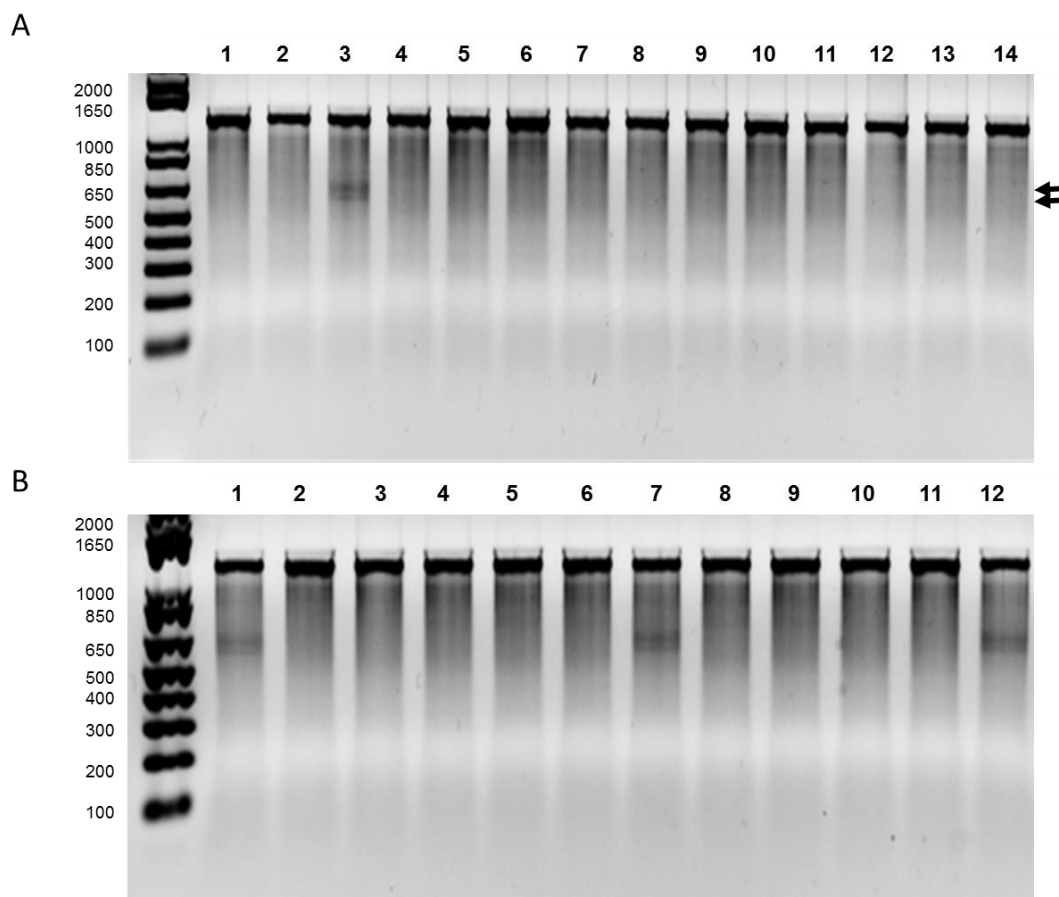


Figure 3. Example data of validation of doubled haploid production in Tef. Primers were designed to amplify a 1400 bp region of the Tef genome. Screening for loss of heterozygosity was first performed with putative doubled haploid plants. The top band in all of the numbered lanes of both panels represents undigested PCR product. Panel A: Lanes 1, and 2 are parents. Lane 3 represents a synthetic F1 hybrid created by mixing equal amounts of genomic DNA from each parent. Two cleavage products are produced indicating heterozygosity in the F1 (marked by arrows). Lanes 4-14 represent putative DH plants that are confirmed as being homozygous for the amplified genomic region, since no heteroduplex cleavage product does appear. Experiments were then performed to determine the parental origin of the allele in putative DH plants. Panel B shows mixtures of samples with parent 1 (lane 1 of panel A). The first lane of panel B is the positive control involving a mixture of both parents. Lanes 2 -12 are putative DH plants mixed with parent 1. This shows that two plants (lanes 7 & 12) inherited the allele from parent 2 (forming cleaved heteroduplexes with the one from parent 1), while all others inherited the allele from parent 1 so that no heteroduplex can be formed and cleaved). The reciprocal experiment was done by mixing samples with parent 2 (not shown).

## 14.5. Notes

1. Most Taq polymerases should be suitable for this method. Compare your favourite Taq with the least expensive you can purchase. If the two produce similar results, use the cheaper version. Some Taq polymerases, like Takara ExTaq come supplied with dNTPs and buffer containing MgCl<sub>2</sub>.
2. Primer pairs should be selected to different genomic regions that are not genetically linked. The total number of primer pairs needed depends on the percentage of primer pairs



where heterozygosity is observed. Primers are typically designed with the program Primer3 (Rozen and Skaletsky 2000) with a melting temperature of 70°C. See note 8.

3. The optimal amount of nuclease varies depending on the activity of the batch being used. Analysis of activity can be performed by screening mutant DNA (mixed with wild-type DNA) and titrating the amount of enzyme to produce clear bands on the gel (see chapter X of this book for a protocol on enzyme preparation and activity optimization).
4. Caution. Ethidium bromide is hazardous. Wear gloves and avoid contamination. Consult Material Safety Data Sheet (MSDS) for proper handling and disposal procedures.
5. Avoid loading dyes containing Bromophenol blue or other dyes that migrate in molecular weight range where you expect to observe DNA fragments. The presence of loading dyes can reduce the intensity of bands.
6. The optimum experimental design includes DNA from the progenitor plant along with DNA from each putatively DH plant produced from the progenitor. The only useful primer pairs will be those where heterozygosity is discovered in the progenitor plant. For example, if making doubled haploids from an F1 hybrid, material from the F1 hybrid is ideally evaluated along with DNA from the parents that were used to make the F1. If the F1 harbours heterozygous SNPs at a particular locus, and the putative DH plants do not, this is evidence that the plant is DH (see note 9 for more on data analysis and interpretation). In this example, if F1 material is not available, a synthetic F1 sample can be prepared by mixing an equal concentration of DNA from each parent in a 1:1 ratio prior to PCR amplification. Screening the parental material alone is informative to learn if parents are heterozygous in any interrogated regions. It is not ideal to screen only putative DH material as it is difficult to estimate the probability that plants are truly DH rather than being homozygous because progenitor material was homozygous at that locus.
7. The optimal amount of genomic DNA to be used should be determined empirically. A PCR product yield of 10 ng/μL is typically sufficient. The amount of genomic DNA needed to produce this amount of product can be roughly estimated by size of the genome (Till et al. 2006c). The yield of PCR product should be sufficiently high to produce cleavage products visible by agarose gel electrophoresis (see figures 2 & 3).
8. PCR conditions may need to be optimized. For example, primers with a melting temperature ( $T_m$ ) of 70°C were used to develop this protocol. Higher  $T_m$  primers increase specificity of amplification, but may not be necessary for all species. If lower  $T_m$  primers are used, the annealing temperature must be adjusted accordingly.
9. When analysing data, fragments observed of lower molecular weight than the full-length PCR product are typically the result of cleavage of duplexed DNA at the site of a mismatch due to a nucleotide polymorphism. Truly doubled haploid plants are homozygous and therefore should not show no cleavage products. However, cleavage fragments can also be observed due to a homopolymeric stretch of adenosine residues (Till et al. 2004a). Evaluation of the parental material, which is typically homozygous, is therefore advised. Mixtures of parental material with putative DH plants allow assignment of alleles to a specific parent (see figure 2 and figure 3B). An estimation of the number of suitable primers can be made by calculating the probability that the data results by chance from self-fertilization. For example, the probability that the offspring of a self-fertilized heterozygous F1 is homozygous for a specific locus is 0.5 (assume a Mendelian 1:2:1 ratio). The probability that two genetically unlinked loci are homozygous is  $0.5^2 = 0.25$ . By screening seven primer pairs from unlinked loci that show heterozygosity in the F1,

one can achieve 99% confidence. Such estimations become impossible if parental or F1 material is not available for screening.

## 14.6. Acknowledgments

Funding for this work was provided by the Food and Agriculture Organization of the United Nations and the International Atomic Energy Agency through their Joint FAO/IAEA Programme of Nuclear Techniques in Food and Agriculture. This work is part of IAEA Coordinated Research Project D24012.

## 14.7. Contributors

Bernhard Hofinger, Owen A. Huynh, Joanna Jankowicz-Cieslak, and Bradley J. Till

## 14.8. References

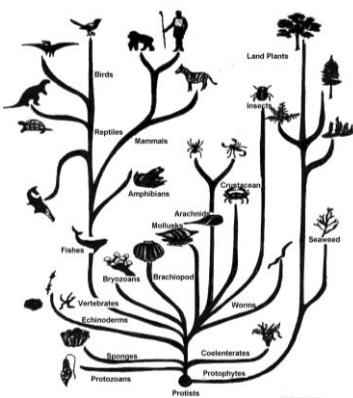
- Gugsa L, Sarial AK, Lörz H, Kumlehn J: Gynogenic plant regeneration from unpollinated flower explants of *Eragrostis tef* (Zuccagni) Trotter. *Plant Cell Reports* 2006, **25**:1287-1293.
- Blakeslee, A. F., J. Belling, M. E. Farnham and A. D. Bergner, 1922 A Haploid Mutant in the Jimson Weed, "Datura Stramonium". *Science* 55: 646-647.
- Burdon, M. G., and J. H. Lees, 1985 Double-strand cleavage at a two-base deletion mismatch in a DNA heteroduplex by nuclease S1. *Biosci Rep* 5: 627-632.
- Chaudhry, M. A., and M. Weinfeld, 1995 Induction of double-strand breaks by S1 nuclease, mung bean nuclease and nuclease P1 in DNA containing abasic sites and nicks. *Nucleic Acids Res* 23: 3805-3809.
- Colbert, T., B. J. Till, R. Tompa, S. Reynolds, M. N. Steine *et al.*, 2001 High-throughput screening for induced point mutations. *Plant Physiol* 126: 480-484.
- Comai, L., K. Young, B. J. Till, S. H. Reynolds, E. A. Greene *et al.*, 2004 Efficient discovery of DNA polymorphisms in natural populations by EcoTilling. *Plant J* 37: 778-786.
- Forster, B. P., E. Heberle-Bors, K. J. Kasha and A. Touraev, 2007 The resurgence of haploids in higher plants. *Trends Plant Sci* 12: 368-375.
- Galeano, C. H., M. Gomez, L. M. Rodriguez and M. W. Blair, 2009 CEL I Nuclease Digestion for SNP Discovery and Marker Development in Common Bean (*Phaseolus vulgaris* L.). *Crop Science* 49: 381-394.
- Garvin, M. R., and A. J. Gharrett, 2007 DEco-TILLING: an inexpensive method for single nucleotide polymorphism discovery that reduces ascertainment bias. *Molecular Ecology Notes* 7: 735-746.

- Geiger, H. H., and G. A. Gordillo, 2009 Doubled haploids in hybrid maize breeding. *Maydica* 54: 485-499.
- Hofinger, B. J., O. A. Huynh, J. Jankowicz-Cieslak, A. Muller, I. Otto *et al.*, 2013 Validation of doubled haploid plants by enzymatic mismatch cleavage. *Plant Methods* 9: 43.
- Howard, J. T., J. Ward, J. N. Watson and K. H. Roux, 1999 Heteroduplex cleavage analysis using S1 nuclease. *Biotechniques* 27: 18-19.
- Kurowska, M., A. Daszkowska-Golec, D. Gruszka, M. Marzec, M. Szurman *et al.*, 2011 TILLING - a shortcut in functional genomics. *Journal of Applied Genetics* 52: 371-390.
- Maluszynski, M., K. J. Kasha, B. P. Forster and I. Szarejko, 2003 *Doubled haploid production in cop plants. A manual*. Kluwer Academic Publishers, Dordrecht, Boston, London.
- McCallum, C. M., L. Comai, E. A. Greene and S. Henikoff, 2000 Targeted screening for induced mutations. *Nat Biotechnol* 18: 455-457.
- Perry, J. A., T. L. Wang, T. J. Welham, S. Gardner, J. M. Pike *et al.*, 2003 A TILLING reverse genetics tool and a web-accessible collection of mutants of the legume *Lotus japonicus*. *Plant Physiol* 131: 866-871.
- Ravi, M., and S. W. Chan, 2010 Haploid plants produced by centromere-mediated genome elimination. *Nature* 464: 615-618.
- Rozen, S., and H. Skaletsky, 2000 Primer3 on the WWW for general users and for biologist programmers. *Methods Mol Biol* 132: 365-386.
- Sato, Y., K. Shirasawa, Y. Takahashi, M. Nishimura and T. Nishio, 2006 Mutant Selection from Progeny of Gamma-ray-irradiated Rice by DNA Heteroduplex Cleavage using Brassica Petiole Extract. *Breeding Science* 56: 179-183.
- Slade, A. J., S. I. Fuerstenberg, D. Loeffler, M. N. Steine and D. Facciotti, 2005 A reverse genetic, nontransgenic approach to wheat crop improvement by TILLING. *Nat Biotechnol* 23: 75-81.
- Sokurenko, E. V., V. Tchesnokova, A. T. Yeung, C. A. Oleykowski, E. Trintchina *et al.*, 2001 Detection of simple mutations and polymorphisms in large genomic regions. *Nucleic Acids Res* 29: E111.
- Till, B. J., 2014 Mining Genetic Resources via Ecotilling, pp. 349-365 in *Genomics of Plant Genetic Resources*, edited by R. Tuberosa, A. Graner and E. Frison. Springer Netherlands.
- Till, B. J., C. Burtner, L. Comai and S. Henikoff, 2004a Mismatch cleavage by single-strand specific nucleases. *Nucleic Acids Res* 32: 2632-2641.
- Till, B. J., J. Jankowicz-Cieslak, L. Sagi, O. A. Huynh, H. Utsushi *et al.*, 2010 Discovery of nucleotide polymorphisms in the *Musa* gene pool by Ecotilling. *Theor Appl Genet* 121: 1381-1389.
- Till, B. J., S. H. Reynolds, E. A. Greene, C. A. Codomo, L. C. Enns *et al.*, 2003 Large-scale discovery of induced point mutations with high-throughput TILLING. *Genome Res* 13: 524-530.
- Till, B. J., S. H. Reynolds, C. Weil, N. Springer, C. Burtner *et al.*, 2004b Discovery of induced point mutations in maize genes by TILLING. *BMC Plant Biol* 4: 12.
- Till, B. J., T. Zerr, E. Bowers, E. A. Greene, L. Comai *et al.*, 2006a High-throughput discovery of rare human nucleotide polymorphisms by Ecotilling. *Nucleic Acids Res* 34: e99.

- Till, B. J., T. Zerr, L. Comai and S. Henikoff, 2006b A protocol for TILLING and Ecotilling in plants and animals. *Nat Protoc* 1: 2465-2477.
- Till, B. J., T. Zerr, L. Comai and S. Henikoff, 2006c A protocol for TILLING and Ecotilling in plants and animals. *Nature Protocols* 1: 2465-2477.
- Zerr, T., and S. Henikoff, 2005 Automated band mapping in electrophoretic gel images using background information. *Nucleic Acids Res* 33: 2806-2

## 15. MULTIVARIATE ANALYSIS – PHYLOGENETICS AND PRINCIPAL COMPONENT ANALYSIS

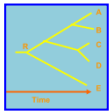
### 15.1. Phylogenetics



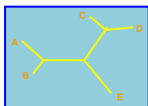
Phylogenetics in the plant kingdom is based on genetic information from accessions. The entities whose affinities are studied are called operational taxonomic units (OTUs, anything from a population to a phylum, including sequence variation and other polymorphisms). Phylogenetics studies the evolutionary relatedness among OTUs using genetic information and is mostly based on genetic distances calculations. The results of these calculations are often synoptically presented as a phylogenetic tree (rooted) or dendrogram (unrooted). There are many methods using different models and assumptions on which the genetic distances calculations are based and ultimately the phylogenetic tree. It is important to understand from the outset what model and *a priori*

assumptions to apply in order to be able to infer valuable information from the raw data to be mined.

There are two different tree types that might be constructed, based on two different purposes in analysing the raw data:



Rooted trees serve to unfold an evolutionary path



Un-rooted trees (dendrograms) are used to visualize relationships

A multitude of tree reconstruction algorithms are available. These can be roughly classified into 4 methods:

- Distance Matrix, based on pairwise evolutionary distances (*e.g.* UPGMA, Neighbour Joining)
- Maximum Parsimony, based on the shortest pathway to the present character state
- Maximum Likelihood, based on choosing the tree with the largest ML value of the character state presented
- Invariants, based on functions of characters that have an expected value of 0 in some trees and non-zero expectation in other trees.

## 15.2. Inferring phylogeny from pairwise distances: construction of a distance tree using clustering with the unweighted pair group method with arithmetic mean (UPGMA).

There are mainly two multivariate methods widely used for pattern analyses of DNA genotypes in biology: principal component analysis (PCA) (Flury 1988) and cluster analysis (Everitt 1992). PCA and cluster analysis seek to uncover hidden or cryptic patterns among objects (*e.g.*, individuals, genetic stocks, or populations) on which two or more independent variables (phenotypic or genotypic characters) have been measured.

- Typical phenotypic variables are morphological traits (*e.g.*, flower petal length and width).
- Typical genotypic variables are DNA marker genotypes or allele sequences. A variety of DNA markers can be employed for genotyping or DNA fingerprinting.

PCA and cluster analysis seek to project multivariate phenotypic or genotypic measurements in lower dimensional spaces so that the underlying patterns or structures can be described and visually displayed. The ‘genetic’ patterns among a set of OTUs (entities, genetic materials) usually cannot be directly discerned from DNA fingerprints (raw multivariate data); however, patterns among the OTUs can nearly always be ‘extracted’ by PCA or cluster analyses of pairwise genetic distance matrices.

Originally developed for constructing taxonomic phenograms, *i.e.* trees that reflect the phenotypic similarities between OTUs, UPGMA is the simplest method of tree construction, if the rates of evolution are approximately constant among the different lineages. For this purpose the number of observed nucleotide or amino-acid substitutions can be used.

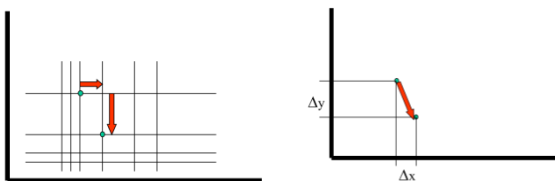
## 15.3. Distance measures

Distance measures are based on topology paths in n-dimensional space. As an example in a two dimensional space we might consider the following:

Travel in a grid *versus* shortest direct distance

$$d = \min(\sum x_i + \sum y_j)$$

$$d = \text{SQR}(\Delta x^2 + \Delta y^2)$$



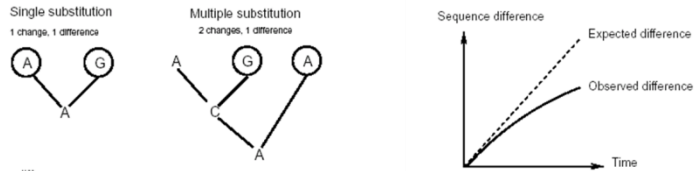
In the context of plant production and protection, the choice of genetic distance estimators depends on what we want to do, what we want to see, what precision of their estimations is needed and the conditions of their applications (in terms of type of markers, genetic structure of the cultivars/accessions/individuals, diversity of reference collections, breeding programmes *etc.*). This defines the dimensions and topologies of the space we are exploring and the paths in this space. Let us construct the following set-up to illustrate the utmost importance of the choice of a genetic distance estimator (*i.e.* it should not be chosen uniquely given the availability of a computer programme): a “naïve” measure of genetic similarity or

measure of genetic distance is the Hamming distance where  $d_0$  = proportion of sites at which two sequences differ:

Sorghum TGTATCGCTC...  
Sugarcane TGTGTCGCTC...

Sorghum TGTATCGCTC...  
Rice AGTCTCGTTC...

Sugarcane TGTGTCGCTC...  
Rice AGTCTCGTTC...



The Hamming Distance is a poor measure of the actual number of evolutionary changes, as a site can undergo repeated substitutions. It might be appropriate for short periods and/or parental inferences.

In order to define a genetic distance estimator, we have to assay the genetic similarities of the entities we are studying. Let these entities be dominant markers (present-absent characters): the genetic similarity between the  $i$ th and  $j$ th entity is  $s_{ij}$ . As such, genetic similarity coefficients are symmetric ( $s_{ij} = s_{ji}$ ), positive and bound by  $1 (0 \leq s_{ij} \leq 1)$ . Two individuals are completely identical, when  $s_{ij} = 1$  and completely different when  $s_{ij} = 0$ . Genotypic scores and counts for a binary variable (dominant marker):

entity i	entity j	count	condition
present (1)	present (1)	$a$ ( $n_{11}$ )	positive match
present (1)	absent (0)	$b$ ( $n_{10}$ )	mismatch
absent (0)	present (1)	$c$ ( $n_{01}$ )	mismatch
absent (0)	absent (0)	$d$ ( $n_{00}$ )	negative match

The two most widely used similarity measures for binary data are the simple matching coefficient and Jaccard's coefficient.

- The simple matching coefficient is the ratio of the sum of matches to the sum of matches and mismatches:

$$s_{ij} = \frac{a + d}{a + b + c + d}$$

- Jaccard's coefficient is the ratio of positive matches to the sum of positive matches and mismatches:

$$s_{ij} = \frac{a}{a + b + c}$$

Based on defined genetic similarity coefficients, genetic distance measures can be inferred. The Euclidean genetic distance between the  $i$ th and  $j$ th entity is:  $d_{ij} = \sqrt{2(1 - s_{ij})}$ , if the genetic similarity matrix is positive semi-definite (Gower 1971). Both simple matching coefficient and Jaccard's coefficient matrices are positive semi-definite.

In linear algebra, a positive-definite matrix is a Hermitian matrix which in many ways is analogous to a positive real number. The notion is closely related to a positive-definite symmetric bilinear form. In mathematics, a definite bilinear form is a bilinear form  $B$  such that  $B(x, x)$  has a fixed sign (positive or negative) when  $x$  is not 0.

To give a formal definition: let  $K$  be one of the fields  $R$  (real numbers) or  $C$  (complex numbers). Suppose that  $V$  is a vector space over  $K$ , and  $B : V \times V \rightarrow K$  is a bilinear form which is Hermitian in the sense that  $B(x, y)$  is always the complex conjugate of  $B(y, x)$ . Then  $B$  is called positive definite if  $B(x, x) > 0$  for every nonzero  $x$  in  $V$ . If  $B(x, x) \geq 0$  for all  $x$ ,  $B$  is said to be positive semidefinite.

A Hermitian matrix (or self-adjoint matrix) is a square matrix with complex entries which is equal to its own conjugate transpose  $a_{i,j} = a_{j,i}^*$  — that is, the element in the  $i$ th row and  $j$ th column is equal to the complex conjugate of the element in the  $j$ th row and  $i$ th column, for all indices  $i$  and  $j$ . Or written with the conjugate transpose:  $A = A^\dagger$

For example,  $\begin{bmatrix} 3 & 2+i \\ 2-i & 1 \end{bmatrix}$  is a Hermitian matrix. For all non-zero  $x \in R^n$  (or, equivalently, all non-zero  $x \in C^n$ ), it is called positive-semi-definite if  $x^* M x \geq 0$ .

The three most common distance estimators which are computed throughout the majority of the literature for different purposes are: the Jaccard's distance (J) (1908), the Nei & Li's distance (NL) (1979) and the Sokal & Michener's distance (SM) (1958):

$$J_{xy} = 1 - (n_{11} / (n_{11} + n_{10} + n_{01})) \quad [1]$$

$$NL_{xy} = 1 - ((2 \times n_{11}) / ((2 \times n_{11}) + n_{10} + n_{01})) \quad [2]$$

$$SM_{xy} = 1 - ((n_{11} + n_{00}) / (n_{11} + n_{10} + n_{01} + n_{00})) \quad [3]$$

where  $n_{11}$  is the number of bands shared by the individuals (cultivars, clones accessions etc.)  $x$  and  $y$  tested (i.e. positive matching between pairs),  $n_{10}$  is the number of bands present in  $x$  and absent in  $y$ ,  $n_{01}$  the number of bands present in  $y$  and absent in  $x$ , and  $n_{00}$  the number of bands absent both in  $x$  and  $y$  (i.e. negative matching). In addition, one may also, using the inverse of the PIC (polymorphism information content of a certain marker), compute a weighted Jaccard's distance (WJ) to take into account the frequency of each marker in the calculation of the distance.

$$PIC = 1 - \sum_{i=1}^n P_i^2 - \sum_{i=1}^n \sum_{j=i+1}^n P_i^2 P_j^2 \quad [4]$$

$P_i$  = frequency of allele  $i$  from 1 to  $n$

This formula produces an indicator of how many alleles a certain marker has and how much these alleles divide evenly. For example if a marker has few alleles, or if the marker has many alleles but only one of them is frequent, the PIC will be low. Obviously:

$$1 = \sum_{i=1}^n P_i \quad [5]$$

The Nei & Li genetic distance estimator was developed for the analysis of restriction site polymorphisms, and is the estimator proposed by Dice (1945) in the pre-molecular era:



$D_{ij} = 2N_{ij}/(N_i + N_j)$ , where  $N_{ij}$  is the number of restriction sites or restriction fragments shared by  $i$  and  $j$  ( $= n_{11}$ ),  $N_i$  is the number of restriction fragments in  $i$  ( $n_{11} + n_{10}$ ), and  $N_j$  is the number of restriction fragments in  $j$  ( $n_{11} + n_{01}$ ). This estimator excludes negative matches.

The simple matching coefficient and Jaccard's coefficient differ in how negative matches (0-0 matches or d counts) are handled. The problem of whether to include or exclude negative matches only arises for present-absent characters (binary or categorical variables), e.g., binary genetic markers with null alleles.

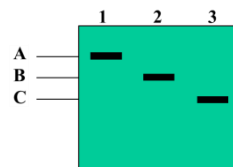
The question as to whether two individuals are similar when they both lack a character does not always have a simple answer. This topic has been hotly debated, particularly in taxonomic circles (Romesburg 1984; Sneath and Sokal 1973). When one allele is absent (null) and the other is present and both alleles are observed among the entities sampled, Dudley (1993) argued that 0-0 matches should be included because the absence of an allele in two entities measures similarity. This may or may not be true. Two individuals, for example, may lack an AFLP band; however, the mutations that abolished the AFLP band in the two individuals could be different (mutation in the restriction sites = elimination of sites, insertion between restriction sites = band too long to amplify, deletion between restriction sites = smaller band appearing but too small to be scored, translocation = reshuffling restriction sites), in which case the two individuals carry different null alleles and the 0-0 score is incorrect. But the probability of these events locus by locus depends on the frequency of these events, and the probability of loss of band due to different mutation events decreases with increasing relatedness. In fact, including 0-0 matches increases homoplasy: loci identical by state but not identical by descent. Thus, when estimated from multiallelic markers, genetic similarities may be upwardly biased by including negative matches, particularly when one or more alleles are rare.

Negative matches should be excluded for multiallelic, co-dominant markers with no null alleles, otherwise, similarities are overestimated. In the following, an example illustrating this will be detailed:

Suppose three lines are genotyped for a locus with three codominant alleles and each line is homozygous for a different allele

Entity	Allele 1	Allele 2	Allele 3	Genotype
1	1	0	0	1
2	0	1	0	2
3	0	0	1	3

(1 = present, 0 = absent)



Now, Gower (1971) proposed a similarity measure for cases where mixed variable types are measured (e.g., mixtures of binary, ordinal, categorical, and continuous variables). This coefficient can be used, for example, to combine dominant (binary) and multiallelic, co-dominant (categorical) DNA markers or discrete genotypic and continuous phenotypic variables and is one of several similarity measures used in genetic pattern analysis. Gower's coefficient and Jaccard's coefficient are the same when the former is estimated from binary

variables and negative matches are excluded. We can use this to illustrate whether negative matches should be included or not.

Gower's coefficient is:

$$s_{ij} = \frac{\sum_{k=1}^m w_{ijk} \times s_{ijk}}{\sum_{k=1}^m w_{ijk}}$$

where the similarity between the  $i$ th and  $j$ th entity measured on the  $k$ th variable is  $s_{ijk}$ , the weight for the  $k$ th variable measured on the  $i$ th and  $j$ th entity is  $w_{ijk}$ ,  $i = 1, 2, \dots, n$ ,  $j = 1, 2, \dots, n$ ,  $n$  is the number of entities,  $k = 1, 2, \dots, m$ , and  $m$  is the number of variables (DNA fragments or bands). The variable weight is either 0 or 1 and is used to include or exclude negative matches for binary or categorical variables (genetic markers). when  $k$  is unknown for one or both entities.

In our example, if we exclude 0-0 matches:

Outcome	Entity $i$	Entity $j$	$s_{ijk}$	$w_{ijk}$
if positive match	1	1	1	1
if mismatch $i - j$	1	0	0	1
if mismatch $i - j$	0	1	0	1
if negative match	0	0	1	0

$s_{12} = ((0 \times 1) + (0 \times 1) + (1 \times 0))/(1 + 1 + 0) = 0/2 = 0$   
 $s_{13} = ((0 \times 1) + (1 \times 0) + (0 \times 1))/(1 + 0 + 1) = 0/2 = 0$   
 $s_{23} = ((1 \times 0) + (0 \times 1) + (0 \times 1))/(0 + 1 + 1) = 0/2 = 0$

Now, if we include 0-0 matches:

Outcome	Entity $i$	Entity $j$	$s_{ijk}$	$w_{ijk}$
if positive match	1	1	1	1
if mismatch $i - j$	1	0	0	1
if mismatch $i - j$	0	1	0	1
if negative match	0	0	1	1

$s_{12} = ((0 \times 1) + (0 \times 1) + (1 \times 1))/(1 + 1 + 1) = 1/3$   
 $s_{13} = ((0 \times 1) + (1 \times 1) + (0 \times 1))/(1 + 1 + 1) = 1/3$   
 $s_{23} = ((1 \times 1) + (0 \times 1) + (0 \times 1))/(1 + 1 + 1) = 1/3$

The genetic similarities among the lines (considering the one locus only) are 0.00; however, if negative matches are included, then the genetic similarities are 0.33.

Obviously this is a "demonstration by the absurd": we have a population of 3 entities, genotyping is based on 1 co-dominant locus, there are only 3 alleles in our population, allele frequency of all the alleles is identical in our population, and we are sure that there is no null allele, further all individuals are homozygotes. Obviously genetic similarities should be 0. In our thought experiment, to include 0-0 matches is wrong.

Unfortunately, in "real" life, matters are not so easy. But our example shows, what questions we have to answer before deciding which model to use: heterozygosity, *a priori* knowledge of the population (structure, phylogeny), allelism (number, frequencies, null-alleles), marker system (dominant/co-dominant). Unfortunately, some of these data cannot be assessed. A fruitful approach, in my opinion, is to compare the results of different models and look for consistencies/differences, which contradict our *a priori* expectations and trying to find an explanation to these puzzles.

In some cases, the simple coefficients of correlation between these four genetic distances (J, NL, SM and WJ) may be calculated, *e.g.* to test whether there is an effect due to the choice of the distance. If the correlation is high for the six pairwise comparisons (*e.g.* over 0.9), then one might not bother about the biology, reproduction system (vegetatively *versus*. sexually propagated, auto/allogamous), ploidy, heterozygosity or population structure. One has not to forget that genetic diversity analysis is not just "number crunching": it is the knowledge of the plant biology and the characteristics of the used marker system(s) which prompts the choice, eventually the construction, of a mathematical model to analyse the data.

For example: the choice of the euclidean distance leading to Jaccard or Dice-indeces is *a priori* a model to consider when using RAPD markers. The Dice index (Jaccard, euclidean distance) is more robust against artefactual bands, but takes into account only common present bands. Now AFLP is more reproducible than RAPD, and absent bands are very significant indeed, and an algorithm such as the "simple matching algorithm", or an algorithm of Sokal and co-workers is more appropriate.

So when confronted with analysing genetic diversity, one should start by acknowledging the biological characteristics of the plant and the general taxonomy (genera, species *e.g.*) of the individuals/accessions in the study (assess the *a priori* structure of the genetic diversity of a collection of individuals, phenotyping). Then look into the characteristics of the marker system(s) used: dominant *vs.* co-dominant, PIC, reproducibility (confidence in reading the pattern, power of resolution of the analysis system, for example). This will prompt a choice of different mathematical models applicable to the problem, or even more interestingly exclude some choices.

In general: the choice of the Dice-index is at least worth a tentative first order approximation to genetic diversity analyses to sketch a rough outline of genetic diversity of the population studied. To confirm/refine this working draft (compare/oppose the *a priori* structure of genetic diversity to the one obtained using the Dice-index), one might have to use co-dominant markers to assess ploidy, heterozygosity. This might bring new insights furthering data re-analyses using more appropriate algorithms, adapted to the plant biology and/or marker characteristics, to get a better modelisation of diversity.

The sampling distributions of genetic distance estimators are not known; thus, parametric methods for estimating sampling variances and constructing confidence intervals have not been developed; however, bootstrapping or other resampling methods can be used to estimate sampling variances. Bootstrapping is done by randomly sampling data with replacement to produce individual samples from which the parameters are estimated. Suppose  $n$  individuals

were sampled from a population to estimate allele frequencies. Bootstrapping would be done by drawing  $b$  bootstrap samples of  $n$  individuals with replacement and producing  $b$  allele frequency estimates from which mean allele frequencies and sampling variance are estimated.

When constructing dendrograms bootstrapping generates multiple data sets (usually 100 random resampling iterations with replacement are sufficient, format of seed number being  $[4n+1]$ ) and adds statistical significance to the branching points in the dendrograms, which are good starting points for discussions in an article. Sometimes PCA (principal component analysis) eigenvector decomposition into major axes for 2D representation of clustering give a better synoptic background to discussions than dendrograms.

#### 15.4. Some reflexions on the comparison between genetic distances.

NL can be easily expressed as an increasing function of J ( $NL = J / [2 - J]$ ), which means that one is to expect them to be very highly correlated and lead to identical rankings of genetic distances. If this expectation is not met, this is very significant and needs to be investigated

In comparison, a high correlation between J and SM is not obvious. The difference between these distances (formula [1] and [3]) come from negative matches which are taken into account in the denominator of SM distance.

Peltier *et al.* (1994), supported that in the case of intra-specific studies, an allelic relation exists between presence and absence of a band and a negative matching is an indication of similarity and might lead to the same kind of results with SM and J.

In addition, if the weighting of Jaccard (WJ) distance by the inverse of the PIC provides similar relationships between cultivars/accessions/individuals to Jaccard ones, this might be due to the structure of the marker frequency between individuals tested. But WJ leads to take the most different individuals further away from each other, enhancing differences and might clarify

#### 15.5. What genetic distance estimator to choose for essential derivation?

In the framework of plant production and protection, the choice of the genetic distance is crucial for determining the level of relatedness between cultivars/accessions. For the distinctness and without any genetic consideration, J and NL are independent of the samples because only bands present in x and/or in y are considered. For SM, negative matches are counted and if a new cultivar/accession carries a new band absent in the previously registered ones, this becomes a new negative matching for these cultivars and the distance will change. For pragmatic reasons, the stability of genetic distance is a very attractive quality for breeders because a distance between two cultivars is constant when the number of cultivars in the reference collection increases.

But on the other side, the disadvantage of  $J$  results in the difficulty of finding statistical distribution of this distance which is important to calculate a confidence interval. This difficulty comes from the denominator, which is not a constant but a random variable. It is easier to work with euclidian distances like SM. They can be modelled as a binomial variable and their statistical properties are well known (Dillmann *et al.* 1997).

## 15.6. Genetic distances between populations

Genetic distance measures between populations are a generalization from the distance measures we have seen above.

Nei's genetic distance between the  $i$ th and  $j$ th population, using the notation of Weir (1996), is

$$D_{ij} = -\ln(d_{ij}) = -\ln \left[ \frac{\sum_l \sum_u p_{lu_i} p_{lu_j}}{\sqrt{\sum_l \sum_u p_{lu_i}^2 \sum_l \sum_u p_{lu_j}^2}} \right],$$

where  $p_{lu_i}$  is the frequency of allele  $A_u$  for locus  $l$  in the  $i$ th population and  $p_{lu_j}$  is the frequency of allele  $A_u$  for locus  $l$  in the  $j$ th population.

Nei's genetic identity between the  $i$ th and  $j$ th population, corrected for sampling bias (Nei 1978), is

$$D_{ij} = -\ln \left[ \frac{(2n-1) \sum_l \sum_u p_{lu_i} p_{lu_j}}{\sqrt{\sum_l \left( 2n \sum_u p_{lu_i}^2 - 1 \right) \sum_l \left( 2n \sum_u p_{lu_j}^2 - 1 \right)}} \right],$$

where  $n$  is the number of individuals sampled within each population.

Hillis (1984) proposed a genetic distance estimator to overcome the problem of Nei's genetic distance estimator producing greatly different estimates when polymorphisms within populations vary. The Hillis genetic distance estimator is

$$D_{ij} = -\ln \left[ \sum_l \left( \frac{\sum_u p_{lu_i} p_{lu_j}}{\sqrt{\sum_u p_{lu_i}^2 \sum_u p_{lu_j}^2}} \right) / m \right],$$

where  $p_{lu_i}$  is the frequency of allele  $A_u$  for locus  $l$  in the  $i$ th population,  $p_{lu_j}$  is the frequency of allele  $A_u$  for locus  $l$  in the  $j$ th population,  $l = 1, 2, \dots, m$ , and  $m$  is the number of loci.

Roger's genetic distance (1972) between the  $i$ th and  $j$ th population is defined by

$$D_{ij} = \frac{1}{m} \sum_l \left[ \frac{1}{2} \sum_u (p_{lu_i} - p_{lu_j})^2 \right],$$

where  $p_{lu_i}$  is the frequency of allele  $A_u$  for locus  $l$  in the  $i$ th population,  $p_{lu_j}$  is the frequency of allele  $A_u$  for locus  $l$  in the  $j$ th population,  $l = 1, 2, \dots, m$ , and  $m$  is the number of loci.

The genetic distance estimators proposed by Nei (1972, 1978) and Rogers (1972) are affected by within population heterozygosity (Swofford *et al.* 1996). Cavalli-Sforza and Edwards (1967) proposed an estimator that overcomes this problem. The arc distance estimator of Cavalli-Sforza and Edwards is:

$$D_{ij} = \sqrt{\frac{1}{m} \sum_l \left[ \left( 2 \cos^{-1} \sum_u \sqrt{p_{lu_i} p_{lu_j}} \right) / \pi \right]^2},$$

where  $p_{lu_i}$  is the frequency of allele  $A_u$  for locus  $l$  in the  $i$ th population,  $p_{lu_j}$  is the frequency of allele  $A_u$  for locus  $l$  in the  $j$ th population,  $l = 1, 2, \dots, m$ , and  $m$  is the number of loci.

Populations are conceptualised as existing as points in an  $m$ -dimensional Euclidean space which are specified by  $m$  allele frequencies (i.e.  $m$  equals the total number of alleles in both populations). The distance is the angle between these points (chord):

$$d_{chord}(x_1, x_2) = \sqrt{2 \left[ 1 - \frac{\sum_{j=1}^p x_{1j} x_{2j}}{\sqrt{\sum_{j=1}^p x_{1j}^2 \sum_{j=1}^p x_{2j}^2}} \right]}$$

where  $x_i$  and  $y_i$  are the frequencies of the  $i$ th allele in populations  $X$  and  $Y$

- If no alleles are shared between populations  $i$  and  $j$ , then  $D_{ij}=1$ , "regardless of the variability within either population" (Swofford *et al.* 1996), a property lacking in the estimators of Nei (1972, 1978) and Rogers (1972).

- The angular transformation of allele frequencies seeks to eliminate the adverse effects of different allele frequency ranges.

Nei's genetic distance estimators are based on the following assumptions: Infinite-Alleles Model, all loci have same rate of neutral mutation, mutation-genetic drift equilibrium, stable/constant effective population size ( $N_e$ ), linear in time

Cavali-Sforza's genetic distance estimator assumes genetic drift only (no mutation), accommodates changes in population size, is linear in sum of  $1/N_e$  over time

## 15.7. Protocol: tree reconstruction

UPGMA employs a sequential clustering algorithm, in which local topological relationships are identified in order of similarity, and the phylogenetic tree is built in a stepwise manner. We first identify from among all the OTUs the two OTUs that are most similar to each other and then treat these as a new single OTU. Such an OTU is referred to as a composite OTU. Subsequently from among the new group of OTUs we identify the pair with the highest similarity, and so on, until we are left with only two OTUs.

### Definition

Distance  $d_{AB}$  between clusters A, B from individual distances  $d_{ab}$  :

$$d_{AB} = \frac{1}{|A||B|} \sum_{a \in A} \sum_{b \in B} d_{ab}$$

The distance between a simple OTU and a composite OTU is the average of the distances between the simple OTU and the constituent simple OTUs of the composite OTU. Then a new distance matrix is recalculated using the newly calculated distances and the whole cycle is being repeated.

### Algorithm

#### Initialisation

- Assign each sequence  $i$  to its own cluster  $C_i$ . Define one leaf for each sequence, and place at height zero.

#### Iteration

- Determine the two clusters  $i, j$  for which  $d_{ij}$  is minimal.
- Define a new cluster  $C_k = C_i \cup C_j$
- Define a new node  $k$  with daughter nodes  $i$  and  $j$ , and place it at height  $d_{ij}/2$ .
- Add  $k$  to the current clusters and remove  $i$  and  $j$ .

#### Termination

- When only two clusters  $i, j$  remain, place the root at height  $d_{ij}/2$ .

Following the first clustering A and B are considered as a single composite OTU (A,B) and we now calculate the new distance matrix as follows:

$$\text{dist}(A,B),C = (\text{dist}AC + \text{dist}BC) / 2$$

$$\text{dist}(A,B),D = (\text{dist}AD + \text{dist}BD) / 2$$

$$\text{dist}(A,B),E = (\text{dist}AE + \text{dist}BE) / 2$$

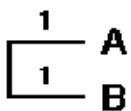
$$\text{dist}(A,B),F = (\text{dist}AF + \text{dist}BF) / 2$$

and so on.

*Example*

Suppose we have the following distance matrix giving the pair wise evolutionary distances of 6 OTUs:

	A	B	C	D	E
B	2				
C	4	4			
D	6	6	6		
E	6	6	6	4	
F	8	8	8	8	8



**First cycle**

We now cluster the pair of OTUs with the smallest distance, being A and B, that are separated by a distance of 2. The branching point is positioned at a distance of  $2 / 2 = 1$  substitution. We thus construct a sub-tree as follows:

Following the first clustering A and B are considered as a single composite OTU (A,B) and we now calculate the new distance matrix as follows:

$$\text{dist}(A,B),C = (\text{dist}AC + \text{dist}BC) / 2 = 4$$

$$\text{dist}(A,B),D = (\text{dist}AD + \text{dist}BD) / 2 = 6$$

$$\text{dist}(A,B),E = (\text{dist}AE + \text{dist}BE) / 2 = 6$$

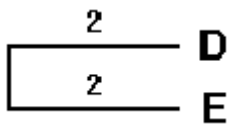
$$\text{dist}(A,B),F = (\text{dist}AF + \text{dist}BF) / 2 = 8$$

In other words the distance between a simple OTU and a composite OTU is the average of the distances between the simple OTU and the constituent simple OTUs of the composite OTU. Then a new distance matrix is recalculated using the newly calculated distances and the whole cycle is being repeated:



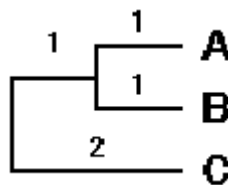
**Second cycle**

	A,B	C	D	E
C	4			
D	6	6		
E	6	6	4	
F	8	8	8	8



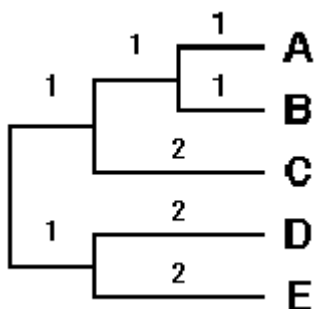
**Third cycle**

	A,B	C	D,E
C	4		
D,E	6	6	
F	8	8	8



**Fourth cycle**

	AB,C	D,E
D,E	6	
F	8	8



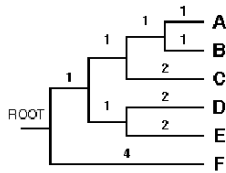
**Fifth cycle**

The final step consists of clustering the last OTU, F, with the composite OTU.

	ABC,DE
F	<b>8</b>

Although this method leads essentially to an unrooted tree, UPGMA assumes equal rates of mutation along all the branches, as the model of evolution used. The theoretical root, therefore, must be equidistant from all OTUs. We can here thus apply the method of mid-point rooting. The root of the entire tree is then positioned at  $\text{dist}(ABCDE),F / 2 = 4$ .

The final tree as inferred by using the UPGMA:



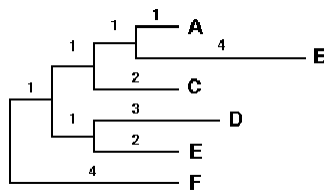
So now we have reconstructed the phylogenetic tree using the UPGMA method. However, there are some **pitfalls**:

- UPGMA clustering is *very sensitive* to unequal evolutionary rates. This means that when one of the OTUs has incorporated more mutations over time than the other OTU, one may end up with a tree that has the wrong topology.
- Clustering works only if the data are *ultrametric*
- Ultrametric distances are defined by the satisfaction of the '*three-point condition*'.

**What is the three-point condition?**

For any three taxa:  $\text{dist} AC \leq \max(\text{dist}AB, \text{dist}BC)$  or in words: the two greatest distances are equal, or UPGMA assumes that the evolutionary rate is the same for all branches

If the assumption of rate constancy among lineages does not hold UPGMA may give an erroneous topology. This is illustrated in the following example; suppose that you have the following relationship:



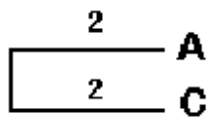
Since the divergence of A and B, B has accumulated mutations at a much higher rate than A. The Three-point criterion is violated! e.g.  $\text{dist}BD \leq \max(\text{dist}BA, \text{dist}AD)$  or,

$10 \leq \max(5,7) = \text{False}$

The reconstruction of the evolutionary history uses the following distance matrix:

	A	B	C	D	E
B	5				
C	<b>4</b>	7			

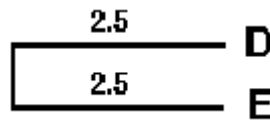
D	7	10	7		
E	6	9	6	5	
F	8	11	8	9	8



We now cluster the pair of OTUs with the smallest distance, being A and C, that are separated a distance of 4. The branching point is positioned at a distance of  $4 / 2 = 2$  substitutions. We thus construct a sub-tree as follows:

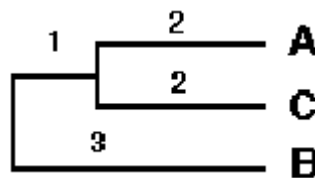
### Second cycle

	A,C	B	D	E
B	4			
D	7	10		
E	6	9	5	
F	8	11	8	9



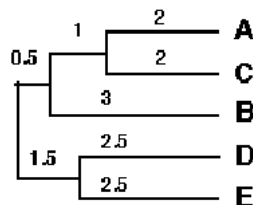
### Third cycle

	A,C	B	D,E
B	6		
D,E	6.5	9.5	
F	8	11	8.5



### Fourth cycle

	AC,B	D,E
D,E	8	
F	9.5	9.5

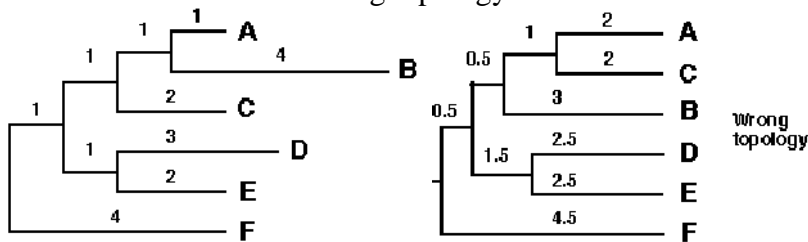


**Fifth cycle**

The final step consists of clustering the last OTU, F, with the composite OTU, ABCDE.

	ABC,DE
F	9

When the original, correct, tree and the final tree are compared it is obvious that we end up with a tree that has the wrong topology.

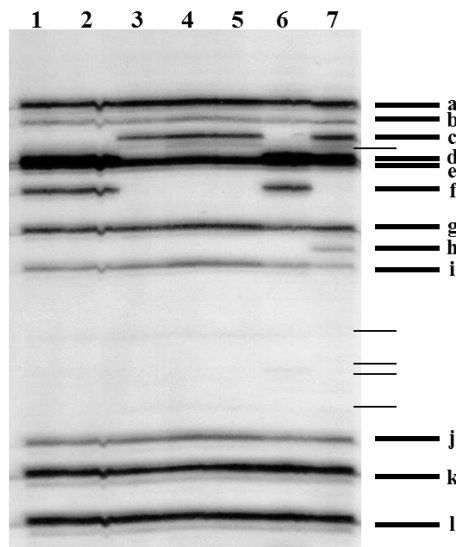


**Conclusion:** The unequal rates of mutation have led to a completely different tree topology.

**15.8. UPGMA exercise**

**UPGMA exercise**

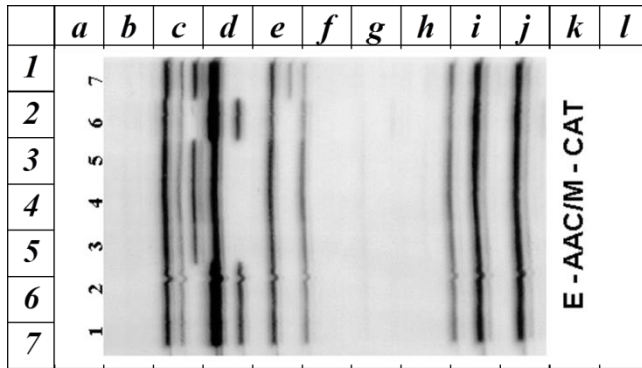
**AFLP results**



Lane identifications:

- 1= reference 1
- 2= line 01
- 3= line 97
- 4= line 02
- 5= line 09
- 6= line 95
- 7= reference 2

**E - AAC/M - CAT**



Accessions 2 to 6 were obtained by mutation induction from supposedly accession 1. Accession 7 is a control.

Verify whether accessions 2 to 6 have been derived from accession 1.

$$s_{ij} = (n_{11} + n_{00}) / (n_{11} + n_{00} + n_{10} + n_{01})$$

$$d_{ij} = \sqrt{2 * (1 - s_{ij})}$$

	a	b	c	d	e	f	g	h	i	j	k	l
1	1	1	0	1	1	1	1	0	1	1	1	1
2	1	1	0	1	1	1	1	0	1	1	1	1
3	1	1	1	0	1	0	1	0	1	1	1	1
4	1	1	1	0	1	0	1	0	1	1	1	1
5	1	1	1	0	1	0	1	0	1	1	1	1
6	1	1	0	1	1	1	1	0	1	1	1	1
7	1	1	1	1	1	0	1	1	1	1	1	1

1:2	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$											
1:3	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$	2:3	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$									
1:4	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$	2:4	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$	3:4	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$							
1:5	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$	2:5	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$	3:5	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$	4:5	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$					
1:6	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$	2:6	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$	3:6	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$	4:6	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$	5:6	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$			
1:7	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$	2:7	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$	3:7	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$	4:7	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$	5:7	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$	6:7	$n_{11} =$ $n_{00} =$ $n_{10} =$ $n_{01} =$	

$$m = \frac{1}{2} [(N-1)^2 + (N-1)]$$

The choice of  $s_{ij}$  and  $d_{ij}$  is given by the problem, (verify relation to parent, AFLP)

Possible simplification based on identity of rows 1, 2 & 6 and rows 3, 4 & 5

	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>	<i>h</i>	<i>i</i>	<i>j</i>	<i>k</i>	<i>l</i>
<b>1</b>	1	1	0	1	1	1	1	0	1	1	1	1
<b>2</b>	1	1	0	1	1	1	1	0	1	1	1	1
<b>3</b>	1	1	1	0	1	0	1	0	1	1	1	1
<b>4</b>	1	1	1	0	1	0	1	0	1	1	1	1
<b>5</b>	1	1	1	0	1	0	1	0	1	1	1	1
<b>6</b>	1	1	0	1	1	1	1	0	1	1	1	1
<b>7</b>	1	1	1	1	1	0	1	1	1	1	1	1

	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>	<i>h</i>	<i>i</i>	<i>j</i>	<i>k</i>	<i>l</i>
<b>1,2,6</b>	1	1	0	1	1	1	1	0	1	1	1	1
<b>3,4,5</b>	1	1	1	0	1	0	1	0	1	1	1	1
<b>7</b>	1	1	1	1	1	0	1	1	1	1	1	1

$$s_{ij} = (n_{11} + n_{00}) / (n_{11} + n_{00} + n_{10} + n_{01})$$

$$d_{ij} = \sqrt{2 * (1 - s_{ij})}$$

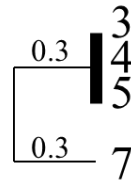
1,2,6 : 3,4,5	$n_{11}=8$ $n_{00}=1$ $n_{10}=2$ $n_{01}=1$	$s_{ij} = 9/12$ $d_{ij} = \sqrt{6/12} = 0.7$	<b>1,2,6</b>	<b>3,4,5</b>	<b>7</b>
1,2,6 : 7	$n_{11}=9$ $n_{00}=0$ $n_{10}=1$ $n_{01}=2$	$s_{ij} = 9/12$ $d_{ij} = \sqrt{6/12} = 0.7$	<b>1,2,6</b>	<b>0</b>	
3,4,5 : 7	$n_{11}=9$ $n_{00}=1$ $n_{10}=0$ $n_{01}=2$	$s_{ij} = 10/12$ $d_{ij} = \sqrt{4/12} = 0.6$	<b>3,4,5</b>	<b>0,7</b>	
			<b>7</b>	<b>0,7</b>	<b>0,6</b>

three-point condition:

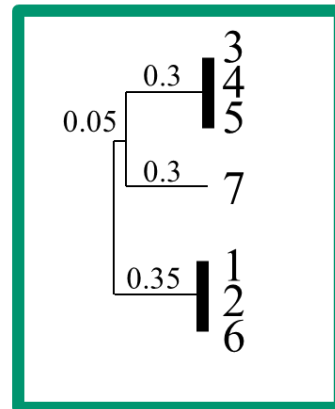
$$\text{dist}(1,2,6:7) \leq \text{Max}[\text{dist}(1,2,6:3,4,5), \text{dist}(3,4,5:7)]$$

$$0.7 \leq \text{Max}(0.7, 0.6) !$$

	1,2,6	3,4,5	7
1,2,6	0		
3,4,5	0.7	0	
7	0.7	0.6	0



	1,2,6	3,4,5/7
1,2,6	0	
3,4,5/7	0.7	0



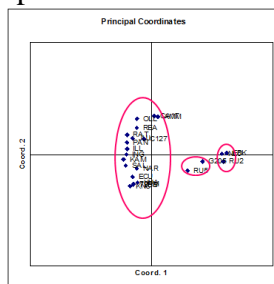
**Conclusion:**

Mutants 3, 4 & 5 are more related to the control 7 than to the putative parent 1. Possible explanations:

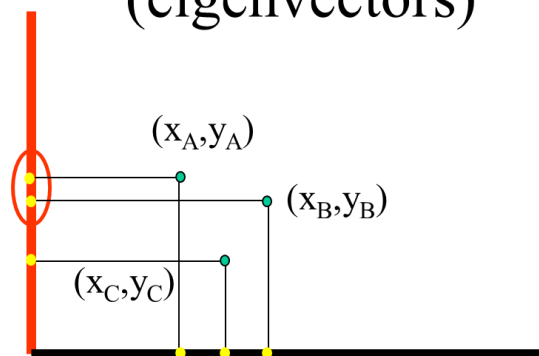
- Mislabelling of part of the M<sub>0</sub> and/or M<sub>1</sub>
- Outcrossing during M<sub>1</sub> selfing

## 15.9. Principal Component Analysis (PCA)

If a multivariate dataset is represented as a set of coordinates in an n-dimensional data space (1 axis per variable), PCA can reduce the dimensionality of the transformed data and supply a lower-dimensional projection when viewed from its most informative viewpoint, using only the first few principal components. For a seemingly random distribution of data points in the n-dimensional results space, PCA starts with finding the analytical plane by slicing the results space into lower dimensional representations of uncorrelated parameters (eigenvectors).



## PCA (eigenvectors)



### • finding the analytical plane

In mathematical terms, PCA is a procedure to transform a set of potentially correlated observations into a set of uncorrelated data points: principal components (in number less than or equal to the original variables). This orthogonal transformation is defined in such a way that the first principal component accounts for as much of the variability in the data as possible (maximum variance), and each succeeding component in turn has the highest variance possible under the constraint that it is uncorrelated with (orthogonal to) the preceding components. Principal components are guaranteed to be independent only if the data set is jointly normally distributed.

PCA is the simplest of the true eigenvector-based multivariate analyses. It might be visualized as uncovering the internal structure of the data in a way which best explains their variance. Sensitive to the relative scaling of the original variables, it can be done by eigenvalue decomposition of a data covariance matrix or singular value decomposition of a data matrix, usually after mean centring the data for each attribute. The results of a PCA are usually discussed in terms of component scores (the transformed variable values corresponding to a particular data point) and loadings (the weight by which each standardized original variable is to be multiplied to get the component score). PCA is closely related to factor analysis; and some statistical packages deliberately merge the two techniques. True factor analysis makes



different assumptions about the underlying structure and solves eigenvectors of a slightly different matrix.

In linear algebra, an orthogonal matrix, is a square matrix with real entries whose columns and rows are orthogonal unit vectors. This means, that a matrix  $Q$  is orthogonal if its transpose is equal to its inverse:  $Q^T = Q^{-1}$ , and thus it follows that  $Q^T Q = Q Q^T = I$  ( $I$  being the identity matrix). An orthogonal matrix  $Q$  is thus square, invertible, unitary ( $Q^{-1} = Q^*$ ), and normal ( $Q^* Q = Q Q^*$ ). As a linear transformation, an orthogonal matrix preserves the dot product of vectors, and therefore acts as an isometry of Euclidean space, such as a rotation or reflection, thus, it is a unitary transformation.

The eigenvectors of a square matrix are the non-zero vectors that, after being multiplied by the matrix, remain parallel to the original vector. For each eigenvector, the corresponding eigenvalue is the factor by which the eigenvector is scaled when multiplied by the matrix. The prefix eigen- is adopted from the German word "eigen" for "own" in the sense of a characteristic description. In mathematical terms: if  $A$  is a square matrix, a non-zero vector  $v$  is an eigenvector of  $A$  if there is a scalar  $\lambda$  (lambda) such that  $Av = \lambda v$

The scalar  $\lambda$  (lambda) is said to be the eigenvalue of  $A$  corresponding to  $v$ . An eigenspace of  $A$  is the set of all eigenvectors with the same eigenvalue together with the zero vector, which however, is not an eigenvector.

## 15.9.1. Considerations and references

### Planning experiments and analyses

#### Which entities should be sampled?

There are no formal statistical rules for deciding this, so empirical testing is needed. When selecting among a large number of potential entities (*e.g.*, germplasm accessions) or when resources are limiting (which they nearly always are), geographical or ancestral origin, morphological phenotypes, or other phenotypic or historical criteria can often be used to select accessions to represent a gene pool or a specific subset of a gene pool. The genetic material chosen for study depends on economic resources, the nature, scale, scope, and goals of the study, and *a priori* knowledge of genetic relationships. Closely related genetic materials, for example, need not be sampled unless there is a compelling biological or economic reason to do so. The 'ideal' sample of genetic material for studying a particular question is profoundly affected by the nature and genetic origin (if known) of the genetic material.

The goal of a DNA fingerprinting study might be to classify every entity belonging to a particular biological or economic class of entities, *e.g.*, a seed company might fingerprint and classify every inbred line and hybrid they own and every hybrid sold by their competitors for the purpose of protecting intellectual property. Many crop plant gene pools are comprised of hundreds or even thousands of germplasm accessions. Depending on the mating biology and breeding systems of the species, accessions could be comprised of outcrossing wild populations (*e.g.*, genetically heterogeneous, segregating populations), mixtures of inbred genotypes, or inbred lines. How genetically heterogeneous accessions are sampled depends on the goal of the study and economic resources.

Another goal of a DNA fingerprinting study might be to assess the minimum set of accessions that comprise an ideal or so-called core set. The purpose of a core set, in theory, is to produce maximum information from a minimum sample of genetic materials. The practical aims might be to eliminate redundant accessions and streamline the maintenance of genetic diversity in a seed or gene bank.

Similar concepts can be applied to surveys of genetic diversity, *e.g.*, the 'optimum' set of genetic materials for assessing the utility of a sample of genetic markers or, more broadly, for classifying new genetic materials or genetic materials of unknown ancestry or origin.

#### What is the best sampling strategy?

The mating biology and breeding system of the species dictate the sampling strategy. The gene pools of many plant species, *e.g.*, maize (*Zea mays* L.) and sunflower (*Helianthus annuus* L.), are comprised of partially or 'fully' inbred genetic stocks, in addition to heterogeneous, segregating populations (natural or experimental). The gene pools of humans, most animal species, and many plant species, more or less domesticated and/or wild types, are comprised of heterogeneous, segregating populations.

The optimum genetic and statistical sampling strategies may be difficult to specify, are nearly always constrained by economic factors, and depend on the nature of the statistical analysis and scope of inference. When analyses are performed on segregating populations, a sufficient number of individuals must be sampled within each population to accurately estimate gene

and genotype frequencies. Weir (1996) proposed sampling over loci for random model analyses and over individuals for fixed model analyses. The line between fixed and random is often blurred. Basically, if the scope of inference is across a species or across other strata where broad inferences are to be made, then random models are used. If the scope of inference is a fixed set of populations or inbred lines, then fixed models are used.

If the goal of the study is to survey allelic diversity among a sample of populations (chosen for some biological or commercial reason), then extensive within-population sampling may not be necessary.

If the goal is to accurately describe genetic patterns among populations, measure linkage disequilibrium or gene flow, or protect intellectual property (e.g., an open-pollinated or synthetic cultivars in crop plants), then individuals within populations must be sampled to accurately estimate gene and genotype frequencies and perhaps to find rare alleles and genotypes.

### **What types of variables should be measured?**

Although we are primarily concentrating on the analysis of genotypic measurements (e.g., DNA marker genotypes), phenotypic measurements should not be overlooked and can be combined with genotypic measurements in analyses of genetic patterns. Special similarity measures can be used to combine phenotypic and genotypic measurements or a 'conceptual synthesis' of patterns can be produced from separate analyses performed on phenotypic and genotypic variables. The choice of variables is usually more complicated for phenotypic than genotypic variables, because the former are heterogeneous, whereas the latter are homogeneous (when a single marker system is employed) in the conceptual sense, however, the information supplied by individual genetic markers can vary. If DNA fingerprints are to be produced, then the types of variables measured are dictated (i) by the types of markers developed for the species, (ii) whether the DNA markers are dominant or co-dominant, (iii) by the homology of DNA fragments across individuals or populations, (iv) by economic factors, (v) by the reproducibility and robustness of the DNA marker system (genotyping errors). The ideal genetic marker is highly polymorphic, co-dominant, locus-specific, robust, and highly reproducible.

### **How many variables should be measured?**

There are no formal statistical rules for deciding how many genetic markers are needed to accurately classify accessions, describe genetic patterns, or accurately estimate genetic distances and phenograms.

- Smith *et al.* (1991) used 200 RFLP markers dispersed across the maize genome to fingerprint 11 inbred lines (the genetic distance matrix was comprised of 55 elements). They estimated distance matrices by sampling 5 to 200 RFLP markers in increments of five (e.g., five, 10, 15, ..., 200). They concluded that accuracy was sufficient with 100 or more markers.
- Bernardo (1993) concluded that 250 or more marker loci were needed to produce precise estimates of coefficients of co-ancestry.

The number of genetic markers used in an analysis may be dictated by non-statistical factors. The outcome of the analysis might be one of the criteria used to select genetic markers for future analyses.

Ideally, genetic markers for protecting intellectual property and classifying unknown genetic materials should be highly polymorphic and dispersed across the genome.

**Should analyses be performed on raw multivariate data or genetic similarities?**

- Typically, multivariate analyses of DNA genotypes (fingerprints) are performed on genetic similarity or distance matrices among entities rather than on raw multivariate data matrices.
- PCA of raw DNA genotypes, although not widely done, can be used to assess the importance of individual genetic markers by comparing principal component coefficients, i.e., individual elements of characteristic vectors (eigenvectors).

**15.10. References**

<sup>(1)</sup> [<http://www.icp.ucl.ac.be/~opperd/private/upgma.html>]

- Bernardo R.** 1993. Estimation of coefficient of coancestry using molecular markers in maize. *Theor. Appl. Genet.* 85: 1055-1062.
- Cavalli-Sforza L.L. and Edwards A.W.F.** 1967. Phylogenetic analysis: models and estimation procedures. *Am. J. Hum. Genet.* 19: 233-257.
- Dice L.R.** 1945. Measures of the amount of ecological association between species. *Ecology* 26: 297-302.
- Dillmann C., Charcosset A., Goffinet B., Smith J.S.C. and Dattée Y.** 1997. Best linear estimator of the molecular genetic distance between inbred lines. In: Krajewski P, Kaczmarek Z (eds) *Advances in biometrical genetics. Proceedings of the tenth meeting of the EUCARPIA section biometrics in plant breeding, 14-16 may 1997, Poznan*, pp 105-110
- Dudley J. W.** 1993. Molecular markers in plant improvement: Manipulation of genes affecting qualitative traits. *Crop Science* (33):660-668 & **Munn R. and Dudley J.** 1995. A PC computer program to generate a dissimilarity matrix for cluster analysis. *Crop Sci.* 35:925-927.
- Everitt B.S.** 1992. *Cluster analysis.* Oxford Univ. Press, New York.
- Excoffier L., Smouse P.E. and Quattro J.M.** 1992. Analysis of molecular variance inferred from metric distances among DNA haplotypes: application to human mitochondrial DNA restriction data. *Genetics* 131:479-491
- Flury B.** 1988. *Common principal components and related multivariate methods.* Wiley, New York.
- Gower J.C.** 1971. A general coefficient of similarity and some of its properties. *Biometrics* 27: 857-872.
- Hamming R.W.** 1950. Error detecting and error correcting codes. *Bell System Technical Journal* 29 (2): 147-160
- Hillis D.M.** 1984. Misuse and modification of Nei's genetic distance. *Syst. Zool.* 33: 238-240.
- Hillis D.M., Moritz C., and Mable B.K.** 1996. *Molecular systematics.* Sinauer, Sunderland, Massachusetts.
- Jaccard P.** 1908. Nouvelles recherches sur la distribution florale. *Bull Soc Vaud Sci Nat* 44: 223-270
- Nei M.** 1972. Genetic distance between populations. *Am. Nat.* 106: 283-292.

- Nei M.** 1978. Estimation of average heterozygosity and genetic distance from a small number of individuals. *Genetics* 89: 583-590.
- Nei M. and Li W.-H.** 1979. Mathematical model for studying genetic variation in terms of restriction endonucleases. *Proc. Natl. Acad. Sci.* 76: 5269-5273.
- Peltier D., Chacon H., Tersac M., Caraux G., Dulieu H. and Bervillé A.** 1995. Utilisation des RAPD pour la construction de phénogrammes et de phylogrammes chez *Petunia*. In: *Techniques et utilisations des marqueurs moléculaires. Coll Les colloques INRA*
- Rogers J.S.** 1972. Measures of genetic similarity and genetic distance. *Univ. Texas Publ.* 7213: 145-153.
- Romesburg H.Ch.** 1990. *Cluster Analysis for Researchers.* Florida, Krieger Publishing Co. (original edition 1984).
- Smith O.S., Smith J.S.C., Bowen S.L. and Tenborg R.A.** 1991. Numbers of RFLP probes necessary to show associations between lines. *Maize Genet. Newsltr.* 65: 66.
- Sneath P.H.A. and Sokal R.P.** 1973. *Numerical taxonomy.* San Francisco, Freeman
- Sokal R.P. and Michener C.D.** 1958. A statistical method for evaluating systematic relationships. *Univ Kansas Sci Bull* 38: 1409-1438
- Swofford D.L., Olsen G.J., Waddell P.J. and Hillis D.M.** 1996. Phylogenetic inference. pp. 407-514. Hillis, D.M., C. Moritz, and B.K. Mable (ed.). *Molecular systematics.* Sinauer, Sunderland, Massachusetts.
- Weir B.S.** 1996. *Genetic data analysis.* Sinauer, Sunderland, Massachusetts.

## 16. POPULATION GENETICS

Population genetics is that branch of genetics that attempts to describe how the frequency of the alleles (of genes) changes over time. To study frequency changes, populations rather than individuals are analysed. The scope of this module however is not to provide an in-depth resource on this branch of science, rather it is aimed at guiding the researcher in a stepwise format through the collection (including coding), analyses and arriving at valid inferences on data for allelic frequencies of molecular markers.

The data coding schemes begin with a random example of a dominant marker gel data. Whether the bands come from RAPD's, ISSR, and AFLP's or similar, does not affect the way data is coded, and more importantly, how it is analysed. What matters, is whether or not *we observe* a given band.

Next, co-dominant markers are dealt with as they are close to the notion of a diploid species where each individual carries  $n$  *maternally* and  $n$  *paternally* inherited gametes for a total ploidy of  $2n$ . Of course, codominant data can be obtained in tetraploid or hexaploid individuals also, as will be demonstrated. The exercises will start with microsatellite data from a population sample. It is important to note however that all these coding systems can be used also for allozyme data. Different coding schemes will be analysed, some 'tricks' with using spread sheets and highlights on what can, and what cannot be done with each coding system will also be shown.

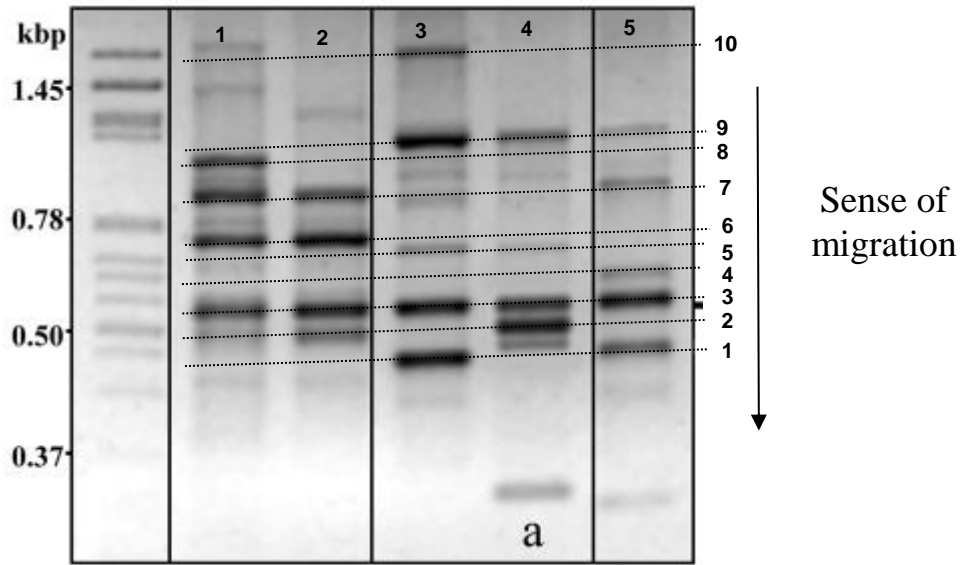
After reviewing how data can be coded, the next step will involve going through the basic concepts of population diversity, population structure, and population divergence. This last part of this module is the basis of phylogenetic studies, although for this manual, only phenetic analyses will be shown.

To conclude this brief introduction to population genetics, two non-exhaustive lists of references and of web-resources of relevance to the study of the subject are provided. Finally, a list of key concepts and equations are provided to complete the definitions given in the text.

### 16.1. Reading and coding genetic data

#### 16.1.1. Presence/absence coding of dominant data

The most commonly used way for coding genotypes or genetic marker data is by doing a matrix of presence/absence of bands, usually with 1's and 0's. This type of markers is easy to read, provided the number of bands is reasonable and clear. Band intensity, is an issue, and interpretations may change from person to person.



**Figure 16-1.** Typical dominant data gel, consisting of 5 lanes, and at least 10 well identifiable bands. Bands are scored 1, if present, zero, otherwise. Table 1 shows one reading of this gel into a spread sheet program (interpretation may vary from person to person, or from day to day).

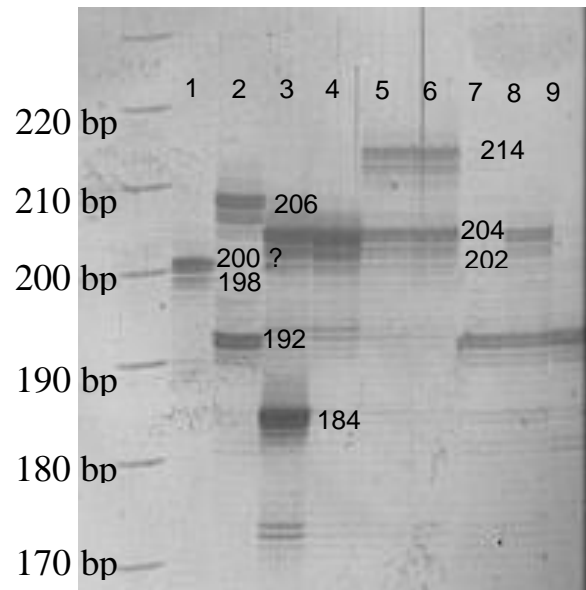
**Table 16.1–1.** Basic transcription of a dominant marker gel into a spread sheet. Data are organized by columns (fields: id, b1, b10) and individuals are rows (records).

	A	B	C	D	E	F	G	H	I	J	K
1	id	b1	b2	b3	b4	b5	b6	b7	b8	b9	b10
2	1	0	0	1	0	0	1	1	1	0	0
3	2	0	1	1	0	0	1	1	0	0	0
4	3	1	0	1	0	1	1	0	0	1	1
5	4	0	1	1	0	0	0	0	0	1	0
6	5	1	0	1	1	0	0	1	0	1	0
7											
8											

As will be seen later, this coding is not complete for analysis with corresponding software, but is a good starting point. Score bands are highlighted grey for clarity purposes.

### 16.1.2. Allele size coding for microsatellites

**Figure 13.1** shows a typical microsatellite data with 7 alleles in 9 individuals (the number of alleles may change according to the person that reads the gel!). This marker is codominant, because we can see that individuals can bear two alleles at the same time. In principle, each product is originated in the two homologous parts for that particular locus, and if the two alleles are the same, a darker, single band should be seen. Figure 13.2 and Table 13.2 show a first interpretation of this gel in a codominant fashion, upon which inbreeding  $f$  or  $f_{is}$  can be computed as well as other statistics (see chapters 2 and so forth).



**Figure 16-2.** Test gel of *Quercus humboldtii* (Andean oak, Colombia) showing 9 individuals (Fernandez, unpublished data). This gel presents many of the typical features of microsatellites: many alleles, stuttering bands, more than two “main” bands, and ambiguity of allele size. A sequencer will also give you results of the type **202.14** bp that the researcher needs to round. Rounding is necessary at this stage or at later steps as most programmes only accept integer numbers.

**Table 16.1–2.** Same data from example gel using a regular spread sheet programme. Note individuals appear in *rows* (records), and particular data (fields) are in *columns*. Note that individuals 7 and 9 are coded as homozygotes and not as one allele with missing data. Some programs deal with “null” alleles, i.e., false homozygotes due to PCR problems, and in that case, the notation would indicate one un-observed allele.

	1	2	3	4
1	population	individual	locus_1_allele1	locus_1_allele2
2	1	1	198	200
3	1	2	192	206
4	1	3	184	204
5	1	4	202	202
6	1	5	204	214
7	1	6	204	214
8	1	7	192	192
9	1	8	192	204
10	1	9	192	192
11	2	1		
12	2	2		
13	2	3		
14				
15				



### 16.1.3. Categorical coding

A second interpretation of this gel, would be simply naming the alleles with letters or numbers (preferred coding) from 1 to 8; this is what is usually called “categorical” or “allelic states” coding of alleles that in this case disregards the size information present in the microsatellites (bp’s). We will see that the size information is important for genetic distances such as Delta  $\mu^2$  and others, but that allelic state is sufficient for genetic distances such as Nei’s standard genetic distance, widely used for allozyme data. Figure 13.3 and Table 13.3 shows the coding in “categorical” or “allelic states” for the same gel.

### 16.1.4. Presence/absence coding of co-dominant data

Yes, you are reading right. A third coding scheme is the popular one that uses 0’s (zeros) and 1’s (ones), usually called “presence/absence” coding that we just saw for Dominant data in the first section (13.1.1). Often times, we are not interested in evolutionary models and/or samples do not come from random samples from natural populations. We may have accessions coming from different countries or regions within countries collected simply because they present an interesting trait: nice fruits, long spikes, little cyanide, etc. This coding is required for traditional statistics such as Principal Components Analysis (PCA) and related multivariate techniques, with the advantage that genetic data can be combined with morphological data for grouping purposes. Table 4 shows the presence/absence coding for the same example gel.

**Important Note:** You may notice that this coding is not exclusively for diploids. In fact, tetraploids or hexaploids can be handled this way. Simply, there can be more than two bands per individual, and the notion of heterozygotes diffuses and becomes secondary.

It is clear that for allozyme data, or morphological data known to be co-dominant (white, lilac and purple flowers in *Lynanathus*, for example), “presence/absence” are perfectly applicable. At this point, we would lose the diploid information so estimation of inbreeding (the parameter  $f_{is}$  that measures the probability that two alleles within an individual are the same) cannot be computed. This coding, however, is highly popular for analysing accessions because if you will, it is “model” free, and as seen from Table 4, we can include in the same database different kinds of data, and potentially in the same analysis (fruit data color could be changed to 1, 2 and 3 *etc.* to run all in the same analysis, but all depends on the programme used).

**Table 16.1–3.** Example of Co-dominant data coded as presence absence of bands. First, the total number of alleles is counted, and the corresponding number of bands is defined, being 8 in our case. Note that for homozygote individuals (we are dealing with diploid data) there is controversy about the scoring. In the example below the individuals 7 and 9 were coded as 1 for allele 1, but some people think we should give them twice as much weight (i.e., two copies are there!) so the genotype should be “2” instead of “1”. This is no longer “presence/absence” strictly, but results change little in practice.

	1	2	3	4	5	6	7	8	9	10	11
1	population	individual	band_01	band_02	band_03	band_04	band_05	band_06	band_07	band_08	fruit_color
2	1	1	0	0	1	1	0	0	0	0	red
3	1	2	0	1	0	0	0	0	1	0	red
4	1	3	1	0	0	0	0	1	0	0	crimson
5	1	4	0	0	0	0	1	1	0	0	green
6	1	5	0	0	0	0	0	1	0	1	red
7	1	6	0	0	0	0	0	1	0	1	red
8	1	7	0	1	0	0	0	0	0	0	crimson
9	1	8	0	1	0	0	0	0	0	1	red
10	1	9	0	1	0	0	0	0	0	0	red
11	2	1									
12	2	2									
13	2	3									
14											

### 16.1.5. Formatting dominant data as co-dominant

As strange as it may sound, we can code dominant data as codominant for using a codominant-based data analysis software. Some functions may not work, and the measure of inbreeding will be totally false, but genetic distances, using shared allele distances can be computed. In this case, we would code as follows:

- 22 for the presence of a band
- 11 for the absence of a band
- (-99 if missing data are allowed... not easy to know for dominant markers!)

The file should look similar to that in Table 3 (categorical or allelic state coding) before we transform it in its final form, as shown in Table 5. We can no longer use zeros, because in this context, zeros are usually reserved for missing data!

**Table 16.1–4.** Dominant data (Figure 13.1) coded as codominant *i.e.*, 2-alleles per band.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V
1	pop	id	b1	b2	b3	b4	b5	b6	b7	b8	b9	b10										
2	1	1	1	1	1	1	2	2	1	1	1	1	2	2	2	2	2	2	1	1	1	1
3	1	2	1	1	2	2	2	2	1	1	1	1	2	2	2	2	1	1	1	1	1	1
4	1	3	2	2	1	1	2	2	1	1	2	2	2	2	1	1	1	1	2	2	2	2
5	1	4	1	1	2	2	2	2	1	1	1	1	1	1	1	1	1	1	2	2	1	1
6	1	5	2	2	1	1	2	2	2	2	1	1	1	1	2	2	1	1	2	2	1	1
7																						

### 16.1.6. Notes of formatting diploid data with spread sheets

Many programmes for analysing diploid data have the bad habit (among many) of using fixed length characters for each marker. For example, our first individual with genotype 198 / 200, may need to be coded as “198200” in a single string of characters. Moreover, the same genotype in categorical coding 3 / 4 may need to be coded “0304” in a so-called **two-allele** coding, or “003004” in a **three-allele** coding. This is particularly true for the programmes Fstat and GenePop on the web. By the way, other programs may need coding as 198.200, or 198, 200, *etc.* but in general, they are handled automatically by some software (see below).

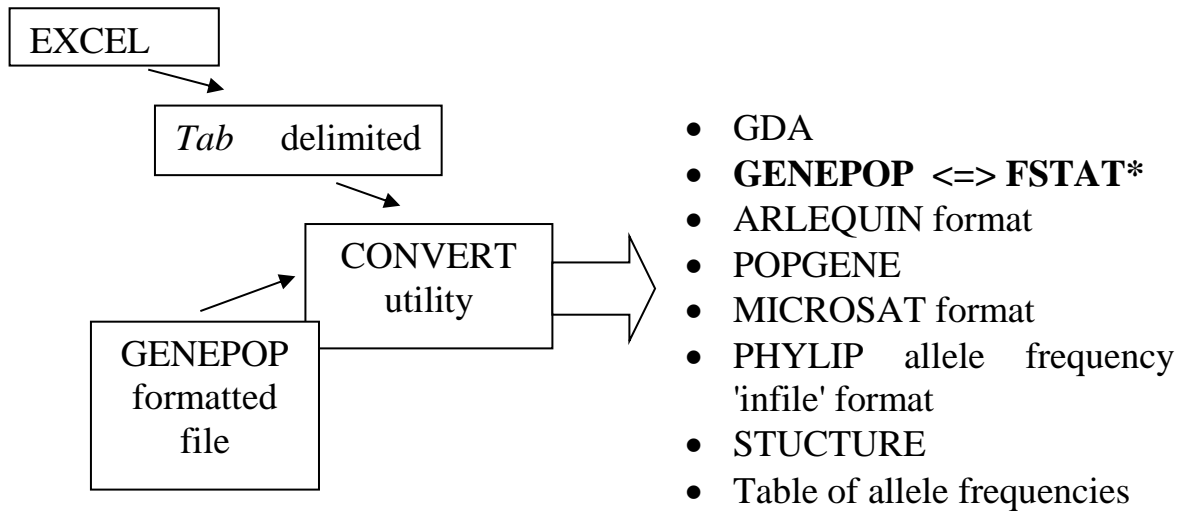
Spread sheet programmes as OpenOffice Calc or Excel handle text conversions with the CONCAT string function that can be seen in the example below.

**Table 16.1–5.** Example of our size type coding where two columns (one for each possible allele) have been collapsed and “concatenated” in a single text. This one is from a French version of the software and the name of the function changes a bit from language to language. For OpenOffice in English, the function is: =CONCATENATE(A1;B1), and they are accessible from the  $f_x$  button, **string** functions.

	A	B	C	D	E	F
1	population	individual	locus_1_allele1	locus_1_allele2		locus1
2	1	1	198	200		198200
3	1	2	192	206		192206
4	1	3	184	204		184204
5	1	4	202	202		202202
6	1	5	204	214		204214
7	1	6	204	214		204214
8	1	7	192	192		192192
9	1	8	192	204		192204
10	1	9	192	192		192192
11	2	1				
12	2	2				
13	2	3				

### 16.1.7. Transforming data types using software

As already noted, there is not a universal data type, but some conversions can be done with available software, at least for some applications. For many programmes, there is no way around and data files must be coded manually. A small utility that we will use is the software CONVERT (Glaubitz 2003). This software can translate from a rather simple data file, to several other programmes, as shown in Figures 13.4 and 13.5:



**Figure 16-3.** Flow chart showing the different data translation paths possible with the CONVERT utility software. Not all possibilities are here, but at least these programmes are glued together. Note, however, that these programs are almost exclusive for diploid codominant data, but some tricks can be done as explained in section 13.1.6. **FSTAT** is marked with an asterisk as is the one that we are going to use for most of the analyses, as explained in the next section.

### 16.1.8. The FSTAT data file

As we will use this programme mostly throughout the exercises let us explain briefly the data structure need.

For running FSTAT, it is first necessary to create an input file named FILENAME.DAT (where FILENAME is anything between 1 and 256 characters) containing the genotypic data, coded numerically, either with a 1, a 2 or a 3-digit number per allele. The file must have the following format:

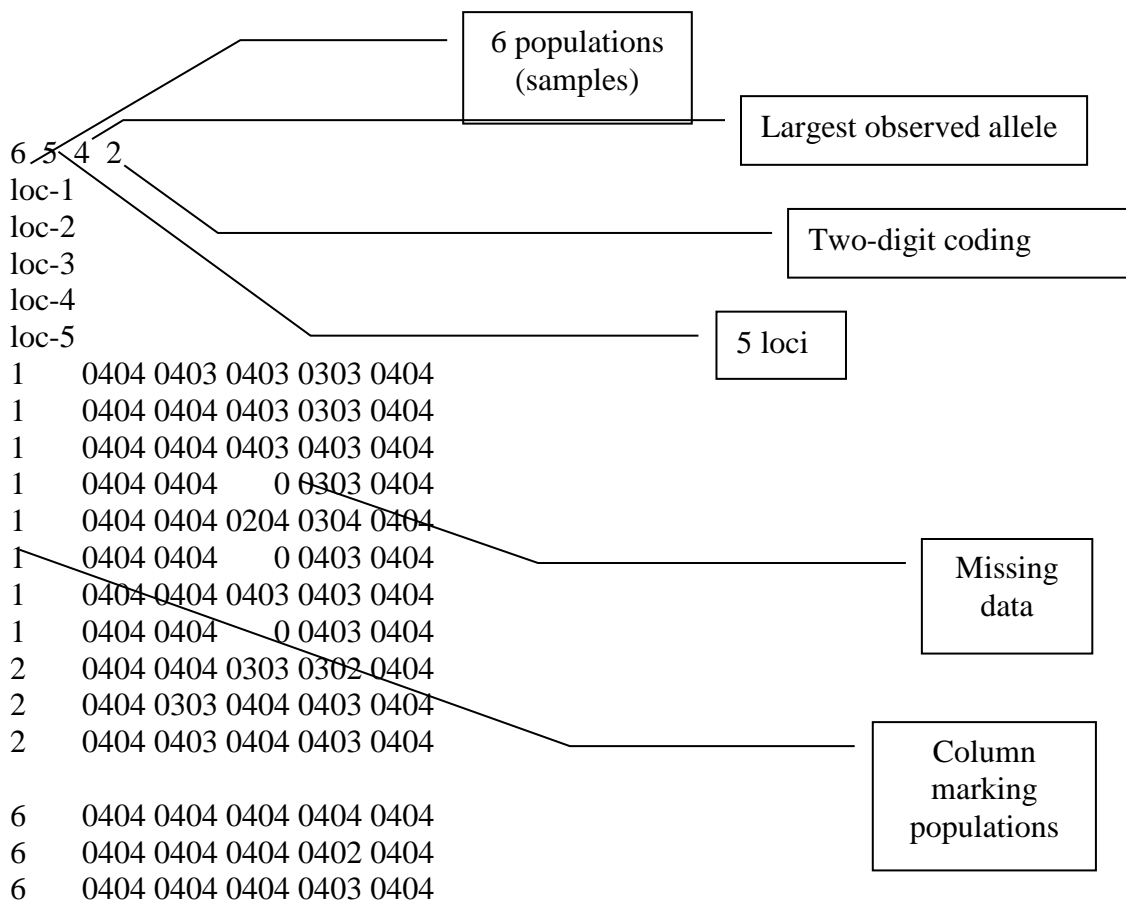
- The first line contains 4 numbers:
  1. the number of populations (here called samples)  $\leq 200$
  2. the number of loci  $\leq 100$
  3. the highest number used to label an allele  $\leq 999$ , and a
  4. data coding type: 1 if the code for alleles is a one digit number (1-9), a 2 if code for alleles is a 2 digit number (01-99) or a 3 if code for alleles is a 3 digit number (001-999).

These 4 numbers need to be separated by any number of spaces.

- Next, the name of the loci are written, one per line, and finally, the main data with first a number for each population followed by the different genotypes, each row for each individual.

- Missing data is encoded as zeros.

A data file for six populations, five loci, 4 alleles maximum and 2-digit allele coding would look then as:



## 16.2. Genetic diversity

Gene or genetic diversity is perhaps the central notion and motivation for conducting research in natural resources and crop improvement. If there were no biodiversity, we wouldn't have a job, and more importantly, we would probably not exist.

Evolution, or the change of heritable characters across generations (in the case of genes, it is simply the change of allele frequencies and genotype frequencies in time) can only occur if there is enough genetic variability upon which, natural and artificial selection can act. Hence, measuring genetic diversity is paramount in population genetics, and we will see that we use several complimentary approaches. First, we will see the descriptive statistics.

**Allelic Richness:** The first measure of genetic diversity is the number of alleles at a locus (see glossary for definitions), usually denoted  $A$ . The more allelic variants are found in a population, the more variable it is.

**Rare Alleles:** Often, we would like to mark a difference between the number of common and rare alleles. One way is to define a threshold of considering all alleles with frequencies below 0.05 as rare. These rare alleles are then considered important, and if they are unique or private to the population, we would stress them in our results. It is somewhat less used today.

**Effective Alleles:** Another way of estimating the number of alleles that contribute more to the diversity is by means of the effective number of alleles, denoted  $A_e$ . This measure uses the frequency of alleles to estimate the number of alleles *if* they were at the same frequency or at the maximum possible diversity, using the formula:  $A_e = \frac{1}{2! \sum_i p_i^2}$ , where  $p_i$  represents the frequency of each allele. This number can be seen also as how many numbers of individuals need to be sample before we repeat an allele. For example, typical results for microsatellite data include  $A = 10$ , and  $A_e = 3.8$  (for example) meaning that we observed 10 alleles, but that 4 are common, and six are rare. Note that here rare is not exactly as in the previous definition, but simply that contribute less to the general diversity.

**Polymorphic Bands:** For dominant marker data, a straight forward measure of diversity is the percentage of polymorphic bands, which is simply the proportion of bands that present presence/absence variability. Usually they are counted with the 0.05 criterion.

**Observed Heterozygosity:** For diploid individuals (and polyploidy in general) this is a key measure obtained when using Co-dominant data. It is simply the proportion of individuals per population that have different recognizable alleles at a given locus and it is denote as  $H_o$  or  $h_o$ , being the former more

used for an average of many populations and the latter for a single population measure.

**Expected Heterozygosity:** This is the actual measure of genetic or gene diversity. It represents the probability that two alleles in a locus are different, and is usually denoted  $H$ ,  $H_e$  or  $h_e$ . It is also known as Nei's genetic diversity as most of the gene diversity theory has been proposed by M. Nei in the 1970's. In general, it is computed as follows, although there are some variations to account for sample size or levels of inbreeding:

$$H_e = 1 - \sum_i p_i^2$$

where  $p_i$  represents again the frequency of each allele. The  $p^2$  term represents the probability of sampling twice the same allele, or probability of homozygosity. Then, one minus this probability computed for all present alleles, gives us the probability of sampling two different alleles at a locus. It will be seen next, that this measure is calculated with respect to an ideal, or reference population that may, or may not have similar values as the observed heterozygosity. These deviations are considered next.

**Shannon Index Diversity:** The equivalent to the gene diversity, but this time cast in information theory, is the Shannon index borrowed from community ecology. Bands can be counted as we count species in lake and a global value can be calculated for a population as:

$$H = - \sum_i^l p_i \ln p_i$$

Sometimes we see this index estimated for co-dominant data. One drawback of this measure is that is not bounded, so values vary from population to population and comparisons are difficult, not as for  $H_e$  whose values are between 0 and 1.

**Inbreeding:** Inbreeding is both the process of reproduction between related individuals, and the result of this type of reproduction. The coefficient of inbreeding, denoted  $F_{is}$  or  $f_{is}$  or simply  $f$ , is a measure of consanguinity, and estimates the probability that within a locus from a given individuals, both alleles are the same, and more importantly, have originated from the same ancestor. It is measured as:

$$F_{is} = (H_e - H_o)/H_e = 1 - H_o/H_e$$

As evident from the above formula, the inbreeding coefficient measures a departure of genotype frequencies from a reference population (a so called Hardy-Weinberg population). When both are the same, or  $H_o = H_e$ , the inbreeding coefficient is 0, and we would say that no significant departures from HW were observed.

Significant deviations from HW, i.e.,  $f_{is}$  significantly greater than zero, can arise for a number of reasons that are not mutually exclusive, mainly:

- *Small population* size that entails the loss of heterozygotes just by chance (genetic drift) and increases the probability of mating with related individuals;
- *Non-random mating* that favours the replication of the same genotypes in the population;
- *Selfing* (plants and certain snails), which is a form of non-random mating
- *Lack of external gene flow*, without migration, alleles will be fixed just by chance in small, isolated populations.

Testing for significant inbreeding is performed with different tests (i.e., fisher's exact tests), but many programmes rely in permutation tests to find a numerical solution for it. For example, FSTAT reshuffles alleles within loci to create a null distribution of possible  $f_{is}$  values from the data, and then compares if the observed value is at one or the other extreme of this distribution that is centred approximately at zero. If the observed  $f_{is}$  is in one of the extremes that contain 2.5 % of the simulated data (a 5% two-sided test), we would conclude that the  $f_{is}$  is true value greater than zero, and not a random result.

### 16.3. Genetic structure

In section 11.2, we saw a series of descriptive genetic diversity parameters, that summarizing are:  $A$ ,  $A_e$ ,  $H_o$ ,  $H_e$ , and  $f_{is}$ . When we have two or three populations, comparisons are feasible, but things can be more complicated for more samples. Moreover, we could begin to loose information, even with few populations, because the measures of inbreeding, for example, are performed with population-specific data that does not tell us anything about the relative value of diversity, or inbreeding of *all* populations.

As a definition, genetic structure refers to the *non*-random distribution of genetic diversity in space and time.

#### 16.3.1. Nei's population genetics parameters: $G_{st}$ family

Casting our question in terms of  $H$ 's or genetic diversities only, we might ask how is the total genetic diversity related to the average sup-population diversity? In other words, has the total population more information than that existing in a single population? Or, are all populations the same?

To answering these questions, Nei developed in 1972 a synthetic parameter called  $G_{st}$ . This parameter takes the value of zero, if all sub-populations contain the same information as the total population, and greater than zero and up to one (rarely achieved), if any of the sub-populations contains levels of diversity that are not distributed at random among the sup-populations.

Its computation is rather straight forward and follows the equation:

$$G_{st} = (H_t - H_s)/H_t = 1 - H_s/H_t$$

Where  $H_t$  is the total population diversity (computed from the average allele frequencies from all subpopulations) and  $H_s$  is the average within population diversity computed for each single population. It is clear that if both values are the same,  $G_{st}$  approaches zero. If not, if  $H_t$  is much larger than  $H_s$ , we would say that the distribution of genetic diversity is not random, or is structured.

#### 16.3.2. Sewall Wright's F-statistics

If instead of thinking of diversity, but inbreeding, or better correlation of alleles within *Individuals*, *Subpopulations* and the *Total* population, a set of relationships can be deduced for the different levels at which genes occur (individuals, subpopulations and the total population, of course). Thus, the inbreeding coefficient that we saw earlier for a single population can be "scaled" to different levels of population organization and different inbreeding coefficient can be used. Thus, we can ask ourselves about of:

- the correlations of gametes within individuals relative to the subpopulation, or  $F_{IS}$ ;
- the correlations of gametes within individuals relative to the total population, or  $F_{IT}$ ;



- the correlations of gametes within subpopulations relative to the total population, or  $F_{ST}$ .

If any of these correlations is  $\gg 0$ , it means that the probability of finding two identical alleles is stronger in the subunit (individual or subpopulation) than in the reference population (subpopulations and total population). Note that in principle, all these values are between zero and one, closeness to one meaning fixation of alleles at the particular scale. Note also that capital letters have been used to distinguish these parameters from single-population parameters. They are related by the expression:

$$(1 - F_{IT}) = (1 - F_{IS}) \dots (1 - F_{SR}) (1 - F_{..}) \dots (1 - F_{IS})$$

Where  $F_{SR}$  and  $F_{..}$  have been introduced between  $F_{IS}$  and  $F_{ST}$  to denote that population structure can be more complex and include regions, watersheds, *etc.*

The two most common used statistics are  $F_{IS}$  and  $F_{ST}$ , but  $F_{IT}$  has been overshadowed by the rest. Note also that for Nei's  $G$ -statistics, there are equivalent  $G_{is}$ ,  $G_{it}$ , but are less and less used.

$F_{st}$  is commonly regarded as *the* population structure parameter that if significantly greater than zero indicates that diversity (or inbreeding) is not randomly distributed. Several other parameters, however, have been proposed by different authors and the list grows almost every year. We will highlight some of the most used:

- Weir's and Cockerham's  $\Theta$  (theta), also now as the co-ancestry coefficient. Reputedly more robust to sampling variation than the basic  $F_{ST}$ .
- Excoffier's et al.  $\Phi$  (*Phi*)-statistics, that are analogous to  $F_{ST}$ , but based on variance components analyses.
- $R_{ST}$  (with its estimator  $\rho$  (*rho*)) that uses the actual microsatellite size to estimate the genetic structure parameter. Note, if microsatellites are coded as allelic states, we would be estimating *Phi*-statistics.
- $N_{ST}$ , analogous to the others, but for sequencing data (seldom used, more of a theoretical value).

## 16.4. Population and individual divergence and phylogenetic trees

So far, we have seen that a complete description of genetic diversity entails first, the estimation of various descriptive parameters for each subpopulation, and then, the use of synthetic values that will allow us to tell if genetic diversity is distributed at random or not (*i.e.*,  $F_{st} \gg 0$ ). However, can we tell apart which population(s) is actually producing this structure? Which populations are more divergent than others, and in which direction?

These questions are then answered by using a divergence analysis based on genetic distances. Strictly speaking, unless we use particular methods that can validate a direction of evolutionary changes (uses of out-groups, identification of ancestral characters or states, *etc.*) we would be doing phenetic analysis. This means that we are able to pinpoint out the

separation of populations, or individuals, but we cannot know which end of the “phylogenetic” tree precedes the rest. In crop improvement, however, this is not usually a big problem as groups are arbitrarily chosen and what matters is what is different from the others.

Similarly as for genetic structure (see section 3), there exist several ways of estimating individual or population genetic distances, but the procedure is always the same:

- Define a distance metric.
- Calculate distances among groups or among individuals (results are usually stored in a pairwise matrix of genetic distances whose diagonal is zero). If possible, bootstrap loci or individuals (*i.e.*, resample information to validate observed results) to get a support for the branches of the tree.
- Visualize the resulting distance using a particular algorithm.

In our case, the two most used algorithm for visualizing distances among groups are UPGMA (Un-weighted Pair Group Method with Arithmetic Mean) and Neighbor-joining. The former is the simplest method of tree construction. It was originally developed for constructing taxonomic phenograms, *i.e.* trees that reflect the phenotypic similarities between species, but it can also be used to construct phylogenetic trees if the rates of evolution are approximately constant among the different lineages. The latter, Neighbor-joining (Saitou and Nei, 1987) is a method that is related to the cluster method but does not require data whose lineages have diverged by equal amounts.

Common genetic distances include:

- Nei’s genetic distance (Nei, 1972);
- Cavalli-Sforza chord measure (Cavalli-Sforza and Edwards, 1967)
- Reynolds, Weir, and Cockerham’s genetic distance (1983).

These types of analyses are well handled by the set of program PHYLIP, and also by POPULATIONS, although any software that can produce a distance matrix will be useful for producing a tree. Testing of the branches and tree structure, however, is a delicate task and is mostly the domain of phylogenetics instead of population genetics, although the two fields overlap.

## 16.5. Web resources and software – non-exhaustive

FSTAT:

<http://www2.unil.ch/popgen/softwares/fstat.htm>

**Pros:** General purpose diploid analysis software with not so difficult data file. Nice interface, very good help files and handles most of the necessary analyses. Output files are also good, almost ready to use.

**Cons:** doesn’t perform nested Fst analyses. Does not report per population  $H_o$ (!).

GenePop on the Web:

<http://wbiomed.curtin.edu.au/genepop/>

**Pros:** Frequently updated, includes many tests for the significance of inbreeding, available everywhere through the web.

**Cons:** doesn't perform nested Fst analyses either. Output tables are awful and confusing. Ho is reported not as a fraction, but as the count (observed and expected) of heterozygote individuals.

Arlequin:

**Pros:** so far, the most comprehensive software devoted for population genetics. Does handle nested Fst (or hierarchical AMOVA's). Excellent manual that serves as a summary of population genetic methods, highly recommended!

**Cons:** one of the worst data file format ever! This has been circumvented by the automatic translation by other software, to certain limits. Interface apparently simple, but results are mixed with original data files, becoming confusing after many runs.

AFLPsurv:

<http://www.ulb.ac.be/sciences/lagev/aflp-surv.html>

**Pros:** I have yet to see a dominant marker program that convinces me, but this is a workable one. Includes many genetic distances and calculates genetic diversity.

**Cons:** Bootstrapping for individuals is restricted as it is population oriented software.

PHYLIP:

<http://evolution.genetics.washington.edu/phylip.html>

**Pros:** this is a collection of programs, and is somewhat the dean of phylogenetic analyses. Has been overshadowed by PAUP, but as free software is a good starting point, and although methods are somewhat outdated, the implementation is serious.

**Cons:** as said, somewhat outdated, but good for most applications.

TreeView:

<http://taxonomy.zoology.gla.ac.uk/rod/treeview.html>

**Pros:** small and effective program for visualizing trees constructed in the PHYLIP format (i.e., out files from NEIGHBOUR, for example).

**Cons:** large trees appear sometimes not so well, no possibility of editing trees.

Populations:

**Pros:** very good collection of genetic distances for codominant markers. It can deal with dominant marker data if we use the 22-11 coding. Produces tree files directly observable with TreeView and accepts GenePop data files.

**Cons:** often times it crashes unexpectedly possibly because of missing data or repeated individual names within populations.

RstCalc:

<http://helios.bto.ed.ac.uk/evolgen/rst/rst.html>

**Pros:** good programme for estimating Rst.

**Cons:** Data file is not difficult, but could be simpler. It does not handle nested Rst.

CONVERT:

<http://www.agriculture.purdue.edu/fnr/html/faculty/Rhodes/Students%20and%20Staff/glaubitz/software.htm>

**Pros:** little programme that uses a simple excel file that can be translated into other software, including GenePop and Arlequin.

**Cons:** does not support FSTAT, so passing through GenePop is necessary.

## COOK BOOK FORMATTING POPULATION GENETIC DATA

### (1) Step 1:

#### Scoring the data

Record data in excel file and transform as necessary. For the programme populations, to be used in our demonstration, a 2-digit formatting is required.

- For dominant markers, ISSR, AFLP, IRAP or others scored as present or absent, *i.e.* 1 or 0, transform as follows:  
Manually select all data input (taking care not to select the names of the individuals, populations or loci)  
First, replace all '1' with '22'  
Second, replace all '0' with '11'  
(At this point it is helpful to check for missing data)
- For codominant markers e.g. SSR data are already scored as 2 digits so no need for transformation
- For mixture of dominant and codominant markers, transform the dominant to codominant by scoring as 2-digit

### (2) Step 2:

#### Formatting the data for populations programme

- Insert a new row between the header row (*i.e.* A, B, C, ...) and first row such that newly inserted row becomes row no. 1 and then do the following in the new row (*i.e.* Row No. 1):
  - First column: type in the number of populations or samples
  - Second column: type in the number of loci or markers
  - Third column: type in the highest number used to label an allele
  - Fourth column: data coding type [*1 if the code for alleles is one digit number (1-9); 2 if code for alleles is a 2 digit number (01-99) or a 3 if code for alleles is a 3 digit number (001-999)*]
- Insert another row between now rows 2 and 3 and do the following:
  - In the first column, type '**pop**'

### (3). Step 3:

#### Formatting the data as a "tab delimited text file"

- Select all entries by highlighting (starting from cell A1X1 to the end of the data entries)
- File > Save as > text (tab delimited) (\*.txt) > OK > Yes.
- Save on same disk and folder as the Populations.exe file (To run the programme, the text file (.txt) must be in same folder as 'Populations.exe'.

### (4). Step 4:

#### Formatting in NOTEPAD:

- Open NOTEPAD
- From File menu, locate the saved .txt file, open file

- Put cursor in front of first locus and hit backspace so that it is now in the first column, second row
- Highlight all entries by select all in Edit menu
- Cut (the entries)
- Paste in Word
- Select All
- Edit > Replace > In “*find what*” box type “^t” and in “*replace with*” box, hit the space bar once. Select ‘replace all’ option. All the tabs have been replaced. (It helps to have the paragraph icon on in order to see that there are only single spaces).
- Delete the dots (after the figures in the first row and after ‘pop’ and insert comma each sample name. Make sure that there are no spaces within a sample name.
- Select all entries
- Cut (the entries)
- Paste again in NOTEPAD
- Put the cursor in front of the first data in each row and hit backspace (the space between the ‘comma’ after the sample name and the first score is deleted)
- Save (Use a simple file name – one word).
- Save the **.txt file** in the same folder as the Programme, ‘Populations.exe’.

#### (5). Step 5:

##### Running the programme

- Open program and choose sequentially by entering the corresponding numbers and hitting ‘Enter’:
- Compute individuals distance + tree (when data has only one population) – No. 1
- Type the exact name of .txt file from last saving in the space provided. The ‘.txt’ extension must be included in the name. The name is also case sensitive.
- Phylogenetic tree of individuals with bootstraps on locus – No. 3
- Nei’s standard genetic distance, Ds (1972) – No. 2
- UPGMA – No. 1
- 10000
- Enter desired name for output file with ‘.tre’ extension
- Wait for the programme to finish running. The output file with the ‘.tre’ extension is now deposited in the same folder as the programme, ‘Populations.exe’
- Double click on the output file with the ‘.tre’ extension in order to see the resulting dendrogram.

## 16.6. References

- Cavalli-Sforza, L. L.; Edwards, A. W. F., 1967: Phylogenetic analysis: models and estimation procedures. *Am. J. Hum. Genet.* 18, 233-257.
- Chakraborty R and Danker-Hopfe H, 1991. Analysis of population structure: A comparative study of different estimators of Wright's fixation indices. In 'Statistical Methods in Biological and Medical Sciences.' Ed C.R. Rao and R. Chakraborty, Elsevier Science Publishers.

- Cockerham CC, 1969. Variance of gene frequencies. *Evolution*. 23:72-84.
- Cockerham CC, 1973. Analysis of gene frequencies. *Genetics*. 74:679-700.
- Cockerham CC and Weir BS, 1993. Estimation of gene-flow from F-statistics. *Evolution*. 47:855-863.
- El Mousadik A and Petit RJ, 1996. High level of genetic differentiation for allelic richness among populations of the argan tree [*Argania spinosa* (L.) Skeels] endemic to Morocco. *Theor. Appl. Genet.* 92:832-839.
- Excoffier L 2001. Analysis of population subdivision. In *Handbook of statistical genetics*, Balding, Bishop & Cannings (Eds) Wiley & Sons, Ltd.
- Fisher R, 1954. *Statistical Methods for Research Workers*. 12th Edition, Oliver & Boyd, Edinburgh. 356pp.
- Goodman SJ, 1997. Rst Calc: a collection of computer programs for calculating estimates of genetic differentiation from microsatellite data and a determining their significance. *Molecular Ecology* 6: 881-885.
- Glaubitz, J.C. (submitted) CONVERT: A user-friendly program to reformat diploid genotypic data. *Molecular Ecology*.
- For commonly used population genetic software packages. *Molecular Ecology Notes*.
- Goudet J, 1995. FSTAT (vers. 1.2): a computer program to calculate F-statistics. *J. Hered.* 86: 485-486.
- Goudet J, Raymond M, Demeus T and Rousset F, 1996. Testing differentiation in diploid populations. *Genetics*. 144:1933-1940.
- Hartl DL, Clark AG (1997) *Principles of Population Genetics*. Third Edition. Sinauer Associates.
- Nei, M. (1972) Genetic distance between populations. *Am. Nat.* 106:283-292.
- Nei M, 1973. Analysis of gene diversity in subdivided populations. *Proc. Natl. Acad. Sci. USA*. 70:3321-3323.
- Nei M, 1988. *Molecular Evolutionary Genetics*. Columbia University Press, New York.
- Nei M and Chesser RK, 1983. Estimation of fixation indices and gene diversities. *Ann. Hum. Genet.* 47:253-259.
- Pamilo P, 1984. Genotypic correlation and regression in social groups: multiple alleles, multiple loci and subdivided populations. *Genetics*. 107:307-320.
- Petit RJ, El Mousadik, A and Pons O, 1998. Identifying populations for conservation on the basis of genetic markers. *Conservation Biology*. 12:844-855.
- Queller DC and Goodnight KF, 1989. Estimating relatedness using genetic markers. *Evolution*. 43:258-275.
- Raymond M. & Rousset F, 1995. GENEPOP (version 1.2): population genetics software for exact tests and ecumenicism. *Journal of Heredity*, 86, 248-249.
- Raymond M and Rousset F, 1995. An exact test for population differentiation. *Evolution*. 49:1280-1283.
- Reynolds J, Weir BS and Cockerham CC, 1983. Estimation of the coancestry coefficient: Basis for a short-term genetic distance. *Genetics*. 105:767-779.
- Rousset F, 1996. Equilibrium Values of Measures of Population Subdivision For Stepwise Mutation Processes. *Genetics* 142:1357-1362.
- Rousset F, 1997. Genetic differentiation and estimation of gene flow from F-statistics under isolation by distance. *Genetics* 145:1219-1228.

- Saitou, N and M Nei, 1987. The neighbor-joining method: a new method for reconstructing phylogenetic trees. *Mol Biol Evol* 4: 406-425.
- Slatkin M, 1993. Isolation by distance in equilibrium and non-equilibrium populations. *Evolution* 47:264-279.
- Slatkin M, 1995. A measure of population subdivision based on microsatellite allele frequency. *Genetics*. 139:457-462.
- Slatkin M and Barton NH, 1989. A comparison of three methods for estimating average levels of gene flow. *Evolution* 43:1349-1368.
- Sokal RR and Rohlf FJ, 1981. *Biometry*. 2nd Edition. Freeman & Co.
- Weir BS and Cockerham CC, 1984. Estimating F-statistics for the analysis of population structure. *Evolution* 38:1358-1370.
- Weir BS, 1996. *Genetic data analysis II*. Sinauer Publ., Sunderland, MA.
- Whitlock MC and McCauley D, 1999. Indirect measures of gene flow and migration:  $F_{st} \approx 1/(4Nm+1)$ . *Heredity*. 82: 117-125.
- Wright S, 1969. *Evolution and the genetics of populations*. Vol. 2. The theory of gene frequencies. University of Chicago Press.

## 16.7. Some key concepts

**Alleles:** All possible forms of a gene.

**Gene:** A unit of inheritance, a non-recombining segment of DNA. A given location on a chromosome

**Genotype:** The combination of the two homologous alleles carried on the two chromosomes of a diploid individual at a given locus.

**Haplotype:** A particular combination of alleles at different loci on a chromosome.

**Heterozygosity:** The probability of an individual to have two different alleles at a given locus (the probability of being heterozygote).

**Homozygosity:** The probability of an individual to be homozygote at a given locus.

**Homozygote:** The fact that an individual has two identical alleles at a given locus.

**Locus:** A given location on a chromosome, a non-recombining segment of a chromosome (usually interchanged with gene)

**Phenotype:** The visible (physical) state of an individual. The relation between the genotype and the phenotype can be complex and will usually depend on the degree of dominance and the interaction of different alleles at a single or multiple loci.

**Polymorphism:** the fact that there exist different alleles at a given locus in a population.

**Population:** A group of interbreeding individuals living together in time and space. It is usually a subdivision of a species.

**Sample:** A collection of individuals or of genes drawn from a population.



## 16.8. Equations

$$p = P + H/2 \quad \hat{H}_E = \frac{2n}{2n-1} \left( 1 - \sum_{i=1}^k p_i^2 \right) \quad H_e = 2pq - \frac{2Spq}{2-S}$$

$$\frac{Q_f}{Q} = \frac{q^2 + fpq}{q^2} = 1 + \frac{fp}{q} \quad H = 2pq(1-f) \quad q_2 = \frac{1}{2}(q_f + q_m)$$

$$\text{var}(p') = \frac{pq}{2N} \quad H_t = 2p_0q_0 + A^t(Q_0 - 2p_0q_0)$$

$$P(X=r) = \frac{2N!}{r!(2N-r)!} p^r (1-p)^{2N-r} \quad f = \sum_{i=1}^m \left(\frac{1}{2}\right)^{N_i} (1 + f_{CA(i)})$$

$$f(t) = 1 - \left(1 - \frac{1}{2N}\right)^t (1 - f(0)) \quad H(t) = \left(1 - \frac{1}{2N}\right)^t H(0)$$

$$t \approx -\ln\left(\frac{H(t)}{H(0)}\right) 2N \quad \text{var}(p(t)) = p_0q_0 \left(1 - \left(1 - \frac{1}{2N}\right)^t\right)$$

$$D = \frac{t}{2N} \approx -\ln\left(1 - \frac{\text{var}(p)}{\bar{p}\bar{q}}\right) \quad N_e = \frac{4N_m N_f}{N_f + N_m} \quad N_e = \frac{N\bar{k} - 1}{\bar{k} - 1 + \frac{\text{var}(k)}{\bar{k}}}$$

$$\frac{1}{N_e} = \frac{1}{t} \sum_{i=0}^{t-1} \frac{1}{N(i)} \quad N_e = \frac{N}{1+f} \quad N_e = 4\pi\sigma^2 d$$

$$\bar{H}_S = 2\bar{p}(1-\bar{p}) - 2\text{var}(p) \quad F_{ST} = \frac{\text{var}(p)}{\bar{p}(1-\bar{p})} \quad (1-x)^y \approx e^{-xy} \text{ if } x < 0$$

$$D = \frac{t}{2N} \approx -\ln(1 - F_{ST}) \quad (1 - F_{IT}) = (1 - F_{IS})(1 - F_{ST}) \quad \bar{F}_{IT} = \frac{\bar{H}_T - \bar{H}_0}{\bar{H}_T}$$

$$\bar{F}_{ST} = \frac{\bar{H}_T - \bar{H}_S}{\bar{H}_T} \quad \bar{F}_{IS} = \frac{\bar{H}_S - \bar{H}_0}{\bar{H}_S} \quad F_{ST} \approx \frac{1}{4Nm + 1}$$

$$F_{ST} \approx \frac{1}{4Nm \left( \frac{k}{k-1} \right)^2 + 1}$$

$$N_e \approx \frac{kN}{1 - F_{ST}}$$

$$p' = p \frac{w_1}{\bar{w}}$$

$$q' = q \frac{\bar{w}_2}{\bar{w}}$$

$$\bar{w} = p^2 w_{11} + 2pq w_{12} + q^2 w_{22}$$

$$\bar{w}_2 = q w_{22} + p w_{12}$$

$$\Delta q = q \frac{\bar{w}_2 - \bar{w}}{\bar{w}}$$

$$q_e = \frac{u}{u + v}$$

$$p_t = (1 - u)^t p_0$$

$$u(p) = p_0 = \frac{1}{2N}$$

$$u(q) = 1 - u(p) = 1 - \frac{1}{2N}$$

$$T_1(p) = 4N_e$$

$$T_0(p) = 2 \frac{N_e}{N} \ln(2N)$$

$$F_e \approx \frac{1}{4Nu + 1} = \frac{1}{\theta + 1}$$

$$H_e = \frac{\theta}{\theta + 1}$$

$$u\left(\frac{1}{2N}\right) \approx 2s$$

$$T_1(p) = \frac{2}{s} \ln(2N)$$

$$r = \frac{1}{u}$$

$$r = \frac{1}{4Nus}$$

$$E(k|n, \theta) = 1 + \theta \sum_{i=1}^{n-1} \frac{1}{\theta + i}$$

$$K = -\frac{3}{4} \ln\left(1 - \frac{4}{3}d\right)$$

## 17. APPENDICES

### 17.1. General DNA extraction techniques

#### 17.1.1. Phenol/chloroform extraction

NOTE: Wear gloves, goggles, and lab coat at all times for safety and to prevent contamination of your preparations.

Removes protein from DNA preparations. Advisable for example if  $A_{260\text{nm}}$ :  $A_{280\text{nm}}$  (from the spectrophotometer readings) of the DNA are below 1.6. Phenol extraction requires subsequent ethanol precipitation of the DNA.

Phenol: freshly distilled and equilibrated with 20 % 0.5 M Tris-Base. Prepare a mixture of phenol/chloroform/isoamylalcohol (PCI) (25:24:1).

NOTE: Use caution as phenol is toxic.

1. The DNA sample is mixed with an equal volume of PCI, vortexed, and centrifuged for about 5 minutes. Remove the upper aqueous phase avoiding contamination with protein from interphase and transfer it to a fresh reaction tube.
2. Remaining traces of phenol in the aqueous phase are extracted with 1 volume of chloroform/isoamylalcohol (24:1). Vortex and centrifuge for 5 minutes. Transfer the upper phase carefully to a fresh reaction tube.

#### 17.1.2. Ethanol precipitation

NOTE: Wear glasses at all time for safety.

1. Determine volume of the sample, add 0.1 volume 3 M sodium acetate and 2.5 volumes cold ethanol (96%). Mix well and leave at  $-20^{\circ}\text{C}$  for 2 hours.
2. Centrifuge for 15 minutes (in microcentrifuge at  $>12,000$  rpm), preferably at  $4^{\circ}\text{C}$ .
3. Carefully remove ethanol and wash pellet with cold 70% ethanol to remove salt from the sample – centrifuge for 5 minutes.
4. Dry DNA pellet in vacuum centrifuge or air dry in flow bench.
5. Dissolve DNA in TE buffer or sterile double distilled  $\text{H}_2\text{O}$  (dd $\text{H}_2\text{O}$ ).

### 17.1.3. Solutions

- 1.5 x CTAB extraction buffer (1 liter):

CTAB	15.0 g
1 M Tris (pH 8.0)	75 ml
0.5 M EDTA	30 ml
NaCl	61.425 g
ddH <sub>2</sub> O	to 1 litre

- 10% CTAB (1 litre)

CTAB	100 g
NaCl (0.7M)	40.95 g
ddH <sub>2</sub> O	to 1 litre

- β-mercaptoethanol,

- Chloroform:isoamylalcohol (24:1),

- Isopropanol,

- Ethanol 96% and 70%

- sodium acetate (3 M)

- TE buffer

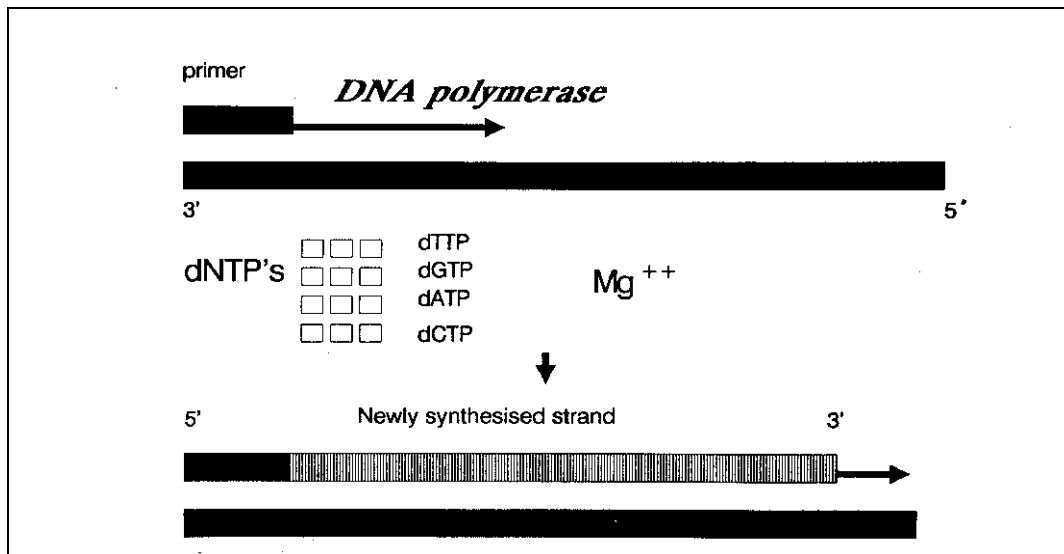
10 mM Tris HCl

1 mM EDTA (pH 8.0)

### 17.2. Polymerase chain reaction protocol

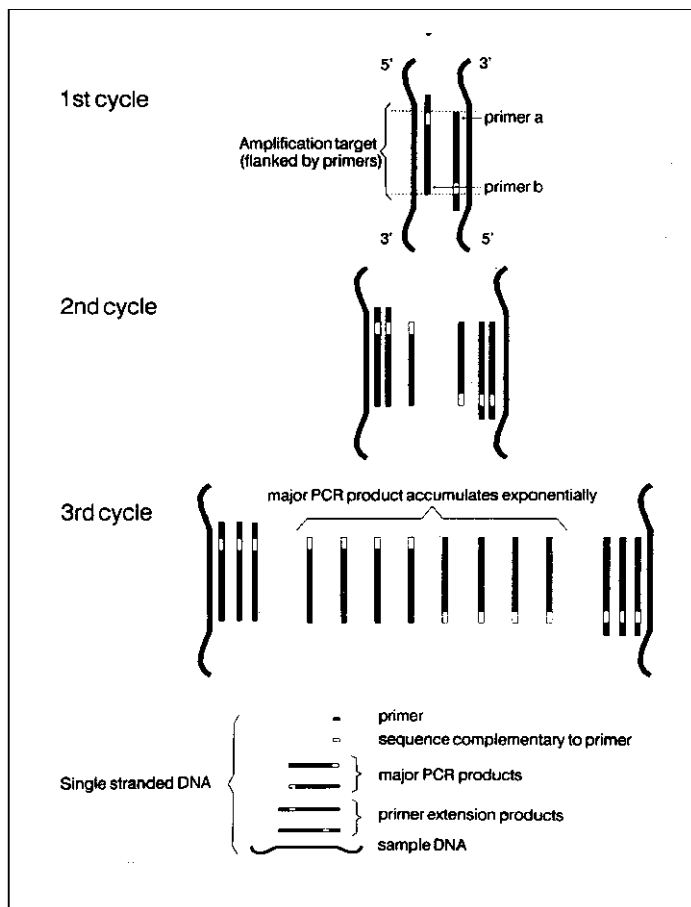
The polymerase chain reaction (PCR) is basically a technique for *in vitro* amplification of specific DNA sequences by the simultaneous primer extension of complementary DNA strands. The principle of primer extension is illustrated in Figure A.2.1 for one DNA strand. The primer binds to its complementary sequence of the single stranded target DNA and the polymerase extends the primer in 5' - 3' direction by using the complementary DNA as a template. For a PCR reaction, two primers are used, one binding to the “lower” strand (forward primer) and one binding to the “upper” strand (reverse primer). Thus, the requirements for the reaction are: template DNA, oligonucleotide primers, DNA polymerase, deoxynucleotides (to provide both energy and nucleosides for DNA synthesis), and a buffer containing magnesium ions. In general the DNA sequence of both ends of the region to be amplified must be known to be able to synthesize proper primer oligonucleotides. The PCR reaction is a cyclic process, which is repeated 25 to 35 times. One cycle consists of three basic steps with characteristic reaction temperatures:

1. Denaturation of the double stranded DNA to make the template accessible for the primers and the DNA polymerase (94°C, 30 seconds).
2. Annealing of primers to complementary sequence on template (between 45 and 60°C, depending on the primer sequences, 30 seconds).
3. Extension of primers by DNA-polymerase (72°C - the optimum temperature of *Taq* DNA-polymerase -, 1 minute per kilobase of template to be amplified).



**Figure A.2.1.** Primer extension. DNA polymerase extends a primer by using a complementary strand as a template (McPherson *et al.*, 1991).

By multiple repetition of this cycle the number of template molecules increases. This result in exponential amplification of the DNA sequence that is bordered by the two primers used (Figure A.2.2).



**Figure A.2.2.** Schematic diagram of PCR. By using primer pairs ‘a’ and ‘b’ (short black lines) annealed to complementary strands of DNA (long black lines), two new strands (shaded lines) are synthesized by primer extension. If the process is repeated, both the sample DNA and the newly synthesized strands can serve as templates, leading to an exponential increase of product which has its ends defined by the position of the primers (McPherson *et al.*, 1991).

Successful performance of a PCR experiment is dependent on a number of different factors; some of them have to be determined empirically.

- The selection of the primers is a very important step. They should be long enough to be specific, not anneal against themselves by folding (avoid palindromic sequences), nor should the forward primer anneal with the reverse primer. Furthermore the G/C content of the primers should be similar and they should have similar melting temperatures ( $T_m$ ). Several computer programs are available on the Internet to help to find the best primer pairs for a given sequence. Try the addresses below- submit the DNA sequence and some required parameters and you will get a list of possible primers:

<http://www-genome.wi.mit.edu/cgi-bin/primer/primer3> [www.cgi](http://www.cgi)  
<http://genome-www2.stanford.edu/cgi-bin/SGD/web-primer>  
<http://www.nwfsc.noaa.gov/protocols/oligoTMcalc.html>

- The annealing temperature must be determined empirically and is dependent from the  $T_m$ 's of the primers. A rule of thumb (Wallace rule) provides a first order approximation for  $T_m$  of oligonucleotides that have 20 bases or less:

$$T_m = 2^\circ\text{C} (A + T) + 4^\circ\text{C} (C + G)$$

The annealing temperature is a few degrees lower than  $T_m$ .

- PCR is extremely sensitive! Thus contamination of samples and solutions with minimal amounts of foreign DNA, or the wrong PCR programme can result in unspecific PCR products. Always include controls without template DNA in order to check if there is any contamination in your nucleotides, primers, etc.

A typical PCR experiment is given in the table below. In the FAO/IAEA course, PCR was demonstrated by amplifying a 1050 bp sequence of the rice retrotransposon Tos 17 accession number D88394:

Forward Primer 1 (100 pmol/ $\mu\text{l}$ ):

Reverse Primer 2 (100 pmol/ $\mu\text{l}$ ):

Reaction volume: 50  $\mu\text{l}$

Stock solutions	$\mu\text{l}$	Final conc./amount
10 x PCR buffer (15 mM $\text{MgCl}_2$ )	5.0 $\mu\text{l}$	1 x PCR buffer (1.5 mM $\text{MgCl}_2$ )
Primer 1 (100 pmol/ $\mu\text{l}$ )	0.5 $\mu\text{l}$	1 pmol
Primer 2 (100 pmol/ $\mu\text{l}$ )	0.5 $\mu\text{l}$	1 pmol
dNTP mix (10 mM)	1 $\mu\text{l}$	0.2 mM
DNA template (100 ng/ $\mu\text{l}$ )	1 $\mu\text{l}$	100 ng
<i>Taq</i> DNA Polymerase (5 U/ $\mu\text{l}$ )	0.5 $\mu\text{l}$	2.5 U
$\text{H}_2\text{O}$	41.5 $\mu\text{l}$	-

NOTE: It is very important to prepare a master mix corresponding to the number of desired samples that contains all the reagents except for the template DNA. Mix well and add the appropriate amount of the master solution to single reaction vials containing the individual template DNA samples you wish analysed. This procedure significantly reduces the number of pipetting steps, avoids errors derived from pipetting small amounts of liquid, and finally ensures that every tube contains the same concentrations of reagents.

For amplification of the Tos17 sequence the PCR machine was programmed as follows:

<i>Step 1</i>	Initial denaturation	94°C	(4:00 minutes)
<i>Step 2</i>	Denaturation	94°C	(0:30 minute)
<i>Step 3</i>	Primer annealing	56°C	(0:30 minute)
<i>Step 4</i>	Primer extension	72°C	(1:10 minutes)
<i>Step 5</i>	Cycling	Repeat steps 2-4	29 times
<i>Step 6</i>	Final extension	72°C	(6:00 minutes)
<i>Step 7</i>	Hold	4°C	(hold)

NOTE: The PCR programme can vary from primer to primer set and species to species with the annealing temperature being the most variable step.

### 17.2.1. References

McPherson, M., P. Quirke, and G Taylor, 1991. PCR: A Practical Approach. Oxford University Press, New York.



### 17.3. Plant genome database contact information

Table 17.3–1 Taken from an IAEA-TECDOC on “Radioactively Labelled DNA Probes For Crop Improvement” VIENNA SEPTEMBER 6-8, 1999).

DATABASE	CROPS	CURATOR	E-MAIL ADDRESS	DATABASE ADDRESS
AAtdB	<i>Arabidopsis</i>	David Flanders	flanders@genome.stanford.edu	<a href="http://genome-www.stanford.edu/Arabidopsis/">http://genome-www.stanford.edu/Arabidopsis/</a>
Alfagenes	Alfalfa ( <i>Medicago sativa</i> )	Daniel Z. Skinner	Dzolek@ksu.ksu.edu	<a href="http://naaic.org/">http://naaic.org/</a>
Bean Genes	<i>Phaseolus</i> and <i>Vigna</i>	Phil McClean	mcclean@beangenes.cws.ndsu.nodak.edu	<a href="http://probe.nalusda.gov:8300/cgi-bin/browse/beangenes">http://probe.nalusda.gov:8300/cgi-bin/browse/beangenes</a>
ChlamyDB	<i>Chlamydomonas reinhardtii</i>	Elizabeth H. Harris	chlamy@acpub.duke.edu	<a href="http://probe.nalusda.gov:8300/cgi-bin/browse/chlamydb">http://probe.nalusda.gov:8300/cgi-bin/browse/chlamydb</a>
CoolGenes	Cool season food legumes	Fred Muehlbauer	muehlbau@wsu.edu	<a href="http://probe.nalusda.gov:8300/cgi-bin/browse/coolgenes">http://probe.nalusda.gov:8300/cgi-bin/browse/coolgenes</a>
CottonDB	<i>Gossypium</i> species	Sridhar Madhavan	msridhar@tamu.edu	<a href="http://probe.nalusda.gov:8300/cgi-bin/browse/cottondb">http://probe.nalusda.gov:8300/cgi-bin/browse/cottondb</a>
GrainGenes	Wheat, barley, rye and relatives	Olin Anderson	oandersn@pw.usda.gov	<a href="http://probe.nalusda.gov:8300/cgi-bin/browse/graingenes">http://probe.nalusda.gov:8300/cgi-bin/browse/graingenes</a>
MaizeDB	Maize	Mary Polacco	maryp@teosinte.agron.missouri.edu	<a href="http://www.agron.missouri.edu/">http://www.agron.missouri.edu/</a>
MilletGenes	Pearl millet	Matthew Couchman	Matthew.Couchman@bbsrc.ac.uk	<a href="http://jiiio5.jic.bbsrc.ac.uk:8000/cgi-bin/ace/search/millet">http://jiiio5.jic.bbsrc.ac.uk:8000/cgi-bin/ace/search/millet</a>
PathoGenes	Fungal pathogens of small-grain cereals	Henriette Giese	h.giese@risoe.dk	<a href="http://probe.nalusda.gov:8300/cgi-bin/browse/pathogenes">http://probe.nalusda.gov:8300/cgi-bin/browse/pathogenes</a>
RiceGenes	Rice	Susan McCouch	srm4@cornell.edu	<a href="http://genome.cornell.edu/rice/">http://genome.cornell.edu/rice/</a>
RiceGenome Project	Rice			<a href="http://www.staff.or.jp">http://www.staff.or.jp</a>
SolGenes	<i>Solanaceae</i>	Molly Kyle	mmk9@cornell.edu	<a href="http://genome.cornell.edu/solgenes/welcome.html">http://genome.cornell.edu/solgenes/welcome.html</a>
SorghumDB	<i>Sorghum bicolor</i>	Russel Kohel/Bob Klein	nus6389@tam2000.tamu.edu	<a href="http://probe.nalusda.gov:8300/cgi-bin/browse/sorghumdb">http://probe.nalusda.gov:8300/cgi-bin/browse/sorghumdb</a>
Soybase	Soybeans	David Grant	dgrant@iastate.edu	<a href="http://129.186.26.94/">http://129.186.26.94/</a>
TreeGenes	Forest trees	Kim Marshall	kam@s27w007.pswfs.gov	<a href="http://dendrome.ucdavis.edu/index.html">http://dendrome.ucdavis.edu/index.html</a>
National Center for Genome Resources	Various			<a href="http://www.ncgr.org/">http://www.ncgr.org/</a>

#### 17.4. Acronyms of chemicals and buffers

<b>AMPPD</b>	4-Methoxy-4-(3-phosphatephenyl)spirol(1,2-dioxetan-3,2'-adamantan)
<b>BCIP</b>	5-Bromo-4-chloro-3-indolyl phosphate
<b>CSPD<sup>®</sup></b>	Chemiluminescence substrate (a registered trademark of Tropix Inc., USA)
<b>CTAB</b>	Hexadecyltrimethylammonium bromide
<b>ddH<sub>2</sub>O</b>	Double distilled water
<b>DIG</b>	Digoxygenin
<b>N<sub>2</sub> liquid</b>	Liquid nitrogen
<b>NBT</b>	Nitro blue tetrazolium
<b>PCI</b>	Phenol/chloroform/isoamylalcohol (25:24:1)
<b>SDS</b>	Sodium dodecyl sulphate
<b>SSC</b>	Saline-sodium citrate buffer
<b>TBE</b>	Tris-borate-EDTA buffer
<b>TE</b>	Tris-EDTA buffer
<b>TEMED</b>	N,N,N',N'-tetramethylenediamine
<b>TRIS</b>	[Tris(hydroxymethyl)aminomethane]

## NOTES





