

UPOV/INF/17/2 Draft 1

Original: English

Date: August 29, 2018

DRAFT
(REVISION)

GUIDELINES FOR DNA-PROFILING: MOLECULAR MARKER SELECTION AND DATABASE CONSTRUCTION (“BMT GUIDELINES”)

Document prepared by the Office of the Union

*to be considered by
the Working Group on Biochemical and Molecular Techniques, and DNA-Profiling in Particular (BMT)
at its seventeenth session, to be held in Montevideo, Uruguay,
from September 10 to 13, 2018*

Disclaimer: this document does not represent UPOV policies or guidance

Note for Draft version

Footnotes will be retained in published document.

Endnotes and **highlighted boxes** are background information when considering this draft and will not appear in the final, published document.

~~Strikethrough~~ (highlighted in grey) indicates deletion from the text of document UPOV/INF/17/1.

Underlining (highlighted in grey) indicates insertion to the text of document UPOV/INF/17/1.

TABLE OF CONTENTS

A.	INTRODUCTION	3
B.	GENERAL PRINCIPLES	3
1.	Selection of a Molecular Marker Methodology	3
2.	Selection of Molecular Markers.....	4
2.1	<i>General Criteria</i>	4
2.2	<i>Criteria for specific types of molecular markers</i>	6
3.	Access to the Technology.....	9
4.	Material to be Analyzed	10
4.1	<i>Source of plant material</i>	10
4.2	<i>Type of plant material</i>	11
4.3	<i>Sample size</i>	11
4.4	<i>DNA reference sample</i>	11
5.	Standardization of Analytical Protocols.....	12
5.1	<i>Introduction</i>	12
5.2	<i>Quality criteria</i>	12
5.3	<i>Evaluation Phase</i>	13
5.4	<i>Scoring of molecular data</i>	14
6.	Databases	16
6.1	<i>Type of database</i>	17
6.2	<i>Database model</i>	17
6.3	<i>Data Dictionary</i>	18
6.4	<i>Table Relationship</i>	19
6.5	<i>Transfer of data to the database</i>	20
6.6	<i>Data access / ownership</i>	20
6.7	<i>Data analysis</i>	20
6.8	<i>Validating the database</i>	20
6.9	<i>Data Exchange</i>	20
7.	Summary	24
GLOSSARY		26
	Microsatellites, or Simple Sequence Repeats (SSRs).....	26
	Single Nucleotide Polymorphisms (SNPs).....	26
	Cleaved Amplified Polymorphic Sequences (CAPS)	26
	Sequence-Characterized Amplified Regions (SCARs)	26
	Pig-tailing.....	26
	Null Allele.....	27
	Stutter Bands.....	27

A. INTRODUCTION

The purpose of this document (BMT Guidelines) is to provide guidance ~~for developing harmonized methodologies to standardize criteria for the use of DNA based markers~~ⁱ with the aim of generating high quality molecular data for a range of applications. The BMT Guidelines are also intended to address the construction of databases containing molecular profiles of plant varieties, possibly produced in different laboratories using different technologies. In addition, the aim is to set high demands on the quality of the markers and on the desire for generating reproducible data using these markers in situations where equipment and/or reaction chemicals might change. Specific precautions need to be taken to ensure quality entry into a database.

B. GENERAL PRINCIPLES

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To add the following text:

Molecular markers sets and subsequently databases developmental process can be subdivided into 6 different phases:

1. Selection of molecular markers
2. Selection of detection method
3. Evaluation of the selected markers set and detection method (fit for purpose validation of the marker set and technological validation of the method)
4. Harmonization and validation of the method
5. Construction of the database.
6. Management of the database

This document describes these different phases in more detail. It is considered that these phases are independent on the stage of development of genotyping technologies and future improvements in high-throughput sequencing.

1. Selection of a Molecular Marker Methodology

1.1 Important criteria for choosing a methodology are:

- (a) reproducibility of data production between laboratories and detection platforms (different types of equipment);
- (b) repeatability over time;
- (c) discrimination power;
- (d) possibilities for databasing; and
- (e) accessibility of methodology.

1.2 As improvements in technology and new equipment become available, it is important for the continued sustainability of databases that the interpretation of the data produced are independent of the equipment used to produce them. This is, for example, the case with DNA sequencing data. Initially, radioactively labeled primers and sequencing gels were used to produce such data, whereas this can now be done using fluorescent dyes followed by separation on high throughput, largely automated, capillary gel electrophoresis systems, real time based techniques and next generation sequencingⁱ.

1.3 Despite these differences, the data produced with the various techniques are consistent with each other and independent of the techniques used to produce them. This can also apply to data produced using, e.g. DNA microsatellites (simple sequence repeats, SSR) or Single Nucleotide Polymorphisms (SNPs). This repeatability and reproducibility is important in the construction, operation and longevity of databases and is very important in generating a centrally maintained database, populated with verified data from a range of sources.

1.4 The molecular techniques readily applicable for variety profiling are constrained by the requirement for the data to be repeatable, reproducible and consistent. Thus, while various multi-locus DNA profiling techniques have been successfully used for research, co-dominance cannot easily be recorded in many of them, and the reproducibility of complex banding patterns between laboratories using different equipment can be problematic.

1.5 These factors present difficulties in the context of variety profiling. Consequently, this document focuses on considerations and recommendations with regard to the well-defined and researched uses of SSRs (microsatellites) and, ~~for the future,~~ⁱ to sequencing information (i.e. single nucleotide polymorphisms, SNPs). Other techniques which rely on DNA sequence information, such as cleaved amplified polymorphic sequences (CAPS) and sequence-characterized amplified regions (SCARs) may also fulfill the above criteria but their use in DNA profiling of plant varieties has not yet been explored.

Comments by Spainⁱⁱⁱ

The document still applies in all the related to SSRs, but it should be broadened to other techniques, mainly SNPs but also those derived from NGS, as their use is now much more extended that at the moment of the writing of the document, in 2010.

Some of the principles can be extended to any other DNA based technology (point 3, access to the technology; 4, material to be analyzed; 5, standardization of analytical protocols, excepting quality criteria, only useful for technologies base on PCR; or 6, databases).

We understand that to do this task it would be useful to delegate in a small working group, that could prepare a review of the new techniques that must be included.

Comments by ESA^{iv}

Several of the marker systems described in the document are already quite "old-fashioned". We are not really in favour of setting up marker databases on the basis of SSR's. The SSR technique has indeed some advantages, but it is an expensive technique with low throughput and it is highly dependent from the equipment that is used.

The use of SNP-markers has a preference, but it will be difficult to choose a SNP technique that all parties involved will apply. The choice for a specific SNP technology also strongly depends from the number of markers that needs to be analysed. For genotyping of bigger number of markers, the use of SNP-chips or sequence-based genotyping technology could be more appropriate, but no reference is made to these technologies at all.

Is this ["but their use in DNA profiling of plant varieties has not yet been explored"] still valid?

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To delete this section 1. "Selection of a Molecular Marker Methodology".

2. Selection of Molecular Markers

Joint comments from the European Union, France and the Netherlandsⁱⁱ

The title of this section to read "1. Phase 1: SELECTION OF MOLECULAR MARKERS".

2.1 General Criteria

The following general criteria for choosing a specific marker or set of markers are intended to be appropriate for molecular markers irrespective of the use of the markers, although it is recognized that specific uses may impose certain additional criteria:

- (a) useful level of polymorphism;

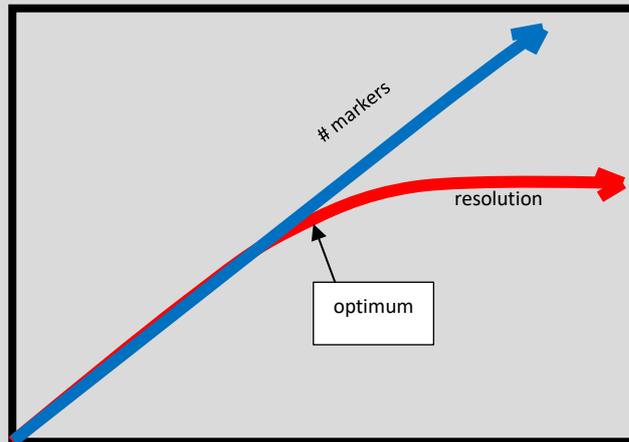
Comments by ESAⁱⁱⁱ

What does "useful" mean? Is there meanwhile any data available that can be used to frame the term in a more precise way?

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To add the following text and figure:

A balance between the number of markers and the resolution or discriminative power should be reached (Cf. figure). The appropriate number of markers should be defined to reach the required discriminative power (optimum). This number is marker-type specific, species-dependent and can also be dependent on the purpose of the genetic analysis. The quality of the markers used is also to take into account: the more high quality markers are used, the more 'error-tolerant' the method will be (e.g. the impact of a single false allele score has a limited effect when the number of markers is high). However, a high number of markers is not a guarantee for a better analysis. Markers with a high error-rate are better left out, since they could hamper the quality of the analysis. Thus, the optimum number of markers should be determined also in respect to error-rate.



- (b) repeatability within, and reproducibility between, laboratories in terms of scoring data;
(c) known distribution of the markers throughout the genome (i.e. map position), which whilst not being essential, is useful information and helps to avoid the selection of markers that may be linked; and

Joint comments from the European Union, France and the Netherlandsⁱⁱ

Section 2.1 (c) to read "Good coverage of the genome. Knowing the position of the selected markers on the genome (i.e. map position) is not essential but allow avoiding the selection of markers that may be linked together."

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To add a new section 1.1 (d) as follows:

- (d) Markers publically accessible, commercial, and/or newly developed depending on group/crop/species;
- Derive molecular markers from reliable public resources (e.g. peer reviewed publications and public databases NCBI, EMBL). This is the easiest and cheapest approach.
 - Commercial molecular marker sets to screen and select a suitable marker set from there (e.g. species-specific chips and arrays). A special attention should be paid to this source of markers as information on position, sequences... for each marker will highly depend on the providers and some might be missing;
 - Generate own sequence data. Sequence data will be generated from the selected varieties that fulfil the requirements mentioned above. However, in general to reduce time and cost, only a subset of varieties will be chosen to detect polymorphism and select the different markers. There are several options to obtain sequence information ranging from sequencing just a part of the genome to sequencing complete genomes of the selected varieties. The investigation needed in money and effort is dependent on the complexity of the genome of the particular species and the sequence data available in the public domain (e.g. reference genome). Many species contain a high level of genetic (botanic) diversity. There is no need to obtain sequence data of the complete genome; the data for only a small part of the genome may be sufficient to develop a suitable marker set for genotyping, depending on the application (e.g. genome-wide sequence capture, transcriptome sequencing (only the coding part of the genome). Presence or absence of a published reference genome is relevant to allow position determination of markers. It will be a matter of time before for a reference genome is available for all economically valuable plant species.

- (d) the avoidance, as far as possible, of markers with “null” alleles (i.e. an allele whose effect is an absence of a PCR product at the molecular level), which again is not essential, but advisable.

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To add new sections 1.1 (f) to (k) as follows:

- (f) Allowance of easy, objective and indisputable scoring of marker profiles. These markers are preferred over complex marker profiles that are sensitive to interpretation. Clear black and white answers also allows for easier harmonization;
- (g) Co-dominant markers are preferred over dominant markers as they have a higher discriminative power;
- (h) Avoidance of linkage disequilibrium;
- (i) Durability of the marker. When a marker is located in a genomic area that is not subject to selection by breeders, there is a better chance that the marker will be informative in a durable way;
- (j) Markers located in coding and/or in non-coding regions and potentially epigenetic markers; and
- (k) The use of molecular markers is species-specific and should take into account the features of propagation of the species.

2.2 *Criteria for specific types of molecular markers*

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To delete section 2.2.

2.2.1 Microsatellite Markers

Comments by ESAⁱⁱⁱ

We are not really in favour of setting up marker databases on the basis of SSR's. The SSR technique has indeed some advantages, but it is an expensive technique with low throughput and it is highly dependent from the equipment that is used. Furthermore a technique like SSR – but also the others – has the weakness when used in polyploid species that they have to be able to distinguish between the genomes.

2.2.1.1 The analysis of simple sequence repeats (SSRs or microsatellites: see Glossary) using the polymerase chain reaction (PCR) is now widely used and has several advantages.

2.2.1.2 SSR markers are expressed co-dominantly, are generally easy to score (record) and can readily be mapped. They have been used and analyzed in different laboratories, and under specific experimental conditions are generally robust and repeatable. In addition, they can be analyzed using automated, high throughput, non-radioactive DNA sequencers, based either on gel electrophoresis or capillary electrophoresis, and several can be analyzed simultaneously (multiplexing).

Proposal by Argentinaⁱ

[...] to add NGS as a technique for SSR scoring.

Comments by Ecuador^v

It is undesirable that the sets of varieties must be defined and included in all analyses and using the same methodology. Also, equipment and the suppliers of materials must be the same to avoid obtaining varying results.

2.2.1.3 For effective microsatellite analysis, selecting high quality markers is essential. This includes a consideration of, *inter alia*:

- (a) the degree of “stuttering” (production of a series of one or more bands, differing by 1 repeat unit in size);
- (b) (n+1) peaks; Taq-polymerase often adds 1 bp to the end of a fragment. This can be prevented by using “pigtailed” primers (see Glossary);
- (c) the size of the amplification product;
- (d) effective separation between the various alleles in suitable detection systems;
- (e) reliable and reproducible scoring of the alleles in different detection systems;
- (f) the level of polymorphism between varieties (note that this requires analysis of a significant number of varieties);
- (g) avoidance of linkage.

2.2.1.4 For scoring SSRs in different laboratories and using different detection equipment, it is crucial that reference alleles (i.e. sets of varieties) are defined and included in all analyses. These reference alleles are necessary because molecular weight standards behave differently in the various detection systems currently available and are therefore not appropriate for allele identification.

2.2.1.5 Primers used in a particular laboratory should be synthesized by an assured supplier, to reduce the possibility of different DNA profiles as a result of using primers synthesized through different sources.

2.2.2 Single nucleotide polymorphism (SNP)

Single nucleotide polymorphisms (SNPs: see Glossary) can be detected via DNA sequencing, a routine technique which generally shows very high levels of repeatability over time and reproducibility between laboratories. ~~However, detection of specific SNPs can be carried out with a range of techniques, many of which are not yet routine.~~ By their nature, SNPs have only two allelic states in diploid plants, although this may vary in polyploids where there will be dosage effects. The simple makeup of SNPs makes the scoring of SNPs relatively straightforward and reliable. It also means that a large number of markers may need to be analyzed, either singly or in multiplexes, to allow the efficient and effective profiling of a particular genotype.

Comments by ESAⁱⁱⁱ

Is this ["are not yet routine"] still valid?

The choice for a specific SNP technology also strongly depends from the number of markers that needs to be analysed. For genotyping of bigger number of markers, the use of SNP-chips or sequence-based genotyping technology could be more appropriate, but no reference is made to these technologies at all. The choice for a technology also depends on the purpose of the database (management of reference collection f.e.).

What about an extra paragraph about genotyping by sequencing? Since more and more reference genomes are available this technology might become more important and it might be useful to address the different sequencing methods and their suitability

The only system/technology that makes any sense is to use full DNA sequences, as anything less would result in being selective, i.e. only really screen/analyse the parts of the genome where the chosen markers are placed. And then how to decide on which parts of the genome is important?

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To add new sections 1.2 and 1.3 as follows:

1.2 Flexibility and adaptability of a marker set

The discriminatory power of a marker set needs to be regularly assessed due to the evolution of the variety collections. Markers may need to be added or discarded depending on the modification of the genetics of varieties. In addition, New Breeding Techniques (NBT) and their resulting products may also require to use specific markers to detect the edited sites in the genome (e.g. additional characteristics could be evaluated these markers provided that a direct correlation between the edited sites and the phenotype has been established).

1.3 Requirements on the molecular profiles

1.3.1 Markers scattered all along the genome are used for the evaluation of distances/similarities between varieties through molecular distances and/or allelic frequencies. Application of this markers set is an assessment of the 'genetic background'.

1.3.2 In addition, markers that correlate with defined morphological qualitative traits, fulfilling the UPOV model 1 criteria, can complement the genetic description.

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To add a new section 2 "Phase 2: SELECTION OF THE DETECTION METHOD" with the following text:

As a prerequisite, whatever the source of material, the method for sampling and DNA extraction should be standardized and documented.

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To add a new section 2.1 "Genotyping methods - general criteria" with the following subsections 2.1.1 and 2.1.2:

2.1.1 Important criteria for choosing a genotyping methods that generate high quality molecular data are:

- Mandatory criteria:
 - (a) Reproducibility of data production within and between laboratories and detection platforms (different types of equipment).
 - (b) Repeatability over time
 - (c) Discrimination power of the method
 - (d) Interpretation of the data produced is independent of the equipment
- Optional criteria
 - (a) Possibilities for databasing
 - (b) Accessibility of methodology
 - (c) Suitable for automation
 - (d) Suitable for multiplexing
 - (e) Applicable for both diploid species and polyploidy species
 - (f) Cost effective; costs, number of samples and number of markers are in balance.

2.1.2 As improvements in technology and new equipment become available, it is important for the continued sustainability of databases that the interpretation of the data produced is independent of the technology and equipment used to produce them. This repeatability and reproducibility is important in the construction, operation and longevity of databases and is very important in generating a centrally maintained database, populated with verified data from a range of sources.

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To add a new section 2.2 "Recommendations for the choice of the method" with the following texts:

- (a) Methods that are simple to perform (limited steps in the protocol) are preferred over methods with a complex protocol that are time and labour consuming.
- (b) Methods that allow easy, objective and indisputable scoring of marker profiles are preferred over methods that produce complex marker profiles that are sensitive for interpretation (e.g. wide range of intensities of the bands).
- (c) Methods that are robust, not sensitive to subtle changes in the protocol or condition, but stable performance in time and conditions are preferred over methods that are sensitive to environmental conditions that are difficult to control.
- (d) Methods that are flexible (vary in the number of samples or the number of markers) are preferred over methods that have a fixed set-up.
- (e) Methods that are open source are preferred over methods that are completely or partly protected by IP rights or by confidential information.
- (f) Methods that are independent of a specific machine or specific chemistry or specific supplier are preferred over methods that require a specific machine, chemistry or supplier that have a monopoly in the market. Methods without dependence on particular partners or products are preferred.
- (g) Methods that detect molecular markers in a co-dominant way are preferred over methods that detect markers in a dominant way.
- (h) Methods that allow multiplexing are preferred over methods that detect only one marker in one assay.
- (i) Methods that are suitable for automation are preferred.

3. Access to the Technology

Some molecular markers and materials are publicly available. However, a large investment is likely to be necessary to obtain, for example, high quality SSR markers and consequently markers and other methods and materials may be covered by intellectual property rights. UPOV has developed guidance for the use of products or methodologies which are the subject of intellectual property rights and this guidance should be followed for the purposes of these guidelines. It is recommended that matters concerning intellectual property rights should be addressed at the start of any developmental work.

Comments by ESAⁱⁱⁱ

We have some doubts regarding choosing markers on the basis of publications: the markers must have been developed for the purpose at stake.

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To renumber the section 3. for a new section 2.3.

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To add a new section 2.4 “Future perspective on technological development” with the following texts and table:

Genotyping methods develop very fast and new technologies will keep being discovered. High-Throughput sequencing of short reads and now massive sequencing of long reads by nanopore sequencing enable the production of more and more data for a decreasing price per datapoint. As a consequence, the methods for marker set detection will alter in the future and shift from single sample endpoint methods towards whole genome sequences approaches. Irrespective of the technology used to detect the defined marker set, the genotype of a particular variety should not be affected. Both SSR markers and SNP/INDEL markers can be detected by High-Throughput Sequencing. In the (near) future, it could be cost effective to just sequence the whole genome of a plant. Even if all data produced will not be used (depending on the application), if the cost of a whole sequence become cheaper than single end point methods it may become the default method. However genotyping error of this technology need to be evaluated carefully before use.

Strategy	Reference genome	Present cost	Ease of use
Genome reduction - NGS	yes	€	+++
Genome reduction - NGS	no	€€	++
Whole genome - NGS	Yes	€€€	++
Whole genome - NGS	no	€€€€	+

4. Material to be Analyzed

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To move current texts and subsections in section 4 to a new section 5.2 “Requirements of the plant material.

The source and type of the material and how many samples need to be analyzed are the main issues with regard to the material to be analyzed.

4.1 Source of plant material

The plant material to be analyzed should be an authentic, representative sample of the variety and, where possible, should be obtained from the sample of the variety used for examination for the purposes of Plant Breeders’ Rights or for official registration. Use of samples of material submitted for examination for the purposes of Plant Breeders’ Rights or for official registration will require the permission of the relevant authority, breeder and/or maintainer, as appropriate. The plant material from which the samples are taken should be traceable in case some of the samples subsequently prove not to be representative of the variety.

Comments by Ecuador^{iv}

What criteria are used to guarantee an authentic, representative sample of the variety?

What is meant by “where possible, [the plant material] should be obtained from the sample of the variety used for examination for the purposes of plant breeders' rights or for official registration”?

What is the protocol for obtaining the sample and how is permission obtained to take the sample?

Will a bank of germplasm solely from plant material be created? What authority will be the custodian of that information or will the samples be kept by each relevant authority?

The document recommended in the general introduction that the size of the representative samples for determining the number of plants is for plants used in the open field and not for determining the number of plants for sequencing.

4.2 Type of plant material

The type of plant material to be sampled and the procedure for sampling the material for DNA extraction will, to a large extent, depend on the crop or plant species concerned. For example, in seed-propagated varieties, seed may be used as the source of DNA, whereas, in vegetatively propagated varieties, the DNA may be extracted from leaf material. Whatever the source of material, the method for sampling and DNA extraction should be standardized and documented. Furthermore, it should be verified that the sampling and extraction methods produce consistent results by DNA analysis.

4.3 Sample size

It is essential that the samples taken for analysis are representative of the variety. With regard to being representative of the variety, consideration should be given to the features of propagation (see the General Introduction). The size of the sample should be determined taking into account suitable statistical procedures.

4.4 DNA reference sample

It is recommended that a DNA reference sample collection should be created from the plant material sampled according to sections 4.1 to 4.3. This has the benefit that the DNA reference samples can be stored and supplied to other laboratories. The DNA samples should be stored in such a way as to prevent degradation.

Comments by Ecuador^{iv}

While it is understood that each relevant authority must establish a collection of germplasm, in the case of Ecuador, which does not have a collection of living material, it is more feasible to create living material than to have and maintain a collection of DNA reference samples.

The cost of the technology, equipment and materials is high and specialized technical equipment is required. If samples are sent to a general authority like UPOV, foreign germplasm banks used by developed countries will continue to grow; a situation that may affect the food sovereignty and security of countries.

A protocol for collection, maintenance, analysis, quality control and evaluation of samples should be standardized.

Joint comments from the European Union, France and the Netherlandsⁱⁱ

Texts in section 4.4 to read “A DNA reference sample collection may be created from the plant material sampled. The DNA samples should then be stored in such a way as to prevent degradation (e.g. storing it at -80C). The transfer of DNA reference samples to other laboratories will be submitted to the agreement of the owners of the varieties”.

5. Standardization of Analytical Protocols

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To delete current section 5 and replace with a new section 4 “Phase 4: HARMONIZATION AND VALIDATION OF THE MARKER SET AND METHOD”.

5.1 *Introduction*

This document is not intended to provide detailed technical protocols for the production of DNA profiles of varieties. In principle, any suitable analytical methodology can be used, but it is important that the methodology is validated in an appropriate way. This may be via an internationally recognized method of validation, or by developing a performance-based approach. In either case, there are some useful general considerations.

Any method used for genotyping and the construction of databases should be technically simple to perform, reliable and robust, allowing easy and indisputable scoring of marker profiles in different laboratories. This requires a level of standardization, for instance in the selection of markers, reference alleles and allele calling/scoring.

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To delete current section 5.1 and replace with a new section 4.1 “Harmonisation and validation – general criteria” with the following texts:

In order to select suitable markers and produce acceptable laboratory protocols for a given species, a harmonisation process involving more than one laboratory (i.e. an internationally recognized method of validation, e.g. a ring test according to internationally agreed standards) is recommended. This phase will involve the validation of genotyping methods and markers from which a defined set will be selected. This selection is based on performance: markers and methods should be robust and give rise to consistent results and DNA profiles in different laboratories using different equipment and chemistry. The consistence of the markers and methods in different laboratories is evaluated in the harmonisation process. The final choice of a number to be validated will be a balance between costs and the requirement to produce a satisfactory set of agreed markers at the end of the process. The objective is to produce an agreed set of markers that can be reliably and reproducibly analysed, scored and recorded in different laboratories, potentially using different types of equipment and different sources of chemical reagents, etc.

5.2 *Quality criteria*

5.2.1 It is important to consider quality criteria concerning, for example:

- (a) the quality of DNA;
- (b) methods of DNA extraction
- (c) the primer sequences;
- (d) the polymerase to be used in PCR-based methodologies;
- (e) for PCR-based methodologies, the amount/concentration of each PCR component and other components; and
- (f) PCR cycling conditions.

5.2.2 The detailed methodology should be set out in a protocol.

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To delete current section 5.2 and replace with a new section 4.2 "Performance criteria" with the following texts:

It is needed to determine whether the selected marker set is suitable (fit-for-purpose). The accuracy should be measured. To determine the adequacy of a method and DNA marker set several points should be considered:

(a) Discriminative capacity/informativeness

This can be determined by testing a defined collection of varieties – test set. For example variety pairs that are derived/excepted mutants, samples from the same variety but maintained in different places during time. Variety pairs with a known very close relation. Known pedigree. Diversity statistics such as Polymorphism Information Content (PIC)-values, expected heterozygosity (He), Effective Multiplex Ratio (EMR), Marker Index (MI) and/or Resolving power (Rp) can be calculated to illustrate the informativeness of a marker or marker set. The number of markers used should be an excess (exhausted number of markers). The minimal number of marker should be assessed and define so that analysis with a random selection of markers should not lead to different conclusions.

(b) Reproducibility

Once chosen for a particular technology and DNA marker set this should repetitively reveal the same DNA profile for a variety. This is inevitable essential within one laboratory and between laboratories, especially when the DNA profiles are stored in databases.

(c) Repeatability

(d) Robustness

(e) Error-rate

Every technology and every machine or platform has its imperfections and deficiencies. It is crucially important to be able to distinguish the technically induced variation from the real genetic diversity. This can be determined by the analysis of replica and/or duplicate samples (several DNA samples derived from the same plant material, so, identical DNA information, that proceed through laboratory process in parallel). Any deviation in the DNA profile must be a technical error.

5.3 Evaluation Phase

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To delete current section 5.3 and replace with a new section 3 "Phase 3: EVALUATION OF THE SELECTED MARKER SET AND DETECTION METHOD (fit for purpose validation of the marker set and technological validation of the method)" with the following subsections:

3.1 General requirements for molecular marker set development

3.1.1. Selection of the varieties - defining the genetic width of the marker set

The selection of the varieties on which the molecular markers are developed is crucial. An appropriate number of varieties, based on the genetic variability within the species and type of variety concerned, should be selected. The selected varieties should be well characterized (morphologically) and true-to-type. The choice of varieties should reflect the maximum range of diversity within the group/crop/species/type - representative sampling of the particular group/crop/species/type must be guaranteed. In addition, some genetically very similar varieties or lines, some parents and offspring, genetically close but morphologically distinct varieties, some morphologically close varieties with different pedigree should be included, to enable to 'measure' the level of discriminative capacity of the markers and to determine the 'suitability' of the marker set.

3.1.2. Generation of molecular data of selected varieties – defining the genetic depth of the marker set

Primers used in a particular laboratory should be synthesized by an assured supplier, to reduce the possibility of different DNA profiles as a result of using primers synthesized through different sources.

There are several ways to collect the data on the genetic diversity within the particular group/crop/species/type for which a marker set is to be developed.

5.3.1 Introduction

In order to select suitable markers and produce acceptable laboratory protocols for a given species, a preliminary evaluation phase involving more than one laboratory (i.e. an internationally recognized method of validation, e.g. a ring test according to internationally agreed standards) is recommended. This phase should be mainly concerned with selecting a set of markers, and will usually involve the evaluation of existing markers, either published or available via other means. The number of markers to be evaluated will vary and depends on the possibilities presented by different species. The markers should derive from reliable sources (e.g. peer-reviewed publications) and be sourced from assured suppliers. The final choice of a number to be evaluated will be a balance between costs and the requirement to produce a satisfactory set of agreed markers at the end of the process. The objective is to produce an agreed set of markers that can be reliably and reproducibly analyzed, scored and recorded in different laboratories, potentially using different types of equipment and different sources of chemical reagents, etc.

5.3.2 Variety choice

An appropriate number of varieties, based on the genetic variability within the species and type of variety concerned, should be selected as the basis for the evaluation phase. The choice of varieties should reflect an appropriate range of diversity and where possible should include some closely related and some morphologically similar varieties, to enable the level of discrimination in such cases to be assessed.

5.3.3 Interpretation of results

The next evaluation stage should, if possible, include an internationally recognized method of validation to assess the whole methodology in an objective way. Any marker which causes difficulties in any of the laboratories involved in this evaluation phase should be rejected for subsequent use. As most errors in the analysis of large variety collections seem to arise from scoring errors, construction of databases should be based on duplicate samples (e.g. different sub-samples of seed from the same variety), analyzed by more than one laboratory. Since the sub-samples (or DNA extracts from them) can be exchanged in the event of any discrepancy, this approach is very effective in highlighting sampling errors, or those due to heterogeneity within the samples, and eliminates possible laboratory artifacts.

5.4 *Scoring of molecular data*

A protocol for allele/band scoring should be developed in conjunction with the evaluation phase. The protocol should address how to score the following:

- (a) rare alleles (i.e. those at a specific locus which appear with a frequency below an agreed threshold (commonly 5-10%) in a population);
- (b) null alleles (an allele whose effect is an absence of PCR product at the molecular level);
- (c) "faint" bands (i.e. bands where the intensity falls below an agreed threshold of detection, set either empirically or automatically, and the scoring of which may be open to question);
- (d) missing data (i.e. any locus for which there are no data recorded for whatever reason in a variety or varieties);
- (e) monomorphic bands (those alleles/bands which appear in every variety analyzed, i.e. are not polymorphic in a particular variety collection).

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To add a new section 4.3 "Consistence criteria - harmonization of markers and methods in different laboratories Performance criteria" with the following texts (underlined texts are from the current section 5.4 "Scoring of molecular data"):

(a) Well defined collections of varieties/samples to be applied as a reference that represent all alleles. An appropriate number of varieties, based on the genetic variability within the species and type of variety concerned, should be selected as the basis for the validation and harmonisation. The choice of varieties should reflect an appropriate range of diversity and where possible should include some closely related and some morphologically similar varieties, to enable the level of discrimination in such cases to be assessed. For example variety pairs that are derived/excepted mutants, samples from the same variety but maintained in different places during time. Variety pairs with a known very close relation. Known pedigree.

(b) Duplicate samples (e.g. different sub-samples of seed from the same variety), analysed by more than one laboratory. Taking into account the data obtained from these samples may help reduce scoring errors and reinforce the reliability of the data stored in the databases. Since the sub-samples (or DNA extracts from them) can be exchanged in the event of any discrepancy, this approach is very effective in highlighting sampling errors, or those due to heterogeneity within the samples, and eliminates possible laboratory artefacts.

(c) Blind samples

(d) Agreements on the scoring of molecular data. A protocol for allele/band scoring should be developed. The protocol should address how to score the following:

(i) rare alleles (i.e. those at a specific locus which appear with a frequency below an agreed threshold (commonly 5-10%) in a population);

(ii) null alleles (an allele whose effect is an absence of PCR product at the molecular level);

(iii) "faint" bands (i.e. bands where the intensity falls below an agreed threshold of detection, set either empirically or automatically, and the scoring of which may be open to question);

(iv) missing data (i.e. any locus for which there are no data recorded for whatever reason in a variety or varieties); and

(v) monomorphic bands (those alleles/bands which appear in every variety analysed, i.e. are not polymorphic in a particular variety collection).

6. Databases

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To delete current section 6 and replace with a new section 5 “Phase 5: CONSTRUCTION OF A SPECIES-SPECIFIC DATABASE” with the following texts and subsection 5.1:

A database and the data that is stored in a shared database and how it is stored in a database reflects the process of producing the data. The database should store:

- (1) the end results, e.g. the genotype as well as how it was derived both in terms of;
- (2) sequencing library preparation; and
- (3) the computational steps for deriving a genotype.

5.1 Requirements of a database

- (a) The database architecture should be flexible, e.g. allow for storing both flat files as well as compressed archives.
- (b) Contains different tables, separate tables and entries are required for library prep (the wet-lab work), data processing and the genotyping scores.
- (c) Store information at different levels (allele scores / how the allele score was called (the rules or the interpretation rules behind a decision) / (links) to the raw data (tiff files, bam files, xx files that came out of the machine that produced the data that were used for Allele scoring and interpretation).
- (d) For sequencing data, variant call files in VCF or BCF format corresponding to the standard version 4.2 or higher should be used. Header entries should contain the name and version of the different scripts used for both sequence read mapping, read filtering, variant calling and variant filtering in such a way that a competent bioinformatician can repeat the analysis.
- (e) In case of replicate samples, one consensus genotype entry can be computed and stored in case the genotypes of the replicates match. In case of non-matching replicates, the record needs to be flagged or filtered out where appropriate. The rules applied for these cases need to be documented in a publicly accessible code repository that is references from the variant call file. Frequencies could also be used for heterogeneous varieties.
- (f) The database should validate the VCF and or BCF data against relevant specifications.
- (g) The database should have a web front-end that enables easy uploading, downloading and interactive exploration of the data. The systems for storing, analysing and interpreting the data should be build and function separately yet function well in concert.
- (h) Easy to share data, an API is recommended.

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To add a new section 5.3 “Data processing” with the following texts:

The pipeline for processing the data should keep a detailed log of:

- (a) Type and versions of tools;
 - (b) Command line used for the tool;
 - (c) Reproducibility counts;
 - (d) Open source tools are preferred;
 - (e) Sharing is encouraged;
 - (f) Raw alignment data (bam or CRAM files) should be stored where possible;
 - (g) Multi-sample VCF files are not suitable, one VCF file per cultivar must be present;
 - (h) If VCF files are stored, all positions (both variants & non-variants) and their depth should be stored;
 - (i) Both heuristic and probabilistic approached should be considered and compared for genotyping methods;
 - (j) Databases should facilitate input and output of genotype call data in standardized format (VCF or BCF);
 - (k) The data processing pipeline should result in a detailed log file which should be stored in conjunction to the variant call data;
 - (l) If possible, raw data should be stored so that data processing can be repeated with new or updated tools ;
- and
- (m) A p-value or uncertainty for a given allele should be stored.

6.1 Type of database

There are many ways in which molecular data can be stored, therefore, it is important that the database structure is developed to be compatible with all intended uses of the data.

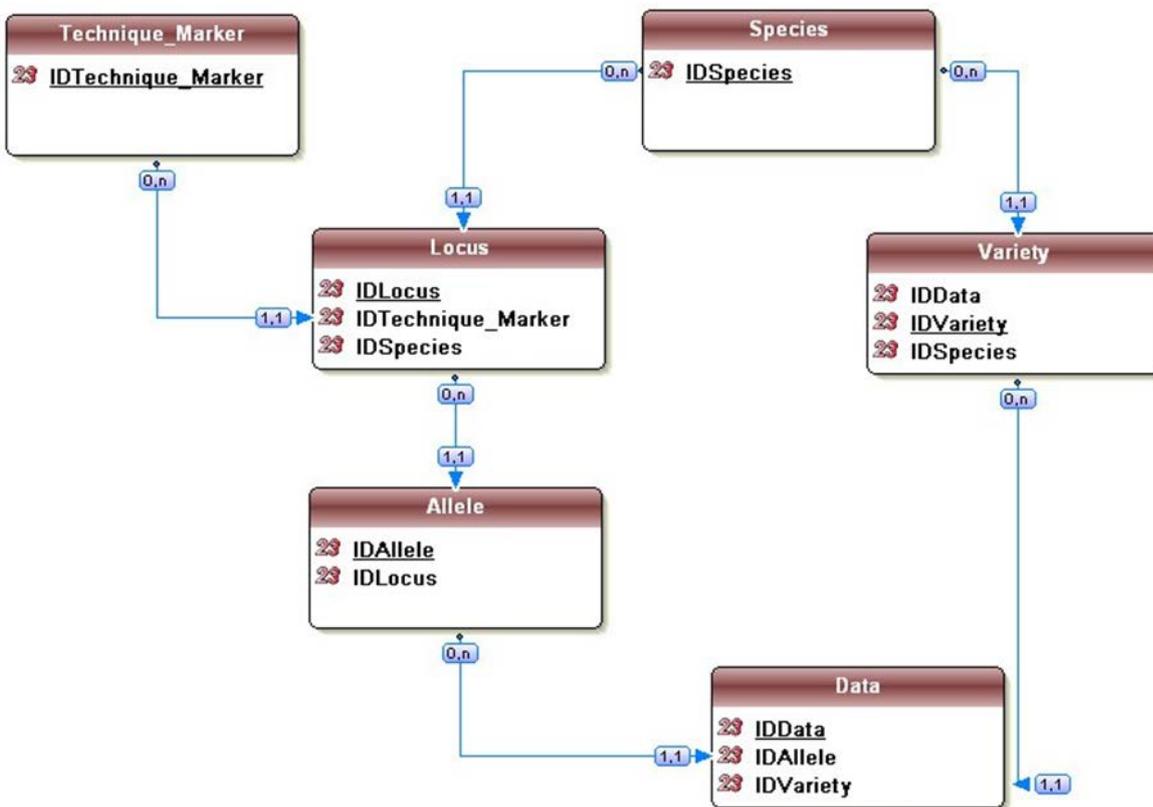
Joint comments from the European Union, France and the Netherlandsⁱⁱ

To renumber section 6.1. for a new section 5.4 and to add the following sentence to the end of the current texts:

For molecular data obtained using next generation sequencing (NGS), the variant call file standard VCFv4.2 is recommended (<https://samtools.github.io/hts-specs/VCFv4.2.pdf>).

6.2 Database model

The database model should be defined by IT database experts in conjunction with the users of the database. As a minimum the database model should contain six core objects: Species; Variety; Technique; Marker; Locus; and Allele.



Proposal by Argentinaⁱ

[Regarding section 6.2, to] change the word "technique" for "scoring system".

If so, to change also point 6.3.1 a).

When saying scoring system, we add a new quality to the data base. A scoring system would be RT PCR, sesquencing, capilar electrophoresis, others.

Joint comments from the European Union, France and the Netherlands ⁱⁱ

To renumber section 6.2. for a new section 5.5 and to add the following sentences to the end of the current texts:

For variants obtained from sequencing data, storing VCF files in a relational or no SQL database is recommended. In this case, each database record for a variant has a defined genome version, chromosome, position, reference allele.

6.3 Data Dictionary

Joint comments from the European Union, France and the Netherlands ⁱⁱ

To renumber section 6.3. for a new section 5.6.

6.3.1 In a database, each of the objects becomes a table in which fields are defined. For example:

(a) Technique/Marker ~~code~~^{vi}:

Marker type^{vi} indicates the code or name of the technique or type of marker used, e.g. SSR, SNP, etc.

(b) Locus ~~code~~^{vi}:

Locus code^{vi} indicates name or code of the locus for the species concerned, e.g. gwm 149, A2, etc.

Joint comments from the European Union, France and the Netherlands ⁱⁱ

The title of section (b) to read "Reference genome position / Locus code:" with the following texts:

Preferably, a genome assembly version, chromosome and position should be provided if a reference genome is available for the species concerned, e.g. SL2.50ch05:63309763 for tomato *Solanum lycopersicum* assembly version 2.50 on chromosome 5 position 63309763. If no reference genome is available or the location is unknown, a name or code of the locus for the species concerned can be used, e.g. gwm 149, A2, etc.

(c) Allele ~~code~~^{vi}:

Allele code^{vi} indicates name or code of the allele of a given locus for the species concerned, e.g. 1, 123, etc.

Example variety: Indicates for each a variety of common knowledge, in which the allele can be observed.^{vi}

Observation method: e.g. polyacrylamide gel electrophoresis (SSR), capillary electrophoresis (SSR), sequencing (SSR), TaqMan (SNP), Kompetitive Allele Specific PCR (SNP), Araay (SNP), High Throughput Sequencing (SNP)^{vi}

Joint comments from the European Union, France and the Netherlands ⁱⁱ

The title of section (c) to read "Genotype:" with the following texts:

For SNP genotypes, the allele composition of the SNP or MNP should be given, e.g. A/T or A/A.

For other technique, genotype indicates the name or code of the allele of a given locus for the species concerned, e.g. 1, 123, etc.

(d) Data value^{vi}:

Data value:^{vi} indicates a data value for a given sample on a given locus-allele, e.g. 0 (absence), 1 (presence), 0.25 (frequency) etc.

Data type: ~~it could be an e.g. Excel file, an XML sequence data, or a value~~^{vi}

Joint comments from the European Union, France and the Netherlandsⁱⁱ

The title of section (d) to read "Allele depths / Data value:" with the following texts:

For SNPs obtained from next generation sequencing data this should indicate the depth of coverage for alleles e.g. 10/20 for an A/T allele in which the A is covered by 10 reads and the T by 20. Otherwise, indicates a data value for a given sample on a given locus-allele, e.g. 0 (absence), 1 (presence), 0.25 (frequency) etc.

(e) Variety:

Variety denomination or breeder's reference:^{vi} the variety is the object for which the data have been obtained.

Grouping type of variety: e.g. Inbred Line or Hybrid^{vi}

(f) Species:

the species is indicated by the botanical name or the national common name, which sometimes also refers to the type of variety (e.g. use, winter/spring type etc.). The use of the UPOV code would avoid problems of synonyms and would, therefore, be beneficial for coordination.

6.3.2 In each table, the number of fields, their name and definition, the possible values and the rules to be followed, need to be defined in the "data dictionary".

6.4 Table Relationship

6.4.1 The links between the tables are an important aspect of the database design. The links between tables can be illustrated as follows:

Table	Link	Table	Description
Woman	0 or 1 to n (0, n)	Child	0: A woman may have no child 1 to n: a woman may have 1 to n children (she is then a mother)
Child	1 to 1 (1,1)	Woman	A given child has only one biological mother

6.4.2 The following table indicates the relationship between the six minimum core objects, as proposed in the database model in Section 6.2:

Table	Link	Table	Description
Technique/marker	0 or 1 to n	Locus	0: A technique/marker can be present in Technique/marker, even if no locus/allele is yet used in the database 1 to n: a given type of marker can provide 1 to n useful loci
Locus	1 to 1	Technique/marker	A given locus is defined within the scope of a given technique/marker
Locus	1 to n	Allele	For each Locus 1, or more than 1, allele can be described
Allele	1 to 1	Locus	A given Allele is defined within the scope of a given Locus
Allele	0 or 1 to n	Data	0: a given Allele can be defined, but without data 1 to n: a given allele can be found in 1 to n data
Data	1 to 1	Allele	data corresponds to a given allele
Variety	0 or	Data	0: the variety has no data

	1 to n		1 to n: the variety has data
Data	1 to 1	Variety	data corresponds to a given variety
Data	1 to 1	Species	data is obtained for a given variety, then for the species of the variety.
Species	0 or 1 to n	Data	0: a species can have no data. 1 to n: a species can have 1 to n data.

Proposal by Argentinaⁱ

Nowadays most labs would use SNP markers, and the data come in an Excel file and then it is analyzed using biostatistical tools, which are common among labs. I do not see the need of a table like 6.4.2.

6.5 *Transfer of data to the database*

To reduce the number of errors in data transfer and transcription, it is advisable to automate transfer of data to databases as much as possible.

6.6 *Data access / ownership*

It is recommended that all matters concerning ownership of data and access to data in the database should be addressed at the beginning of any work.

6.7 *Data analysis*

The purpose for which the data will be analyzed will determine the method of analysis, therefore, no specific recommendations are made within these guidelines.

6.8 *Validating the database*

When the first phase of the database is complete, it is recommended to conduct a 'blind test', i.e. distribute a number of samples to different laboratories and ask them to use the agreed protocol in conjunction with the database to identify them.

6.9 *Data Exchange^{vi}*

For cooperation purposes, the data model should allow:

Scenario 1. Exchange of data produced from a standardized set of markers for a specific crop^{vi}

Scenario 2. Search and view data of selected varieties generated from the same standardized set of markers.^{vi}

Scenario 1^{vi}

In order to exchange data about the marker set used for a specific crop, the following web service can be used:

https://office.org/locus?upov_code={upovcode}&type={marker type}&method={observation method}^{vi}

For example, to obtain marker set information for maize using SSR and CE method, the following URL should be accessed:^{vi}

https://office.org/locus?upov_code=ZEAAA_MAY&type=SSR&method=CE^{vi}

The result would be:

```
{
  "techniqueid": "CN SSR ZEAA MAY CE V 1",
  "locusid": "M01",
  "alleles":
  {
    "alleleid": "238/256",
    "examplevariety":
  
```

1.
"alleleid": "238/271",
"examplevariety":
1.
"alleleid": "246/246",
"examplevariety":
1.
"alleleid": "246/248",
"examplevariety":
1.
"alleleid": "246/250",
"examplevariety":
1.
"alleleid": "246/254",
"examplevariety":
1.
"alleleid": "246/256",
"examplevariety":
1.
"alleleid": "246/260",
"examplevariety":
1.
"alleleid": "246/277",
"examplevariety":
1.
"alleleid": "246/284",
"examplevariety":
1.
"alleleid": "246/288",
"examplevariety":
1.
"alleleid": "248/250",
"examplevariety":
1.
"alleleid": "248/256",
"examplevariety":
1.
"alleleid": "248/271",
"examplevariety":
1.
"alleleid": "248/290",
"examplevariety":
1.
"alleleid": "250/250",
"examplevariety":
1.
"alleleid": "250/252",
"examplevariety":
1.
"alleleid": "250/256",
"examplevariety":
1.
"alleleid": "250/275",
"examplevariety":
1.
"alleleid": "252/256",
"examplevariety":
1.
"alleleid": "252/260",
"examplevariety":
1.
"alleleid": "252/271",
"examplevariety":
1.
"alleleid": "252/273",
"examplevariety":
1.
"alleleid": "252/282",
"examplevariety":

```
1.
{"alleleid": "254/254",
"examplevariety":
1.
{"alleleid": "254/271",
"examplevariety":
1.
{"alleleid": "254/284",
"examplevariety":
1.
{"alleleid": "254/286",
"examplevariety":
1.
{"alleleid": "256/256",
"examplevariety":
1.
{"alleleid": "256/264",
"examplevariety":
1.
{"alleleid": "256/266",
"examplevariety":
1.
{"alleleid": "256/271",
"examplevariety":
1.
{"alleleid": "256/284",
"examplevariety":
1.
{"alleleid": "256/286",
"examplevariety":
1.
{"alleleid": "258/258",
"examplevariety":
1.
{"alleleid": "264/284",
"examplevariety":
1.
{"alleleid": "271/292",
"examplevariety":
1.
1.
{"locusid"="M02",
"alleles": [...]
1.
1.
```

1. vi

[PRO DOMO : Another example for SNP and Sequencing should be provided]

Scenario 2 vi

In order to search and view molecular data of a variety, the following web service can be used:

https://office.org/variety?id={irn}&techniqueid={technique_code} vi

For example,

https://office.org/variety?id=XU_30201800000140 &techniqueid= CN SSR ZEAA MAY CE V 1 vi

The result would be:

```
{ "techniqueid": "CN SSR ZEAA MAY PAGE ",
"varietyid": " XU 30201800000140 ",
"data":
```

```
[  
  "id": "M01",  
  "value": "254/254"  
],  
[  
  "id": "M02",  
  "value": "347/347"  
],  
[  
  "id": "M03",  
  "value": "292/292"  
],  
[  
  "id": "M04",  
  "value": "361/361"  
],  
...  
] vi
```

Joint comments from the European Union, France and the Netherlandsⁱⁱ

To add a new section 6 "Phase 4: DATABASE MANAGEMENT" with the following texts:

The effective management and updating of the database on the long term requires that appropriate agreements between partners are signed at the start of the creation process. These agreements should cover general principles (defining precisely the ownership of the materials and data, conditions of access and use, confidentiality, etc.) and technical principles (describing the types of data, identifiers, role of the partners, rules and planning of updating, etc.). The conditions under which the database could be open to additional partners wishing to contribute to its feeding after it was built needs also to be clearly established.

7. Summary

The following is a summary of the approach recommended for the selection and use of molecular markers to construct central and sustainable databases of DNA profiles of varieties (i.e. databases that can be populated in the future with data from a range of sources, independent of the technology used).

- (a) consider the approach on a crop-by-crop basis;
- (b) agree on an acceptable marker type and source;
- (c) agree on acceptable detection platforms/equipment;
- (d) agree on laboratories to be included in the test;
- (e) agree on quality issues (see section 5.2);
- (f) verify the source of the plant material used (see section 4);
- (g) agree which markers are to be used in a preliminary collaborative evaluation phase, involving more than one laboratory and different detection equipment (see section 2);
- (h) conduct an evaluation (see section 5.3);
- (i) develop a protocol for scoring the molecular data (see section 5.4);
- (j) agree on the plant material/reference set to be analyzed, and the source(s);
- (k) analyze the agreed variety collection, in different laboratories/different detection equipment, using duplicate samples, and exchanging samples/DNA extracts if problems occur;
- (l) use reference varieties/DNA sample/alleles in all analyses;
- (m) verify all stages (including data entry) – automate as much as possible;
- (n) conduct a 'blind test' in different laboratories using the database;
- (o) adopt the procedures for adding new data.
- (p) store the standardized marker set and variety data in a global database^{vi}
- (q) search and retrieve data using the global web services (see section 6.9).^{vi}

Joint comments from the European Union, France and the Netherlands^v

To add a new section C "DEFINITIONS" as follows (proposed texts are underlined):

C. DEFINITIONS

Locus: a position on a chromosome/ a set of homologous chromosomes

Allele: a variant of a locus

Polymorphism: alleles at a particular locus that are different between individual organisms

Marker (genetic marker/DNA marker/ molecular marker): a single piece of DNA or a set of pieces of DNA that mark one or more specific alleles and can be detected through a single assay.

Genotype: the genetic constitution of an individual organism

Genotyping: the process of elucidating the genotype of an individual organism with a biological assay.

DNA profile/DNA fingerprint: a unique pattern of molecular markers, specific for an individual organism and representative for the genotype of this individual organism

Dominant marker: A marker that can mark only one of the possible alleles as either present or absent through a single assay

Co-Dominant marker: A marker that can mark different alleles through a single assay

Locus-specific marker: The position of the marker on the chromosome is exactly known. Prior knowledge on the adjacent DNA sequence is needed to develop the locus-specific assay

Random marker: The position of the marker on the chromosome is NOT known. Prior knowledge on the adjacent DNA sequence of the marker locus is NOT required

Heterozygosity: The state of an individual in which a locus carries at least two alleles

Homozygosity: A state of an individual in which a locus carries only one single allele in the number of copies equal to the ploidy level.

GLOSSARY

Microsatellites, or Simple Sequence Repeats (SSRs)

Microsatellites, or simple sequence repeats (SSRs) are tandemly repeated DNA sequences, usually with a repeat unit of 2-4 base pairs (e.g. GA, CTT and GATA). In many species, multiple alleles have been shown to exist for some microsatellites, due to variations in the copy number of this repeat unit. Microsatellites can be analyzed by PCR using specific primers, a procedure known as the sequence-tagged-site microsatellite (STMS) approach. The alleles (PCR products) can be separated by agarose or polyacrylamide gel electrophoresis. In order to develop sequence-tagged site microsatellites, information about the sequence of the DNA flanking the microsatellite is needed. This information can sometimes be acquired from existing DNA sequence databases, but otherwise has to be obtained empirically.

Single Nucleotide Polymorphisms (SNPs)

Single nucleotide polymorphisms (SNPs) (pronounced “snips”) are DNA sequence variations that occur when a single nucleotide (A,T,C, or G) in the genome sequence is altered. For example a SNP might change the DNA sequence AAGGCTAA to ATGGCTAA. Generally, for a variation to be considered a SNP, it must occur in at least 1% of the population. The potential number of SNP markers is very high, meaning it should be possible to find them in all parts of the genome. SNPs can occur in both coding (gene) and non-coding regions of the genome. The discovery of SNPs involves comparative sequencing of numbers of individuals from a population. More commonly, potential SNPs are identified by comparing aligned sequences from the available sequence databases. Although they can be detected by relatively straightforward PCR + gel electrophoresis, high throughput and micro-array procedures are being developed for automatically scoring hundreds of SNP loci simultaneously.

Joint comments from the European Union, France and the Netherlands^v

To add the following new sentences at the end of this paragraph:

By their nature, SNPs have only two allelic states in diploid plants, although this may vary in polyploids where there will be dosage effects. The simple makeup of SNPs makes the scoring of SNPs relatively straightforward and reliable. It also means that a large number of markers may need to be analyzed, either singly or in multiplexes, to allow the efficient and effective profiling of a particular genotype.

Cleaved Amplified Polymorphic Sequences (CAPS)

Cleaved amplified polymorphic sequences (CAPS) are DNA fragments amplified by PCR using specific 20-25 bp primers, followed by digestion with a restriction endonuclease. Subsequently, length polymorphisms resulting from variation in the occurrence of restriction sites are identified by gel-electrophoresis of the digested products. In comparison with markers such as RFLPs, polymorphisms are more difficult to identify because of the limited size of the amplified fragments (300-1800 bp). CAPS analysis, however, does not require Southern blot hybridization and radioactive detection. CAPS have generally been applied predominantly in gene mapping studies to date.

Sequence-Characterized Amplified Regions (SCARs)

Sequence-characterized amplified regions (SCARs) are DNA fragments amplified by PCR using specific 15-30 bp primers, designed from previously identified polymorphic sequences. By using longer PCR primers, SCARs avoid the problem of low reproducibility. They are also usually co-dominant markers. SCARs are locus specific and have been applied in gene mapping studies and marker assisted selection.

Pig-tailing

In SSR analysis, “pig-tailing” is the addition of a short specific oligonucleotide sequence to the primers used in the PCR, as a way of improving the clarity of the amplification products and reducing artifacts.

Null Allele

In SSR analysis, a “null allele” is an allele at a particular locus whose effect is seen as an absence of a PCR product.

Stutter Bands

In SSR analysis, “stutter bands” is the occurrence of a series of one or more bands, differing by 1 repeat unit in size, following PCR.

[End of document]

ENDNOTES

-
- i Comments from Argentina in reply to UPOV circular E-18/004 of February 15, 2018 (see document BMT/17/10 “Review of document UPOV/INF/17 “Guidelines for DNA-Profiling: Molecular Marker Selection and Database Construction (‘BMT Guidelines’), paragraph 8)
 - ii Joint comments from the European Union, France and the Netherlands Argentina in reply to UPOV circular E-18/004 of February 15, 2018 (see document BMT/17/10 “Review of document UPOV/INF/17 “Guidelines for DNA-Profiling: Molecular Marker Selection and Database Construction (‘BMT Guidelines’), paragraph 8)
 - iii Comments from Spain in reply to UPOV circular E-18/004 of February 15, 2018 (see document BMT/17/10 “Review of document UPOV/INF/17 “Guidelines for DNA-Profiling: Molecular Marker Selection and Database Construction (‘BMT Guidelines’), paragraph 8)
 - iv Comments from the ESA in reply to UPOV circular E-18/004 of February 15, 2018 (see document BMT/17/10 “Review of document UPOV/INF/17 “Guidelines for DNA-Profiling: Molecular Marker Selection and Database Construction (‘BMT Guidelines’), paragraph 8)
 - v Comments from Ecuador in reply to UPOV circular E-18/004 of February 15, 2018 (see document BMT/17/10 “Review of document UPOV/INF/17 “Guidelines for DNA-Profiling: Molecular Marker Selection and Database Construction (‘BMT Guidelines’), paragraph 8)
 - vi Proposals from the Office of the Union (see document BMT/17/10 “Review of document UPOV/INF/17 “Guidelines for DNA-Profiling: Molecular Marker Selection and Database Construction (‘BMT Guidelines’), paragraph 9)