

Criteria for evaluating molecular markers: Comprehensive quality metrics to improve marker-assisted selection

John Damien Platten^{1,*}, Joshua N. Cobb¹, Rochelle E. Zantua¹

¹ International Rice Research Institute, National Road, Los Banos, Philippines

* Corresponding author

Email: j.platten@irri.org

Abstract

Despite strong interest over many years, the usage of quantitative trait loci in plant breeding has often failed to live up to expectations. A key weak point in the utilisation of QTLs is the “quality” of markers used during marker-assisted selection (MAS): unreliable markers result in variable outcomes, leading to a perception that MAS products fail to achieve reliable improvement. Most reports of markers used for MAS focus on markers derived from the mapping population. There are very few studies that examine the reliability of these markers in other genetic backgrounds, and critically, no metrics exist to describe and quantify this reliability. To improve the MAS process, this work proposes five core metrics that fully describe the reliability of a marker. These metrics give a comprehensive and quantitative measure of the ability of a marker to correctly classify germplasm as QTL[+]/[-], particularly against a background of high allelic diversity. Markers that score well on these metrics will have far higher reliability in breeding, and deficiencies in specific metrics give information on circumstances under which a marker may not be reliable. The metrics are applicable across different marker types and platforms, allowing an objective comparison of the performance of different markers irrespective of the platform. Evaluating markers using these metrics demonstrates that trait-specific markers consistently outperform markers designed for other purposes. These metrics also provide a superb set of criteria for designing superior marker systems for a target QTL, enabling the selection of an optimal marker set before committing to design.

Introduction

The world population is expected to top 9 billion people by 2050. To feed this population, it is estimated that agricultural output of cereals alone will need to increase by approximately 1 billion tons [1]. It is widely acknowledged that meeting this growth target will require the integration of new technologies into the breeding process. Many authors have discussed the promise of molecular marker technologies for improving the speed and efficiency of the breeding process, and extensive literature has accumulated on methodologies to incorporate the use of markers into breeding decisions (e.g. [2]). At its core, this marker-assisted selection (MAS) is about two correlations:

Trait ↔ QTL
and
QTL ↔ marker

A QTL is identified as a genetic position (locus) associated with some degree of phenotypic variation in a specific trait. Markers are assayable polymorphisms with some degree of association with a QTL in a specific gene pool. Both sets of correlations may be broad-ranging or narrowly applicable, and the success or reliability of MAS is directly determined by the strength of these correlations. However, since the middle factor (QTL) is almost never tracked *per se*, both correlations are often conflated into the indirect association of the marker with the trait. Since this association is indirect, reliable markers may not always result in reliable improvement of the trait – this depends on the quality of the QTL, and is a topic for a different discussion. However, unreliable markers will always result in unreliable improvement, and it is thus essential to identify what constitutes a reliable marker, and metrics to objectively measure and evaluate this quality.

Current MAS programs use one of several genotyping platforms, depending on requirements for marker density and sample throughput. These platforms range from low-throughput, PCR-based techniques such as the traditional microsatellite/simple-sequence repeats (SSRs) and newer insertion/deletion mutations (indels), to the explosion of high-throughput single-nucleotide polymorphism (SNP) platforms and new sequencing-based methods such as genotyping-by-sequencing (GBS) and amplicon sequencing [3, 4, 5, 6]. Most recent literature has focused on new SNP technologies, but by far the most common systems in use by public sector breeding programs are traditional SSRs. These are widely employed in biparental mapping studies such as QTL mapping and fine-mapping [7, 8, 9, 10, 11, 12, 13, 14], but have also been used in cross-population meta-analyses [12] and allelic diversity assessments [15, 16]. Given their high throughput nature SNPs and GBS are the platforms of choice for strategies requiring high sample volumes and/or marker densities, such as genomic selection and genome-wide association studies (e.g. [17, 18]).

Despite the acknowledged importance of incorporating molecular markers into the breeding process, there has been little discussion on factors influencing the success of such endeavours. Some studies have investigated the utility of SSRs in breeding processes [15, 19, 20, 21], however the only criterion for usefulness that is considered

is genetic linkage with the QTL; other issues such as how reliably the markers classify favourable and unfavourable alleles are rarely examined, and if so done in a cursory manner.

Poor classification ability in markers can lead to many undesirable outcomes. For example, a typical MAS workflow involving SSRs starts with a parental polymorphism survey; SSRs are chosen to introgress a QTL based on their linkage (position) and the fact they are polymorphic between a donor and the chosen recipient line. However, this does not examine whether the chosen recipient already contains the target QTL; the ability of the QTL to improve the recipient is assumed, not tested. SSR markers cannot provide information on whether the chosen recipient already possesses the QTL - this is a circular argument - and there are documented cases where SSRs give misleading indications as to the presence/absence of a QTL. This can be seen for example in [15] and [22], which in both cases confuse varieties with different *Saltol* alleles and distinguish varieties with the same allele (compare for example [23]). Many similar situations have arisen in breeding programs, irrespective of marker platform in use (SNP or gel-based), leading to unreliable outcomes, usually as a result of classifying a variety as lacking a QTL when in fact the QTL is already present. To avoid similar problems in future, some method of characterising markers is required, to help identify markers at risk of giving misleading results and to aid in the design of superior markers as replacements. Most importantly, a set of objective measures should be derived that describe how accurately a marker classifies both QTL[+] and QTL[-] material.

Surprisingly there is no literature available dealing with this subject, which may contribute to the ambiguity surrounding the quality of existing marker systems. The closest parallels in other disciplines are accuracy metrics used to evaluate clinical tests (e.g. [24]), but the concepts do not map directly to each other. For example, medical diagnostics assume a large number of case-[+] and case-[-] datapoints are available, and often deal with quantitative measures such as enzyme or antibody activity levels. By contrast, “datapoints” in the case of molecular markers are characterised varieties with or without the target QTL, which are typically far fewer in number. Combined with this, for many genotyping platforms the genotype is essentially a binary output, making it strictly impossible to distinguish more than two allelic states.

To stimulate discussion in this area, a set of five core and nine supporting metrics is presented along with strategies for their calculation, which attempt to capture the level and type of association between a marker and

its target QTL. These metrics are focused primarily on assessing the classification ability of a marker against a background of high allelic polymorphism such as is found in rice [25], but also assess several other parameters related to the reliability of scoring and usefulness of a marker in breeding. The metrics are then used to evaluate a set of SSR, SNP and indel markers targeting a range of QTLs in rice (*Oryza sativa* L.), to illustrate their application in designing marker systems that give higher confidence for deployment in breeding programs.

Materials and Methods

A summary of proposed marker quality metrics is found in Table 1. The metrics are comprised of five core metrics that quantify the reliability of a marker, and a further nine supporting metrics that allow the determination of the core metrics. They fall into three main categories: (1) Technical metrics, (2) Biological metrics, and (3) Breeding metrics. To evaluate these metrics on existing marker systems in rice several different kinds of markers from multiple platforms were compared. These include 20 SSR markers [3], approximately 4500 SNP markers identified as part of the OryzaSNP project [25] and utilised on fixed genotyping platforms such as the Illumina Infinium chip ([6]; hereinafter the “anonymous SNP panel”), and 137 candidate QTL-specific SNP and indel markers. A list of markers evaluated can be found in S1 Table. Markers were evaluated using one of two main datasets, depending on the metrics to be evaluated: empirically-determined PCR results to compare technical performance and breeding metrics for gel-based SSR and indel markers, and a large genome resequencing dataset to compare biological and breeding metrics for anonymous SNPs, QTL-specific SNPs, and QTL-specific indel markers.

107 **Table 1.** Summary of core and supporting metrics to describe marker quality

Technical metrics		
Version	Support	Identifier designating the version of a marker being examined; particularly important for technical performance metrics.
Call rate	Core	The proportion of samples which give a scorable result. For gel-based markers this is easy to visualise and understand (unless the marker is dominant), but many SNP platforms will call an allele even if the underlying data is substandard.
Clarity	Core	How reliably a sample can be classified as allele A, B or heterozygous. This is a critical parameter, and for each new marker it is worthwhile determining – even commercial SNP platforms do not always perform to the standard the sales brochure will advertise, and each individual marker will perform in different ways. The clarity can be quantified by determining what proportion of known germplasm is correctly classified, and/or how well it correlates with tightly-linked alternative markers.
Biological metrics		
Linkage	Support	Genetic distance of the marker from the QTL peak, it is probably best expressed as a genetic distance gained from one or more mapping populations. Diagnostic markers will by definition have a distance of 0cM.
Position	Support	Position targeted by a marker, for example chromosomal location (preferable), mapping bin or consensus genetic position.
Derived QTL state	Support	Which state appears to be the most derived, and therefore most informative in determining presence of the QTL: Favourable or Unfavourable. This is not a property of a specific marker, but rather the QTL.
Marker Target	Support	A qualitative description of the allele targeted by the marker, the favourable or unfavourable allele.
Favourable allele	Support	Not quality metrics per se but necessary for automated analyses, and helpful for new researchers wanting to use a marker. Depending on the platform these could be a size (143bp, Large/Small), allele (A/C/G/T) or even an allele definition (e.g. for amplicon sequencing).
Unfavourable allele		
False Positive Rate	Core	The proportion of known negative genotypes incorrectly classified as QTL[+]. Assayed as the number of <i>known</i> recipients identified as not having an unfavourable allele(s) of the marker (and thus incorrectly classified as QTL[+]). $\frac{(\# \text{ recipients withOUT unfavourable alleles})}{(\text{Total } \# \text{ known recipients})}$
False Negative Rate	Core	This is the converse of FPR, i.e. the proportion of known QTL[+] genotypes incorrectly classified as QTL[-] due to not having a favourable allele of the marker. $\frac{(\# \text{ donors withOUT favourable alleles})}{(\text{Total } \# \text{ known donors})}$
Breeding metrics		
Breeding program False Positive Rate	Support	Analogous to the FPR and FNR respectively, but assessed on a particular breeding panel. Since the breeding panel likely has a lower level of allelic diversity, markers should score higher on the BpFPR and BpFNR than on the true FPR/FNR. These metrics give more precise information on the behaviour of a marker in the target breeding program, but the association is then specific to that program, so must be reassessed for other programs.
Breeding program False Negative Rate		
Utility	Core	The proportion (percentage) of a prospective breeding pool across which a marker could be used to introgress a QTL. This is equivalent to the proportion of the pool which does NOT carry the donor allele of a marker: $\frac{(\# \text{ cultivars withOUT favourable alleles})}{(\text{Total } \# \text{ cultivars assessed})}$

Technical metrics: Call rate and clarity

Technical metrics such as call rate and clarity must be determined empirically from genotyping data using a specific marker assay (primers, probes, etc.); it is entirely possible for two independent sets of primers targeting the same locus to give vastly different results on these metrics. To this end, a set of 20 SSR and 86 trait-specific indel markers (S1 Table) were empirically evaluated on a set of 121 diverse varieties released by the International Rice Research Institute and others, supplemented with various QTL donor and recipient germplasm characterised as part of QTL mapping exercises by numerous groups. A list of varieties examined is found in S2 Table. Indel marker genotypes were scored as Large/Small, while SSR products were assigned to band size categories as appropriate. Missing and unclear results were flagged as such. The call rate was determined as the percentage of samples giving a visible result (i.e. not “missing”), while clarity was the percentage of samples giving a clearly scorable result (i.e. as opposed to unclear or ambiguous results).

Biological and breeding metrics: false positive rate, false negative rate, and utility

These metrics need to be determined against a background of high allelic diversity. To achieve this, whole-genome resequencing data was obtained for a set of 242 diverse rice accessions, comprising 173 cultivated lines (named, released varieties) and 69 landraces, most of which were chosen for their status as QTL donors or recipients. Much of this data was obtained from the rice 3000 genome dataset [26], supplemented with resequencing of high-value donors and recipients for specific QTLs. A list of varieties examined is found in S3 Table. Raw data (reads) were mapped to the MSU7 rice genome build using bwa, and resulting bam files processed using samtools [27, 28]. A total of 352 anonymous SNPs (i.e. not designed specifically for a given QTL) were chosen from the ~4500 useable features represented on the Infinium SNP chip [6], either within QTL limits (for large QTLs) or within similar distances to the QTL-specific SNP and indel markers. A total of 482 QTL-specific markers, both SNP and indel, were chosen within QTL limits, or within short physical distances of known, cloned genes. Details of the marker positions examined are found in S4 Table. Nucleotide base calls were obtained at all SNP sites – both QTL-specific SNPs and anonymous SNP markers – using standard Samtools/bcftools pipelines [27]. Genotype calls for QTL-specific indel marker positions were determined

manually from the same dataset, as automated variant-calling algorithms did not produce reliable results for indels $> \sim 5$ nt.

All data was consolidated in a MS Access database and information on favourable and unfavourable alleles was recorded for each marker. For anonymous SNPs, which do not have defined favourable or unfavourable alleles, data on accuracy metrics was calculated in two ways: first using the assigned allele A as favourable and B as unfavourable, and secondly classifying the allele with the highest frequency in known donor lines (lowest false negative rate; FNR) as favourable. The latter method was designated as “FNR corrected”. Data was analysed across 42 target QTLs for various disease resistance, abiotic stress, yield and flowering-related traits (S5 Table). As each QTL typically spanned a significant physical distance, summary data was calculated first across all markers of a particular type within a QTL, then averaged across the 42 target QTLs. The final dataset consisted of 352 anonymous SNPs, 251 trait-specific SNPs and 223 trait-specific indels, scored across a common set of 242 genomes for the 42 target QTLs.

Results

Derivation of quality metrics.

In developing a set of metrics to assess the performance of a candidate marker, it is necessary to break down the features of a marker that impact on its reliability. Broadly speaking, markers may vary in three main areas:

1. Technical aspects related to the assay and scoring of the marker;
2. Biological aspects of the association between the marker and its target locus;
3. Practical aspects of the marker’s use in a breeding pool.

These three areas are quite independent of each other; there are many markers which score well in some categories but fail completely in others, and thus for an accurate picture of how a marker behaves, all three areas must be examined. In addition, to adequately characterise a marker in each area, the areas themselves must be broken down into several measurable quality metrics (Table 1). The metrics are divided into five core metrics that quantify the reliability of a marker, and another nine supporting metrics that enable the calculation of the core metrics and deal with more difficult case studies (Fig 1).

160

161 **Fig 1. Overview of the marker quality metrics.** Core metrics capture the most critical information relating to
162 marker performance, accuracy and usefulness. Supporting metrics are needed in the calculation of the core
163 metrics and/or capture other important information, but are not routinely required in making breeding decisions.

164

165 **Technical metrics**

166 Technical metrics relate to how confidently a randomly selected sample from a genotyping job gives an accurate
167 result and can be captured clearly by two core metrics: call rate and clarity. Call rate is the proportion of
168 samples that give a scorable result (as opposed to a “missing” result). Many commercial genotyping platforms
169 already report call rates as a metric of platform performance; typical claims are >99%. Less often reported are
170 estimates of a marker’s clarity. At its simplest, clarity is a subjective opinion on how *clear* the results are, i.e.
171 how reliably genotypes A, B and H can be distinguished. In a more objective sense, an estimate of this could be
172 obtained from how often samples with known genotype are reported to have the correct score, or how often
173 duplicate samples match. Commercial SNP platforms occasionally report statistics on clarity (or repeatability),
174 but without recourse to raw data – which is rarely available – these are difficult to verify. Finally, since every
175 existing technology makes some use of target-specific oligonucleotides, different marker assays on the same
176 platform will vary in their quality on these metrics, even if they target the same position and polymorphism.
177 Thus, versioning of the marker is necessary to allow distinguishing the performance of alternate forms of a
178 marker.

179 **Biological metrics**

180 Biological metrics can be broken down into qualitative metrics that describe the *type* of association between a
181 marker and QTL, as well as quantitative metrics that describe the *level* of association. These are the most
182 important metrics for successful MAS, and also the most complex. The root cause of differences in marker
183 associations stem from the evolution of traits in an organism, and specifically the relative evolutionary timelines
184 in which the causative allele for a QTL and the polymorphism at the proposed candidate marker emerged. An
185 illustration of this is given in Figs 2 and 3. In all cases, irrespective of whether the causal allele is favourable or

unfavourable, the polymorphism most reliable in classifying the donor and recipient lines is one which arose at the same time (and in the same lineage) as the mutation giving rise to the causal, derived allele. This is because the causal mutation and the marker allele are in perfect LD when they emerged and remain in perfect LD in the gene pool consistent with the probability of a recombination event between them.

Fig 2. Illustration of the evolution of a QTL. **A** Starting from an ancestral point, mutations in a particular gene accumulate (numbered black bars, representing mutations in **B**), resulting in new alleles. At some point, a mutation arises which improves a trait, resulting in a donor allele for a QTL (dark/red), distinguishing it from known recipient alleles (light/blue); typically the status of many alleles is unknown (white). **C** Each mutation (1 – 21) is a potential marker, and all are found in the same gene, but some are more informative than the others. Comparing the false positive and false negative rates for each mutation allows the determination of which polymorphism gives the most reliable discrimination between the donor and recipient phenotype classes.

Fig 3. Determining the optimal polymorphism under several evolutionary scenarios. Depending on when the causative mutation arose (arrowed bars; outwards pointing for mutations conferring favourable or inwards facing for unfavourable alleles respectively), there may be only one favourable allele (**A**), a small number of alternative favourable alleles and multiple unfavourable alleles (**B**), or multiple favourable alleles and a few unfavourable (**C**). In **A** and **B**, the derived allele for the QTL is the favourable allele(s); in **C** it is the unfavourable. In all cases the polymorphism which gives the most accurate classification of donor and recipient status is one which arose in the same lineage and at a similar time to the causal, derived mutation.

From this theoretical consideration, it is clear that a number of parameters must be specified in order to accurately describe the association of a marker with its target QTL. Descriptive metrics such as which allele of the QTL (favourable or unfavourable) represents the derived state, the allele (favourable or unfavourable) targeted by a marker, and specifications of marker linkage, favourable and unfavourable alleles, all describe the *type* of association the marker has with its target QTL. Most are easily determined, although determining which QTL allele is derived may require a detailed genomic investigation. Nonetheless, if this can be determined accurately, then markers specific to the derived allele have the greatest chance of also providing a reliable classification across novel allelic diversity, reducing the risk of incorrect classification in future breeding efforts.

Therefore, expending some effort to determine the derived QTL allele before designing large numbers of markers is justified.

The proposed quantitative QC metrics of false positive rate and false negative rate (FPR and FNR) describe the *level* of association between the marker and its target QTL and are arguably the most important of the metrics presented, but also the most difficult to estimate. The FNR is the proportion of known donor lines that are incorrectly classified as QTL[-] by the marker. Since the lines are known to carry favourable alleles of the QTL, classification of any of these lines as QTL[-] thus represents a false-negative call by the marker. This is the converse of the marker's specificity. Markers with a high FNR pose a significant risk of mis-classifying samples as QTL[-] when they do in fact carry a favourable allele; thus breeding material may be discarded which in reality could have been advanced. Markers with a low FNR will correctly identify all samples that possess the QTL[+] state but may still mis-classify samples with an unfavourable (non-donor) alleles as QTL[+], in other words a low FNR does not imply a low FPR.

The FPR is simply the converse and represents the proportion of non-donor (recipient) lines that are incorrectly classified as QTL[+], and is the converse of the marker's sensitivity. In a breeding context, a high FPR means there is a significant risk of investing in and advancing lines based on MAS results that indicate the presence of the QTL[+] allele, but in reality are QTL[-].

It is important to reiterate here that the FPR and FNR are the proportion of lines with *known* QTL[+]/[-] alleles that are correctly classified as such. They are fundamentally linked to the diversity of alleles with known function, which is determined by the effort that has been put into characterising/mapping donor and non-donor diversity. Many markers are chosen for breeding applications based on their linkage with QTL alleles in specific mapping populations where the QTL is discovered. But those markers – while informative in the chosen mapping population – could still score poorly on both FPR and FNR, resulting in poor performance once deployed as MAS targets. The difficulty arises because the marker is being applied to new populations, where at least one of the parents is of unknown status with respect to the QTL. Markers with low FPR and low FNR will faithfully report the presence or absence of the QTL in any sample, irrespective of allelic diversity in any gene pool of interest.

Breeding metrics

Breeding QC metrics describe the relative value of applying a marker in a specific breeding program. These consist of three metrics: Breeding program false positive rate (BpFPR), Breeding program false negative rate (BpFNR) and Utility. BpFPR and BpFNR are equivalent to the FPR and FNR metrics described above, but are specific to particular breeding program in which they are assessed, rather than on the full diversity of known donors and recipients. Since the breeding pool may be expected to have lower allelic diversity than occurs species-wide, and because selection and genetic drift are modifying patterns of LD independently across breeding programs, these rates can be quite different from the true FPR and FNR (the usual expectation would be the breeding program rates to be lower than the true rates, though the opposite could also occur). They may also be different for different breeding programs, and must be assessed independently for each. They will require the determination of donor and recipient lines *within a breeding program*, which will involve collecting phenotype data for each program under investigation. But once gathered, the QC metrics directly quantify the marker's reliability for making breeding decisions in that specific program. It's worth noting however, that this is predicated on the assumption that the assessed panel represents the complete allelic constituency present in a breeding program, and that the breeding strategy focuses on increasing the frequency of favourable haplotypes through recombination in a closed gene pool, minimizing the introduction of novel allelic variation which may introduce marker alleles that are not in LD with the causal polymorphism.

Finally, utility is the proportion of a breeding panel over which a marker could be used to select for its associated QTL[+] allele. This is basically an assessment of the frequency of QTL[+] alleles in the breeding program, easily calculated as the proportion of breeding lines which possess *non-donor* allele(s) of the marker (Fig 4). A program where the marker offers high utility by definition has the QTL[+] allele at low frequency. Note that this is entirely separate from the FPR and FNR: a marker may perfectly classify all material as to its QTL status, but if the QTL is fixed in the breeding program then the marker (QTL) is of little utility in improving the trait. A good example of this in rice would be *sd1*, which due to intense selection pressure for plant height and heavy usage of QTL[+] *sd1* green revolution varieties by breeding programs is fixed in nearly all breeding populations and therefore is not available to manipulate plant height.

267

268 **Fig 4. Utility of a marker.** Alternative markers within a QTL region (markers A – E) each have multiple
 269 alleles (numbered). Alternative alleles of each marker are found at differing frequencies within a breeding pool.
 270 Those markers with a high frequency of the favourable allele in the breeding pool (B, C, D) – and thus low
 271 utility – can only distinguish the donor genotype in a small number of breeding backgrounds. By contrast
 272 markers A and E have high or very high utility, as they are polymorphic with respect to nearly all target genomes
 273 in the breeding pool.

274

275 A summary of all the proposed metrics is presented in Table 1. Several of these metrics are purely descriptive
 276 (derived QTL allele, marker target allele, donor and recipient alleles), but are required to allow the calculation of
 277 the more quantitative parameters; these are called supporting metrics. The quantitative parameters then provide
 278 a detailed assessment of the performance of a marker; these are the core metrics, and provide the best criteria for
 279 assessing markers. Ideal target values and consequences using markers with unfavourable performance scores
 280 using these metrics are explained in Table 2. Of particular note are the core metrics FPR and FNR; poor scores
 281 on these will increase the probability of discarding valuable breeding germplasm, or worse, wasting resources
 282 advancing QTL[-] lines. Assuming good scores on FPR and FNR (or BpFPR and BpFNR), a poor utility value
 283 indicates the marker (and thus the QTL[+]) allele is nearly monomorphic in the breeding program and can only
 284 be used to select for the QTL across a narrow/small proportion of the breeding pool. Finally the derived marker
 285 allele (donor/recipient) should ideally match the derived allele of the QTL; if so, the marker stands a much better
 286 chance of correctly classifying additional unknown or uncharacterised alleles.

Table 2. Ideal target values for key metrics

Metric	Target value	Consequences of unfavourable score
Call rate	100%	High number of samples with missing data (cannot be genotyped); reduction in efficiency of program and possibility of missing high-value germplasm.
Clarity	100%	High number of samples with ambiguous genotype. Possibility of misclassifying material in the donor/recipient categories; propagating germplasm with low value, or discarding germplasm with high value.
Linkage	0cM	High chance of recombination between the marker and actual QTL, leading to a breakdown in the marker-trait association.
Marker target	Derived QTL state	Misclassification of new, unrecognised or uncharacterised allelic diversity. This is of particular use in situations where few donor and recipient lines have been characterised; if the marker target allele matches the derived QTL allele (favourable or unfavourable), the marker is more likely to correctly classify new, uncharacterised alleles.
FPR/BpFPR	0%	Failure to distinguish some unfavourable alleles from some favourable alleles; reports presence of QTL when it is actually absent. May result in use of a line as donor when it is not, or a lack of effort to use a QTL to improve a trait when this would be appropriate. This is in general a more serious failure than low FNR.
FNR/BpFNR	0%	Failure to distinguish some favourable alleles from unfavourable ones; reports absence of QTL when it is present. May result in ineffective MAS (wasted effort) due to recipient lines being classified as QTL[-] and thus the QTL being introgressed, when in fact it is already present.
Utility	100%	Marker can be used to track/introgress QTL across only a narrow range of germplasm. For most populations alternative markers are needed.

Applying quality metrics: Evaluation of existing marker systems

Technical metrics: Indels vs. SSRs.

Since technical metrics relate to the performance of a specific marker assay, they are by necessity empirical and may vary widely between markers even when these have the same biological properties and are run on the same platform. Indeed, wide variation was seen between markers for both call rate and clarity, even within trait-specific indel markers in specific QTL regions such as qDTY4.1 and qNa1L (Fig 5).

Fig 5. Technical metrics for candidate indel markers. Markers were evaluated within the qDTY4.1 (A) and qNa1L (B) QTL regions. Significant variation was seen for different markers in both QTL regions, showing some markers clearly performed better than others.

Technical metrics can be used to assess the relative performance of platforms as well as specific markers. The mean call rate and clarity were compared between a set of SSR and QTL-specific indel markers on a panel of 122 diverse cultivars (Fig 6). Trait-specific indel markers significantly out-performed SSR markers for clarity ($P < 0.05$). They also scored better than SSRs for missing data, though this was not significant ($0.05 < P < 0.12$), and may reflect a greater-than-average contribution from a few specific indels that scored very poorly, due largely to a few markers that covered genomic deletions in *qHTSF4* and *qDTY4.1*.

Fig 6. Comparison of performance of PAGE-based marker systems. The performance of SSRs was compared to that of trait-specific indels. Indels consistently out-performed SSRs, particularly in marker clarity but also in call rate.

Biological metrics: Anonymous vs. Trait-specific markers.

Working with whole-genome resequencing data it is evident that numerous polymorphisms can be easily identified between two varieties. However, these polymorphisms vary widely in their level of association with a target QTL. The level of association (FPR and FNR) for candidate markers throughout a salinity tolerance QTL in rice between 37 and 41Mb on the long arm of chromosome 1 shows wide variation, all through the QTL interval (Fig 7). This shows linkage with a QTL is not sufficient to give reliable selection. Secondly, none of the anonymous SNPs found on the Infinium chip within the QTL region (Fig 7a) scored perfectly on both the FPR and FNR, indicating they all suffer from errors in classifying known varieties. These SNPs have been filtered for those which show polymorphism between known donors and recipients, and corrected to minimise the false negative rate (correctly identifying as many donors as possible). In contrast, while QTL-specific indel and SNP markers (Figs 7b and c) also show variation in their association across the QTL, both classes have several markers achieving ideal scores (0%) on both metrics. Those markers scoring $>0\%$ on either the FPR or

FNR may find niche applications in fine-mapping or specific populations, but those with perfect scores would make the most reliable marker system for breeding purposes. These metrics thus give information allowing the identification and design of optimal marker systems, even for quite large QTL regions.

Fig 7. Use of marker quality metrics to determine optimal markers for a QTL. Comparison of quality metrics for different markers within a QTL region for salinity tolerance, qNa1L, between 37 and 41Mb on the long arm of chromosome 1. Multiple markers from the Illumina Infinium chip (Anonymous SNPs; favourable allele corrected to minimise FNR), QTL-specific SNPs and QTL-specific indels were assessed for their utility (A), false-positive rate (B) and false-negative rate (C) across IRRI germplasm. Anonymous SNPs typically scored poorly on FPR, FNR and Utility, and none scored well on all metrics. QTL-specific SNP and indel markers typically showed low or perfect scores on the FPR, and several markers scored 0% (no misclassified entries) on both FPR and FNR metrics. In addition QTL-specific markers scored far better for utility, indicating wider applicability in breeding. Arrows to the right indicate ideal target values for a new marker.

Improvements in the false positive and negative rates with QTL-specific markers are also seen for other QTL targets. Mean values from 42 QTLs for a range of traits including stress tolerance, grain quality, disease resistance etc. show that anonymous SNPs derived from the 6k diversity set consistently under-perform (Fig 8). Re-assignment of favourable and unfavourable alleles to minimise the false negative rate improves that metric to a level equivalent to those seen for the QTL-specific markers, but no better ($P > 0.05$). Optimising the FNR also improved the FPR, but not to a level equivalent to QTL-specific markers, so anonymous markers still performed worse on average ($P < 0.0001$). This means these anonymous markers, even with “corrected” assignments of favourable allele, have no relative benefit in detecting presence of the QTL, but do have a significant penalty in their FPR, incorrectly classifying lines as QTL[+].

Fig 8. Comparison of mean accuracy metrics for diversity SNPs and QTL-specific SNP and indel markers. Anonymous SNP markers initially have very low scores on both the FPR and FNR. Correcting the assignment of favourable and unfavourable alleles to minimise the FNR improves that score to equivalent levels with the QTL-specific markers, but the FPR remains poor and no benefit is seen for the utility.

Breeding metrics: Utility.

Anonymous SNPs also showed lower average utility values (Fig 8), i.e. the designated favourable allele is present at higher frequencies in elite germplasm, and so the marker is less useful for introgressing a given QTL into a range of elite material. The utility metric is especially useful in the case of diagnostic markers, as it then indicates directly the proportion of elite material that a QTL may improve. Utility values for a range of QTL controlling various yield, grain quality, disease resistance and stress tolerance traits show a wide range in variation (Fig 9). These range from less than 20% for *LTG1*, *qSCT1* and *SCM2*, which appear to be fixed in nearly all *indica* elite material, to 100% for many disease resistance QTL. The latter observation is surprising considering the substantial selective pressure exerted on disease resistance in most breeding programs, and further work to determine its cause seems warranted.

Fig 9. Variation in utility between various QTL for agronomic traits in *indica* breeding germplasm. QTL were selected that have diagnostic markers or markers scoring 0% on both FPR and FNR (and thus could be accurately scored). Wide variation in QTL utilities were seen, from near-fixation (utility ~0%) to absent (utility 100%), but most were rare or absent.

Discussion

Since the 1980s with the advent of SSR markers, it has become almost a mantra that the ideal marker should be highly polymorphic. This is certainly a useful feature for certain applications such as in bi-parental mapping, where high polymorphic information contents (PIC) increase the chances a given marker will be polymorphic between random parents. However, marker-assisted selection places different demands on the markers – the number of alleles displayed by the marker is not relevant, rather the ability to unambiguously discriminate between all donor and recipient material becomes critical. Surprisingly, there is a dearth of literature on designing reliable markers and almost no criteria for judging what makes a marker “good” or “bad”. Most MAS programs use markers identified in QTL mapping populations – typically SSRs, applying them to other genetic backgrounds, and even attempting to use them to determine the

presence of a QTL in diverse germplasm panels. These applications require very stringent false positive and false negative rates, but few examples exist where some validation of these false positive and negative rates has been conducted. Bernardo *et al.* [20] examined the reliability of SSR markers in selecting for stem rust resistance in wheat. Association of the markers with donor and recipient germplasm was analysed, but not quantified as a metric; association was not always good, for example markers for *Sr32* targeted the recipient allele not the (derived/wild) donor allele, and so run the risk of classifying some lines without the QTL as positive. Similarly, markers for other resistance genes variously failed to distinguish some or all known recipients from known donors (poor FPR; many SSRs had this difficulty), while others showed poor separation of donor and recipient alleles, produced major stutter bands, or produced non-target amplicons (poor clarity, figures 1 and 3 in their paper). In another strategy, Mohammadi-Nejad *et al.* [15] performed allele mining of the *Saltol* QTL using SSRs. Of particular note is their table 3, which lists the SSR haplotype and varieties which possess this haplotype. This can be related to alleles of the *HKT1;5* gene [23], which is causal for this QTL [29]. This reveals the SSR markers confounded (failed to distinguish) lines that had different alleles, such as IR64 and Kala Rata. Likewise, the converse was even more common: IR64 / IR29, and Pokkali / Kala Rata / Sadri were placed in separate haplotypes while sharing the same allele at the causal gene. Thus again the SSR genotype, and even haplotype, was not sufficient to reliably classify donor and recipient material for the QTL. In a counter-example, Tian *et al.* [30] designed indel markers for the allelic major rice blast resistance gene *Pi2/Pi9* based on sequence comparisons of parental varieties. The FPR and FNR were not quantified, but the inclusion of multiple reference alleles presumably helped in the design of markers highly specific to the favourable alleles. This marker was then able to demonstrate near-zero occurrence of this gene in a set of Chinese breeding germplasm, where previous marker sets were prone to false positives (see [31]) – incidentally indicating a potentially very high utility, though this was not articulated as a metric. A similar effort was conducted by Scheuermann and Jia [32] using a different approach. The latter marker from Scheuermann and Jia [32] should have the same FPR and FNR as that designed by Tian *et al.* [30], but suffered very low apparent call rates and clarity. In neither case were any metrics similar to FPR, FNR, Utility and Clarity quantified by either group, despite their datasets being sufficient to do so. Had these metrics been quantified, they could easily

404 demonstrate which marker system is better and how reliably these markers could be used in other breeding
 405 programs, thereby greatly enhancing the impact of this work.

406 These examples show the need for a better system for describing the association of a marker with its target QTL.
 407 The fourteen metrics described in Table 1 are a substantial step towards providing such a description.

408 Association of a marker with its target QTL is captured by a range of biological metrics rooted in the preceding
 409 discussion on the evolution of markers. Additional metrics describe parameters relating to reliability of the
 410 genotyping information, and the applicability of a marker in specific breeding situations. The preceding
 411 considerations have shown how the ideal marker – one which reliably identifies all donor and recipient
 412 germplasm – is based on the same polymorphism as gives rise to the QTL phenotype. Such a marker can be
 413 called diagnostic, and requires the identification of the gene *and the mutation* giving rise to a QTL – something
 414 that is very rarely done, even in rice. While ideal, this is difficult and time consuming. Alternative, flanking
 415 markers can still accurately classify observed alleles provided they arose at similar times and in the same lineage
 416 as the causative mutation (i.e. the derived marker allele matches the derived QTL allele; Fig 3). Again, the
 417 metrics in Table 1 provide a means to evaluate and validate candidate markers before committing to design and
 418 implementation (very important for expensive SNP systems) as well as assess performance after implementation.

419 Validating existing marker systems with these metrics illustrates several points. First and foremost, QTL-
 420 specific marker systems consistently out-perform both older SSR and new anonymous SNP systems in most of
 421 these metrics, but most notably in the accuracy metrics FPR and FNR. For many QTL, no anonymous markers
 422 showed the required level of association with the target QTL. Thus QTL-specific markers will give consistently
 423 more reliable results in selection. In addition, accurate markers (scoring 0% on both FPR and FNR) can be used
 424 to determine the proportion of a breeding panel that may benefit from that QTL – the utility. Utility values vary
 425 widely between QTL (Fig 9), which reflects a complex interplay of the QTL's origin and the artificial and
 426 natural selective pressures it has been subjected to in breeding programs. For example, *SCM2* is widely regarded
 427 as a candidate to reduce lodging, a major problem even in semi-dwarf rice. Unfortunately, however, the
 428 characterised donor allele of *SCM2* from Habataki [33] appears identical to that already found in the vast
 429 majority of *indica* breeding germplasm. This means the donor allele is already present in most or all improved

430 *indica* material, and therefore the Utility of *SCM2* (and accurate markers for this gene) is very low in most *indica*
 431 breeding programs. By contrast many of the disease resistance loci have high utility, despite being under strong
 432 selective pressure in breeding programs.

433 Secondly, there is no inherent advantage of SNP genotyping platforms in terms of selection accuracy. Existing
 434 anonymous SNP marker sets such as the 6k Infinium chip were designed to maximise the probability of
 435 polymorphism (discriminatory power) between randomly-chosen varieties [6]. This makes them ideal for
 436 population genetics and as a fixed panel of markers to genotype any random set of parents and progeny.

437 Ironically though, this means they have the least power to track specific QTL, and indeed they perform rather
 438 poorly in accuracy metrics overall. By contrast, QTL-specific markers, whether high-throughput SNP or low-
 439 throughput indel, perform quite well on accuracy metrics – and it is certainly possible to use these metrics to
 440 identify both SNP and indel markers with “perfect” associations with their target QTL (Fig 7). Therefore, the
 441 choice of marker *platform* has less to do with selection accuracy than with the expected sample throughput.

442 However, the choice of best *marker* will be based on the level of association with the target QTL.

443 Thirdly the application of the metrics in evaluating individual markers is easily demonstrated (Figs 5 and 7).
 444 Biological and breeding metrics are mostly useful in distinguishing between candidate markers prior to
 445 committing time and resources to implementing these on a particular genotyping platform; these are about
 446 choosing the optimal target polymorphisms. Yet although a candidate marker may have perfect biological and
 447 breeding metrics, a given assay for that polymorphism may perform very poorly on its technical metrics (call
 448 rate and clarity), such as seen for *qDTY4.1* (Fig 5). For PCR-based systems such as SSRs, indels and some SNP
 449 technologies, much of this variation in the technical metrics is due to inherent issues with primer efficiency.

450 However, while SNP assays also utilise oligonucleotides as either probe or primer sequences, it is unfortunately
 451 rare to see “validation” of the performance of new markers on SNP platforms due to the expense of an assay. In
 452 addition, SNP assays very rarely return the raw data, instead reporting a digested summary –the actual SNP call
 453 –thus glossing over such factors as whether the clustering of fluorescence intensities was unambiguous. It seems
 454 advisable going forward to implement some form of technical and biological replication when validating a new
 455 assay to determine its clarity.

456 Allied to this discussion on technical metrics, it is evident that a given polymorphism may be interrogated by
 457 multiple different “markers”. For example, a SNP could be targeted by one of several gel-based assay systems,
 458 any of the many SNP platforms, or an amplicon/sequencing approach. All are based on the same polymorphism,
 459 but as each platform has its own design quirks and even a single platform may use alternate primer pairs for
 460 amplification, different instances of a marker could have very different scores on technical metrics, despite
 461 targeting the same polymorphism. Therefore, information on the version/instance of a marker under
 462 consideration is needed to distinguish between alternate forms and platforms targeting the same polymorphism.
 463 These quality metrics thus provide a good framework for assessing the accuracy and reliability of any specific
 464 marker. This information can be used to evaluate existing markers, design better markers, and even to compare
 465 performance of different marker types and platforms. Nevertheless, they are by no means complete or perfect.
 466 Two issues worth highlighting concern the Clarity metric and estimation of the FPR and FNR metrics.
 467 The Clarity metric is currently slightly ambiguous; it could refer to either how clearly/reliably genotyping data
 468 (bands on a gel, fluorescence signal clusters on a SNP platform, or other measures) can distinguish between the
 469 allelic states of the marker. It could also refer to how often duplicate samples cluster together – repeatability.
 470 This is of course closely related to the former situation, but is also subtly different. For the sake of simplicity
 471 these are not distinguished here, but further discussion on whether Clarity as described here should be broken
 472 down into two metrics – clarity and repeatability – seems warranted.
 473 Of the metrics presented, the accuracy metrics FPR and FNR are arguably the most important in distinguishing
 474 between candidate markers. A common objection and difficulty in the assessment of these is that they require a
 475 large number of case-[+] and case-[-] data points to accurately estimate, whereas in most cases the number of
 476 defined donor and recipient lines is limited. This is absolutely true – and should inform how QTL mapping and
 477 validation efforts are undertaken – but is not a justification to reject the importance of these metrics. On one
 478 hand, some estimate is better than none, and it should be recognised that all datasets are inadequate, to some
 479 extent. It is thus better to report the statistics, together with data on how accurate these might be – such as the
 480 total number of defined donor and recipient lines available, where a greater number of both implies greater
 481 accuracy in estimation. On the other hand, estimates of FPR and FNR based on inadequate (small) datasets

actually *inflate* their values. In the extreme case of a single known donor and recipient line, any marker polymorphic between these parents will then score 0% on both metrics. An inadequate dataset thus does not lead to a rejection of markers, but rather the opposite: the power to distinguish between them is limited, and so any feature that shows polymorphism will be accepted. In this situation an educated guess as to whether the favourable or unfavourable allele is likely to be derived (i.e. the derived QTL allele). Markers targeting this (i.e. the marker target allele is the same as the derived QTL allele) are then effectively making the assumption that the derived QTL allele is rare in the overall allelic diversity, thus deliberately biasing the error towards false-negatives and away from false positives – maximising the probability of a good FPR at the potential penalty of FNR – and thus biasing risk away from advancing unfavourable genotypes at the penalty of increasing risk of discarding favourable ones.

In summary, the metrics proposed in Table 1 quantify all significant parameters describing a marker’s behaviour when assayed, the level and type of association it displays with its target QTL (both species-wide and in specified breeding panels), and its distribution within a breeding panel. These metrics give a fast, comprehensive and objective means to discriminate between and evaluate alternative markers (e.g. Fig 7), allowing an optimal marker system to be designed. In addition, by including such “housekeeping” metrics as the favourable and unfavourable alleles, it is possible to automate these calculations, providing the possibility to scan genomic datasets for optimal SNP markers programmatically, greatly simplifying the deployment of QTL in breeding. It also allows the development of a marker database with live updating of metrics as new data is added, enabling continual refinement of marker systems. The advantages of adopting a set of metrics are manifold, and it is hoped that the proposed system will assist in developing a new generation of reliable marker systems to improve the efficiency of plant breeding.

Acknowledgements

The authors wish to thank Irish Bagsic, Chenie Zamora and Katreena Titong for technical assistance in various aspects of this work.

References

1. Alexandratos, N and Bruinsma, J. World agriculture towards 2030/2050: the 2012 revision. ESA Working paper No. 12-03. Rome: FAO; 2012.
2. Collard BCY, Mackill DJ. Marker-assisted selection: an approach for plant breeding in the twenty-first century. Phil Trans R Soc Lond B Biol Sci. 2008;363: 557–72.
3. Temnykh S, Park WD, Ayres N, Cartinhour S, Hauck N, Lipovich L, et al. Mapping and genome organization of microsatellite sequences in rice (*Oryza sativa* L.). Theor Appl Genet. 2000;100: 697–712.
4. International Rice Genome Sequencing Project. The map-based sequence of the rice genome. Nature 2005;436: 793–800.
5. Hayashi K, Yoshida H, Ashikawa I. Development of PCR-based allele-specific and InDel marker sets for nine rice blast resistance genes. Theor Appl Genet. 2006;113: 251–60.
6. Thomson MJ. High-Throughput SNP Genotyping to Accelerate Crop Improvement. Plant Breed Biotechnol. 2014;2: 195–212.
7. Bradbury LMT, Fitzgerald TL, Henry RJ, Jin Q, Waters DLE. The gene for fragrance in rice. Plant Biotech J. 2005;3: 363–70.
8. Weng J, Gu S, Wan X, Gao H, Guo T, Su N, et al. Isolation and initial characterization of *GW5*, a major QTL associated with rice grain width and weight. Cell Res. 2008;18: 1199–209.
9. Dixit S, Swamy B, Vikram P, Ahmed H, Sta Cruz M, Amante M, et al. Fine mapping of QTLs for rice grain yield under drought reveals sub-QTLs conferring a response to variable drought severities. Theor Appl Genet. 2012;125: 155–69.
10. Singh R, Singh AK, Sharma TR, Singh A, Singh NK. Fine mapping of grain length QTLs on chromosomes 1 and 7 in Basmati rice (*Oryza sativa* L.). J Plant Biochem Biotechnol. 2012;21: 157–66.
11. Ghimire KH, Quiatchon LA, Vikram P, Swamy BM, Dixit S, Ahmed H, et al. Identification and mapping of a QTL (qDTY1.1) with a consistent effect on grain yield under drought. Field Crops Res. 2012;131: 88–96.
12. Swamy BPM, Ahmed HU, Henry A, Mauleon R, Dixit S, Vikram P, et al. Genetic, Physiological, and Gene Expression Analyses Reveal That Multiple QTL Enhance Yield of Rice Mega-Variety IR64 under Drought. PLoS ONE 2013;8: e62795.
13. Yadaw RB, Dixit S, Raman A, Mishra KK, Vikram P, Swamy BM, et al. A QTL for high grain yield under lowland drought in the background of popular rice variety Sabitri from Nepal. Field Crops Res. 2013;144: 281–87.
14. Yang T, Zhang S, Zhao J, Liu Q, Huang Z, Mao X, et al. Identification and pyramiding of QTLs for cold tolerance at the bud bursting and the seedling stages by use of single segment substitution lines in rice (*Oryza sativa* L.). Mol Breed. 2016;36: 96.
15. Mohammadi-Nejad G, Arzani A, Rezal AM, Singh RK, Gregorio GB. Assessment of rice genotypes for salt tolerance using microsatellite markers associated with the *Saltol* QTL. Afr J Biotechnol. 2008;7: 730–36.
16. Singh AK, Singh PK, Arya M, Singh NK, Singh US. Molecular Screening of Blast Resistance Genes in Rice using SSR Markers. Plant Pathol J. 2015;31: 12–24.
17. Begum H, Spindel JE, Lalusin A, Borromeo T, Gregorio G, Hernandez J, et al. Genome-Wide Association Mapping for Yield and Other Agronomic Traits in an Elite Breeding Population of Tropical Rice (*Oryza sativa*). PLoS ONE 2015;10: e0119873.
18. Spindel J, Wright M, Chen C, Cobb J, Gage J, Harrington S, et al. Bridging the genotyping gap: using genotyping by sequencing (GBS) to add high-density SNP markers and new value to traditional bi-parental mapping and breeding populations. Theor Appl Genet. 2013;126: 2699–716.
19. Babu R, Nair SK, Prasanna BM, Gupta HS. Integrating marker-assisted selection in crop breeding - Prospects and challenges. Current Sci. 2004;87: 607–19.
20. Bernardo R. Genomewide markers as cofactors for precision mapping of quantitative trait loci. Theor Appl Genet. 2013;126: 999–1009.
21. Miah G, Rafii MY, Ismail MR, Puteh AB, Rahim HA, Islam KN, Latif MA. A Review of Microsatellite Markers and Their Applications in Rice Breeding Programs to Improve Blast Disease Resistance. Int J Mol Sci. 2013;14: 22499–528.

22. Ul Haq T, Gorham J, Akhtar J, Akhtar N, Steele K. Dynamic quantitative trait loci for salt stress components on chromosome 1 of rice. *Funct Plant Biol.* 2010;37: 634–45.
23. Platten J, Egdate J, Ismail A. Salinity tolerance, Na⁺ exclusion and allele mining of *HKT1;5* in *Oryza sativa* and *O. glaberrima*: many sources, many genes, one mechanism? *BMC Plant Biol.* 2013;13: 32.
24. Lalkhen AG, McCluskey A. Clinical tests: sensitivity and specificity. *Continuing Education in Anaesthesia, Critical Care and Pain.* 2008;8: 221–3.
25. McNally KL, Childs KL, Bohnert R, Davidson RM, Zhao K, Ulat VJ, et al. Genomewide SNP variation reveals relationships among landraces and modern varieties of rice. *Proc Natl Acad Sci USA.* 2009;106: 12273–8.
26. The 3000 rice genomes project. The 3,000 rice genomes project. *GigaScience.* 2014;3: 7.
27. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence alignment/map (SAM) format and SAMtools. *Bioinformatics.* 2009;25: 2078–9.
28. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler Transform. *Bioinformatics.* 2009;25: 1754–60.
29. Ren Z-H, Gao J-P, Li L-G, Cai X-L, Huang W, Chao D-Y, et al. A rice quantitative trait locus for salt tolerance encodes a sodium transporter. *Nat Genet.* 2005;37: 1141–6.
30. Tian D, Chen Z, Chen Z, Zhou Y, Wang Z, Wang F, Chen S. Allele-specific marker-based assessment revealed that the rice blast resistance genes *Pi2* and *Pi9* have not been widely deployed in Chinese *indica* rice cultivars. *Rice.* 2016;9: 19.
31. Cho Y-C, Kwon S-W, Choi I-S, Lee S-K, Jeon J-S, Oh M-K, et al. Identification of Major Blast Resistance Genes in Korean Rice Varieties (*Oryza sativa* L.) Using Molecular Markers. *J Crop Sci Biotechnol.* 2007;10: 265–76.
32. Sheuermann KK, Jia Y. Identification of a *Pi9*-Containing Rice Germplasm with a Newly Developed Robust Marker. *Phytopathology.* 2016;106: 871–6.
33. Ookawa T, Hobo T, Yano M, Murata K, Ando T, Miura H, et al. New approach for rice improvement using a pleiotropic QTL gene for lodging resistance and yield. *Nat Commun.* 2010;1: 132.

Supporting information

S1 Table. List of indel and SSR markers assessed for technical performance.

S2 Table. List of varieties examined for technical performance evaluation.

S3 Table. List of varieties examined in calculating biological accuracy and breeding metrics.

S4 Table. List of marker positions interrogated for assessing biological accuracy and breeding metrics.

S5 Table. List of QTL examined, with start and end positions.

SUPPORTING METRICS

CORE METRICS

OUTCOMES

Technical

Version

Biological

Linkage

Position

Derived QTL state

Marker target allele

Favourable allele

Unfavourable allele

Breeding

Program-specific FPR

Program-specific FNR

Call rate
Clarity

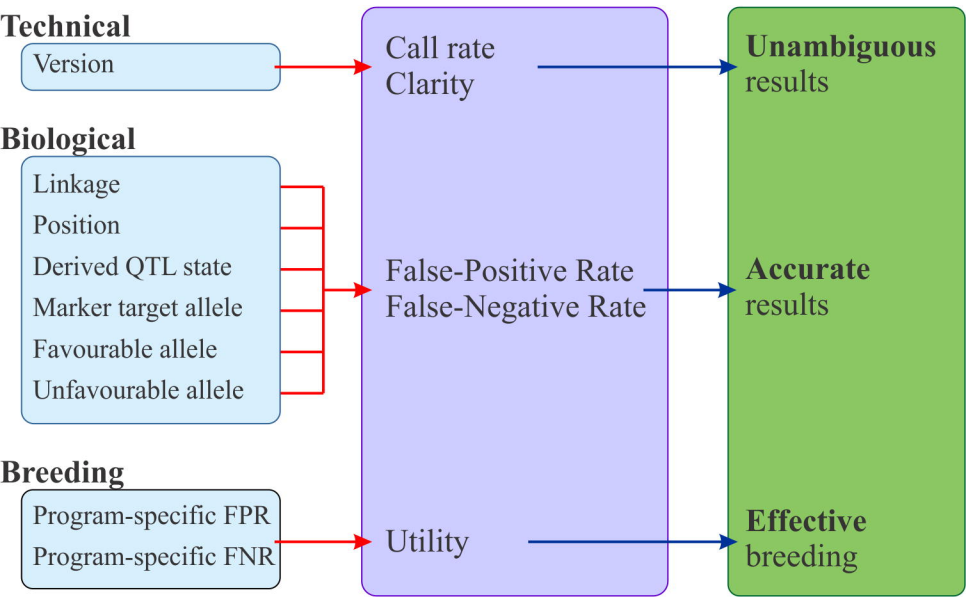
False-Positive Rate
False-Negative Rate

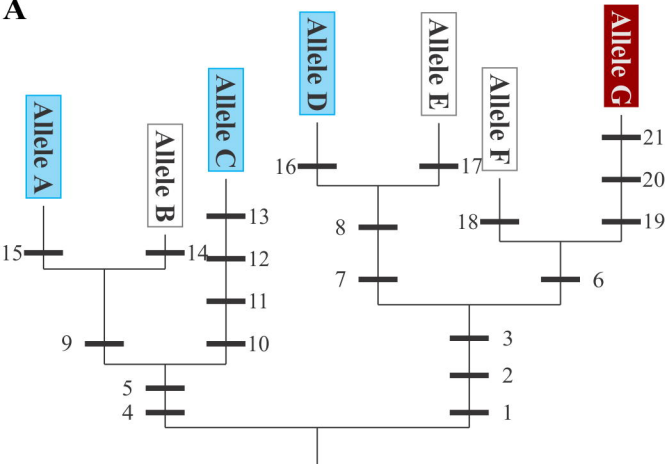
Utility

Unambiguous
results

Accurate
results

Effective
breeding



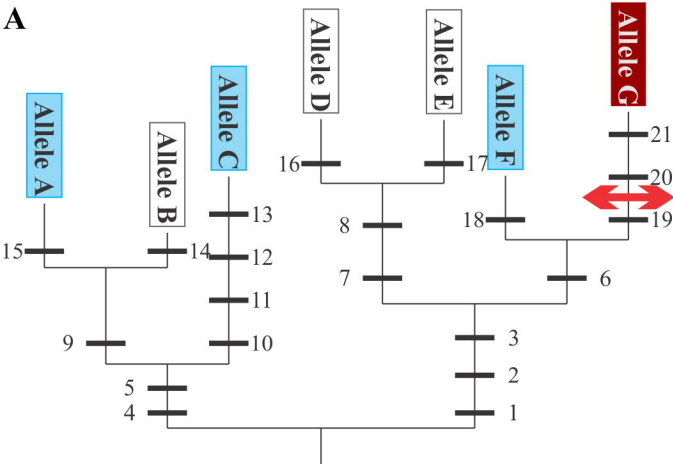
A**B**

Position	Allele						
	A	B	C	D	E	F	G
1 - 3	C	C	C	<u>T</u>	<u>T</u>	<u>T</u>	<u>T</u>
4 - 5	<u>G</u>	<u>G</u>	<u>G</u>	A	A	A	A
6	A	A	A	A	A	<u>G</u>	<u>G</u>
7 - 8	T	T	T	<u>A</u>	<u>A</u>	T	T
9	<u>C</u>	<u>C</u>	G	G	G	G	G
10 - 13	A	A	<u>C</u>	A	A	A	A
14	T	<u>A</u>	T	T	T	T	T
15	<u>G</u>	A	A	A	A	A	A
16	T	T	T	<u>G</u>	T	T	T
17	C	C	C	C	<u>A</u>	C	C
18	G	G	G	G	G	<u>A</u>	G
19 - 21	A	A	A	A	A	A	<u>T</u>

C

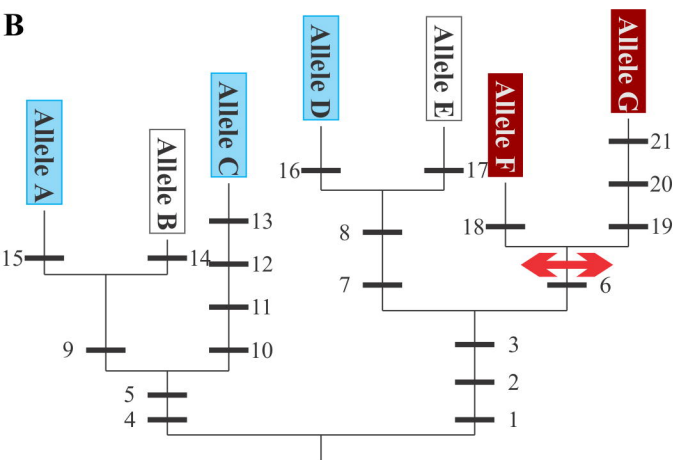
Polymorphism	FPR	FNR	Comments
1 - 3, 4 - 5	33%	0%	Good discriminatory power - good for diversity and variety typing. However conflate donor with one or more recipients.
7, 8	67%	0%	Conflate recipients A and C with donor.
6	0% *	0%	FPR (discrimination of donors) depends on status of allele F. If F is a donor, marker #6 will be accurate; if it is not, then markers 19 - 21 are accurate.
19 - 21	0% *	0%	

A



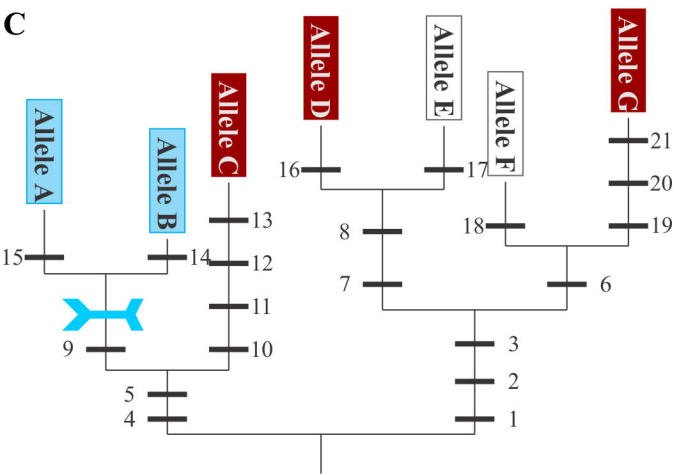
Polymorphism	FPR	FNR	Comments
1 – 3, 4 – 5	33%	0%	Conflates donor with one recipient (F).
18	67%	0%	Conflates most recipients (except F) with donor.
6	33%	0%	Conflates one recipient (F) with donor.
19 – 21	0%	0%	Accurate.

B



Polymorphism	FPR	FNR	Comments
1 – 3, 4 – 5	33%	0%	Conflates donors with one or more recipients.
7,8	67%	0%	Conflates recipients A and C with donors.
18, 19 – 21	0%	50%	Distinguish one donor, but conflate other with recipients.
6	0%	0%	Accurate.

C



Polymorphism	FPR	FNR	Comments
1 – 3, 4 – 5	0%	33%	Conflates donor C with recipients.
6, 7, 8, 19 – 21	0%	67%	Specific to one donor allele, but conflate others with recipients.
14, 15	50%	0%	Conflate one of the recipients with the donors.
9	0%	0%	Accurate.

Donor		Breeding pool (prospective recipients)										# non-donor	Utility	
	Allele													
Marker A	1	2	2	1	2	2	3	2	3	2	2	A	9	90%
Marker B	2	2	1	2	2	2	1	1	2	2	2	B	3	30%
Marker C	1	3	3	1	1	2	1	1	1	2	1	C	4	40%
Marker D	1	1	1	2	1	1	1	1	1	1	1	D	1	10%
Marker E	2	1	1	1	1	1	1	1	1	1	1	E	10	100%

