

Thousands of trait-specific KASP markers designed for diverse breeding applications in rice (Oryza sativa)

Steele, Katherine; Quinton-Tulloch, Mark; Vyas, Darshna; Witcombe, John

G3: Genes, Genomes, Genetics

DOI: 10.1093/g3journal/jkae251

E-pub ahead of print: 01/11/2024

Peer reviewed version

Cyswllt i'r cyhoeddiad / Link to publication

Dyfyniad o'r fersiwn a gyhoeddwyd / Citation for published version (APA): Steele, K., Quinton-Tulloch, M., Vyas, D., & Witcombe, J. (2024). Thousands of trait-specific KASP markers designed for diverse breeding applications in rice (Oryza sativa). G3: Genes, Genomes, Genetics. Advance online publication. https://doi.org/10.1093/g3journal/jkae251

Hawliau Cyffredinol / General rights Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

• Users may download and print one copy of any publication from the public portal for the purpose of private study or research.

- · You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal ?

Take down policy If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

1	Research article: Thousands of trait-specific KASP markers designed for diverse breeding
2	applications in rice (Oryza sativa)
3	
4	Authors: Katherine Steele ¹ , Mark Quinton-Tulloch ^{1,2} , Darshna Vyas ³ , John Witcombe ¹
5	
6	1 School of Natural Sciences, Bangor University, Bangor, Gwynedd, LL57 2UW, UK
7	2 Current address: EMBL-EBI, Wellcome Genome Campus, Hinxton, Cambridgeshire, CB10
8	1SD, UK
9	3 LGC BioSearch Technologies, Units 1 & 2, Trident Industrial Estate, Pindar Road, Hoddesdon,
10	Herts, EN11 0WZ, UK
11	
12	Corresponding Author: Katherine A. Steele, email: k.a.steele@bangor.ac.uk tel: +44(0)1248
13	388655, ORCID 0000-0003-4896-8857
14	
15	Abstract
16	This study aimed to broaden applicability of KASP for Oryza sativa across diverse genotypes
17	through incorporation of ambiguous (degenerate) bases into their primer designs and to validate
18	4000 of them for genotyping applications. A bioinformatics pipeline was used to compare 129
19	rice genomes from 89 countries with the <i>indica</i> reference genome R498 and generate ~1.6
20	million KASP designs for the more common variants between R498 and the other genomes. Of
21	the designs, 98,238 were for predicted functional markers. Up to five KASP each for 1024
22	breeder-selected loci were assayed in a panel of 178 diverse rice varieties, generating 3366
23	validated KASP. The 84% success rate was within the normal range for KASP demonstrating
24	that the ambiguous bases do not compromise efficacy. The 3366-trait-specific marker panel was
25	applied for population structure analysis in the diversity panel and resolved them into four
26	expected groups. Target variations in thirteen of the genome sequences used for designs were
27	compared with the corresponding KASP genotypes of other accessions of the same thirteen
28	varieties in the diversity panel. There was agreement across 12 varieties for 79% of markers. Ten
29	varieties had high agreement (>88%) but a variety selected from a landrace had only 46.5%
30	agreement. Breeders can now search for the validated KASP and >1 million so-far untested
	© The Author(s) 2024. Published by Oxford University Press on behalf of The Genetics Society of America. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

1 designs across three alternative reference genomes (including Niponbare MSU7), search for

2 designs proximal to previously published SSR markers and retrieve the target variations in 129

3 rice genomes plus their genomic locations with +/-25 bp flanking sequences.

4

5 Keywords

6 Genetic resources utilisation, InDel, SNP, trait selection

7

8 Introduction

9 Public and private sector rice breeders require efficient markers for selective breeding to enable
10 global rice production to sustainably increase. In many rice-dependent regions the benefits of
11 genomic markers for rice improvement are still to be fully explored (Chakraborti et al., 2021),
12 yet it has been demonstrated that genomics-derived molecular markers can be effectively
13 integrated into traditional rice breeding programmes (Cobb et al. 2019).

14

15 Rice genome resequencing and bioinformatics have previously been used to identify large 16 numbers of useful genomic DNA variants - single nucleotide polymorphism (SNP) and 17 insertion/deletion (InDel) - that can be of use to breeders (Pariasca-Tanaka, 2015; Cheon et al., 18 2018; Sandhu et al., 2022). The bioinformatics skills necessary to identify suitable assays 19 represent a high technology barrier for some breeders. Searchable databases for genomic variants 20 exist for all major cereals (Thudi et al., 2021). Many rice researchers and breeders use IRRI's 21 Rice SNP-Seek Database (snp-seek.irri.org; Mansueto, et al., 2017) and the Chinese Rice 22 VarMap (ricevarmap.ncpgr.cn; Zhao et al., 2015). Despite this wide availability of variants 23 associated with genes and QTLs (quantitative trait loci), there has been limited uptake of 24 genomic breeding tools by public sector and small-scale rice breeders. Many types of marker 25 technologies have been developed for SNP genotyping, but not all can be adopted readily in 26 existing laboratories. Some marker technologies are less transferable to marker-assisted selection 27 applications than others. Some SNP panels are population specific (Heslot et al., 2013) and 28 markers targeting suitable variants might not be readily identifiable for selection of traits in 29 specific crosses (Makhoul et al., 2020). Such issues hinder adoption of new marker technologies, 30 hence, microsatellite (SSR) markers developed in the 1990s remain popular among rice breeders,

largely due to the readily searchable information in the Gramene markers database
 (archive.gramene.org/markers/; Liang et al., 2008).

3

KASP is PCR-based genotyping technology ideally suited for small- or large-scale genotyping
applications. However, it is not always easy for breeders to locate useful information about
suitable KASP markers for their uses. This is partly because KASP is a patented technology of
LGC BioSearch Technologies (LGC) and primer sequences constitute intellectual property (IP).
This study resolves this limitation by using the KASP genomic locations and +/-25 bp flanking
sequences so that breeders have sufficient information to either order them from LGC or use the
location and sequence information to design their own primers.

11

There are considerable benefits to be gained in moving a marker-assisted breeding programme from SSRs to a KASP-based approach (Steele et al., 2018; Kim et al., 2021). High-throughput KASP offer greater cost-effectiveness than SSRs but have similar levels of flexibility and can be used for population studies (e.g. Shikari et al., 2020), linkage mapping (e.g. Qureshi et al., 2018) and MAS (Kim et al., 2021).

17

18 There are 2055 KASP in LGC's original Rice assay search tool, developed by Generation 19 Challenge Programme (Pariasca-Tanaka et al, 2015). Separately, Lee et al., (2022) developed 20 2565 KASP from the C7AIR SNP array. KASP are increasing in popularity for quantitative trait 21 locus (QTL) analysis and marker-assisted selection (MAS) in a range of agriculturally important 22 species (Cheon et al., 2018; Kaur et al, 2020; Paudel et al 2019; Van Inghelandt et al., 2019; 23 Kante et al., 2018; Zhang et al., 2020; Devran et al., 2019; Zhao et al., 2021). A valid concern for 24 OTL mapping with KASP assays is that they may miss many rare variations or common alleles 25 absent from the samples used to develop the assays (Scott et al., 2020). This can be overcome by 26 sequencing the parents used by breeders for their crosses under study for *de-novo* marker 27 development. However, this step is not practical for many rice breeders, so the goal of this study 28 was to develop off-the-shelf SNP and InDel markers that should be representative across Oryza 29 sativa.

KASP technology was selected for this study following a successful feasibly study with breeders 1 2 from Nepal who incorporated KASP-derived genotyping data into their breeding programmes. 3 The feasibility study only sampled nine genomes to identify variants and design KASP primers 4 (Steele et al., 2018). This study used a comparison of 129 genomes to identify suitable target 5 variants (SNP or InDel) and also identified SNP variation occurring in the flanking regions of 6 each target variant so that this information could inform the incorporation of degenerate bases in 7 primers. The addition of degenerate bases was predicted to extend the efficacy of the resultant 8 KASP assays and thereby broaden their applicability in different varieties. In this study we 9 selected 4000 KASP assays of potential value for precision trait selection and applied them in 10 genotyping an independently obtained diverse rice population. The study aimed to (i) determine 11 if the number or location of ambiguous bases differed between successful and failing designs, (ii) 12 demonstrate the utility of a ~4K KASP panel for resolving population structure (iii) provide a 13 database of information for all the new KASP designs including their proximity to existing SSRs 14 and C6IAR SNPs that can help breeders select them for different applications.

- 15
- 16 **Materials and Methods**
- 17

18 Rice genome data

19

20 The KASP design and bioinformatics filtering steps were done at Bangor University (BU). In 21 total 78 indica genomes and 51 non-indica genomes were used for in-silico KASP design (Figure 22 1). This project incorporated variation from the sequencing data from 118 rice genomes selected 23 and retrieved from the 3,000 Genomes Project (3K RGP, 2014) alongside the paired-end 24 sequencing reads for 11 varieties selected by BU's project partners (nine *indica* rice genomes 25 selected by breeders in Nepal (Steele et al., 2018) and two Indian upland varieties (Kalinga III 26 and Ashoka 200F). File S1 contains the methods used for sequencing these two previously 27 unpublished genomes. The 118 genomes included at least one line from each of the 89 countries 28 of origin in the 3K project and all seven rice varietal groups represented in the 3K RGP dataset. 29 All genome sequences are available in the EBI Sequence Read Archive, accession numbers PRJNA395505 (for Bangor University genomes) and PRJEB6180 (for 3K RGP genomes). 30 31 Tables A and B in File S2 provides further details for all genomes used in this study.

2 Sequence read processing, alignment and variant calling

3

4 Quality trimming of sequencing reads was carried out using Sickle (Joshi and Fass, 2011). An

5 average Phred score of 30 was set as the threshold for the trimming window, with sequences

6 truncated at the position of the first N. Trimmed reads shorter than 20 bp were discarded.

7

8 Trimmed reads were aligned against the Shuhui 498 (R498) *indica* rice reference genome (Du *et*9 *al.*, 2017). Version 2 of the genome sequence was downloaded from MBKbase
10 (<u>http://mbkbase.org/R498/</u>) and this version was used for all subsequent R498 alignments and
11 positions generated in this study. Sequence read alignment was carried out with Bowtie2

(Langmead and Salzberg, 2012). Alignments were only reported if both mates of a read pair
aligned in the expected orientation. A single best alignment was reported for each pair, selected
at random in the case of equally good alternative alignments. All other alignment, scoring and

15 reporting options for Bowtie2 were kept as default.

16

17 Genotype likelihood calculation and variant calling was carried out with SAMtools (Li *et al.*,

18 2009). SNPs or insertions with a read depth of less than five were filtered out. Base coverage

19 was calculated using BEDTools (Quinlan, 2014).

20

21 Bioinformatics filtering for KASP marker design generation

22

23 KASP marker design was carried out utilising the data from 129 rice genomes using custom Perl scripts through a process of sequential filtering. For each SNP or InDel variation identified in 24 25 one of the resequenced rice lines, the following tests were carried out to determine whether the 26 variation was suitable for KASP marker design. Here the term 'target variation' is used to 27 describe an allele detected in a particular rice line that is different to the allele in the R498 28 reference genome and under consideration for KASP marker design, and 'alternative variation' 29 for any other alternative alleles (up to two are possible) at the same genomic site detected in the other rice lines. 30

1. **Removal of rare alleles:** If an alternative variation (ie. alleles other than the R498 allele 1 2 and the most common target variant at that position) was identified in more than 10% of 3 lines, then the target variation was removed. That is, at least 90% of lines had to possess 4 either the reference allele or the most common alternative allele. In the case of multiple 5 variations fulfilling this criterium for any site, only the most common variant was carried 6 forward. Any target variations not demonstrating polymorphism among the resequenced 7 lines, were also discarded (i.e. where all 129 test lines shared an allele that was different 8 from the reference allele).

9

Removal of targets with low base coverage: In the case of the target variation being a
 SNP or In/Del the base coverage was checked at the target site. The resequenced
 genomes utilised variable read depths, and a particular genomic site was considered to
 have low base coverage if the read depth was less than one tenth of the average read
 depth for the genome in question. If the target variation site was identified as having low
 base coverage in more than 10% of the 129 resequenced genomes it was not included.

16

3. Removal of targets with low base coverage in flanking sequences: Tests were then
carried out on the 50 bp either side of the target variation, described here as the 'flanking
sequence' with each base referred to as a 'flanking site'. The same base coverage check
described in check 2 above for each target site at a SNP or insertion was made for each
base position in the flanking sequence, with the target variation being rejected in the case
of a single failure at any base position along the flanking sequence.

23

24 4. **Removal of targets with high variation in flanking sequences:** For each base position 25 in the flanking sequence, all 129 genomes were checked for the presence of variations. If 26 alternative variations were present at a flanking site then the target variation was rejected 27 unless: (a) all the alternative variations in the flanking sequence were insertions or 28 deletions of equal length, there was only one insertion or deletion, no insertion or deletion 29 was within 5 bp of the target variation, and no more than 10 bases were inserted or 30 deleted; (b) all the alternative variations were at a single base position, i.e. were SNPs, 31 and no more than five flanking sites were SNPs.

A KASP design sequence consists of the target variation with the 50 bp flanking either side of it.
Preliminary KASP designs were generated for all target variations that had passed the filtering
tests 1-4 (above), with degenerate nucleotides included for flanking sequence variants identified
among the 129 genomes. SNPs in the flanking sequence were represented using the appropriate
International Union of Pure and Applied Chemistry (IUPAC) nucleotide code, and insertions and
deletions were represented by sequences of Ns (e.g. NNN for a 3 base deletion or insertion).
The KASP design sequences were checked for the presence of repeats by first removing any Ns

from the design sequence and then creating a set of test sequences that represented all the
possible combinations of SNPs in the design sequence. If a tandem repeat consisting of more
than five copies of any one to five nucleotide pattern was detected in any member of the test set,
the target variation was excluded.

14

15 To enable end-users to cross-reference marker positions between both *indica* and *japonica* 16 reference genomes, potential KASP design sequences were aligned against the *indica* rice R498 17 (version 2, Du et al., 2017) reference genome and the japonica rice Nipponbare reference 18 genome (version IRGSP-1.0, International Rice Genome Sequencing Project, 2005; downloaded 19 from https://plants.ensembl.org this version was used for all subsequent Nipponbare alignments 20 and positions generated in this study) reference genomes using BLAST (Camacho et al., 2009). 21 Only design sequences that had a single best alignment in both reference genomes were kept. 22 23 A final check for inclusion was for the GC content within a 55 bp window of the KASP design 24 sequence to be between 35-65%, any that did not meet this criterion were excluded to optimise 25 their reliability in PCR. 26 27 Cross referencing with the historic *indica* reference genome

- 28
- 29 The KASP marker design sequences derived via the above filtering steps were aligned against
- 30 the older *indica* rice 93-11 reference genome assembly (ASM465v1, Yu et al., 2002: this version
- 31 was used for all subsequent 93-11 alignments and positions generated in this study) using

BLAST (Camacho et al., 2009). Any best-hit sequences with less than a 90% identity over the 1 2 full sequence length, or those having multiple best hits, were given an unknown position in the 3 93-11 reference genome. This was done to enable subsequent annotation of gene features, 4 C6IAR SNPs and SSRs to include data from two indica reference genomes. 5 6 Annotating target variation with predicted function 7 8 The gene feature annotations for the same three genomes were downloaded in GFF format: The 9 R498 genome (version2, Du et al., 2017), Nipponbare (November 2018 release of RAP-DB 10 (Sakai et al., 2013)), and 93-11 (2010 release of the BGI RISe Rice Information System (Zhao et 11 al., 2004)). These files were processed with custom Perl scripts to categorise each target 12 variation with respect to location within gene features and record information about predicted 13 effect (e.g. functional/non-functional) for each reference genome. 14 15 Target variations (SNPs or InDels) were classified as either being intergenic, or genic. Genic 16 variations located within protein coding genes were further classified according to their location 17 in the 5[']/3['] UTR regions, introns, or coding sequences. SNPs within coding regions were further 18 categorised according to their predicted effect on the translated amino acid sequence: 19 synonymous, non-synonymous, premature stop codon, stop codon loss. Insertions or deletions 20 overlapping coding regions were classified as either frameshift or non-frameshift and start/stop 21 codon loss mutations were identified. The effect on all isoforms was predicted for any variants 22 located within coding genes with multiple annotated transcript isoforms. 23 24 Determination of C6IAR SNP genomic positions 25 26 This was done so that database users can cross reference KASP designs with SNPs in the Cornell 27 6K Infinium rice array (C6IAR) (Thomson et al., 2017). C6IAR SNPs were aligned in the same 28 three reference genomes using the same criteria described above for annotating target variation 29 with predicted function. Of the 5274 C6IAR SNPs, 94% aligned with a position in at least one

30 reference genome, with 75% aligning to the same chromosome in all three genomes (Table C in

File S2). Of the 4569 (86%) C6IAR SNPs that aligned to R498 autosomes only 2099 (46% of
 these) fulfilled KASP design criteria (Table D in File S2).

- 3
- 4 Determination of SSR markers genomic positions
- 5

6 This was done so that database users can cross-reference between previously published SSR

7 markers and the KASP designs in this study. Forward and reverse pairs of primer sequences for

8 19,475 SSR rice markers were downloaded from <u>www.gramene.org</u> and each primer was aligned

9 using BLAST (Camacho et al., 2009) against the R498 (Du et al., 2017), Nipponbare

10 (International Rice Genome Sequencing Project, 2005), and 93-11 (Yu et al., 2002) reference

11 genome sequences.

12

13 Individual primer alignments were rejected if they had an identity of less than 95% for full sequence length. In the case of multiple best hits for SSR primers, all combinations of primer 14 15 pair alignments were considered. Ninety-eight percent (19,138) of SSR primer pairs used for the 16 analysis fulfilled the criteria for alignment, of which 16,980 (89% of all SSRs considered) 17 aligned with a position in at least one of the three reference genomes (Shuhui 498, 93-11 or 18 Nipponbare). Seventy-three percent of aligned SSRs were positioned on the same chromosome 19 in all three genomes. Then SSRs were given a known position if both left and right primers 20 aligned within 10 kb of each other and when only a single pair of best hit alignments fulfilled 21 these criteria in at least one reference genome (Table E in File S2).

22

23 Selection of 4000 KASP for validation test and population analysis

24

The genomic positions of 1080 breeder-specified target genes or SSR markers previously associated with traits or QTLs were used to identify KASP designs that were situated within 0-19913 bp of a target gene or SSR position (Table F in File S2, where columns B and C, headed 'Marker/gene' and 'Alternative IDs' give the names or codes used in previous publications or databases for target genes or markers). When more than five KASP designs were located within this range, five were selected from them according to predicted functionality, followed by closeness to the target. One-hundred and forty-three targets had fewer than five KASP designs, and for these all targets were included (Table G in File S2). Designs for KASP targeting variations predicted to result in functional mutations were preferentially selected, while designs in very close proximity to others in the set were preferentially removed. This resulted in 5,028 KASP designs selected for their proximity to genes or SSRs commonly targeted in rice breeding programs. This set only included designs that had not been selected for validation in breeding applications being done in the wider project. These designs (the target sequences including +/- 50 bp flanking either side) were submitted to LGC who tested them *in-silico* with their proprietary KrakenTM software for primer design and they rejected 43 designs because they did not pass the criteria for primer production.

removing all that were within 100 bp of another marker with the same predicted genotype for all varieties; (ii) removing any non-functional markers furthest from its SSR target, starting with the targets that have the most markers and continuing until 4000 remained. These 4000 designs (for 1024 loci) were submitted to LGC for synthesis of the corresponding KASP primers for use in genotyping in the rice collection.

17

1

2

3

4

5

6

7

8

9

10

11

There was no selection for C6IAR loci during selection of the 4000 designs and only 20 of these
submitted designs had C6IAR equivalents. There were 3275 KASP designs selected and
submitted in the vicinity of 635 Cornell SSR markers (maximum of five designs per SSR).

22 Development of a diverse rice panel for genotyping

23

24 Diverse rice (O. sativa) genotypes, selected to include a wide range of landraces, modern 25 varieties and advanced breeding lines, were supplied by breeders or researchers from the 26 International Rice Research Institute (IRRI), the National Institute for Biotechnology and 27 Genetic Engineering, Pakistan (NIBGE), the Sheri-e-Kashmir University of Agricultural 28 Sciences and Technology of Kashmir, India (SKUAST), the Nepal Agriculture Research Council 29 (NARC), Anamolbiu PVT, Nepal and the Earlham Institute, UK. Modern varieties or advanced 30 breeding lines (including some for direct seeding in uplands) from Brazil, Bangladesh and 31 Pakistan were sourced from the International Rice Genebank at IRRI. The collection (Table H in

File S2) included samples originating from 16 countries and eight breeding programmes and 1 2 their designations included 132 indica, 22 japonica, 5 boro and 17 basmati (Table I in File S2). 3 Seed samples were grown at Bangor University's Henfaes Research Centre, sown either in May 4 2018 or June 2019. Seeds were sown directly into compost and grown under the glasshouse

- 5 conditions described in Note B in File S1.
- 6

7 Leaf samples for DNA extraction were taken after about 7 weeks growth from a single plant of 8 each line using BioArk plant sampling kits, with 96 sampled in July 2018 and a further 82

9

sampled in August 2019. The 178 rice DNA samples were genotyped with KASP by LGC

10 Biosearch Technologies, Hoddesdon, UK.

11 Genotype data were converted to a numeric matrix (1 = R498 allele; 0 = target variant;

12 heterozygotes were run either coded as the most common allele or as 0.5, and results did not

13 vary) and used for Hierarchical cluster analyses with the FactoMineR and Factoextra libraries in

R (Husson et al., 2020). Distances were calculated using the 'dist' function and the 'euclidean' 14

15 method to give a distance matrix. Clusters were produced from the distance matrix using the

16 method 'average' in the function 'hclust" and plotted using the function 'plot'.

17

18 The Wilcoxon rank-sum test was used to test whether various KASP design properties differed 19 between the designs that produced successful genotyping assays among the 178 rice samples and 20 those that did not. The tests related to the number and location of ambiguous bases representing 21 non-target variations, the number and location of InDel bases (only for KASP designs targeting 22 InDels), and the GC content of the design sequence. The two flanking sequences of KASP 23 designs could include either no ambiguous bases or one or more, up to a maximum of five (according to the filtering step 4 above). The distances in bp between the target SNP or InDel 24 25 and the furthest ambiguous base in either or both flanking sequences were used to test properties relating to the distance to the nth ambiguous base (where n was in the range 0-5). Separate tests 26 27 were done for left flank distance, right flank distance, shortest distance in either flank and 28 longest distance in either flank. Only designs with n or more ambiguous bases in the design the 29 flanking the target sequence were included in the tests relating to the distance between the design 30 target and the *n*th ambiguous base. Only designs with at least *n* ambiguous bases in both flanks 31 were included for tests on properties relating to the longest distance to the *n*th ambiguous base.

1	

2 Database development

4	A 'back-end' database was constructed at BU to act as a repository for the KASP assays
5	designed in this study. It contains the bp location of the target variant for each KASP design in
6	up to three reference genomes (Shuhui498, 93-11 or Nipponbare) and the expected variant
7	genotype at each KASP design for each of the 129 source data genomes. Options were included
8	to enable breeders to search for KASP designs, either within a specified region of the genome or
9	at a specified distance from either a named gene, SSR or a SNP from the C6IAR panel. Each
10	KASP design from this study was assigned a KASP ID number (pKey) for information
11	management. This back-end database was provided to LGC for them to use to update their Rice
12	Assay Search Tool. LGC released a beta version of the search tool which has been tested by the
13	authors and two independent rice breeders. A user's manual was written by BU and LGC (File
14	S3).
15	
16	Results
17	
18	Sequence read alignment to indica reference genome
18 19	Sequence read alignment to <i>indica</i> reference genome
18 19 20	Sequence read alignment to <i>indica</i> reference genome Rates of alignment for 129 rice whole genome sequences with the R498 <i>indica</i> reference genome
18 19 20 21	Sequence read alignment to <i>indica</i> reference genome Rates of alignment for 129 rice whole genome sequences with the R498 <i>indica</i> reference genome (Du <i>et al.</i> , 2017) ranged from 52.4% to 95.4%. Genome coverage was between 81.9% and 97.4%
18 19 20 21 22	Sequence read alignment to <i>indica</i> reference genome Rates of alignment for 129 rice whole genome sequences with the R498 <i>indica</i> reference genome (Du <i>et al.</i> , 2017) ranged from 52.4% to 95.4%. Genome coverage was between 81.9% and 97.4% and average read depth ranged from 6.2 to 90.3 (Table B in File S2). Across all 129 genomes, a
18 19 20 21 22 23	Sequence read alignment to <i>indica</i> reference genome Rates of alignment for 129 rice whole genome sequences with the R498 <i>indica</i> reference genome (Du <i>et al.</i> , 2017) ranged from 52.4% to 95.4%. Genome coverage was between 81.9% and 97.4% and average read depth ranged from 6.2 to 90.3 (Table B in File S2). Across all 129 genomes, a total of 15,140,996 variant sites were identified, with an average of one variant site for every
18 19 20 21 22 23 24	Sequence read alignment to <i>indica</i> reference genome Rates of alignment for 129 rice whole genome sequences with the R498 <i>indica</i> reference genome (Du <i>et al.</i> , 2017) ranged from 52.4% to 95.4%. Genome coverage was between 81.9% and 97.4% and average read depth ranged from 6.2 to 90.3 (Table B in File S2). Across all 129 genomes, a total of 15,140,996 variant sites were identified, with an average of one variant site for every 25.8 bases in the 391 Mb R498 genome. The number of variant sites for each variety against
18 19 20 21 22 23 24 25	Sequence read alignment to <i>indica</i> reference genome Rates of alignment for 129 rice whole genome sequences with the R498 <i>indica</i> reference genome (Du <i>et al.</i> , 2017) ranged from 52.4% to 95.4%. Genome coverage was between 81.9% and 97.4% and average read depth ranged from 6.2 to 90.3 (Table B in File S2). Across all 129 genomes, a total of 15,140,996 variant sites were identified, with an average of one variant site for every 25.8 bases in the 391 Mb R498 genome. The number of variant sites for each variety against R498 ranged from 737,046 to 2,343,000.
18 19 20 21 22 23 24 25 26	Sequence read alignment to <i>indica</i> reference genome Rates of alignment for 129 rice whole genome sequences with the R498 <i>indica</i> reference genome (Du <i>et al.</i> , 2017) ranged from 52.4% to 95.4%. Genome coverage was between 81.9% and 97.4% and average read depth ranged from 6.2 to 90.3 (Table B in File S2). Across all 129 genomes, a total of 15,140,996 variant sites were identified, with an average of one variant site for every 25.8 bases in the 391 Mb R498 genome. The number of variant sites for each variety against R498 ranged from 737,046 to 2,343,000.
18 19 20 21 22 23 24 25 26 27	Sequence read alignment to <i>indica</i> reference genome Rates of alignment for 129 rice whole genome sequences with the R498 <i>indica</i> reference genome (Du <i>et al.</i> , 2017) ranged from 52.4% to 95.4%. Genome coverage was between 81.9% and 97.4% and average read depth ranged from 6.2 to 90.3 (Table B in File S2). Across all 129 genomes, a total of 15,140,996 variant sites were identified, with an average of one variant site for every 25.8 bases in the 391 Mb R498 genome. The number of variant sites for each variety against R498 ranged from 737,046 to 2,343,000.
18 19 20 21 22 23 24 25 26 27 28	Sequence read alignment to <i>indica</i> reference genome Rates of alignment for 129 rice whole genome sequences with the R498 <i>indica</i> reference genome (Du <i>et al.</i> , 2017) ranged from 52.4% to 95.4%. Genome coverage was between 81.9% and 97.4% and average read depth ranged from 6.2 to 90.3 (Table B in File S2). Across all 129 genomes, a total of 15,140,996 variant sites were identified, with an average of one variant site for every 25.8 bases in the 391 Mb R498 genome. The number of variant sites for each variety against R498 ranged from 737,046 to 2,343,000. Novel KASP marker designs
18 19 20 21 22 23 24 25 26 27 28 29	Sequence read alignment to <i>indica</i> reference genome Rates of alignment for 129 rice whole genome sequences with the R498 <i>indica</i> reference genome (Du <i>et al.</i> , 2017) ranged from 52.4% to 95.4%. Genome coverage was between 81.9% and 97.4% and average read depth ranged from 6.2 to 90.3 (Table B in File S2). Across all 129 genomes, a total of 15,140,996 variant sites were identified, with an average of one variant site for every 25.8 bases in the 391 Mb R498 genome. The number of variant sites for each variety against R498 ranged from 737,046 to 2,343,000. Novel KASP marker designs
 18 19 20 21 22 23 24 25 26 27 28 29 30 	Sequence read alignment to <i>indica</i> reference genome Rates of alignment for 129 rice whole genome sequences with the R498 <i>indica</i> reference genome (Du <i>et al.</i> , 2017) ranged from 52.4% to 95.4%. Genome coverage was between 81.9% and 97.4% and average read depth ranged from 6.2 to 90.3 (Table B in File S2). Across all 129 genomes, a total of 15,140,996 variant sites were identified, with an average of one variant site for every 25.8 bases in the 391 Mb R498 genome. The number of variant sites for each variety against R498 ranged from 737,046 to 2,343,000. Novel KASP marker designs Over 1.6 million KASP marker designs were generated <i>in silico</i> by this bioinformatics study. Based on their positions in relation to the gene models of the <i>japonica</i> Nipponbare and <i>indica</i>

1 variations that would cause a change in the expressed proteins (Table 1). The maximum distance

2 between adjacent KASP marker designs was 573 kb for the *indica* reference genome and 270 kb

3 for *japonica*, with the median distance between designs for both being 55 bp (Figure 2). The

4 number of KASP targets predicted to be polymorphic between all 8,256 pairs of possible varietal

5 comparisons ranged between 23,078 (minimum) and 581,415 (maximum), with a median of

- 6 328,256 (Figure 3).
- 7

8 Demonstration of utility of novel KASP panel through genotyping

9

10 Of the 4,000 trait-specific KASP designs tested, 3,371 were passed for wet-lab validation 11 according to the service provider. But more stringent data analysis revealed that 3366 KASP 12 gave successful genotype calls in >90% of samples which resulted in successful KASP (Tables I 13 and J in File S2), hence 3366 were considered as validated. Eleven markers were monomorphic in this set of rice germplasm so data for the remaining 3355 markers were used in cluster 14 15 analysis to reveal separation into four major groups corresponding to indica, japonica, 16 intermediate and aromatic sub-types. This grouping well-reflects the diverse population tested 17 which includes diverse landraces as well as breeding lines and modern varieties, many of which 18 are derived from crosses between sub-types (Figure 4). 19 20 Some 635 designs (15.9%) failed in all samples (Figure Aa in File S1) and for designs with 21 genotype calls below 90%, the lowest call rate was 17.4%. For the successful 3,366 markers, 22 23.5% produced calls in all samples while 81.9% produced calls in >90% of samples (Figure Ab 23 in File S1). The percentage successful allele calls per variety ranged from 72.5% to 83.8%. There 24 was no significant difference in success rate of KASP between different rice sub-groups or 25 countries of origin. 26 27 Comparison between sequenced and genotyped datasets 28 29 Thirteen of the genotyped rice samples had the same names as thirteen of the 129 sequenced

30 genomes used in the marker design process. For each of these named lines the genotype calls for

31 each marker were compared against the sequenced genotypes (Table L in File S2). There were

and "?". In the subset of 13 varieties there were three Bad calls (one each in Anmol Masuli, Lok Tantra and Kalinga III at different loci); two Uncallable (in Loc Tantra and Kalinga III at different loci); and 501 calls of ?, ranging from 20 to 63 alleles called as ? per line.
The percentage matching alleles for line ranged from 97% to as low as 46.5%. However, >77% of markers had calls matching predicted genotypes in all but one line (Chommrong) and 10 of the 13 lines had >88% agreement between the genotype data and sequence data.
There was complete agreement for 1043 (31%) markers and non-agreement for the remaining 69%. Of these, 1610 (48%) did not match in only one line. The number of non-agreements reduced rapidly for additional lines: 14% in two lines for, 5% in three lines and 0.9% in four lines. Only nineteen lines had non-matching calls in five or more lines. Only one marker (R498 locus 10:20882568) had no matches in all 13 lines, but all were called as heterozygotes (possibly an artifact, see discussion).
Success rates for the 13 lines used in design generation were compared against the 165 that were not. They had median success rates of 83.33% and 83.05% respectively and a Wilcoxon rank-sum test showed no significant difference in the distributions of success rates at the 99% confidence level, with a p-value of 0.035.
Effect of ambiguous bases in designs on success rate
Examining the subsets of the 4,000 KASP designs submitted for genotyping showed that KASP designs with a mean of 2.5 ambiguous bases were significantly more likely to fail than those with a mean of 1.86 ambiguous bases (Table 2). Several properties of the distance of ambiguous bases from the target variation also showed significant differences in distribution between the successful and failed markers (Table 2). KASP designs with higher GC had significantly more failures. Wilcoxon rank-sum test results for all 29 properties are shown in Table 2.
14

3389 non-matching calls (7.7%) out of 43,758 datapoints in this subset. Included in non-

matching calls are failed calls which were reported in the dataset as either "Bad", "Uncallable"

1	The overall success rates are likely improved by reducing the cut-off for the maximum number
2	of ambiguous bases permitted in a KASP design. However, this comes at a large cost in the
3	reduced number of potential designs and hence an increased distance between adjacent markers
4	(Table 3).
5	
6	Proximity of validated KASP to widely used Cornell markers
7	
8	All 20 of the KASP designs with a Cornell C6IAR equivalent included in the 4000 designs tested
9	were successfully validated (Table C in File S2). The success rate of KASP designs located near
10	to Cornell's SSRs was 81%, with 620 unvalidated KASP designs in the vicinity of SSR markers
11	listed on the Gramene database. Breeders can use Table K in File S2 to identify selected designs
12	for trait selection that are close to previously published SSR loci.
13	
14	Rice Assay Search Tool
15	
16	Breeders and other end-users can access the database containing details of the ~ 1.6 million
17	KASP assay designs developed in this study through the Rice Assay Search Tool
18	(www.biosearchtech.com/kasp-assay-search) (Further details and search tips are provided in
19	Table M in File S2 and File S3).
20	
21	
22	
23	Discussion
24	Previous studies have also mined genomic variations within the rice 3K RGP (2014) data to
25	identify a large number of target SNPs for use by rice researchers (Alexandrov et al., 2015;
26	Tareke Woldegiorgis et al., 2019). Others have identified KASP for specific traits in rice (e.g.
27	Addison et al. 2020; Angira et al., 2021; Sandhu et al., 2022). To our knowledge no large-scale
28	previous study has designed KASP primers which include ambiguous flanking variation, or
29	specifically selected thousands of KASP for loci that are relevant to breeding programmes.

This study used a bioinformatics pipeline to filter ~15 billion potential target variants detected among 129 publicly accessible rice genomes and remove those within problematic regions as well as any with unsuitable non-target variations flanking the target variation. It used the R498 *indica* reference genome as the baseline for KASP targets, meaning that every KASP we designed assayed for the R498 allele and the most common alternate allele at that target among a sampled population of 129 resequenced genomes. Targets having multiple alternate alleles where they together accounted for >10% of that population were not developed as KASP in this study.

9 This pipeline resulted in ~1.6 million KASP assay designs optimised to include IUPAC
10 nucleotide codes at a maximum of five non-target variations in each region flanking the target
11 SNP or InDel. The number of KASP assays generated compares well with other rice KASP
12 development projects (Cheon et al., 2018). Similar numbers of KASP designs were present in
13 both *indica* and *japonica* genomes although *indica* had more potentially functional markers
14 (Table 1) and slightly larger gaps between markers (Table 3).

15

16 The frequency of polymorphic sites showed a bi-modal pattern when plotted as a histogram for 17 pairwise comparisons between genotypes (Figure 3). A similar bimodal pattern of 18 polymorphisms observed by Alexandrov et al. (2015) was considered to indicate the absence of a 19 proportion of mapped reads in some genomes. There were differences in coverage between the 20 genomes used in this study. However, it was observed that the number of polymorphisms in pairs 21 was associated with how closely related each pair are to each other, with pairs in the peak on the 22 left side of the histogram made up of varieties from the same Oryza sub-group group while those 23 on the right are made up of two varieties from different groups (aus, boro, *indica, japonica* etc.). 24 Pairs with intermediate values are the product of lower polymorphic pairs between groups or 25 higher polymorphic pairs within groups.

26 Limitations of the design pipeline

27

The KASP design algorithm used in this study resulted in fewer designs generated (10.6% of the ~15 million variant sites identified in the sequenced lines) than the 51.9% variation site to KASP design conversion rate reported in our previous study (Steele *et al.*, 2018). This was expected because of the wider range of varieties used in the current study giving many more variant sites in the flanking sequences of target SNPs and InDels. The avoidance of potential markers with
excessive variation in flanking sequences has reduced the overall number of successful KASP
designs for targets that are predicted to be polymorphic at target sites. In practical terms,
particular combinations of parental lines may have large genomic regions lacking detectable
polymorphism using the available assay designs.

6

When there is polymorphism in regions harbouring specific traits, these markers offer precision
and effectiveness for trait selection. However, if breeders' populations do not show
polymorphism with any of these markers in specific regions (e.g. for fine mapping) they can
consider *de-novo* cross-specific marker design generation (without ambiguous bases) which can
be carried out using the KASP design software code provided by Steele et al. (2018).

12

13 Factors affecting success of KASP assays

14

Of the assay designs submitted, 99.2% passed the final *in-silico* step for primer design. Of the subset of 4000 KASP designs developed into 'wet lab' assays, 84% were successfully amplified with alleles called. This rate was only slightly lower than was obtained for KASP designed from only nine genomes without the inclusion of ambiguous bases (Steele et al., 2018) and for KASP designed from previously published rice SNPs by Yang et al (2019). It is higher than the success rate of 71% validated KASP converted from SNPs derived from RNA-seq in maize (Jagtap et al., 2020).

22

23 Of the 3366 validated KASP, 82% gave allele calls in a panel of 178 diverse varieties, providing 24 genotypes for >90% of the panel. No significant differences were observed in the genotyping 25 success rates of 116 not-resequenced lines as compared to 13 genotyped resequenced lines that 26 were used in the assay design. The rates are within the range of other KASP desing studies in 27 rice (e.g. 70% reported by Gouda et al., 2021). Overall, this result indicates that these KASP 28 assays have a high probability of working in a wide range of rice populations and should be 29 considered widely applicable for breeding. The following discussion considers some of the 30 reasons that could lead to failure for genotyping.

In the thirteen varieties used for both KASP design generation and genotyping, the vast majority 1 2 of genotype calls matched those predicted in the design stage, with the percentage of matching 3 genotypes being within the bounds of normal within-variety variation. The notable exception 4 was Chhomrong for which only 46% of genotype calls matched expectations. Clearly there were genetic differences between the two seed lots of Chhomrong used for genotyping in this study 5 6 and for resequencing the 3K RGP (2014). Chhomrong was originally considered a landrace and 7 subsequent selection and purification produced the released variety with the same name (Joshi et 8 al., 2017), which was the source of the sample used here for genotyping. Chhomrong did not 9 have more heterozygote calls than other lines, however the line FL_478 had only 74% agreement 10 with the sequenced version and the disagreement was exacerbated by numerous heterozygous 11 calls in the genotyped sample (Table L in File S2).

12

Nearly 15.8% of the 4000 marker designs submitted for genotyping failed to result in any
genotype calls in this study. Failures might be explained by the extracted DNA quality, the assay
conditions in a particular genotyping run or they could potentially be due to issues related to
using ambiguous bases in the designs.

17

18 From the KASP assay designs submitted for validation, it was possible to infer the aspects most 19 likely to influence success rates. A statistical comparison of various design-related properties 20 suggested that as the number of ambiguous bases increased to accommodate non-target 21 variations, the rate of success decreased (Table 3). The position of ambiguous bases within the 22 design sequence also affects the chances of success, with a greater distance between target 23 variations and ambiguous bases resulting in a higher proportion of successful designs. The 24 distance between ambiguous and target variation did not show significant differences between 25 successful and unsuccessful marker designs, although there were relatively few designs with a high number of ambiguous bases so statistical power was reduced (Table 2). Reducing the 26 27 number of permitted ambiguous bases in the KASP designs led to a relatively small increase in 28 assay success rate (Table 3), but the predicted number of potential designs is reduced 29 substantially, and there are larger gaps in genome coverage (Table 3). If no ambiguous bases are 30 permitted in the flanking sequences, then the predicted success rate increases by only $\sim 5\%$, but

the median distance between markers is increased by more than 10 times, and the number of
potential designs is reduced by 86% (Figure B in File S1).

3

The GC content was also linked to design success rate (Table 2). If the resultant primers are
leading to failed assays due to non-optimal assay conditions for the primer GC content, then it is
possible that adjusting assay conditions could result in successful genotyping.

7

No significant differences were observed between the distributions of the number of inserted or
deleted bases of passing and failing designs, nor in their distance from the target variation (Table
3). This may be because only a single insertion or deletion was permitted in the design algorithm
and, thus, any InDels could be avoided in primer design.

12

13 Although genotype calls with question marks (?) were rare in the diverse population (0.011%), it is noteworthy that they often occurred in several different genotypes at the same locus, 14 15 suggesting they have a biological cause rather than being an artifact. For example, one marker 16 (locus R498 position 2:31326133) had 47 ? calls, 55 homozygous A calls (the alternate allele 17 used in the KASP design) and 76 homozygous T calls (the R498 target allele) in the diverse 18 population (Table J in File S2). In contrast, the same marker had no calls for the alternate allele 19 (A) in a different population largely composed of commercial aromatic rice varieties genotyped 20 by Steele et al. (2020). In that study all genotype calls were either ? or T homozygote. Our 21 working hypothesis is that the ? calls in both studies may denote presence of an alternate ('third') 22 allele that was not included in the designs (either in a homozygote or heterozygote state) in the 23 accessions with ? calls. Further work is needed to test this hypothesis either using sequencing of 24 accessions carrying the ? allele or by querying the genome assemblies of such accessions. An 25 alternative approach could be to produce alternative KASP assays at these loci with primers 26 selected to call the rarer alternate allele instead of the most common one, which was our default 27 strategy during KASP design in this study. Anecdotally, breeders in Nepal, Anmolbiu PVT and 28 NARC, who used such KASP markers for selective breeding have found that the calls for ? 29 segregate as expected in some populations, and often can show an identical pattern of 30 segregation to adjacent, tightly linked, markers with clarity in calls, indicating that data from 31 such markers can be used, in some circumstances and with caution, to inform selection decisions.

Downloaded from https://academic.oup.com/g3journal/advance-article/doi/10.1093/g3journal/jkae251/7863403 by Prifysgol Bangor University user on 04 November 2024

1

2 Efficacy and value for breeding applications

3

The novel panel of 3366 'ambiguous base', trait-specific KASP developed in this study were validated in a panel of diverse rice, demonstrating the efficacy of such designs for genotyping use across a wide range of potential varieties. By sampling multiple KASP for 1024 target loci, we hope to have widened the range of assays available so that breeders can select the ones that are most useful in their crosses.

9

10 Many of the markers located near to specific targets or genes are relevant to multiple breeding 11 programmes. Some of the validated trait-specific markers have been functionally confirmed by 12 breeders to be linked to traits including disease resistance genes Xa5, Piz-t, Pi33 (Arif M., 13 NIBGE Pers. comm.) and QTL for bakanae foot rot resistance (Shikari A., SKUAST, Pers. comm.). Nepalese breeders at NARC and Anamolbiu PVT have used KASP designs from this 14 15 study in marker-assisted backcrossing to successfully incorporate blast and bacterial leaf blight 16 resistance genes in Khumal-4, Sunaulo Sugandha, Sugandha-1, and Anmol Mansuli that are 17 being tested for potential release in Nepal.

18

19 The applied validation of the ~4K KASP panel was demonstrated through their ability to resolve 20 groups in hierarchical cluster analysis (Figure 4). It is noteworthy that members of the 21 intermediate group derived through this analysis were two Vietnamese varieties (Khara Ganga 22 and OM 479 expected to be *indica*) and the approved Basmati variety Pusa Basmati 1. Two 23 varieties originally thought to be japonica (SKAU_D40 and SKAU_D54) were confirmed as 24 indica in this analysis, supporting a similar finding by Shikari et al. (2020) with a different sub-25 set of our KASP designs. Those authors successfully used 114 (of 213 genotyped loci) of our 26 marker designs for genotyping and structural analysis of a 470 line population of Himalayan-27 grown rice. At the same time (in the same LGC project), this sub-set of markers were also 28 genotyped on the USDA minicore collection and Pakistan landraces and also applied 29 successfully for population structure analysis (M. Arif, NIBGE, Pakistan, Pers. comm), 30 supporting the value of the wider set of KASP designs for this application, and also highlighting their potential high-throughput scalability. With co-ordinated teamwork and careful management 31

of resources, a large sub-set of KASP can be used efficiently for screening large populations
 including multiple sets of material from different groups to increase efficiency.

3

In contrast to SSR markers, which can detect multiple alleles at a single locus, KASP only detect
a maximum of two alternate alleles (SNPs or InDels) at each target locus (although in some cases
a repeatable null allele can be identified (with ? calls) that follows expected Mendelian patterns
of inheritance, and thus inferred as a 'third' allele, discussed above). For genetic diversity
studies, estimates suggest that 7-11 times more KASP markers are needed to reveal a similar
amount of diversity (in the form of haplotypes) compared to a single SSR (Hamblin et al. 2007;
Van Inghelandt et al., 2010).

11

12 For breeding applications, Ashfaq et al. (2023) used a sub-set of KASP derived from this study 13 and found that a similar numbers of foreground or background marker loci are required for KASP as compared to SSRs when applied for QTL mapping and haplotype discovery, so long as 14 KASP markers known to be polymorphic in the population were used. The number of assays can 15 16 be scaled up for high-throughput applications such as genomic selection or down for marker-17 assisted backcrossing and panels including more as-yet unvalidated KASP can be selected from 18 the online rice assay search tool. The rice assay search tool links each marker to genome 19 annotation information and contains information about predicted gene functionality as well as 20 alleles in resequenced genomes. The resources in the search tool could be used by researchers to 21 integrate these KASP with other Omics data.

- 22
- 23

24 Data Availability

25 The genome sequences used for KASP design development are publicly available via the EBI

- 26 Sequence Read Archive, accession numbers PRJNA395505 for Bangor University genomes
- 27 (www.ebi.ac.uk/ena/browser/view/PRJNA395505) and PRJEB6180 for 3K RGP genomes (
- 28 www.ebi.ac.uk/ena/browser/view/PRJEB6180). Genomic locations of all validated KASP are
- 29 available in supplemental files. Genomic locations of KASP target variants for all ~1.6 M KASP
- 30 designs generated during this study are available via the BU-LGC_plus rice assay search tool:
- 31 www.biosearchtech.com/kasp-assay-search

2 Acknowledgements

3 The authors are grateful to: Innovate UK Agi-Tech Catalyst for funding this study (Project 4 Number 103711); Asif Shikari, SKUAST (India), Mohammed Arif, NIBGE (Pakistan), Resham 5 Amgai, NARC (Nepal), Jose DeVega, Earlham Institute (UK), and IRRI (Philippines) for 6 providing rice samples; Bangor University's Henfaes Research Centre staff Llinos Hughes and 7 Mark Hughes for glasshouse operations; The LGC Biosearch Technologies team including Rhian 8 Gwillam, Dominique Fauvin for genotyping project management and Laima Barr, Chris Bond 9 and Ariane Mogha for developing database functionality on LGC's server; Agri EPI Centre for 10 overall project management. 11

12 Funding

This public-private partnership was possible due to funding by Innovate UK Agi-Tech Catalyst
(Project Number 103711) with UK government aid (ODA) funds designated to promote the
economic development and welfare of developing countries.

16

17 Author Declarations

18 Conflict of Interest

19 There is no conflict of interest, no personal gain (financial or otherwise) and no vested interest

- 20 by any of the authors. The project was funded by UK public funds designated for Official
- 21 Development Assistance through an Innovate UK grant where all parties signed a collaboration
- 22 agreement which included maintaining the Intellectual Property Rights of individual partners:
- 23 KASP primers fall under this IP agreement.
- 24

25 *Ethics approval and consent to participate.*

26 This research was screened under Bangor University Research Ethics Framework, no issues were27 identified

28

- 29 Consent for publication
- 30 Not Applicable

1	Authors' information(optional)
2	Not Applicable
3	
4	Author Contributions
5	JW, KS and DV conceived the study. All authors contributed to the study design. Material
6	preparation was performed by KS. MQT conducted the bioinformatics analysis and wrote the
7	first draft. KS revised the manuscript, and all authors approved the final manuscript.
8	
9	References
10	
11	3K RGP (2014) The 3,000 Rice Genomes Project. Gigascience 3:7
12	
13	Addison, C.K., Angira, B., Kongchum, M., Harrell, D.L., Baisakh, N., Linscombe, S.D. and
14	Famoso, A.N., 2020. Characterization of haplotype diversity in the BADH2 aroma gene and
15	development of a KASP SNP assay for predicting aroma in US rice. Rice, 13, pp.1-9.
16	
17	Alexandrov, N., Tai, S., Wang, W., Mansueto, L., Palis, K., Fuentes, R.R., Ulat, V.J.,
18	Chebotarov, D., Zhang, G., Li, Z. and Mauleon, R., 2015. SNP-Seek database of SNPs
19	derived from 3000 rice genomes. Nucleic Acids Research, 43(D1), pp.D1023-D1027.
20	
21	Angira, B., Addison, C.K., Cerioli, T., Rebong, D.B., Wang, D.R., Pumplin, N., Ham, J.H.,
22	Oard, J.H., Linscombe, S.D. and Famoso, A.N., 2019. Haplotype characterization of the sd1
23	Semidwarf gene in United States Rice. The Plant Genome, 12(3), p.190010.
24	
25	Ashfaq H, Rani R, Perveen N, Babar AD, Maqsood U, Asif M, Steele KA, Arif M (2023) KASP
26	mapping of QTLs for yield components using a RIL population in Basmati rice (Oryza
27	sativa L.) Euphytica 219:79
28	
29	Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL (2009)
30	Blast+: architecture and applications. BMC Bioinformatics 10:421.

1	Chakraborti, M., Kumar, A., Verma R.L, R A.F., Reshmi Raj KR, Patra B C, Balakrishnan D,
2	Sarkar S, Mandal N.P., Kar M.K., Meher J., R M Sundaram, LV Subba Rao. (2021) Rice
3	breeding in India: eight decades of journey towards enhancing the genetic gain for yield,
4	nutritional quality, and commodity value. ORYZA-An International Journal of
5	Rice, 58:S69-88
6	
7	Cheon K.S., Baek J., Cho Y.I., Jeong Y.M., Lee Y.Y., Oh J., Won Y.J., Kang D.Y., Oh H., Kim
8	S.L., Choi I. (2018) Single nucleotide polymorphism (SNP) discovery and kompetitive
9	allele-specific PCR (KASP) marker development with Korean japonica rice varieties. Plant
10	Breeding and Biotechnology 6: 391-403.
11	
12	Cobb, J.N., Biswas, P.S. and Platten, J.D., 2019. Back to the future: revisiting MAS as a tool for
13	modern plant breeding. Theoretical and Applied Genetics, 132:647-667.
14	
15	Devran, Z. and Kahveci, E., 2019. Development and validation of a user-friendly KASP marker
16	for the Sw-5 locus in tomato. Australasian Plant Pathology, 48:503-507.
17	
18	Du H, Yu Y, Ma Y, Gao Q, Cao Y, Chen Z, Ma B, Qi M, Li Y, Zhao X, Wang J. (2017)
19	Sequencing and de novo assembly of a near complete indica rice genome. Nature
20	Communications. 4: 1-2.
21	
22	Gouda, A.C., Warburton, M.L., Djedatin, G.L. Kpeki S.B., Wambugu P.W., Gnikoua K., and
23	Ndjiondjop M.N. (2021) Development and validation of diagnostic SNP markers for quality
24	control genotyping in a collection of four rice (Oryza) species. Scientific Reports 11, 18617.
25	
26	Heslot, N., Rutkoski, J., Poland, J., Jannink, J.L. and Sorrells, M.E. (2013) Impact of marker
27	ascertainment bias on genomic selection accuracy and estimates of genetic diversity. PloS
28	one, 8: e74612.
29	

1	Hamblin, M.T., Warburton, M.L. and Buckler, E.S. (2007) Empirical comparison of simple
2	sequence repeats and single nucleotide polymorphisms in assessment of maize diversity and
3	relatedness. PloS one, 2:.e1367.
4	
5	Husson, F., Josse, J., Le, S., Mazet, J. and Husson, M.F., 2020. Package 'FactoMineR'.
6	Multivariate Exploratory Data Analysis and Data Mining. Available on CRAN: https://cran.
7	r-proje ct. org/web/packa ges/Facto MineR/Facto MineR. Pdf
8	International Rice Genome Sequencing Project (2005) The map-based sequence of the rice
9	genome. Nature 436:793-800
10	
11	Jagtap, A. B., Vikal, Y., & Johal, G. S. (2020). Genome-wide development and validation of
12	cost-effective KASP marker assays for genetic dissection of heat stress tolerance in
13	maize. International Journal of Molecular Sciences, 21: 7386.
14	
15	Joshi, B.K.; Bhatta, M.R.; Ghimire, K.H.; Khanal, M.; Gurung, S.B.; Dhakal, R.; Sthapit, B.R.
16	(2017). Released and promising crop varieties of mountain agriculture in Nepal (1959-
17	2016). Pokhara, Nepal: LI-BIRD/Bioversity International 207 p. ISBN: 978-92-9255-060-8
18	
19	Joshi NA, Fass JN (2011) Sickle: A sliding-window, adaptive, quality-based trimming tool for
20	FastQ files (Version 1.33) [Software]. Available at https://github.com/najoshi/sickle
21	
22	Kante, M., Rattunde, H.F.W., Nébié, B., Weltzien, E., Haussmann, B.I. and Leiser, W.L., 2018.
23	QTL mapping and validation of fertility restoration in West African sorghum A 1 cytoplasm
24	and identification of a potential causative mutation for Rf 2. Theoretical and Applied
25	Genetics, 131: 2397-2412.
26	
27	Kaur, B., Mavi, G.S., Gill, M.S. and Saini, D.K. (2020) Utilization of KASP technology for
28	wheat improvement. Cereal Research Communications, pp.1-13.
29	

1	Kim, M.S., Yang, J.Y., Yu, J.K., Lee, Y., Park, Y.J., Kang, K.K. and Cho, Y.G. (2021) Breeding
2	of high cooking and eating quality in rice by Marker-Assisted Backcrossing (MABc) using
3	KASP markers. Plants, 10: 804.
4	
5	Langmead B, Salzberg SL (2012) Fast gapped-read alignment with Bowtie 2. Nature Methods 9:
6	357–359
7	
8	Lee, C., Cheon, K.S., Shin, Y., Oh, H., Jeong, Y.M., Jang, H., Park, Y.C., Kim, K.Y., Cho, H.C.,
9	Won, Y.J. and Baek, J. (2022) Development and application of a target capture sequencing
10	SNP-genotyping platform in rice. Genes, 13: 794.
11	
12	
13	Liang, C., Jaiswal, P., Hebbard, C., Avraham, S., Buckler, E.S., Casstevens, T., Hurwitz, B.,
14	McCouch, S., Ni, J., Pujar, A. and Ravenscroft, D. (2007). Gramene: a growing plant
15	comparative genomics resource. Nucleic Acids Research, 36: D947-D953
16	
17	Makhoul, M., Rambla, C., Voss-Fels, K.P., Hickey, L.T., Snowdon, R.J. and Obermeier, C.,
18	(2020) Overcoming polyploidy pitfalls: A user guide for effective SNP conversion into
19	KASP markers in wheat. Theoretical and Applied Genetics, 133: 2413-2430.
20	
21	Mansueto, L., Fuentes, R.R., Borja, F.N., Detras, J., Abriol-Santos, J.M., Chebotarov, D.,
22	Sanciangco, M., Palis, K., Copetti, D., Poliakov, A. and Dubchak, I., (2017) Rice SNP-seek
23	database update: new SNPs, InDels, and queries. Nucleic Acids Research 45: D1075-D1081
24	
25	Pariasca-Tanaka J, Lorieux M, He C, McCouch S, Thomson MJ, Wissuwa M (2015)
26	Development of a SNP genotyping panel for detecting polymorphisms in Oryza
27	glaberrima/O. sativa interspecific crosses. Euphytica 201:67–78.
28	
29	Paudel, L., Clevenger, J. and McGregor, C., 2019. Refining of the egusi locus in watermelon
30	using KASP assays. Scientia Horticulturae, 257: 108665.
31	

8
ž
n
0
DB
Ð
0
fr
Ĕ
1
Ħ
Ö
a
2
þ
Ð
⊒.
0
2
d,
O
9
/L
Q
3
лc
n
a
0
JQ
Sa
n
e
à
Int
0
Ð
b
Ξ
0
<u> </u>
00
3
g
<u>g3i</u>
g3jou
'g3journ
'g3journal.
'g3journal/jk
ˈg3journal/jkae
g3journal/jkae2
g3journal/jkae251
/g3journal/jkae251/7
g3journal/jkae251/78
g3journal/jkae251/7863
g3journal/jkae251/78634.
/g3journal/jkae251/7863403
/g3journal/jkae251/7863403
g3journal/jkae251/7863403 by
g3journal/jkae251/7863403 by P
g3journal/jkae251/7863403 by Prif
g3journal/jkae251/7863403 by Prifys
ˈg3journal/jkae251/7863403 by Prifysgu
g3journal/jkae251/7863403 by Prifysgol
/g3journal/jkae251/7863403 by Prifysgol B₄
g3journal/jkae251/7863403 by Prifysgol Ban
g3journal/jkae251/7863403 by Prifysgol Bange
g3journal/jkae251/7863403 by Prifysgol Bangor
g3journal/jkae251/7863403 by Prifysgol Bangor Ut
g3journal/jkae251/7863403 by Prifysgol Bangor Univ
g3journal/jkae251/7863403 by Prifysgol Bangor Unive
g3journal/jkae251/7863403 by Prifysgol Bangor Universi
g3journal/jkae251/7863403 by Prifysgol Bangor University
g3journal/jkae251/7863403 by Prifysgol Bangor University u
g3journal/jkae251/7863403 by Prifysgol Bangor University use
g3journal/jkae251/7863403 by Prifysgol Bangor University user u
g3journal/jkae251/7863403 by Prifysgol Bangor University user or
g3journal/jkae251/7863403 by Prifysgol Bangor University user on 0
g3journal/jkae251/7863403 by Prifysgol Bangor University user on 04
g3journal/jkae251/7863403 by Prifysgol Bangor University user on 04 Nc
g3journal/jkae251/7863403 by Prifysgol Bangor University user on 04 Nov
g3journal/jkae251/7863403 by Prifysgol Bangor University user on 04 Nover
g3journal/jkae251/7863403 by Prifysgol Bangor University user on 04 Novemb
g3journal/jkae251/7863403 by Prifysgol Bangor University user on 04 Novembei
g3journal/jkae251/7863403 by Prifysgol Bangor University user on 04 November 2
g3journal/jkae251/7863403 by Prifysgol Bangor University user on 04 November 20;
g3journal/jkae251/7863403 by Prifysgol Bangor University user on 04 November 2024

1	Quinlan AR (2014) BEDTools: the Swiss-army tool for genome feature analysis. Current
2	Protocols in Bioinformatics 47:11.12.1-11.12.34
3	
4	Qureshi, N., Bariana, H.S., Zhang, P., McIntosh, R., Bansal, U.K., Wong, D., Hayden, M.J.,
5	Dubcovsky, J. and Shankar, M., 2018. Genetic relationship of stripe rust resistance genes
6	Yr34 and Yr48 in wheat and identification of linked KASP markers. Plant disease, 102: 413-
7	420.
8	
9	Sakai H, Lee SS, Tanaka T, Numa H, Kim J, Kawahara Y, Wakimoto H, Yang CC, Iwamoto M,
10	Abe T, Yamada Y, Muto A, Inokuchi H, Ikemura T, Matsumoto T, Sasaki T, Itoh T (2013)
11	Rice annotation project database (RAP-DB): an integrative and interactive database for rice
12	genomics. Plant Cell Physiology 54:e6
13	
14	Sandhu, N., Singh, J., Singh, G., Sethi, M., Singh, M.P., Pruthi, G., Raigar, O.P., Kaur, R.,
15	Sarao, P.S., Lore, J.S. and Singh, U.M., 2022. Development and validation of a novel core
16	set of KASP markers for the traits improving grain yield and adaptability of rice under direct
17	seeded cultivation conditions. Genomics, 110269.
18	
19	Scott, M.F., Ladejobi, O., Amer, S., Bentley, A.R., Biernaskie, J., Boden, S.A., Clark, M.,
20	Dell'Acqua, M., Dixon, L.E., Filippi, C.V. and Fradgley, N. (2020). Multi-parent
21	populations in crops: a toolbox integrating genomics and genetic mapping with
22	breeding. Heredity 125: 396–416
23	
24	Shikari, A.B., Najeeb, S., Khan, G., Mohidin, F.A., Shah, A.H., Nehvi, F.A., Wani, S.A., Bhat,
25	N.A., Waza, S.A., Rao, L.S. and Steele, K.A. (2020). KASP TM based markers reveal a
26	population sub-structure in temperate rice (Oryza sativa L.) germplasm and local landraces
27	grown in the Kashmir valley, north-western Himalayas. Genetic Resources and Crop
28	Evolution 68: 821-834.
29	

1	Steele KA, Quinton-Tulloch MJ, Amgai RB, Dhakal R, Khatiwada SP, Vyas D, Heine M,
2	Witcombe JR (2018) Accelerating public sector breeding with high-density KASP markers
3	derived from whole genome sequencing of <i>indica</i> rice. Molecular Breeding 38(4):38
4	
5	Steele, K.A., Quinton-Tulloch, M., Burns, M. and Nader, W., (2020). Developing KASP
6	Markers for Identification of Basmati Rice Varieties. Food Analytical Methods 14: 663-673
7	
8	Tareke Woldegiorgis, S., Wang, S., He, Y., Xu, Z., Chen, L., Tao, H., Zhang, Y., Zou, Y.,
9	Harrison, A., Zhang, L. and Ai, Y. (2019) Rice stress-resistant SNP database. Rice 12: 1-12.
10	
11	Thomson, M.J., Singh, N., Dwiyanti, M.S., Wang, D.R., Wright, M.H., Perez, F.A., DeClerck,
12	G., Chin, J.H., Malitic-Layaoen, G.A., Juanillas, V.M. and Dilla-Ermita, C.J., 2017. Large-
13	scale deployment of a rice 6 K SNP array for genetics and breeding applications. Rice 10: 1-
14	13.
15	
16	Thudi, M., Palakurthi, R., Schnable, J.C., Chitikineni, A., Dreisigacker, S., Mace, E., Srivastava,
17	R.K., Satyavathi, C.T., Odeny, D., Tiwari, V.K. and Lam, H.M. (2021) Genomic resources
18	in plant breeding for sustainable agriculture. Journal of Plant Physiology, 257: 153351.
19	
20	Van Inghelandt, D., Melchinger, A. E., Lebreton, C., & Stich, B. (2010). Population structure
21	and genetic diversity in a commercial maize breeding program assessed with SSR and SNP
22	markers. TAG. Theoretical and Applied Genetics. 120: 1289–1299.
23	
24	Van Inghelandt, D., Frey, F.P., Ries, D. and Stich, B. (2019) QTL mapping and genome-wide
25	prediction of heat tolerance in multiple connected populations of temperate maize. Scientific
26	reports, 9: 1-16.
27	
28	Virk, D., Singh, D., Prasad, S., Gangwar J.S. and Witcombe J.R. (2003) Collaborative and
29	consultative participatory plant breeding of rice for the rainfed uplands of eastern
30	India. Euphytica 132: 95–108
31	

1	Yang, G., Chen, S., Chen, L., Sun, K., Huang, C., Zhou, D., Huang, Y., Wang, J., Liu, Y., Wang,
2	H. and Chen, Z. (2019) Development of a core SNP arrays based on the KASP method for
3	molecular breeding of rice. Rice 12: 1-18.
4	
5	Yu J, Hu S, Wang J, Wong GK, Li S, Liu B, Deng Y, Dai L, Zhou Y, Zhang X, Cao M, Liu J,
6	Sun J, Tang J, Chen Y, Huang X, Lin W, Ye C, Tong W, Cong L, Geng J, Han Y, Li L, Li
7	W, Hu G, Huang X, Li W, Li J, Liu Z, Li L, Liu J, Qi Q, Liu J, Li L, Li T, Wang X, Lu H,
8	Wu T, Zhu M, Ni P, Han H, Dong W, Ren X, Feng X, Cui P, Li X, Wang H, Xu X, Zhai W,
9	Xu Z, Zhang J, He S, Zhang J, Xu J, Zhang K, Zheng X, Dong J, Zeng W, Tao L, Ye J, Tan
10	J, Ren X, Chen X, He J, Liu D, Tian W, Tian C, Xia H, Bao Q, Li G, Gao H, Cao T, Wang
11	J, Zhao W, Li P, Chen W, Wang X, Zhang Y, Hu J, Wang J, Liu S, Yang J, Zhang G, Xiong
12	Y, Li Z, Mao L, Zhou C, Zhu Z, Chen R, Hao B, Zheng W, Chen S, Guo W, Li G, Liu S,
13	Tao M, Wang J, Zhu L, Yuan L, Yang H (2002) A draft sequence of the rice genome (Oryza
14	sativa l. ssp. Indica). Science 296(5565):79-92
15	
16	Zhang, G., Li, J., Zhang, J., Liang, X., Wang, T. and Yin, S., (2020) A high-density SNP-based
17	genetic map and several economic traits-related loci in Pelteobagrus vachelli. BMC
18	Genomics, 21: 1-17.
19	
20	Zhao, Y., Chen, W., Cui, Y., Sang, X., Lu, J., Jing, H., Wang, W., Zhao, P. and Wang, H.,
21	(2021) Detection of candidate genes and development of KASP markers for Verticillium
22	wilt resistance by combining genome-wide association study, QTL-seq and transcriptome
23	sequencing in cotton. Theoretical and Applied Genetics 134:1063-1084.
24	Zhao W, Wang J, He X, Huang X, Jiao Y, Dai M, Wei S, Fu J, Chen Y, Ren X, Zhang Y, Ni P,
25	Zhang J, Li S, Wang J, Wong G, Zhao H, Yu J, Yang H, Wang J (2004) BGI-RIS: an
26	integrated information resource and comparative analysis workbench for rice genomics.
27	Nucleic Acids Research 32:D377-382
28	Zhao H, Yao W, Ouyang Y, Yang W, Wang G, Lian X, Xing Y, Chen L, Xie W. (2015)
29	RiceVarMap: a comprehensive database of rice genomic variations. Nucleic Acids Research
30	43:D1018-22.
31	Tables

1	
2	
3	Table 1. Total number and number of predicted functional KASP designs generated per
4	chromosome

	Indica (Shuhui 498)		Japonica (Nipponbare)	
Chromosome	Total	Functional	Total	Functional
1	190,895	12,780	191,156	10,252
2	166,771	10,565	166,949	8,134
3	171,281	9,696	170,739	7,128
4	130,900	8,752	130,880	6,573
5	141,145	7,210	141,455	5,284
6	144,269	8,444	144,479	6,573
7	130,922	7,731	130,778	6,122
8	122,189	7,147	122,194	5,613
9	96,906	5,951	96,910	4,491
10	100,483	5,524	99,552	4,355
11	116,812	8,119	116,167	6,359
12	96,040	6,319	95,971	4,830
Chloroplast	0	0	n/a	n/a
Mitochondrion	5	0	n/a	n/a
Unanchored contigs	n/a	n/a	1,388	0
Total	1,608,618	98,238	1,608,618	75,714

1	Table 2. Results of Will	lcoxon rank-sum tests	(W) of difference	between properties of KASP
---	--------------------------	-----------------------	-------------------	----------------------------

2 designs resulting in successful and failing assays.

KASP design property	W	p-value	Mean	
			Success	Failure
No. of indel bases	2187	0.689	3.67	3.38
Distance of indel from target variation	2285	0.984	31.24	31.62
No. of ambiguous bases	1327001	1.20 x 10 ⁻¹⁹ ***	1.86	2.5
No. of ambiguous bases in left flank	1243608.5	7.19 x 10 ⁻¹⁰ ***	0.92	1.21
No. of ambiguous bases in right flank	1263164	5.09 x 10 ⁻¹² ***	0.94	1.29
No. of ambiguous bases in flank with least	1242962.5	1.48 x 10 ⁻¹¹ ***	0.43	0.62
No. of ambiguous bases in flank with most	1317236	1.01 x 10 ⁻¹⁸ ***	1.43	1.87
Distance to 1 st ambiguous base in left flank	371206.5	0.075	20.22	18.58
Distance to 1 st ambiguous base in right flank	372326.5	4.9 x 10 ⁻⁵ ***	20.32	17.34
Shortest distance to 1 st ambiguous base in either flank	364947	3.63 x 10 ⁻⁶ ***	15.44	12.32
Longest distance to 1 st ambiguous base in either flank	164961.5	7.39 x 10 ⁻⁴ ***	27.98	25.26
Distance to 2 nd ambiguous base in left flank	90414	0.365	29.61	28.84
Distance to 2 nd ambiguous base in right flank	93176	0.004**	29.98	27.41
Shortest distance to 2 nd ambiguous base in either flank	92249	0.002**	27.65	24.94
Longest distance to 2 nd ambiguous base in either flank	9470	0.014*	39.4	36.9
Distance to 3 rd ambiguous base in left flank	16438	0.937	34.91	34.87
Distance to 3 nd ambiguous base in right flank	13871.5	0.082	35.34	33.28
Shortest distance to 3 rd ambiguous base in either flank	13879	0.083	35.28	33.16
Distance to 4 th ambiguous base in left flank	741.5	0.007**	39.52	34.81
Distance to 4 th ambiguous base in right flank	1173.5	0.844	38.77	37.86
Shortest distance to 4 th ambiguous base in either flank	1152	0.734	38.77	37.59
Distance to 5 th ambiguous base in left flank	12.5	0.435	40.67	45
Distance to 5 th ambiguous base in right flank	39	0.456	38	42
Shortest distance to 5 th ambiguous base in either flank	39	0.456	38	42
%GC content	1192550	9.40 x 10 ⁻⁵ ***	44.07	45.82
Left flank %GC content	1210513	4.74 x 10 ⁻⁶ ***	44.44	46.4
Right flank %GC content	1165491	3.71 x 10 ⁻³ ***	43.61	45.19

Lowest flank %GC content	1184005	3.30 x 10 ⁻⁴ ***	40.07	41.56
Highest flank %GC content	1197391	4.35 x 10 ⁻⁵ ***	44.07	45.82

P values <0.05, <0.01 and <0.001 are denoted by one, two or three asterisks (at 99% confidence level)

4

Table 3. Effect of reducing the cut-off for the maximum number of bases permitted in the

4 flanking sequences of KASP designs on number of available designs and predicted success rates.

5 Distances between markers are based on their position in the Shuhui 498 *indica* reference

6 genome. Predicted success rates are calculated from the subsets that fulfil the cut-off criteria out

7 of the 4,000 markers submitted for genotyping in 178 rice lines.

8

ccessful
says
3.76
5.16
6.36
7.92
7.44
3 3 3 3 3

1 Figure Legends

2

3	Figure 1. Steps used for incorporation of rice diversity during <i>in-silico</i> design of KASP markers
4	that should discriminate between a R498 reference allele and the most common alternate allele in
5	a population of 129 diverse genomes
6	
7	Figure 2. Distribution of distances between adjacent KASP designs aligning in the Shuhui 498
8	indica (Du et al., 2017) and Nipponbare japonica (International Rice Genome Sequencing
9	Project, 2005) reference genomes. Horizontal lines represent the 1 st quartile, median, and 3 rd
10	quartile
11	
12	Figure 3. Number of polymorphic sites identified in the 8,256 possible pairwise crosses of the
13	129 diverse sequenced rice lines used for KASP design generation
14	
15	Figure 4. Hierarchical cluster analysis of 178 genotypes with the 3355 polymorphic KASP
16	markers. The groups from the PCA are indicated on the y axis (* = indica group, ***=japonica
17	group, ***** = intermediate group (Int), ******= aromatic group)
18	



Figure 1 159x230 mm (x DPI) Downloaded from https://academic.oup.com/g3journal/advance-article/doi/10.1093/g3journal/jkae251/7863403 by Prifysgol Bangor University user on 04 November 2024







Figure 3 159x101 mm (x DPI)



Figure 4 159x170 mm (x DPI)