

RESEARCH ARTICLE

Genetic diversity and core collection construction of glutinous rice landraces in 'He' cultivation zone of Guizhou using SSR sequencing

Wenhui Yang¹, Mingyi Mao¹, Jianquan Qin¹, Jamal Nasar¹, Zongdong Pan², Lijie Zhou^{1*}, Quanzhi Zhao^{1*}

1 Institute of Rice Industry Technology Research, College of Agriculture, Guizhou University, Guiyang, Guizhou, China, **2** Academy of Agricultural Sciences, Qiandongnan Miao and Dong Autonomous Prefecture, Guizhou, China

* lzhou@gzu.edu.cn (LZ); qzzhao@gzu.edu.cn (QZ)



OPEN ACCESS

Citation: Yang W, Mao M, Qin J, Nasar J, Pan Z, Zhou L, et al. (2026) Genetic diversity and core collection construction of glutinous rice landraces in 'He' cultivation zone of Guizhou using SSR sequencing. PLoS One 21(3): e0343623. <https://doi.org/10.1371/journal.pone.0343623>

Editor: Muhammad Abdul Rehman Rashid, Government College University Faisalabad, PAKISTAN

Received: October 18, 2025

Accepted: February 9, 2026

Published: March 10, 2026

Copyright: © 2026 Yang et al. 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 use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data availability statement: All raw sequencing data files are available from the Genome Sequence Archive (GSA) database (accession number: CRA027773; URL: <https://ngdc.cncb.ac.cn/gsa/browse/CRA027773>).

Abstract

Kam Sweet Rice (KSR), a distinctive group of glutinous rice landraces, has evolved over millennia through agro-ecological adaptation by the Dong ethnic group in the 'He' cultivation zone of Southeast Guizhou, China. This study examined the genetic diversity of 388 glutinous rice landraces from the region, comprising 325 KSR and 63 non-KSR varieties, using Simple Sequence Repeat (SSR) sequencing. Results revealed that non-KSR germplasm exhibited significantly higher genetic diversity than KSR germplasm. Collectively, diversity patterns were strongly shaped by the numerical predominance of genetically similar KSR germplasms, resulting in an uneven distribution of genetic diversity between KSR and non-KSR groups. Five strategies were applied to construct and evaluate core collections (see Methods for full details). Among them, the simulated annealing algorithm (SA)-based Allelic Richness Maximization Strategy (SANA) (20% sampling intensity) demonstrated superior performance in preserving genetic diversity, except for the number of alleles (N_a) and observed heterozygosity (H_o), where the Modified Heuristic Sampling (M-HS) strategy (13.66% sampling intensity) performed better at lower sampling intensities. By optimizing both approaches, a core collection of 65 germplasms was established, capturing 90.86% of alleles and retaining key genetic parameters. This core set effectively represents the genetic diversity of the entire collection, providing a strong foundation for future germplasm innovation and utilization.

Introduction

Glutinous rice (*Oryza sativa*), known for its sticky texture, plays an important role in the dietary and cultural traditions in East and Southeast Asia, contributing to what

Funding: This work was supported by the National Natural Science Foundation of China (Grant No. 32260444), Congjiang Terraced Wetland Ecosystem Observation and Research Station of Guizhou Province (Grant No. Qian-Ke-He YWZ [2024]004), Innovative Talent Team in Rice Crop Science and Technology in Karst Mountainous Areas of Guizhou Province (Grant No. Qian-Ke-He-Platform-talent BQW [2024]001), and Foundation of Guizhou University (Grant No. Guidarenjihezi [2015]09). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

has been termed the "Glutinous Rice Cultural Zone" [1]. The southern mountainous region of Qiandongnan Miao and Dong Autonomous Prefecture(QDN) in Guizhou province, China, represents a significant area for glutinous rice cultivation. Local ethnic minorities have traditionally depended on glutinous rice as a staple food, developing a distinctive cultural tradition centered on this crop. The Dong people, in particular, have developed and maintained ecological rice varieties known as 'He', which have adapted to the region's mountainous climate through a combination of natural domestication and over a thousand years of artificial selection [2,3].

The cultivation of 'He' rice is primarily concentrated in the border region of Guizhou, Hunan, and Guangxi provinces, with the largest planting areas located in Congjiang, Liping, and Rongjiang counties of Guizhou. This distribution formed a distinct 'He' cultivation zone, where glutinous varieties dominate [4]. Local Dong people named this glutinous rice "Oux Yag" in their native language, and is internationally referred to as "Kam Sweet Rice" [5]. Kam Sweet Rice (KSR) is prized for its favorable agronomic and culinary traits, including a rich, mellow aroma, high stickiness, cold tolerance, and resistance to pests and diseases [3]. Later, the Food and Agriculture Organization of the United Nations (FAO) has recognized KSR as a specialty rice of global significance [6]. Besides the wide range of KSR varieties, the 'He' rice Cultivation Zone in Guizhou is also a home to diverse array of non-KSR glutinous rice landraces. These two groups differ not only in name but also in both harvesting and post-harvest processing. For example, KSR varieties are difficult to thresh under natural conditions and require traditional milling methods, specifically, the use of mortar and pestle, to effectively remove husk. In contrast the non-KSR varieties are relatively easy to thresh using conventional methods [7].

Genetic resources are a vital component of biodiversity, forming the foundation for agricultural productivity, resilience, and adaptability. They play a crucial role in promoting sustainable agricultural development and ensuring global food security [8]. To improve the management and utilization of large germplasm resources, Frankel et al. [9] proposed the concept of a core collection, a subset designed to capture the genetic diversity of the entire collection using the smallest number of representative samples. Brown et al. [10] demonstrated that for germplasm collection exceeds 3,000 accessions, sampling 10% of the total can retain around 70% of the overall genetic diversity. Typically, core collections are developed using 5% to 30% of the original collection, a range shown to be effective for preserving genetic variation and improving the efficiency of resource use [11,12]. The optimal sampling proportion depends on factors such as the total population size, the level of genetic diversity, and the specific sampling strategy used [13]. Various strategies have been developed to construct core collections based on different objectives. For example, Maximization (M strategy) [14], Allele Coverage (CV) [15], and simulated annealing algorithm (SA)-based strategies have been adopted: one maximizing allelic richness (SANA) for allele representation, and the other maximizing genetic diversity (SAGD) for genetic diversity optimization [16]. Additionally, the A-NE strategy is used to optimize the average genetic distance between germplasm and the closest core set entry [17].

These strategies with appropriate weighting can help achieve both comprehensive allele representation and high genetic diversity in core collections [18,19].

Accurate assessment of genetic diversity is critical for the effective construction of core collections [20]. Compared to traditional phenotypic assessment, molecular marker-based approaches offer distinct advantages. They are not influenced by environmental conditions and more accurately reflect the inherent genetic variation within germplasm, providing a reliable basis for core collection development [21]. Among molecular markers, Simple Sequence Repeats (SSRs), short tandem repeat sequences widely distributed throughout plant genomes, are particularly valuable due to their high polymorphism, co-dominant inheritance, reproducibility, and ease of use [22,23]. However, traditional SSR detection techniques, which rely mainly on electrophoresis, are limited in resolution and accuracy because they only measure fragment length and lack detailed sequence information [24]. The Advent of Next-Generation Sequencing (NGS) technology has enabled the development of SSR sequencing, which integrates SSR markers with high-throughput sequencing. This approach allows for precise detection of SSR repeat units, allele frequency, and relative abundance, thereby significantly improving resolution and accuracy [25,26]. Unlike traditional methods, SSR sequencing enhances the detection of allelic variation and reveals greater levels of genetic polymorphism [27]. This technology has been successfully applied in genetic studies of various crop germplasm resources, including cucumber [26], radish [28], and *Camellia oleifera* [27].

In this study, we analyzed and compared the genetic diversity of Kam Sweet Rice (KSR) and non-KSR glutinous rice collected from the 'He' cultivation zone Qiandongnan (QDN) region. Based on SSR sequencing, we constructed a core collection of glutinous rice using different strategies. These findings provide important insights into the genetic structure of glutinous rice in the 'He' cultivation zone and offer a valuable reference for developing core collections using integrated strategies frameworks.

Materials and methods

Rice materials

A total of 388 glutinous rice germplasm samples, including 325 KSR and 63 non-KSR glutinous rice landraces, were collected from the 'He' cultivation zone by the Qiandongnan Miao and Dong Autonomous Prefecture Academy of Agricultural Sciences in Guizhou Province, China (S1 Table). Among them, 223 accessions originated from Liping County, 150 from Congjiang County, and 11 from Rongjiang County in Guizhou Province, while four accessions were collected from Sanjiang County, Guangxi Province. Fig 1 shows the geographic distribution of the sampling sites at the provincial and county levels. In addition, four representative *indica* (9311, Minghui 63, IR50, Huanghuazhan) and *japonica* varieties (Nipponbare, Daohuaxiang, Wuyujing 3, Koshihikari) were used as controls for clustering to inaugurate *indica-japonica* differentiation.

DNA extraction and quality control

For each sample, ten seeds from a single panicle of each landrace were germinated under controlled indoor conditions. Fresh leaf tissue (20 mg) was collected from each sample, ground in liquid nitrogen, and genomic DNA was extracted using the DNA secure Plant Kit (Tiangen Biotech, China). DNA concentration and purity were evaluated using a NanoDrop 2000 spectrophotometer (Thermo Scientific, USA), with acceptable samples showing OD260/280 values between 1.7 and 1.9, and OD260/230 values above 2.0. DNA integrity was confirmed via agarose gel electrophoresis.

SSR primer screening and genotyping

Initially, 348 SSR markers distributed across the 12 rice chromosomes [29,30] were assessed for suitability in SSR sequencing. Selection criteria for optimal markers includes: (1) repeat units of at least 3bp; (2) avoidance of homopolymer-rich sequences (e.g., only GC or AT); (3) no nearby SSR loci within the flanking region; and (4) fewer than

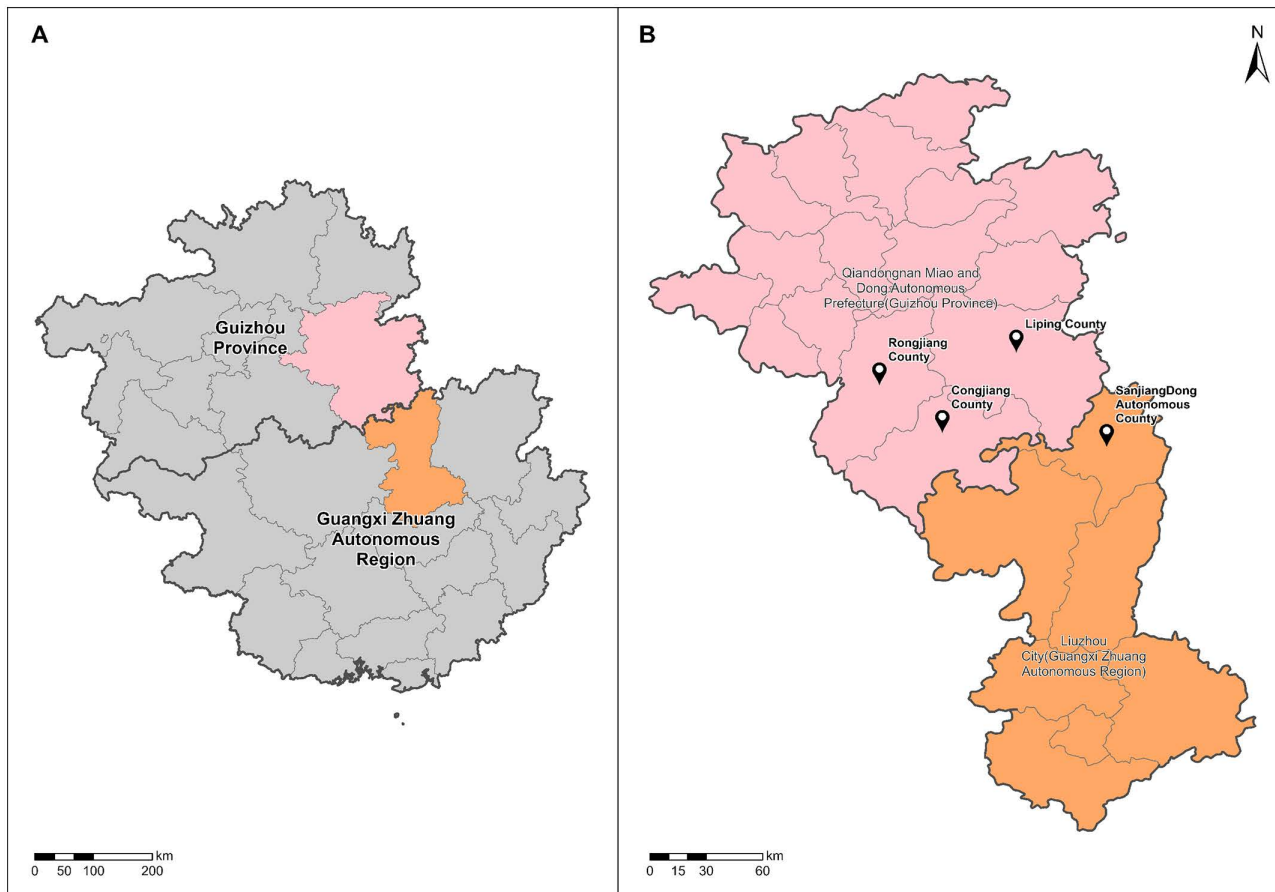


Fig 1. Geographic distribution of sampling counties for glutinous rice germplasm in the 'He' cultivation zone of Southwest China. (A) Location of Qiongdongnan Miao and Dong Autonomous Prefecture (pink, Guizhou Province) and Liuzhou City (orange, Guangxi Province). (B) Distribution of the four sampling counties, including Liping, Congjiang, and Rongjiang Counties in Qiongdongnan Prefecture (Guizhou Province), and Sanjiang County in Liuzhou City (Guangxi Province). Maps were created with ArcGIS Pro 3.6.0. The administrative boundary data were obtained from the GADM database (<https://gadm.org/>).

<https://doi.org/10.1371/journal.pone.0343623.g001>

10 repeat units. Based on preliminary bioinformatics analysis, 80 SSR markers were deemed suitable. Further screening prioritized marker polymorphism and even chromosomal distribution, resulting in the selection of 37 SSR primer pairs (see [S2 Table](#)).

Genotyping was employed using SSR sequencing technology on the Illumina HiSeq 4000 platform (paired-end 2 × 150 bp) [27] by Genesky Biotechnologies Inc. (Shanghai, China). The 37 SSR markers were grouped into two multiplex PCR panels (18 or 19 markers per panel). Amplicons from each panel were pooled based on fragment intensity and count, with standard and GC-rich PCR systems processed in parallel. Following initial amplification, Index PCR added sample-specific barcodes. Indexed amplicons were pooled, gel-purified, and verified using an Agilent 2100 Bioanalyzer (Illumina, San Diego, CA, USA). Raw sequencing reads were first assessed using FastQC to evaluate base quality scores and nucleotide composition. Paired-end reads were merged using FLASH, and only successfully merged sequences were retained for subsequent analyses. The merged reads were then aligned to primer-captured reference sequences

using BLASTn, and only high-confidence, locus-specific reads were defined as effective reads and used for genotyping. For each sample and SSR locus, the number of effective reads was calculated to ensure adequate sequencing depth. To further evaluate the specificity and efficiency of target locus enrichment, merged reads were aligned to the reference genome using BLASTn, and the proportion of reads mapping to target regions was calculated. Potential genotyping errors, including allele dropout, PCR stutter interference, and null alleles, were assessed indirectly based on read depth distribution, locus-specific missing rates, and allelic read count balance within heterozygous genotypes. Loci or samples showing extremely low effective read counts, excessive missing data, or highly unbalanced allelic read proportions were excluded from downstream analyses to reduce the influence of unreliable genotypes. Allele frequencies and allele counts were then calculated from the filtered dataset to generate the final SSR genotyping matrix for subsequent population genetic analyses.

Construction of core collections

Five sampling strategies were used to construct preliminary core collections based on SSR data: (1) Modified Heuristic Sampling (M-HS) Strategy (PowerCore v1.0) [31]: Uses a modified heuristic algorithm to reduce redundancy and capture 100% of allele diversity without manual sampling ratio adjustment, (2) SANA Strategy (PowerMarker v3.25) [16]: Simulated annealing algorithm (SA) that maximizes allele richness, (3) A-NE Strategy and (4) E-NE Strategy (Core Hunter v3.0) [17,32]: Optimize average (A-NE) and extreme (E-NE) genetic distances, using the Modified Rogers distance, and (5) CV Strategy (Core Hunter v3.0) [15]: Focuses on maximizing allele coverage with minimal sample size [33]. While the M-HS strategy determines the sampling ratio automatically, the other four strategies were evaluated at fixed sampling proportions of 5%, 10%, 15%, 20%, 25%, and 30%.

Data analysis

Genetic diversity parameters, including the number of alleles (N_a), effective alleles (N_e), Shannon's information index (I), observed heterozygosity (H_o) and expected heterozygosity (H_e) were analyzed using GenAIEx v.6.51b2 software [34]. The major allele frequency (MAF), polymorphism information content (PIC), and fixation index (F) were analyzed using PowerMarker v.3.25 software [16]. The Shapiro-Wilk test was used to evaluate the normality of the genetic diversity parameters. The results indicated that all genetic diversity parameters significantly ($P < 0.01$) deviated from normality in both the entire germplasm collection and KSR. Therefore, the Mann-Whitney U test (SPSS version 27; IBM, Armonk, NY, USA) was used to assess differences between KSR and non-KSR groups, and between the core collection and the full germplasm set. Retention rate was calculated as the mean of locus-wise parameter values in the core collection divided by the mean of locus-wise parameter values in the whole collection (37 loci). Population genetics parameters, including intra-population inbreeding coefficient (F_{is}), total inbreeding coefficient (F_{it}), population differentiation index (F_{st}), and gene flow (Nm), were calculated using Popgene v1.32 [35]. AMOVA (Analysis of Molecular Variance) was conducted using Arlequin v3.5.2.2 [36].

Clustering analysis based on Nei's genetic distance was performed using the UPGMA method in PowerMarker v3.25, and visualized using Evolview [37]. Population structure analysis was carried out using STRUCTURE v2.3.4 [38], with a burn-in of 10,000 iterations and 50,000 MCMC replications. The number of subpopulations (K) ranged from 1 to 11, with 10 replicates per K . STRUCTURE HARVESTER v0.6.91 [39] was used to determine the optimal K using the ΔK method. Q -values were averaged across replicates using CLUMPP v1.1.2 [40]. Individuals with $Q \geq 0.6$ were assigned to a specific group; those with all Q -values < 0.6 were designated as admixed or unassigned. To evaluate core collection representativeness, six genetic parameters (N_a , N_e , I , H_o , H_e , and PIC) were compared between core subsets and the full collection. Principal Coordinate Analysis (PCoA) was also conducted using GenAIEx v6.51b2 to visualize genetic relationships and confirm the genetic representativeness of the core collections.

Results

Genetic diversity based on SSR sequencing

A genetic diversity analysis was conducted on 388 glutinous rice germplasms from the 'He' cultivation zone using 37 SSR markers derived from SSR sequencing (Table 1). In total 186 alleles were identified across all loci. The number of alleles per locus (N_a) ranged from 2 to 16, with an average 5.0. This number for N_e ranged from 1.0536 to 6.4317 (mean = 1.8453). The index I , reflecting allele distribution uniformity and evenness, ranged from 0.1554 to 2.1573, averaging 0.6970. MAF values ranged from 0.2964 to 0.9741 (mean = 0.7663), higher values indicating the dominance of certain alleles across the population. H_o ranged from 0.0026 to 0.5902, averaging 0.0406, while H_e ranged from 0.0509 to 0.8445 (mean = 0.3442). PIC values ranged from 0.0505 to 0.8291, averaging 0.3166. Overall, the N_a and I values reflect allelic variation across loci, whereas the relatively high MAF and low H_o (and comparatively lower H_e) indicate skewed allele-frequency distributions and low heterozygosity, as expected in a predominantly self-pollinating crop such as rice.

A comparative analysis of genetic diversity between KSR and non-KSR germplasms is summarized in Table 1, S3 and S4 Tables. There was no significant difference in N_a between the KSR and non-KSR groups (Table 1). However, non-KSR germplasms showed significantly higher N_e (1.3463 [1.3087–1.9385] vs. 1.2268 [1.1905–1.9059]; $P=0.0001$), I (0.4938 [0.4007–0.8150] vs. 0.3851 [0.2972–0.7444]; $P=0.0002$), and H_e ($P<0.001$), along with significantly lower MAF (0.8544 [0.6671–0.8634] vs. 0.8969 [0.6841–0.9123]; $P=0.0002$). Notably, H_o was lower in both groups. The number of highly polymorphic sites ($PIC \geq 0.50$) were greater in non-KSR (12) than in KSR (7). Moderately polymorphic sites ($0.25 \leq PIC < 0.50$) were 19 in non-KSR and 6 in KSR (S3 and S4 Tables). These loci accounted for 83.78% and 35.14% of all markers in non-KSR and KSR, respectively. Overall, non-KSR accessions exhibited significantly higher with-group diversity (N_e , I , H_e and PIC) than KSR, indicating that panel-level diversity estimates are strongly influenced by the numerical predominance and lower within-group diversity of KSR accessions in this regional collection.

Genetic population structure

UPGMA clustering analysis divided 388 glutinous rice germplasms and eight reference varieties (4 *indica* and 4 *japonica*) into two major genetic clusters corresponding to the *indica* and *japonica* groups (Fig 2A). One cluster comprises all 4 *japonica* reference varieties and 335 landraces, while the other contains 4 *indica* reference varieties and 53 landraces. *Japonica* glutinous rice dominates the 'He' cultivation zone, comprising 86.3% of the total. Differences were observed

Table 1. Genetic diversity analysis of 388 glutinous rice germplasm using SSR sequencing.

Parameters ^a	Overall ^b (n = 388)	KSR ^c (n = 325)	non-KSR ^c (n = 63)	P ^d
N_a	5.0 (2-16)	4(2-6)	4(3-6)	0.6749
N_e	1.8453(1.0536-6.4317)	1.2268(1.1905-1.9059)	1.3463(1.3087-1.9385)	0.0001
I	0.6970(0.1554-2.1573)	0.3851(0.2972-0.7444)	0.4938(0.4007-0.8150)	0.0002
MAF	0.7663(0.2964-0.9741)	0.8969(0.6841-0.9123)	0.8544(0.6671-0.8634)	0.0000
H_o	0.0406(0.0026-0.5902)	0.0123(0.0092-0.0248)	0.0129(0.0103-0.0289)	0.0278
H_e	0.3442(0.0509-0.8445)	0.1849(0.1600-0.4753)	0.2572(0.2359-0.4841)	0.0001
PIC	0.3166(0.0505-0.8291)	0.1722(0.1472-0.3826)	0.2383(0.2081-0.3937)	0.0001

^a N_a : Number of alleles; N_e : Effective number of alleles; I : Shannon's information index; MAF : Major allele frequency; H_o : Observed heterozygosity; H_e : Expected heterozygosity; PIC : Polymorphic information content.

^b Mean(Minimum-Maximum).

^c Median(interquartile range).

^d P from Mann-Whitney U tests comparing locus-wise parameter values(37 loci) between KSR and non-KSR groups.

<https://doi.org/10.1371/journal.pone.0343623.t001>

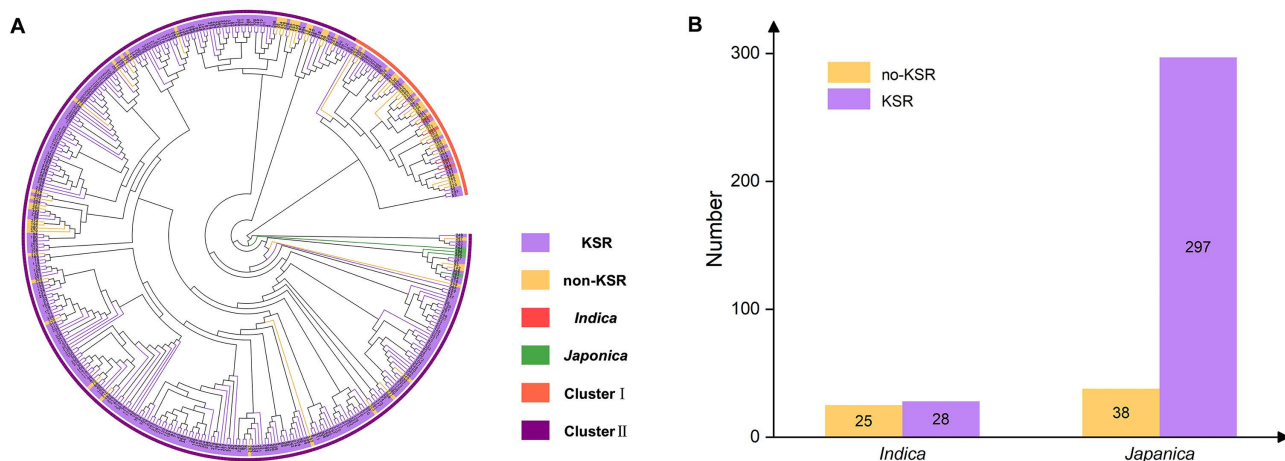


Fig 2. UPGMA clustering of glutinous rice germplasms (4 indica and 4 japonica reference controls). (A) UPGMA clustering diagram of 388 glutinous rice germplasms alongside 8 *indica/japonica* references. (B) The *indica/japonica* composition in KSR and non-KSR germplasms.

<https://doi.org/10.1371/journal.pone.0343623.g002>

between KSR and non-KSR germplasm types. Among non-KSR samples, 25 were indica and 38 were japonica out of 63. In contrast, 91.4% of KSR varieties (297 out of 325) belonged to the japonica group (Fig 2B), in agreement with previous findings [2].

STRUCTURE analysis further confirmed the population structure. Structure Harvester analysis revealed ΔK value trends as K increased from 1 to 11, with a significant ΔK peak occurring at K=2 (Fig 3A), indicating that the 388 glutinous rice germplasm resources could be categorized into two distinct subpopulations based on genetic structure. The hierarchical genetic structure analysis with K=2 resulted in one subpopulation of 332 germplasms and another of 53 germplasms (Fig 3B). Among these germplasms, Niumaohu, Dongsui, and Gaicaohe (with Q-values < 0.6; S5 Table) exhibited admixed ancestry and were not assigned to either subpopulation. Comparing results from the Bayesian model-based STRUCTURE analysis with UPGMA clustering analysis based on Nei's genetic distance revealed that 385 rice germplasms showed completely consistent groupings, except for three KSR germplasms (Niumaohu, Dongsui, and Gaicaohe).

Construction of core collections

Five strategies such as M-HS, SANA, A-NE, E-NE, and CV, were used to construct primary core collection from 388 glutinous rice germplasms based on SSR sequencing. The M-HS strategy automatically selected 53 accessions (13.66%) as core subset. The other sampling strategies such as SANA, A-NE, E-NE, and CV generated subset at six sampling intensities (i.e., 5%, 10%, 15%, 20%, 25%, and 30%), respectively, resulting in core sets of 19–116 accessions (Table 2).

Comparison of different methods for core collection construction

To evaluate the representativeness, genetic diversity parameters (i.e., N_a , N_e , I , H_e , H_o , and PIC) were compared between 25 subsets and the entire collection (Table 2). The M-HS strategy core subset did not show any significant differences in N_a and H_o , but did differ significantly ($P < 0.01$) in N_e , I , H_e and PIC at 13.66% sampling intensity. For the SANA strategy, subsets at 5% and 10% sampling intensities showed significant differences ($P < 0.01$) in N_a and/or H_o . However, subsets with 15% to 30% sampling intensities exhibited no significant ($P > 0.05$) differences across all six genetic diversity parameters. In contrast, core subsets generated via A-NE and E-NE strategies at various sampling intensities showed significant difference from the entire collection. Among CV strategy subsets, only the 30% sampling intensity subset showed no significant difference.

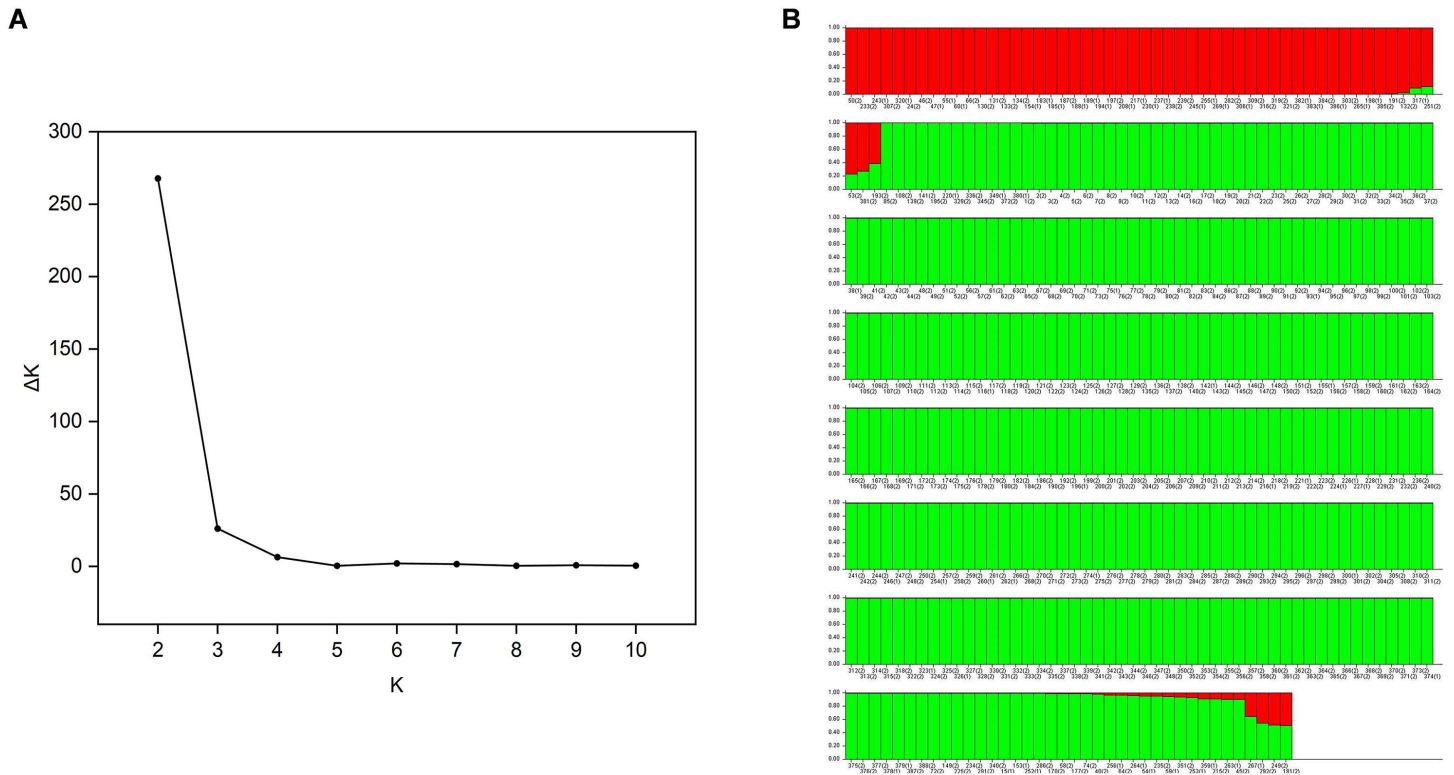


Fig 3. Population genetic structure analysis of 388 glutinous rice germplasms. (A) Estimation of the optimal K value based on ΔK analysis. (B) Inferred population structure of 388 samples at $K=2$, where each vertical bar represents an individual sample. The numbers in parentheses indicate that 1 denotes non-KSR and 2 denotes KSR. Colors indicate genetic clusters (green for subpopulation I and red for subpopulation II), with segment lengths proportional to membership probabilities (Q-values).

<https://doi.org/10.1371/journal.pone.0343623.g003>

Genetic diversity retention rates for representative core subsets are shown in [Table 3](#). The CV strategy core subset (30% sampling intensity) demonstrated retention rates of 100% (Na), 106.54% (Ne), 108.39% (I), 115.72% (Ho), 108.89% (He), and 108.44% (PIC), with the higher sampling intensity more effectively capturing the diversity concentrated in the entire collection. The SANA strategy core subset (20% sampling intensity) showed over 100% genetic diversity for Ne , I , He , and PIC , but lower retention Na (85.48%) and Ho (96.13%). The M-HS strategy achieved 100% Na and 115.04% Ho retention at 13.66% sampling intensity, showing no significant difference in these parameters compared to the entire collection ([Table 3](#)).

Although several single-strategy core subsets showed no significant differences from the entire collection, they generally required relatively high sampling intensities to achieve overall representativeness (e.g., SANA at $\geq 15\%$ and CV at 30%). In contrast, the M-HS strategy achieved complete allelic richness (Na) retention at a much smaller sampling intensity (13.66%) but still differed from the entire collection in other diversity parameters. Therefore, no single strategy simultaneously satisfied high Na coverage, overall representativeness across diversity parameters, and a practical core size. To address this trade-off, we used the above comparative results to guide a results-driven optimization ([Fig 4](#)). Based on the observed differences among strategies in allelic richness retention, overall diversity representativeness, and sampling efficiency, we prioritized candidate solutions that integrated complementary strengths across these evaluation criteria. Multiple alternative core sets were then compared within a reduced candidate space to identify statistically acceptable and

Table 2. Comparison of Genetic diversity parameters between primary core subsets compared with the entire collection.

Method	Sampling Ratio	<i>Na</i> ^a	<i>Ne</i> ^a	<i>I</i> ^a	<i>Ho</i> ^a	<i>He</i> ^a	<i>PIC</i> ^a
Whole collection	100.00%	4	1.3463	0.4938	0.0129	0.2572	0.2383
M-HS	13.66%	4	1.9181**	0.7742**	0.0192	0.4786**	0.3904**
SANA	5.00%	2**	1.4979	0.5367	0.0000**	0.3324	0.2772
	10.00%	3	1.4077	0.5144	0.0000**	0.2896	0.2554
	15.00%	3	1.3484	0.5156	0.0175	0.2584	0.2329
	20.00%	3	1.4316	0.6115	0.0130	0.3015	0.2854
	25.00%	3	1.3525	0.5280	0.0206	0.2607	0.2493
	30.00%	3	1.3226	0.4842	0.0087	0.2439	0.2247
A-NE	5.00%	2**	1.3623	0.4362	0.0000*	0.2659	0.2306
	10.00%	3	1.5321	0.5315	0.0263*	0.3473	0.2870
	15.00%	3	1.7185**	0.6622*	0.0172	0.4181**	0.3428**
	20.00%	3	1.8410**	0.7441**	0.0260*	0.4568**	0.3755**
	25.00%	4	1.8363**	0.7624**	0.0309**	0.4554**	0.3891**
	30.00%	4	1.9093**	0.7929**	0.0259**	0.4763**	0.4032**
E-NE	5.00%	3*	2.0570**	0.8570**	0.0526	0.5139**	0.4254**
	10.00%	3	2.0870**	0.8629**	0.0263	0.5208**	0.4453**
	15.00%	3	2.0867**	0.8622**	0.0345**	0.5208**	0.4450**
	20.00%	4	2.1048**	0.8676**	0.0263**	0.5249**	0.4342**
	25.00%	4	2.0118**	0.8344**	0.0309**	0.5029**	0.4217**
	30.00%	4	1.8715**	0.7808**	0.0330**	0.4657**	0.3998**
CV	5.00%	3	2.2192**	0.8676**	0.0000	0.5494**	0.4479**
	10.00%	4	1.9295**	0.7737**	0.0270*	0.4817**	0.4003**
	15.00%	4	1.7133**	0.6679*	0.0192**	0.4163**	0.3459*
	20.00%	4	1.7535*	0.7189*	0.0260*	0.4297*	0.3674*
	25.00%	4	1.5988*	0.6462	0.0235**	0.3745*	0.3239*
	30.00%	4	1.4394	0.5804	0.0174	0.3053	0.2728

Abbreviations for all parameters are consistent with those listed in Table 1.

^aData in the table represent median values across (37 SSR loci). *P<0.05 and **P<0.01 indicate significance based on Mann-Whitney U tests comparing locus-wise parameter values between primary core subsets and the entire collection. As only medians are shown for brevity, significant differences may occur even when medians are similar.

Bold values indicate that there are no significant differences between the core subsets and the entire collection across the six parameters.

<https://doi.org/10.1371/journal.pone.0343623.t002>

sampling-efficient solutions. Candidate core sets were evaluated against the entire collection using Mann-Whitney U tests for *Na*, *Ne*, *I*, *Ho*, *He*, and *PIC*, together with diversity retention rates. Core sets showing no significant differences across all parameters were considered acceptable, and when multiple candidates met this criterion, lower sampling intensity was prioritized while maintaining high *Na* retention. The final selection was further interpreted in light of population composition to avoid adding redundant accessions with limited contribution to overall diversity.

Optimized construction strategy for the final core collection

Based on the optimization framework described above, we combined two preliminary subsets to form an enriched candidate pool for final core selection: 120 germplasms in total, selected from the M-HS (53) and SANA (77) strategies. Ten germplasms (Heironghe, Goubadang, Zaogaonuo, Heimahe, Goulieshibe, Huangmaozaoan, Liuyuegu, Zhuyanuo-2, Nuopangu, and Shuba) were common to between the two, while 110 germplasms were unique. From the established ten core sets, Core set 1 was generated using the M-HS strategy based on the 110 unique accessions, yielding 49 unique

Table 3. Genetic diversity retention rates for representative core germplasm subsets.

Method	Sampling Ratio	<i>Na</i>	<i>Ne</i>	<i>I</i>	<i>Ho</i>	<i>He</i>	<i>PIC</i>
CV	30.00%	100.00%	106.54%	108.39%	115.72%	108.89%	108.44%
SANA	15.00%	80.11%	101.23%	98.64%	108.92%	100.39%	99.89%
	20.00%	85.48%	103.21%	105.33%	96.13%	108.65%	107.43%
	25.00%	84.95%	98.42%	97.56%	91.53%	98.75%	98.38%
	30.00%	85.48%	98.51%	97.11%	95.81%	97.53%	97.68%
M-HS	13.66%	100.00%	NA ^a	NA ^a	115.04%	NA ^a	NA ^a

Abbreviations of all parameters are consistent with those listed in Table 1.

^aNA indicates values excluded from analysis due to lack of meaningful interpretation, as these parameters indicates significant differences ($p < 0.01$) from the entire collection presented in Table 1.

Bold values indicate the superior core subsets based on genetic diversity retention rates.

<https://doi.org/10.1371/journal.pone.0343623.t003>

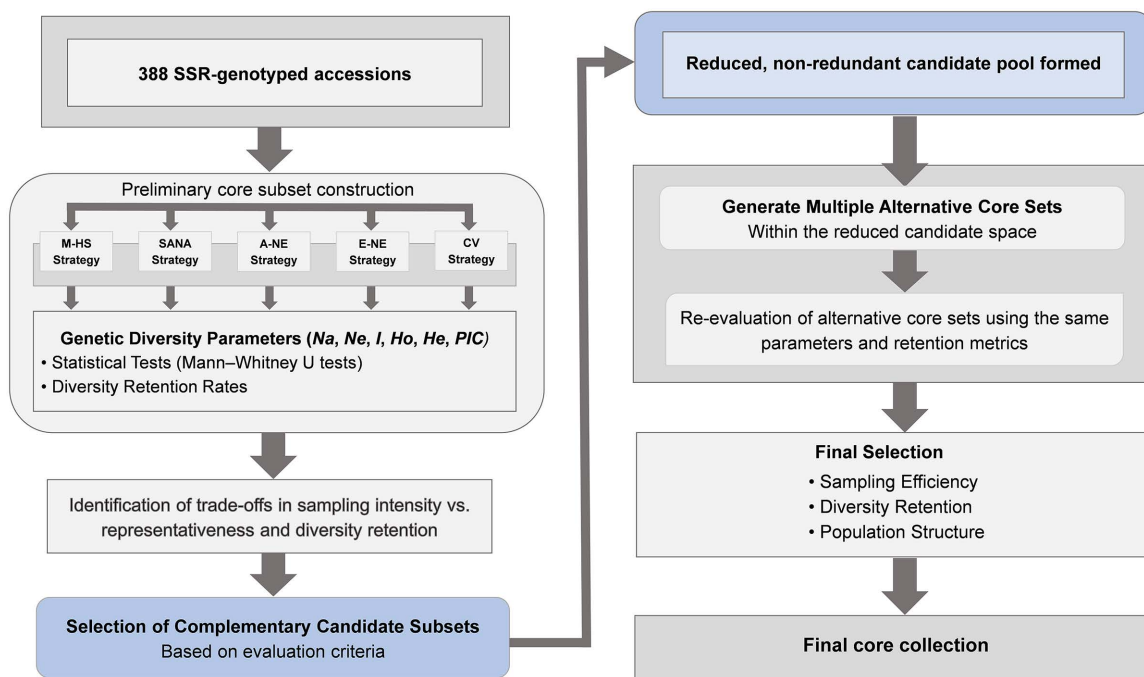


Fig 4. Results-driven workflow for identifying the final core germplasm collection.

<https://doi.org/10.1371/journal.pone.0343623.g004>

accessions, to which the 10 shared accessions were added (total $n = 59$). Core Sets 2–10 were developed using the SANA strategy with incremental sampling intensities (10%–90% of the 110 unique accessions), with the same 10 shared accessions added back to each set (Table 4).

Diversity parameters for all core sets are summarized in Table 4. Among the ten optimized core sets, only Core sets 6 and 9 showed no significant differences ($P > 0.05$; Mann–Whitney U test) from the entire collection across all six genetic diversity parameters (*Na*, *Ne*, *I*, *Ho*, *He*, and *PIC*) and were therefore considered acceptable candidates. Both core sets also showed allele coverage above 90% and retention rates above 100% for *Ne*, *I*, *Ho*, *He*, and *PIC* (Table 5) outperforming most preliminary core subsets (Tables 3 and 5). The two acceptable candidates differed in sampling efficiency and *Na* retention: Core set 6 achieved 16.75% sampling intensity with 90.86% *Na* retention, whereas Core set 9 required a higher

sampling intensity(25.26%) but retained 97.85% of *Na*. Considering that non-KSR accessions exhibited significantly higher within-group diversity (*Ne*, *I*, *He* and *PIC*) than KSR, and that panel-level diversity estimates in this regional collection are strongly influenced by the numerical predominance and lower within-group diversity of KSR accessions, increasing sampling intensity primarily tends to add redundant KSR genotypes, thereby offering limited additional contribution to overall diversity representation. Therefore, among the two acceptable candidates, we selected the more sampling-efficient Core set 6 as the final core collection.

PCoA analysis further confirmed that the core sets 6 and 9 reflected the genetic structure of the entire glutinous rice population, with even distribution across principle coordinates (Fig 5). Based on sampling efficiency, genetic diversity parameters, and population structure representation, Core set 6 comprising 65 germplasms (S6 Table), was determined to be the optimal core collection for the 'He' cultivation zone.

Discussion

Genetic resources are the foundation of crop improvement, with genetic variation directly affecting breeding outcomes. Historic breakthroughs in rice breeding, such as the "Green Revolution" and hybrid rice development, stemmed from

Table 4. Genetic diversity parameters for different optimized core sets compared with the entire collection.

Core set code	Optimization strategy	Core set size	Sampling Ratio	<i>Na</i>	<i>Ne</i>	<i>I</i>	<i>Ho</i>	<i>He</i>	<i>PIC</i>
Whole collection		388	100.00%	4	1.3463	0.4938	0.0129	0.2572	0.2383
Core set 1	M-HS ^a	49+10 ^b	15.21%	4	1.9840**	0.8026**	0.0172	0.4960**	0.4028**
Core set 2	SANA ^a & 10%	11+10 ^b	5.41%	3	1.9417**	0.6829**	0.0000*	0.4850**	0.3698**
Core set 3	SANA ^a & 20%	22+10 ^b	8.25%	3	1.8737**	0.8047**	0.0313	0.4663**	0.4156**
Core set 4	SANA ^a & 30%	33+10 ^b	11.08%	3	1.7018**	0.7128*	0.0233	0.4124**	0.3669*
Core set 5	SANA ^a & 40%	44+10 ^b	13.92%	3	1.7071**	0.7517*	0.0185	0.4142**	0.3800*
Core set 6	SANA ^a & 50%	55+10 ^b	16.75%	3	1.5305	0.6838	0.0154	0.3466	0.3087
Core set 7	SANA ^a & 60%	66+10 ^b	19.59%	4	1.6370*	0.7129*	0.0133	0.3891*	0.3557*
Core set 8	SANA ^a & 70%	77+10 ^b	22.42%	4	1.5780*	0.6781*	0.0230	0.3663*	0.3490*
Core set 9	SANA ^a & 80%	88+10 ^b	25.26%	4	1.5209	0.6242	0.0108	0.3425	0.3016
Core set 10	SANA ^a & 90%	99+10 ^b	28.09%	4	1.6183*	0.6949*	0.0183	0.3821*	0.3410*

Abbreviations for all parameters are consistent with those listed in Table 1.

^a M-HS strategy and SANA strategy were applied to 110 samples not shared between the two core subsets constructed using the SANA strategy with a 20% sampling intensity and the M-HS strategy.

^b values before and after '+' represent number obtained from each optimization strategy based on 110 samples and the common number for both core subsets using the SANA strategy with a 20% sampling intensity and the M-HS strategy.

*P < 0.05 and **P < 0.01 indicate significance based on Mann-Whitney U tests comparing locus-wise parameter values between the optimized core sets and the entire collection.

Values in bold indicate that the optimized core sets did not differ significantly from the entire collection across six parameters.

<https://doi.org/10.1371/journal.pone.0343623.t004>

Table 5. Genetic diversity retention rates of core set 6 and 9.

Core Set code	Sample size	Sampling Ratio	<i>Na</i>	<i>Ne</i>	<i>I</i>	<i>Ho</i>	<i>He</i>	<i>PIC</i>
Core set 6	65	16.75%	90.86%	117.61%	116.10%	102.63%	118.42%	117.32%
Core set 9	98	25.26%	97.85%	113.50%	115.41%	104.16%	116.57%	115.72%

Abbreviations for all parameters are consistent with those listed in Table 1.

<https://doi.org/10.1371/journal.pone.0343623.t005>

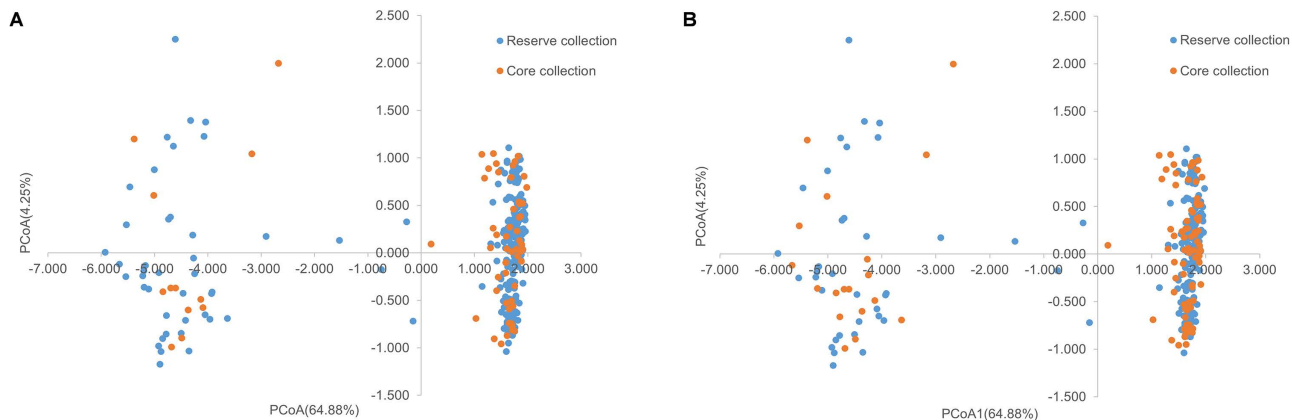


Fig 5. PCoA Analysis of Core sets 6(A) and 9(B).

<https://doi.org/10.1371/journal.pone.0343623.g005>

discovery and utilization of key genetic resources [41,42]. Despite their importance, the vast volume of genetic resources poses challenges for their effective conservation and utilization [43]. Preserving and harnessing effective genetic diversity remains essential for maximizing their breeding potential.

Kam Sweet Rice (KSR), primarily grown in the mountainous regions of southeastern Guizhou, represent a unique group of glutinous rice landraces. Beyond its value as a genetic resource, KSR is a culturally embedded landrace closely tied to the cultural identity and traditional livelihoods of the Dong people. This cultural embeddedness supports its continued on-farm cultivation and maintenance [2,3]. Introduced during the Dong people migration, KSR possess a distinct genetic background compared to other Guizhou landraces [2]. In this study, 325 KSR and 63 non-KSR rice landraces from the 'He' cultivation zone were collected and preserved by the Qiandongnan Agriculture Science Institute [44] representing the genetic diversity characteristic of this region. The genetic diversity observed in this study is consistent with previous SSR-based analyses of glutinous rice landraces in southwest China, while also revealing clear regional and population-specific differences. For instance, a SSR-based study on glutinous rice landraces and promoted cultivated varieties in Tengchong, Yunnan reported an average of 4.6 alleles per locus, a mean I of 0.398 and PIC values spanning 0.26–0.84, providing a useful regional reference for diversity estimates [45]. In contrast, a large-scale analysis of KSR from southeast Guizhou reported low within-group diversity and heterozygosity and suggested that KSR may represent a highly stabilized landrace group with a distinct genetic background [46]. Our results are consistent with these observations: non-KSR landraces in the 'He' cultivation zone exhibited higher diversity than KSR, whereas KSR accessions were more genetically similar and numerically predominant. Clustering and genetic structure analyses classified both KSR and non-KSR varieties into two subpopulations: *japonica* and *indica*, without significant differentiation between the two groups. This supports the view that natural selection, more than artificial selection, drives the main divergence between rice subspecies [47]. KSR cultivation is deeply influenced by traditional dietary culture of ethnic minorities in southeastern Guizhou [48], and continuous selection and exchange among farmers have helped shaped rice diversity [49]. Notably, H_o was extremely low across loci in both KSR and non-KSR landraces, which is expected for predominantly self-pollinating rice. This pattern may be further reinforced by on-farm seed saving/purification and localized cultivation that limit gene flow, as farmer management can affect heterozygosity in landraces [50]. For KSR, recent genomic evidence also indicates domestication-related bottlenecks and reductions in effective population size, plausibly linked to major agricultural transitions in southeastern Guizhou (the "glutinous-to-nonglutinous" shift and the adoption of hybrid rice) [2].

Genetic differentiation analysis of KSR and non-KSR revealed significant inbreeding within populations ($F_{is}=0.9141$), moderate genetic differentiation between populations ($F_{st}=0.0737$), and frequent gene exchange between the two populations ($Nm=37.7232$) (S7 Table). AMOVA showed that 16.57% of the genetic variations was between the KSR and non-KSR populations, while 73.43% accounted among individuals within populations (S8 Table), indicating that most diversity is intra-population.

KSR exhibited lower allele richness and frequency than non-KSR (Table 1), with dominance of specific alleles likely due to the long-term ethnic selection [48]. Non-KSR landraces had higher proportion of highly/moderately polymorphic loci (83.78%) compared to KSR (35.14%), indicating greater genetic diversity, consistent with previous findings [51].

Using 37 SSR markers, we constructed a core collection of 65 germplasm resources (16.75% of total), including 51 KSR and 14 non-KSR samples (S6 Table). Five strategies (M-HS, SANA, A-NE, E-NE, and CV) were evaluated for a core collection construction of 388 glutinous rice resources. Since single genetic parameters cannot fully represent diversity, comprehensive evaluation using N_e , H_e and other indices is essential [52,53].

Individual strategies proved insufficient for constructing a core collection, either due to the significant variations in genetic diversity or excessive sampling proportions. In general, increasing sampling intensity improves allele coverage but simultaneously increases redundancy among accessions, thereby reducing the efficiency of the core collection for practical utilization [10,11,54]. For example, while the CV strategy at a 30% sampling generated representative subset, the proportion was relatively high given the lower genetic diversity of the entire collection. Previous studies have suggested that an optimal sampling proportion for core collections typically ranges between 5% and 20%, providing a practical balance between manageable sample size and effective preservation of genetic variability for breeding and utilization purposes [10,11]. Moreover, minimizing redundancy while retaining maximal genetic diversity has been emphasized as a central principle in core collection construction [20]. Although a larger subset may achieve higher allele coverage, this gain often comes at the cost of a substantial increase in collection size. In this study, while a higher sampling ratio (25.26%) resulted in improved allele coverage, Core set 6 (16.75%) retained more than 90% of the alleles while exhibiting superior or comparable performance across multiple genetic diversity parameters, including N_e and distance-based measures. This indicates a more favorable balance between representativeness and redundancy for a working core collection. The final core collection was therefore optimized through the integration of multiple strategies, retaining over 90% of alleles and showing genetic diversity indices (N_e , I , H_o , H_e , and PIC) comparable to those of the entire germplasm set. This finding is consistent with previous reports demonstrating that combining complementary strategies enhances both efficiency and accuracy in core collection construction [19]. Overall, the optimized core collection of 65 of glutinous rice accessions effectively captured the genetic diversity of the 'He' cultivation zone. It offers a valuable foundation for future breeding efforts, enabling genetic improvement with minimal but representative germplasm resources. The core collection can support pre-breeding applications, such as prioritized multi-environment phenotyping and the selection of genetically diverse parents. Recent studies have identified agronomically relevant loci/genes and resistance-related haplotypes in KSR [2,46]. These findings underscore KSR's potential for breeding climate-resilient and high-yield rice varieties. However, trait-marker relationships are often population- and environment-dependent. Therefore, systematic phenotypic evaluation and validation of candidate alleles/markers in the core collection will be an essential next step.

Conclusion

In this study, we analyzed the genetic diversity and population structure of 388 glutinous rice germplasms from the 'He' cultivation zone using SSR sequencing data from 37 SSR loci. The result revealed an uneven distribution of genetic diversity between KSR and non-KSR landraces, with KSR accessions being more genetically similar and numerically predominant, which strongly shaped the diversity pattern at the whole-collection level. Based on these findings, a core collection was constructed using optimized sampling. The SANA method, applied at a 16.75% sampling intensity

yielded a core collection of 65 samples. This core collection retained a high level of genetic diversity, with retention rate 90.86% (N_a), 117.61% (N_e), 116.10% (I), 102.63% (H_o), 118.42% (H_e), and 117.32% (PIC). No significant differences were observed between the core set and the entire collection, confirming the effectiveness of the strategy. Principal coordinate analysis also validated its representativeness, showing an even distribution across the entire collection. These results demonstrated the feasibility of constructing a compact yet representative core collection through integrated strategies. This provides a strong foundation for future research and targeted utilization of glutinous rice genetic resources in Guizhou.

Supporting information

S1 Table. Information of 388 glutinous rice landraces.

(XLSX)

S2 Table. Information of 37 SSR markers used in SSR sequencing.

(XLSX)

S3 Table. Genetic diversity analysis of 325 KSR germplasms.

(XLSX)

S4 Table. Genetic diversity analysis of 63 non-KSR germplasms.

(XLSX)

S5 Table. The Q-values of 388 germplasms in two subpopulations.

(XLSX)

S6 Table. Composition of the final core collection.

(XLSX)

S7 Table. Genetic differentiation parameters between KSR and non-KSR.

(XLSX)

S8 Table. Molecular varisance analysis (AMOVA) between KSR and non-KSR.

(XLSX)

Acknowledgments

We are grateful to the Academy of Agricultural Sciences, Qiandongnan Miao and Dong Autonomous Prefecture, Guizhou, China for providing the glutinous rice landraces used in this study.

Author contributions

Conceptualization: Lijie Zhou, Quanzhi Zhao.

Data curation: Wenhui Yang, Mingyi Mao.

Formal analysis: Wenhui Yang, Lijie Zhou.

Funding acquisition: Lijie Zhou.

Investigation: Wenhui Yang, Mingyi Mao.

Methodology: Lijie Zhou, Quanzhi Zhao.

Project administration: Lijie Zhou, Quanzhi Zhao.

Resources: Zongdong Pan.

Software: Wenhui Yang, Lijie Zhou.

Supervision: Lijie Zhou, Quanzhi Zhao.

Validation: Wenhui Yang, Mingyi Mao, Lijie Zhou, Quanzhi Zhao.

Visualization: Wenhui Yang, Lijie Zhou.

Writing – original draft: Wenhui Yang, Lijie Zhou.

Writing – review & editing: Jianquan Qin, Jamal Nasar, Quanzhi Zhao.

References

1. You XL, Zeng XS. History of rice culture in China. Shanghai: Shanghai People's Publishing House. 2010.
2. Liu C, Wang T, Chen H, Ma X, Jiao C, Cui D, et al. Genomic footprints of Kam Sweet Rice domestication indicate possible migration routes of the Dong people in China and provide resources for future rice breeding. *Mol Plant*. 2023;16(2):415–31. <https://doi.org/10.1016/j.molp.2022.12.020> PMID: [36578210](https://pubmed.ncbi.nlm.nih.gov/36578210/)
3. Wang Y, Jiao A, Chen H, Ma X, Cui D, Han B, et al. Status and factors influencing on-farm conservation of Kam Sweet Rice (*Oryza sativa* L.) genetic resources in southeast Guizhou Province, China. *J Ethnobiol Ethnomed*. 2018;14(1):76. <https://doi.org/10.1186/s13002-018-0256-1> PMID: [30497534](https://pubmed.ncbi.nlm.nih.gov/30497534/)
4. Lai JH. Preliminary study on the 'He' resources in Congjiang. *Crops Germplasm Resources*. 1988;3:11–2. <https://doi.org/10.19462/j.cnki.1671-895x.1988.03.004>
5. Stone R. Intellectual property. Chinese province crafts pioneering law to thwart biopiracy. *Science*. 2008;320(5877):732–3. <https://doi.org/10.1126/science.320.5877.732> PMID: [18467564](https://pubmed.ncbi.nlm.nih.gov/18467564/)
6. Chaudhary RC, Tran DV, Duffy R. *Speciality Rices of the World: Breeding, Production, and Marketing*. Rome: Science Publishers. 2001.
7. Ma HL. The production problem of He and He area. Guiyang: Department of Science, Technology and Education, Guizhou Provincial Department of Agriculture. 1979.
8. Khoury CK, Brush S, Costich DE, Curry HA, de Haan S, Engels JMM, et al. Crop genetic erosion: understanding and responding to loss of crop diversity. *New Phytol*. 2022;233(1):84–118. <https://doi.org/10.1111/nph.17733> PMID: [34515358](https://pubmed.ncbi.nlm.nih.gov/34515358/)
9. Frankel OH. Genetic perspectives of germplasm conservation. In: Arber W, Llimensee K, Peacock WJ, Starlinger P. Genetic manipulation: impact on man and society. Cambridge: Cambridge University Press. 1984. 161–70.
10. Brown AHD. Core collections: a practical approach to genetic resources management. *Genome*. 1989;31(2):818–24. <https://doi.org/10.1139/g89-144>
11. Van Hintum TJL, Brown AHD, Spillane C, Hodgkin T. Core collections of plant genetic resources. Rome, Italy: International Plant Genetic Resources Institute. 2000.
12. Yonezawa K, Nomura T, Morishima H. Sampling strategies for use in stratified germplasm collections. In: Hodgkin T, Brown AHD, van Hintum TJL, Morales EAV. Core collections of plant genetic resources. West Sussex: IPGRI Wiley. 1995. 35–53.
13. Li ZC, Zhang HL, Cao YS, Qiu ZE, Wei XH, Tang SX, et al. Studies on the sampling strategy for primary rice core collection of Chinese ingenious rice. *Acta Agronomica Sinica*. 2003;(1):20–4. <https://doi.org/10.3321/j.issn:0496-3490.2003.01.004>
14. Schoen DJ, Brown AH. Conservation of allelic richness in wild crop relatives is aided by assessment of genetic markers. *Proc Natl Acad Sci U S A*. 1993;90(22):10623–7. <https://doi.org/10.1073/pnas.90.22.10623> PMID: [8248153](https://pubmed.ncbi.nlm.nih.gov/8248153/)
15. De Beukelaer H, Davenport GF, Fack V. Core Hunter 3: flexible core subset selection. *BMC Bioinformatics*. 2018;19(1):203. <https://doi.org/10.1186/s12859-018-2209-z> PMID: [29855322](https://pubmed.ncbi.nlm.nih.gov/29855322/)
16. Liu K, Muse SV. PowerMarker: an integrated analysis environment for genetic marker analysis. *Bioinformatics*. 2005;21(9):2128–9. <https://doi.org/10.1093/bioinformatics/bti282> PMID: [15705655](https://pubmed.ncbi.nlm.nih.gov/15705655/)
17. Odong TL, van Heerwaarden J, Jansen J, van Hintum TJL, van Eeuwijk FA. Statistical Techniques for Defining Reference Sets of Accessions and Microsatellite Markers. *Crop Science*. 2011;51(6):2401–11. <https://doi.org/10.2135/cropsci2011.02.0095>
18. Boccacci P, Aramini M, Ordidge M, van Hintum TJL, Marinoni DT, Valentini N, et al. Comparison of selection methods for the establishment of a core collection using SSR markers for hazelnut (*Corylus avellana* L.) accessions from European germplasm repositories. *Tree Genetics & Genomes*. 2021;17(6). <https://doi.org/10.1007/s11295-021-01526-7>
19. Wu H, Duan A, Wang X, Chen Z, Zhang X, He G, et al. Construction of a Core Collection of Germplasms from Chinese Fir Seed Orchards. *Forests*. 2023;14(2):305. <https://doi.org/10.3390/f14020305>
20. Gu R, Fan S, Wei S, Li J, Zheng S, Liu G. Developments on Core Collections of Plant Genetic Resources: Do We Know Enough? *Forests*. 2023;14(5):926. <https://doi.org/10.3390/f14050926>

21. Salgotra RK, Chauhan BS. Genetic Diversity, Conservation, and Utilization of Plant Genetic Resources. *Genes (Basel)*. 2023;14(1):174. <https://doi.org/10.3390/genes14010174> PMID: 36672915
22. Merritt BJ, Culley TM, Avanesyan A, Stokes R, Brzyski J. An empirical review: Characteristics of plant microsatellite markers that confer higher levels of genetic variation. *Appl Plant Sci*. 2015;3(8):apps.1500025. <https://doi.org/10.3732/apps.1500025> PMID: 26312192
23. Varshney RK, Graner A, Sorrells ME. Genic microsatellite markers in plants: features and applications. *Trends Biotechnol*. 2005;23(1):48–55. <https://doi.org/10.1016/j.tibtech.2004.11.005> PMID: 15629858
24. Barthe S, Gugerli F, Barkley NA, Maggia L, Cardi C, Scotti I. Always look on both sides: phylogenetic information conveyed by simple sequence repeat allele sequences. *PLoS One*. 2012;7(7):e40699. <https://doi.org/10.1371/journal.pone.0040699> PMID: 22808236
25. Šarhanová P, Pfanzelt S, Brandt R, Himmelbach A, Blattner FR. SSR-seq: Genotyping of microsatellites using next-generation sequencing reveals higher level of polymorphism as compared to traditional fragment size scoring. *Ecol Evol*. 2018;8(22):10817–33. <https://doi.org/10.1002/ece3.4533> PMID: 30519409
26. Yang J, Zhang J, Han R, Zhang F, Mao A, Luo J, et al. Target SSR-Seq: A Novel SSR Genotyping Technology Associate With Perfect SSRs in Genetic Analysis of Cucumber Varieties. *Front Plant Sci*. 2019;10:531. <https://doi.org/10.3389/fpls.2019.00531> PMID: 31105728
27. Cui X, Li C, Qin S, Huang Z, Gan B, Jiang Z, et al. High-throughput sequencing-based microsatellite genotyping for polyploids to resolve allele dosage uncertainty and improve analyses of genetic diversity, structure and differentiation: A case study of the hexaploid *Camellia oleifera*. *Mol Ecol Resour*. 2022;22(1):199–211. <https://doi.org/10.1111/1755-0998.13469> PMID: 34260828
28. Li X, Wang J, Qiu Y, Wang H, Wang P, Zhang X, et al. SSR-Sequencing Reveals the Inter- and Intraspecific Genetic Variation and Phylogenetic Relationships among an Extensive Collection of Radish (*Raphanus*) Germplasm Resources. *Biology (Basel)*. 2021;10(12):1250. <https://doi.org/10.3390/biology10121250> PMID: 34943165
29. McCouch SR, Teytelman L, Xu Y, Lobos KB, Clare K, Walton M, et al. Development and mapping of 2240 new SSR markers for rice (*Oryza sativa* L.). *DNA Res*. 2002;9(6):199–207. <https://doi.org/10.1093/dnares/9.6.199> PMID: 12597276
30. International Rice Genome Sequencing Project. The map-based sequence of the rice genome. *Nature*. 2005;436(7052):793–800. <https://doi.org/10.1038/nature03895> PMID: 16100779
31. Kim K-W, Chung H-K, Cho G-T, Ma K-H, Chandrabalan D, Gwag J-G, et al. PowerCore: a program applying the advanced M strategy with a heuristic search for establishing core sets. *Bioinformatics*. 2007;23(16):2155–62. <https://doi.org/10.1093/bioinformatics/btm313> PMID: 17586551
32. Odong TL, Jansen J, van Eeuwijk FA, van Hintum T.J.L. Quality of core collections for effective utilisation of genetic resources review, discussion and interpretation. *Theor Appl Genet*. 2013;126(2):289–305. <https://doi.org/10.1007/s00122-012-1971-y> PMID: 22983567
33. Thachuk C, Crossa J, Franco J, Dreisigacker S, Warburton M, Davenport GF. Core Hunter: an algorithm for sampling genetic resources based on multiple genetic measures. *BMC Bioinformatics*. 2009;10:243. <https://doi.org/10.1186/1471-2105-10-243> PMID: 19660135
34. Peakall R, Smouse PE. GenAEx 6.5: genetic analysis in Excel. Population genetic software for teaching and research—an update. *Bioinformatics*. 2012;28(19):2537–9. <https://doi.org/10.1093/bioinformatics/bts460> PMID: 22820204
35. Yeh FC, Yang R, Boyle T. POPGENE Version 1.32 Microsoft Windows-based freeware for population genetic analysis. Edmonton: University of Alberta. 1999.
36. Excoffier L, Lischer HEL. Arlequin suite ver 3.5: a new series of programs to perform population genetics analyses under Linux and Windows. *Mol Ecol Resour*. 2010;10(3):564–7. <https://doi.org/10.1111/j.1755-0998.2010.02847.x> PMID: 21565059
37. Subramanian B, Gao S, Lercher MJ, Hu S, Chen W-H. Evolview v3: a webserver for visualization, annotation, and management of phylogenetic trees. *Nucleic Acids Res*. 2019;47(W1):W270–5. <https://doi.org/10.1093/nar/gkz357> PMID: 31114888
38. Pritchard JK, Stephens M, Donnelly P. Inference of population structure using multilocus genotype data. *Genetics*. 2000;155(2):945–59. <https://doi.org/10.1093/genetics/155.2.945> PMID: 10835412
39. Earl DA, vonHoldt BM. STRUCTURE HARVESTER: a website and program for visualizing STRUCTURE output and implementing the Evanno method. *Conservation Genet Resour*. 2011;4(2):359–61. <https://doi.org/10.1007/s12686-011-9548-7>
40. Jakobsson M, Rosenberg NA. CLUMPP: a cluster matching and permutation program for dealing with label switching and multimodality in analysis of population structure. *Bioinformatics*. 2007;23(14):1801–6. <https://doi.org/10.1093/bioinformatics/btm233> PMID: 17485429
41. Li S, Yang D, Zhu Y. Characterization and Use of Male Sterility in Hybrid Rice Breeding. *JIPB*. 2007;49(6):791–804. <https://doi.org/10.1111/j.1744-7909.2007.00513.x>
42. Spielmeier W, Ellis MH, Chandler PM. Semidwarf (sd-1), “green revolution” rice, contains a defective gibberellin 20-oxidase gene. *Proc Natl Acad Sci U S A*. 2002;99(13):9043–8. <https://doi.org/10.1073/pnas.132266399> PMID: 12077303
43. Han P, Tian X, Wang Y, Huang C, Ma Y, Zhou X, et al. Construction of a core germplasm bank of upland cotton (*Gossypium hirsutum* L.) based on phenotype, genotype and favorable alleles. *Genet Resour Crop Evol*. 2022;69(7):2399–411. <https://doi.org/10.1007/s10722-022-01379-6>
44. Liu YH, Chen JX, Zheng GY. Current Status and Development Strategies for Kam Sweet Rice Industry in Southeast Guizhou, China. *Tillage and Cultivation*. 2011;(06):6-7 9. <https://doi.org/10.13605/j.cnki.52-1065/s.2011.06.004>
45. Li S, Zhang SS, Dong YW, Shi DL, Shi YD. Genetic diversity analysis of rice varieties in Tengchong, Yunnan based on SSR markers. *Crops*. 2019;(05):15–21. <https://doi.org/10.16035/j.issn.1001-7283.2019.05.003>

46. Liu C, Cui D, Jiao A, Ma X, Li X, Han B, et al. Kam Sweet Rice (*Oryza sativa* L.) Is a Special Ecotypic Rice in Southeast Guizhou, China as Revealed by Genetic Diversity Analysis. *Front Plant Sci.* 2022;13:830556. <https://doi.org/10.3389/fpls.2022.830556> PMID: [35330871](https://pubmed.ncbi.nlm.nih.gov/35330871/)
47. Wang X, Wang W, Tai S, Li M, Gao Q, Hu Z, et al. Selective and comparative genome architecture of Asian cultivated rice (*Oryza sativa* L.) attributed to domestication and modern breeding. *J Adv Res.* 2022;42:1–16. <https://doi.org/10.1016/j.jare.2022.08.004> PMID: [35988902](https://pubmed.ncbi.nlm.nih.gov/35988902/)
48. Liu C, Wang Y, Jiao A, Ma X, Cui D, Li X, et al. Effects of Traditional Ethnic Minority Food Culture on Genetic Diversity in Rice Landraces in Guizhou Province, China. *Agronomy.* 2022;12(10):2308. <https://doi.org/10.3390/agronomy12102308>
49. Lu B-R. Diversity of rice genetic resources and its utilization and conservation. *Biodiversity Science.* 1998;06(1):63–72. <https://doi.org/10.17520/biods.1998011>
50. Deu M, Sagnard F, Chantereau J, Calatayud C, Vigouroux Y, Pham JL, et al. Spatio-temporal dynamics of genetic diversity in *Sorghum bicolor* in Niger. *Theor Appl Genet.* 2010;120(7):1301–13. <https://doi.org/10.1007/s00122-009-1257-1> PMID: [20062963](https://pubmed.ncbi.nlm.nih.gov/20062963/)
51. Yang SL, Yang SH, Li T, Guan YW, Yang WH, Pan ZD. Genetic diversity analysis and core germplasm construction of Guizhou Kam sweet rice based on SNP chip. *Seed.* 2024;43(07):9–16. <https://doi.org/10.16590/j.cnki.1001-4705.2024.07.009>
52. Pan H, Deng L, Zhu K, Shi D, Wang F, Cui G. Evaluation of genetic diversity and population structure of *Annamocarya sinensis* using SCoT markers. *PLoS One.* 2024;19(9):e0309283. <https://doi.org/10.1371/journal.pone.0309283> PMID: [39231174](https://pubmed.ncbi.nlm.nih.gov/39231174/)
53. Li X, Cui L, Zhang L, Huang Y, Zhang S, Chen W, et al. Genetic Diversity Analysis and Core Germplasm Collection Construction of Radish Cultivars Based on Structure Variation Markers. *Int J Mol Sci.* 2023;24(3):2554. <https://doi.org/10.3390/ijms24032554> PMID: [36768875](https://pubmed.ncbi.nlm.nih.gov/36768875/)
54. Fu Y, Li S, Ma B, Liu C, Qi Y, Pang C. Genetic Diversity Analysis and Core Collection Construction of Ancient *Sophora japonica* L. Using SSR Markers. *Int J Mol Sci.* 2024;25(23):12776. <https://doi.org/10.3390/ijms252312776> PMID: [39684487](https://pubmed.ncbi.nlm.nih.gov/39684487/)