

SINE jumping contributes to large-scale polymorphisms in the pig genomes

Cai Chen

College of Animal Science and Technology, Yangzhou University

Enrico D'Alessandro

Department of Veterinary Sciences, University of Messina

Eduard Murani

Leibniz Institute for Farm Animal Biology, Dummerstorf

Yao Zheng

College of Animal Science and Technology, Yangzhou University

Domenico Giosa

Department of Clinical and Experimental Medicine, University Hospital of Messina

Naisu Yang

College of Animal Science and Technology, Yangzhou University

Xiaoyan Wang

College of Animal Science and Technology, Yangzhou University

Bo Gao

College of Animal Science and Technology, Yangzhou University

Kui Li

Institute of Animal Science, Chinese Academy of Agriculture Science

Klaus Wimmers

Leibniz institute for Farm Animal Biology, Dummerstorf

Chengyi Song (✉ cysong@yzu.edu.cn)

Yangzhou University <https://orcid.org/0000-0002-0488-4718>

Research

Keywords: retrotransposon, insertion polymorphism, RIP, SINE, pig, molecular marker

Posted Date: April 5th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-352249/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Abstract

Background: Molecular markers based on retrotransposon insertion polymorphisms (RIPs) have been developed and are widely used in plants and animals. Short interspersed nuclear elements (SINEs) exert wide impacts on gene activity and even on phenotypes. However, SINE RIP profiles in livestock remain largely unknown, and not be revealed in pigs.

Results: Our data revealed that SINEA1 displayed the most polymorphic insertions (22.5% intragenic and 26.5% intergenic), followed by SINEA2 (10.5% intragenic and 9% intergenic) and SINEA3 (12.5% intragenic and 5.0% intergenic). We developed a genome-wide SINE RIP mining protocol and obtained a large number of SINE RIPs (36,284), with over 80% accuracy and an even distribution in chromosomes (14.5/Mb), and 74.34% of SINE RIPs generated by SINEA1 element. Over 65% of pig SINE RIPs overlap with genes, with significant enrichment in the first and second introns of protein-coding and long non-coding RNA genes. Nearly half of the RIPs are common in these pig breeds. Sixteen SINE RIPs were applied for population genetic analysis in 23 pig breeds, the phylogeny tree and cluster analysis were generally consistent with the geographical distributions of native pig breeds in China.

Conclusions: Our analysis revealed that SINEA1–3 elements, particularly SINEA1, are high polymorphic across different pig breeds, and generate large-scale structural variations in the pig genomes. And over 35,000 SINE RIP markers were obtained. These data indicate that young SINE elements play important roles in creating new genetic variations and shaping the evolution of pig genome, and also provide strong evidences to support the great potential of SINE RIPs as genetic markers, which can be used for population genetic analysis and quantitative trait locus (QTL) mapping in pig.

Background

Retrotransposons—a heterogeneous group of genetic sequences that have the ability to be transcribed into RNA, reverse-transcribed into DNA, and inserted into a new site in a genome—account for 30% to 50% of mammalian genomes and thus represent major genomic parasites of mammals [1]. Accordingly, they play key roles in the structural organization of the genome, in the orchestration of biological processes, and even in the diversity and evolution of species. Retrotransposons are classified into three main groups: long interspersed nuclear elements (LINEs), short interspersed nuclear elements (SINEs), and long terminal repeats (LTRs), including endogenous retroviruses (ERVs) [2]. In general, the retrotransposon landscape of mammal genomes is dominated by LINEs and SINEs, followed by LTR retrotransposons. In most mammals, the predominant LINE is L1, which has a length of 6–8 kb, with a structure comprising of a 5′ untranslated region (UTR), two open reading frames (ORFs), a nuclear chaperone protein (ORF1), and a reverse transcriptase (ORF2). A very short third protein (ORF0) has recently been described in primate L1 elements but its function is still unknown [3]. LINEs represent ~10% to 20% of mammalian genomes, while LTR retrotransposons, with a length of 7–9 kb and a structure comprising two LTRs, GAG, POL, and ENV, are present in mammal genomes with a moderately high copy number and can occupy between ~4% and 10% of the genome [1].

Although SINEs, which are transcribed by RNA polymerase III, only occupy up to ~10% of mammalian genomes, they display an extremely high occurrence rate and/or high copy number in genomes because they usually appear as short fragments 150–300 bp long. SINEs typically have three parts: a 5' head, a body, and a 3' tail; they are non-autonomous retrotransposons that retrotranspose by hijacking the reverse transcriptases (RTs) and endonucleases of their partner LINEs [4]. Because LINEs and LTRs are large fragments (7–9 kb) and are believed to have a greater ability to disrupt genes and genomes than the shorter SINEs (150–300 bp), they are evolutionarily purged from genomes at a greater rate. Thus, SINEs are believed to be more tolerable for hosts, can co-evolve with host genomes, and can exert a wider impact on the shaping of genes and on genome evolution [5].

SINEs have been found to insert frequently in gene regions, suggesting that they might play important roles in regulating gene activity. Approximately 38% of SINE insertions overlap with transcribed regions in wheat, and 30% of SINE insertions overlap with genes in Solanaceae [6,7]; 65.69% of the transposons found in introns were SINEs in bovine genomes [8]. Around 85%–90% of mouse and human protein-coding genes contain transposon sequences in their introns [9], while in pigs, nearly 50% of retrotransposons are inserted into over 80% of protein-coding and long non-coding RNA (lncRNA) genes, with SINEs representing the highest insertion frequency compared with LINE and LTR retrotransposons [10]. It has been suggested that SINEs can shape gene and genome evolution by offering exons, splicing sites, and start and stop codons, thus creating novel genes [11,12]. SINE insertions can play roles in gene regulation by diverse mechanisms: by acting on the promoters, enhancers [13], or transcription factor binding sites [14] of corresponding genes. SINE transposons can also contribute to epigenetic regulation; in fact, SINEs possess a high GC content, which makes them hot spots for DNA methylation, a well-known mechanism related to transcriptional repression [15–17]. Furthermore, SINEs can activate miRNAs by acting as promoters for miRNA synthesis, or as miRNA-binding sites in target mRNAs [18–20], or by regulating gene expression from SINE transcripts [21,22]. When SINEs accumulate in 3' UTRs, they influence mRNA degradation by Staufen-mediated mRNA decay [23]. When a SINE is inserted into a lncRNA, it can promote translation of partially overlapping sense protein-coding mRNAs (designated a SINEUP), leading to regulate the expression of the target gene [24,25].

Recently, retrotransposon-based markers, such as retrotransposon-based insertion polymorphisms (RBIPs) [26], inter-retrotransposon amplified polymorphisms (IRAPs), and retrotransposon–microsatellite amplified polymorphisms (REMAPs) [27], have been developed and widely used in studies of genetic diversity, phylogeny, genetic mapping, and cultivar identification in plants [26–29]. It is commonly accepted that retrotransposon insertion polymorphism (RIP) markers have high prevalence in genomes and are more informative and polymorphic compared with other marker systems [26,28,30,31]. RIP markers have also been developed for several domesticated animals, including sheep, deer, and chicken. A number of related endogenous retroviruses (ERVs) have been used as genetic markers to study the history of domestic sheep and have provided valuable insights into their history [32]. The ERV insertional polymorphism data in cervids have also been shown to have greater resolution than microsatellites in the detection of geographic clustering of related deer [33,34]. GGERV10, which is the youngest ERV family in chicken, might have contributed to recent genomic variations in different genetic populations and has

been suggested as a molecular marker for chicken breed identification [29]. Our previous study on pig mobilome annotation revealed that most (80%) of the protein-coding and lncRNA genes contain retrotransposon insertions in pig genomes, and retrotransposons tend to be enriched in lncRNAs, with nearly half of protein-coding genes generating chimeric transcripts with retrotransposons [10]. Furthermore, this indicated that SINEs are the most widespread retrotransposons in the pig genome, accounting for about 11% with over 1 million copies [10]. These data indicate that SINE RIP markers may be important tools for studying biodiversity and genetics, and even for molecular breeding in domestic animals. In particular, several SINE insertions causing phenotype changes have been reported in pigs, horses, and dogs [35–40]. SINEA1–3 represents the youngest subfamily of pig-specific SINEs and reflects the most recent expansion activity during the last 10 million years [10]. The ‘copy and paste’ life cycle of replicative transposition produces new genome insertions without excising the original element. Thus, the movement and accumulation of SINEs has been a major force shaping the genes and genomes of most organisms. The retrotransposon junction can be mapped using high-throughput next-generation sequencing and retrieved by bioinformatics analysis. However, genome-wide SINE RIPs are rarely reported in the study of genetics and breeding of livestock, including pig.

Here, we developed a genome-wide SINE RIP mining protocol and performed genome-level screening by using the assembled pig genomes deposited in the NCBI database (see below); the resulting RIPs were further verified by polymerase chain reaction (PCR) amplification. We also evaluated the genomic coverage and breed distribution of these SINE RIPs, their insertion bias, and applied them for population genetic analysis and for evaluating domestication processes in Chinese pig breeds. We obtained a novel set of highly informative RIP markers, with a wide distribution (average 14.5 SINE RIPs per 1 Mb) and high coverage (36,284) in the pig genome, which display great potential as genetic markers for application in phylogeny and genetic diversity studies as well as in quantitative trait locus (QTL) mapping to benefit the conservation and utilization of local pig genetic resources and modern molecular breeding.

Results

Young SINE retrotransposon insertions are highly polymorphic in the pig genomes

Three pig-specific SINE families (SINEA, SINEB, and SINEC), with different evolutionary histories, were identified in a previous study showing that SINEA represents the youngest family with some of its subfamilies still displaying activity in the last 10 million years [10]. Eleven subfamilies of SINEA (A1–A11) were identified previously, and they display high sequence similarity, but with minor differences: SINEA1–SINEA3 have six specific nucleotides, SINEA1 and SINEA2 have two specific nucleotides, while SINEA1 contains the longest polyA sequence (Additional file 1: Fig. S1), which is unique and different from other subfamilies and might act in their transposition activities. Insertion age analysis revealed that SINEA1–SINEA3 displayed activity 2 million years ago (Mya); the activity of SINEA4 was hard to detect in the last 2 million years, while the activity of other subfamilies (SINEA5–SINEA11), SINEB, and SINEC was

totally extinct in this period (Fig. 1A and Additional file 1: Fig. S2). Overall, SINEA1 showed dominant current activity (<2 Mya), followed by SINEA2, while SINEA3 exhibited very weak current activity, indicating that these subfamilies, particularly SINEA1, might still jump and contribute to genomic variations in pigs.

To investigate the jumping activity of these SINE elements, 1,400 SINE insertions distributed in the intragenic regions, and 1,400 in the intergenic regions in the reference genome from seven SINE subfamilies (SINEA1–SINEA4, SINEB2, SINEB6, and SINEC4), representing different insertion ages, were selected randomly for polymorphism prediction by local BLAST searching as described in the methodology. The predicted polymorphic ratio varied significantly across subfamilies, as expected. SINEA1 showed the highest polymorphic ratios at 22.50% and 26.50% in intragenic and intergenic regions, respectively. The SINEA2 and SINEA3 subfamilies showed polymorphism rates ranging from 5.00% to 12.50%, while other subfamilies displayed very low insertion polymorphism rates (<2%) (Fig. 1B and Additional file 1: Table S1). Furthermore, 25 predicted polymorphic and 25 non-polymorphic insertions between a non-reference (Meishan) genome and the reference (Duroc) genome were used to evaluate the accuracy of local BLAST searching, by PCR (Fig. 1C). The accuracies of finding polymorphic and non-polymorphic insertions were 88.00% (22/25) and 84.00% (21/25), respectively (Fig. 1D and Additional file 1: Table S2), indicating that the local *BLAST* protocol for SINE RIP prediction is highly reliable. These findings confirmed that SINEA1–SINEA3 are still active and can jump into the pig genome, and proved that SINEA1, the youngest element, is very active, and tends to generate highly polymorphic insertions.

Development of the genome-wide SINE RIP screening protocol

To identify SINE RIPs in all assembled pig genomes (15 non-reference and one reference) we developed a genome-wide SINE RIP mining protocol, summarized in Fig. 2A and described in detail in the methodology. A total of approximately 100,000 SINEA1–SINEA3 insertions in each genome were mapped by RepeatMasker. On average, more than 95% of these insertions in the non-reference genomes were mapped successfully to the reference genome. Based on the comparison of non-reference and reference genomic SINE insertion positions, we obtained 263,837 putative SINE RIPs from all genomes, which were submitted to local BLAST searching and checked manually for each RIP (Additional file 1: Table S3). The ambiguous SINE RIPs were discarded based on their alignment patterns (Fig. 2A), and 94,074 SINE RIPs remained for further analysis (Additional file 1: Table S3).

Because the assembly levels of non-reference genomes were lower than the reference genome, the gaps in the non-reference genomes could result in a false positive estimation for the SINE RIP deletion allele. Therefore, we discarded those predicted SINE RIP deletion alleles that were detected only in one non-reference genome, and verified those present in two, three, and four non-reference genomes using PCR (Fig. 2B). As expected, we found a high rate of false positives when the SINE deletion alleles occurred

only in two or three non-reference genomes, with accuracies of SINE RIP prediction of only 32.14% and 37.50%, respectively, so these sites were removed from further analysis. However, the accuracy (81.25%) was significantly improved when SINE RIP deletions were detected in four non-reference genomes. The SINE RIP insertion alleles identified in one, two, 14, or 15 non-reference genomes were also verified by PCR, and all of them showed high accuracy (>80%) (Fig. 2C; Additional file 1: Table S4). These data indicate that the SINE deletion alleles identified in more than three non-reference genomes and all SINE RIP insertions (one or more non-reference genomes) were at least 80% accurate.

Large-scale RIPs generated by SINE jumping in the pig genomes

After removing the inaccurate and redundant RIPs, a final total of 36,284 SINE RIPs were obtained at the genome level (Table 1). Then, 230 SINE RIPs were selected randomly for PCR verification, and 185 RIPs were confirmed as positive, 30 RIPs were false positives, and 15 RIPs were uncertain (Fig. 3A), resulting in an accuracy of predicting SINE RIPs of >80% (Fig. 3A, Additional file 1: Table S5). Thus, our genome-wide SINE RIP screening protocol was reliable. Overall, 74.34%, 20.21%, and 5.45% SINE RIPs came from the SINEA1, SINEA2, and SINEA3 subfamilies, respectively, which generally corresponds to their age distributions in the genome (Fig. 3B). Furthermore, SINE RIPs were evenly distributed on each chromosome, with an average of 14.5 (range 11.28–21.63) SINE RIPs in each 1 Mb window (Fig. 3C, Additional file 1: Table S6). While chromosomes 10, 11, 12, 17, and 18 tended to be slightly enriched for SINE RIPs (>18 RIPs/Mb, Fig. 3C), which is generally consistent with the retrotransposon coverage on each chromosome (Fig. 3D), chromosomes 1, 13, and X showed a tendency to be slightly depleted of SINE RIPs (<13 RIPs/Mb, Fig. 3C). The Y chromosome was excluded from analysis because of its multiple repeats, which resulted in difficulties in sequencing and assembly, with too many gaps remaining.

Over 65% of SINE RIPs overlapping with genes

By calculating the genomic positions of each SINE RIP with the biogenic regions, 66.08% of the SINE RIPs (21,596/32,684) overlapped with the genic regions (NCBI annotated genes and NONCODE annotated lncRNA genes), which represent 23.09% of the total genes. In all, 51.36% of the SINE RIPs (16,787/32,684) were found to be overlapping with protein-coding genes, which account for 29.78% (6,154/20,666) of the total, and most of them (99.09%) are in introns (16,635/16,787). While 13.59% SINE RIPs (4,443/32,684) overlap with the lncRNA genes, which account for 17.30% (2,504/14,477) of the total lncRNA genes, most of them (96.89%) were found to be overlapping with introns (4,305/4,443) as well (Table 2). Furthermore, significant biases of SINE RIPs in the biogenic locations of lncRNA and protein-coding genes and their transcripts were observed. SINE RIPs tended to be enriched in the first and second introns of the protein-coding and lncRNA genes compared with other introns and their flanking sequences (Fig. 3E). In addition, a total of 260 SINE RIPs were identified in the exon regions of the protein-

coding genes. These SINE RIPs appear to be significantly enriched in the 3' UTRs (151/260) of mRNAs compared with 5' UTRs (98/260) and CDS (8/260) (Fig. 3F).

Nearly half of all SINE RIPs are common in pig genomes

For the 36,284 SINE RIPs, approximately 10,000 (6,612–12,703) of them appeared as insertion alleles, while the rest of them were identified as deletion alleles in each breed's genome (Fig. 4A). Deletion or insertion alleles of the predicted SINE RIPs detected in >12 or <4 breed genomes were designated as rare RIPs. In contrast, deletion or insertion SINE alleles present in 4–12 genomes were considered to be common RIPs. Based on this classification, we identified 16,694 common RIPs, representing 46.01% of all SINE RIPs identified (Table 1), resulting in highly polymorphic sequences in most breeds and with great potential for genetic analysis and QTL mapping. In addition, a pairwise comparison of SINE RIPs across the assembled genomes revealed that, on average, 11,482 differential alleles (range 7,532–14,751) were observed between genomes (Additional file 1: Fig. S3). Comparison across the commercial pig breed genomes (Duroc, Landrace, Large White, Pietrain, Hampshire, Berkshire) revealed that they exhibited relatively few alleles that differed between genomes, representing about 8,000 SINE RIP alleles, ranging from 7,817 between Berkshire and Hampshire pigs to 9,044 between Duroc and Hampshire (Fig. 4B). By contrast, the Chinese native pigs displayed more SINE RIP alleles that differed between breeds, with an average of 11,103 (range 9,721–12,622) (Fig. 4C). Comparison of the most important commercial pig breeds (Duroc, Landrace, and Large White) revealed that 23,189 RIP loci shared the same alleles, with each genome containing about 4,000 (range 4,051–4,793) breed-specific RIP alleles (Fig. 4D).

PCA and cluster analysis of the SINE RIPs

Cluster analysis showed the presence of two main groups of pig breeds, in fact, all Western pigs: Large White, Landrace, Duroc, Pietrain, Hampshire, Berkshire, Duroc, and the cross-breeds form a clade that is well separated from the one comprising all Chinese pigs, including Rongchang, Jianghua, Meishan, Bamei, Tibet, Bama, Wuzhishan, and Göttingen pigs which contained Asian pig genetic material (Fig. 5A). As expected, the SINE RIP-based clusters were also well supported by PCA (Fig. 5B), in which both clusters are separated horizontally in accord with the direction of maximal variance.

Analysis of the population structure and genetic diversity of some Chinese native pigs based on SINE RIP molecular markers

To evaluate the potential application of SINE RIPs in population genetic analysis, 16 SINE RIPs were selected to detect polymorphisms in 22 native Chinese pig breeds and in one native Italian pig breed. The

PCR analysis revealed that all the markers were polymorphic and biallelic. Detection of SINE RIPs in each breed and their primers are summarized in Additional file 1: Table S7 and Additional file 2.

The N_e statistic per locus ranged between 1.537 (REF-16266) and 2.000 (ESA1-16), with a mean across loci of 1.765. The expected heterozygosity was higher than the observed heterozygosity at most loci. Observed and expected heterozygosity values ranged from 0.166 (DR-68328) to 0.468 (REF-3992) and from 0.350 (REF-16266) to 0.500 (ESA1-16) with overall means of 0.354 ± 0.088 and 0.423 ± 0.055 , respectively. While the PIC values, which can reveal the usefulness of a marker in diversity analysis of a breed, are moderately informative for all 16 SINE RIPs (PIC 0.25–0.5), with an overall mean of 0.335 ± 0.031 , ranging from 0.288 to 0.375, the negative F_{IS} values (-0.106 ± 0.153), ranging from -0.315 to 0.328 , indicated a low value of inbreeding of each breed detected. The F_{ST} values ranged from 0.117 (REF-14902) to 0.369 (ESA1-33), with a mean F_{ST} value of 0.252 for all loci, indicating that 74.8% of the genetic variation was caused by differences between individuals and 25.2% arose from differentiation between breeds. Agreement with Hardy–Weinberg equilibrium was tested by loci within breeds at $P < 0.05$. For all loci combined, on average about one-third of the breed–loci combinations did not comply with Hardy–Weinberg equilibrium (Additional file 1: Table S7).

The UPGMA method was used to construct a phylogenetic tree (Fig. 6A) based on Nei's unbiased genetic distance. This clearly shows three clusters that generally correspond to their geographic locations (Fig. 6B), especially for southern Chinese breeds (Bamaxiang, Wuzhishan, Dahuabai, and Lantang) and most pig breeds of central China (Qingping, Hanjiang Black, Shaziling, Tongcheng, Lepinghua, Ningxiang, Erhualian, Laiwu Black, Dapulian, Dingyuan, and Mingguang Small Ear), with the exception of Bamei, Wei, and Anqinliubai. Bamei is a northern Chinese breed, but clustered with the southern Chinese pigs, while Wei and Anqinliubai were separated from their original geographical location (central China) and clustered with the northern Chinese pigs (Mashen and Dongbei Min) and the Italian pig breed (Nero Siciliano pig), which also has the highest genetic distance from Chinese pig breeds.

Discussion

SINE RIPs have great potential as genetic markers

Young retrotransposons, which are very recently evolved elements and still retain jumping activity, have been exploited widely in tagging for gene function annotation [41,42], and as molecular markers for evolution and population genetic studies [28,43] in plants and humans [44,45]. A comprehensive profile of genomic RIPs is critical for the development and application of molecular markers in evolutionary and population genetic studies. However, until now, genome-wide RIP profiles for only a few animal species have been well defined, such as for *Drosophila melanogaster*, *Caenorhabditis elegans*, and *Homo sapiens* [46–48]. Previous studies suggested that every 20 human genomes contain one new Alu insertion [49], which may be an underestimate, as 60,743 RIPs in rice and 16,192 RIPs in human genomes were identified by analysing re-sequencing data [50–53]. In mice, 695 polymorphic ERVs were identified by comparing the genomic sequences of four common mouse strains [54]. Wei and Kirkness [55] identified

at least 10,000 polymorphic SINEC_Cf loci in different dog breeds, and Sara *et al.* detected an average of 10,423 polymorphic loci. In all, the libraries identify 81,747 putative polymorphic SINEs from 62 dogs representing 59 breeds [56]. Our previous study revealed that the pig genome harbours multiple young retrotransposon families/subfamilies, and we have demonstrated that some of them can generate polymorphic insertions [10,36], suggesting that these retrotransposons, particularly SINE retrotransposons, which represent the most abundant and widest distribution in the pig genome [10,57], have great potential for the development of genetic markers.

Here, we first evaluated the evolutionary dynamics of different subfamilies of SINEA and compared polymorphism frequencies across these subfamilies. Our data suggest that SINEA1–3 are the youngest subfamilies, possessing a functional transposition activity in the last 2 million years. SINEA1 displayed the most polymorphic insertions (22.5% intragenic and 26.5% intergenic), followed by SINEA2 (10.5% intragenic and 9% intergenic) and SINEA3 (12.5% intragenic and 5.0% intergenic), which is consistent with the SINEC_Cf data in dogs, where an average 9% polymorphism rate [56] and 8% of polymorphic SINEC_Cfs was reported for the Boxer reference genome [58]. However, the current activities of SINEA4 and other subfamilies are very limited and have very low polymorphic insertions (<2%). These data indicate that SINEA1–3 elements, particularly SINEA1, are major mutators of the pig genome and play important roles in generating new variations in individuals, in population differentiation, and in genomic evolution. The contributions of SINEA1 to the formation of local pig and commercial breeds are also worthy of further evaluation.

We conducted large-scale SINE RIP mining in the pig genome by developing a comprehensive screening protocol using the 16 assembled pig genome sequences deposited in the NCBI database. By comparing the SINEA1–SINEA3 insertion differences across these genomes based on this protocol, we identified 36,284 SINE RIPs. The density of these (14.5 SINE RIPs/Mb) is similar to the single nucleotide polymorphisms (SNPs) represented on the widely used Illumina CAUPorcine 50 K SNP microarray. Furthermore, beside the high density of SINE RIPs, our data also show that these are evenly distributed in the pig genome, strongly supporting them as promising molecular markers for genetic analysis. Here, to improve the accuracy of SINE RIP prediction (>80%), we applied a strict standard for multiple key steps of the protocol as described in the methodology and confirmed the reliability of the prediction by PCR evaluation. In addition, a large number of pig breed genomes have been re-sequenced, and the sequences have been deposited in the NCBI database, so SINE RIP mining using the re-sequenced data is expected to increase the number of SINE RIPs significantly and is worth further exploration.

Most SINE RIPs might be involved in gene regulation

It is commonly accepted that retrotransposons contribute extensively to the diversification of gene function by shaping gene structure or by altering gene activity. Our previous study revealed that about 80% of protein-coding and lncRNA genes contain retrotransposon insertions in pigs [10], and similar annotations were also observed for the bovine, mouse, and human genomes [59–61]. Over 120 cases of

genetic diseases have been reported to be associated with retrotransposon insertions in humans [62]. Furthermore, retrotransposons can regulate gene expression by affecting chromatin structure, gene transcription, pre-mRNA processing, or aspects of mRNA metabolism (for a review see [63]). These data suggest that most retrotransposon insertions can alter the activities of nearby genes. Here, our data demonstrated again that most SINE RIPs (over 65%) overlap with genic regions, and 29.78% of protein-coding genes and 17.30% of lncRNA genes contain SINE RIPs. Furthermore, we found that SINE RIPs tend to be enriched in the first and second introns, where gene regulatory elements also tend to be enriched and play an important role in transcription [64], suggesting that most SINE RIPs might be involved in the regulation of gene activities. In contrast, SINE RIPs are significantly depleted in the exons of protein-coding and lncRNA genes: thus, only 260 and 157 were detected, respectively. Additionally, SINE RIPs appear to be significantly enriched in the 3' UTRs (151/260) of mRNAs, which is generally consistent with the insertion preferences of SINEs in the pig, mouse, and human transcripts [10,65,66]. These data indicate again that SINE RIP markers may have larger impacts on gene activities and higher application values in research on population genetics, QTL mapping, and molecular breeding than other types of genetic markers.

Application of SINE RIPs in population genetic analysis

DNA-based molecular markers such as microsatellites and SNPs are very powerful methods for distinguishing between animal genotypes and have been used extensively in the genetic analysis of pigs [67–73]. SNPs are usually biallelic as co-dominant markers, and less informative compared with that of highly polymorphic microsatellites, but this can be compensated for by employing large numbers of markers (e.g., SNP chips) or WGS [74,75]. Microsatellite markers are co-dominant, multi-allelic, highly polymorphic, relatively evenly spaced throughout genomes, and require low quality template DNA input (10–100 ng); but they are time-consuming and expensive to develop, and require technical expertise or fluorescently labelled primers for simple sequence repeats (SSR) analysis and high-resolution agarose or polyacrylamide gel separation [76–81]. By contrast, SINE RIPs are biallelic, co-dominant, highly polymorphic, give accurate and reproducible results, and exhibit high coverage and an even distribution among mammal genomes, suggesting great potential as genetic markers.

We applied 16 SINE RIPs in 23 pig breeds for population genetic analysis. As expected, the PIC and observed and expected heterozygosity values estimated by the SINE RIPs, as important parameters of genetic diversity, were lower than predictions based on microsatellite markers [67–73,82] but similar to estimates based on SNPs [79–81,83]. This is probably because microsatellite markers are multi-allelic, while SINE RIPs and SNPs are biallelic. In addition, the individual-level and population-level allele frequency of each type of genetic marker has an important impact on its application. The rare and low-frequency genetic variants are routinely excluded from genome-wide association studies (GWAS) studies because when an allele is present in a few individuals, the statistical analysis used to draw correlations between traits and alleles is not powerful enough to obtain significant results [84,85]. Genetic variants presenting only in very few populations/breeds also have significantly limited application value in animal

genetics and breeding. Here, by excluding the rare and low-frequency alleles of SINE RIPs (alleles present in >12 or <4 assembled genomes), we found that about 50% SINE RIPs are common, indicating that most of these RIP loci are polymorphic in these breeds, and applicable in population genetic analysis.

Mitochondrial DNA sequences, microsatellite markers, and SNP markers have been used to trace the domestication and origins of European and Asian domestic pigs [86–88]. Over the past decade, regions in China, including the Mekong River basin, the downstream region of the Yangtze River, the upper stream region of the Yangtze River, the Tibetan highlands and the lower region of the Yellow River [89–93] have been suggested as regions from which wild boar might have contributed to the domestic pig gene pool, and which may have represented independent centres for pig domestication. Here, heatmaps of the clusters related to breed comparison were also well supported by PCA with the whole array of SINE RIPs (Fig. 5A, B), consistent with the results of pig evolutionary research and geographical distribution. Clustering using 16 SINE RIPs in 23 breeds is generally consistent with the geographical distributions of Chinese pig breeds, However, a few breeds do not match completely with their geographical distributions. This discrepancy can be explained by gene flow between these regions or breeds; alternatively, these breeds might not originate locally but were imported historically, which is worth further study.

Conclusion

Our data suggest that SINEA1–3 are the youngest subfamilies in pig genome. SINEA1 displayed the most polymorphic insertions (22.5% intragenic and 26.5% intergenic), followed by SINEA2 (10.5% intragenic and 9% intergenic) and SINEA3 (12.5% intragenic and 5.0% intergenic), These data indicate that SINEA1–3 elements, particularly SINEA1, are major mutators of the pig genome and play important roles in generating new variations in individuals, in population differentiation, and in genomic evolution. Then we developed a genome-wide SINE RIP mining protocol to mining the young SINE insertion polymorphic sites and obtained a large number of SINE RIPs (36,284), with over 80% accuracy and an even distribution in chromosomes. Nearly half of the RIPs are common in these pig breeds. Over 65% of pig SINE RIPs overlap with genes, with significant enrichment in the first and second introns of protein-coding and long non-coding RNA genes. Sixteen SINE RIPs were successfully applied for population genetic analysis in 23 pig breeds. Our experiments have demonstrated the efficiency of the SINE RIP mining protocol and provide evidence to support their potential as genetic markers in pigs as well as in other livestock.

Methods

Assembled genomes and gene annotation files used

Sixteen assembled pig genomes: Duroc, Landrace, Large White, Pietrain, Berkshire, Hampshire, cross-bred (Large White ´ Landrace ´ Duroc), two lines of Göttingen minipigs (Göttingen minipig, Ellegaard Göttingen minipig), Wuzhishan, Tibetan, Rongchang, Meishan, Bamei, Bama, and Jinhua were used for genome-wide screening of SINE RIPs and were obtained from the NCBI whole-genome sequencing (WGS) database (<https://www.ncbi.nlm.nih.gov/assembly/>). These assembled genomes had an average

sequencing depth of 108.80'. The Duroc is the reference genome (susScr11.1) used for the pig, and the other 15 genomes were re-sequencing genomes obtained by next-generation sequencing technology, which are called non-reference genomes here. Five of them are commercial pig breeds (Duroc, Landrace, Large White, Pietrain, Berkshire, and Hampshire), seven of them (Wuzhishan, Tibetan, Rongchang, Meishan, Bamei, Bama, and Jinhua) are Chinese native pig breeds, and five of them (Göttingen, Ellegaard Göttingen minipig, Wuzhishan, Tibetan, and Bama) are miniature pigs. Detailed information about these genomes is shown in Additional file 1: Table S8.

The file on lncRNA gene annotation was downloaded from the NONCODE database (<http://www.noncode.org/download.php>). The Bed format file of lncRNA genes, which represents 17,811 such genes corresponding to Sscrofa10.2, were converted to Sscrofa11.1 by LiftOver (<http://genome.ucsc.edu/cgi-bin/hgLiftOver>), and finally, the coordinates of 14,477 lncRNA genes were obtained. The coordinates of protein-coding genes (20,674) and exons, the mRNAs (63,568) of protein-coding genes, and the 5' UTR, 3' UTR, and coding sequence (CDS) information of protein-coding genes were retrieved from the annotation of Sscrofa11.1 in the NCBI database (https://ftp.ncbi.nlm.nih.gov/genomes/all/annotation_releases/9823/106/).

Insertion age estimation and multiple alignments of SINEs

The reference genome (susScr11.1) was masked using RepeatMasker [94] (version 4.0.9, -nolow) with the custom repeat library [10]. Then, the diversity (K) value of each subfamily in SINEA was calculated using the calcDivergenceFromAlign.pl tool in the RepeatMasker program. The ages of SINE subfamilies were then calculated according to the formula $T = K / 2r$ ($r = 2.2 \times 10^{-9}$ substitutions/site/year) [95]. Multiple alignments were constructed from the reference sequences of the SINEA subfamilies using ClustalX2 [96] (version 2.0) with default parameters.

Genome-wide SINE RIP screening protocol

A protocol for the genome-wide screening for SINE RIPs based on the 16 assembled pig genomes was established in this study, of which the main process is shown in Fig. 2A and divided into four main steps.

Step 1. Screening SINE insertions in the genomes. The custom library (including all SINE subfamilies, DNA, LINE, and LTR repeats) was built in advance [10], and used to mask the 15 non-reference and reference genomes using RepeatMasker (-nolow, -lib custom library). Then, insertions masked by three young SINE subfamilies (SINEA1, SINEA2, and SINEA3) with a length of 100–330 bp and mask score > 1000 were kept for further analysis. The 200 bp upstream or downstream flanking sequences of these insertions were extracted using the bedtools [97] (version 2.27.1) *flank* and bedtools *getfasta* commands.

Step 2. Mapping to the reference genome. The flanking sequences of these SINE insertions in the non-reference genomes were mapped to the reference genome using Blat [98] (-minIdentity = 90, -minScore =

180). The mapping results were filtered by a length of 180–220 bp, and insertions with flanking sequences mapping to more than one genomic position were also removed. For insertions that failed to be mapped against the reference genome by the upstream 200 bp flanking sequence, the 200 bp downstream flanking sequences were mapped in the same way, then the results of these two sets were merged. Thereby, each insertion's information corresponding to the reference genome was obtained from each non-reference genome.

Step 3. Identification of SINE RIPs. The differential insertions, designated as putative SINE insertion polymorphisms between the non-reference and reference genomes, were obtained using a bedtools window (-w 50, -v). The SINE insertions from non-reference genomes that did not fall into the same window (SINE insertion site and 50 bp flanking region) as in the reference genome were considered to be putative SINE RIPs.

Step 4. Verification of SINE RIPs by local BLAST and PCR. The putative SINE RIPs were manually verified by local BLAST [99]. The sequences, including the 200 bp flanking sequences and the SINE sequence of each putative SINE RIP, were extracted using bedtools getfasta, and aligned using a local BLAST platform (blastn -task megablast -evalue 1.0e-5 -max_target_seqs 1 -max_hsp 1) between the non-reference and reference genomes. After the alignment, those putative SINE RIPs exhibiting the expected alignment patterns between genomes were kept for further analysis (Fig. 2). The SINE RIPs from all genomes were merged with bedtools merge (-s, -d 10) and redundancies were removed; 403 of the predicted SINE RIPs were selected for accuracy evaluation using PCR amplification.

PCR verification

Twelve domestic pig breeds (Large White, Landrace, Duroc, Meishan, Erhualian, Sujiang, Fengjing, Diannan small-ear, Wuzhishan, Bama, Tibetan and Nero Siciliano) were used for PCR verification of SINE RIP polymorphisms. The Sicilian black pigs were from Italy and other breeds were from China (Additional file 3). From each pig breed, three individual DNA samples were pooled. DNA was isolated from ear samples using MiniBEST Universal Genomic DNA Extraction kits (TaKaRa, Dalian, China). The primer pairs were designed for the up- and downstream flanking regions of RIPs and spanned the SINE insertions. PCR amplifications were carried out in a total volume of 20 µL, containing 40 ng of genomic DNA, 2 µL Taq Master Mix buffer (Vazyme, Nanjing, China) and 10 pmol of each primer. PCR amplifications were carried out using the following method: an initial denaturation at 94 °C for 3 min; 30 cycles at 94 °C for 30 s; 58 °C for 20 s; 72 °C for 30 s; and a final extension of 10 min at 72 °C. Finally, 7 µL of PCR products and 5 µL of DL2000 molecular weight markers were detected by electrophoresis using 1.0% agarose gels in 1 × TAE buffer with a constant voltage of 130 V for 30 min. Gels were stained with ethidium bromide and visualized with ultraviolet fluorescence.

Intersection analysis

The distribution of these SINE RIPs in the genome and their relationship with genes and biogenic regions were analysed. Only the overlapping sequences of SINE with gene or biogenic regions above 25 bp were considered for further analysis. Some SINE RIPs interacted with more than one biogenic region or gene, so were counted more than once.

Principal component analysis (PCA) and cluster analysis of the SINE RIPs

Based on the SINE RIPs identified in this study, the R statistics package (version 3.6.3) were used to generate a presence/absence matrix and performed the PCA analysis. On the same dataset, heatmaps and cluster analysis were computed by the use of the R package pheatmap tool (version 1.0.12) [100], using the “Euclidean” distance method for clustering.

Genetic diversity and population structure analysis

Sixteen SINE RIPs from 16 chromosomes (Additional file 2) and 585 individuals from 23 breeds (Additional file 3) were selected for genetic diversity and population structure analysis. PCR amplification and detection were performed as described above. The genetic parameters—allele/genotype frequency, effective allele number (N_e), observed heterozygosity (H_o), expected heterozygosity (H_e), Wright’s F-statistics (F_{IT} , F_{IS} , F_{ST}) and Nei’s genetic distance was analysed with Popgen32 (version 1.32)(https://sites.ualberta.ca/~fyeh/popgene_info.html) and polymorphic information content (PIC) was calculated as based on the PCR results. Finally, a phylogenetic tree was constructed using the unweighted pair group method with arithmetic mean (UPGMA) method with Popgen32.

Declarations

Ethics approval

All treatments and protocols involving animals in this study were strictly done in accordance with the guidelines of the Animal Experiment Ethics Committee of Yangzhou University (approval number: NSFC2020-dkxy-02).

Consent for publication

Not applicable

Availability of data and materials

All data needed to evaluate the conclusions in this paper are present either in the main text or the supplementary materials.

Competing interests

The authors declare that have no competing interests.

Funding

This research was funded by the National Natural Science Foundation of China [grant numbers 31872977 and 32002146]; the Independent Innovation Fund Project of Agricultural Science and Technology in Jiangsu Province [CX(19)2016], the Priority Academic Program Development of Jiangsu Higher Education Institutions, China Postdoctoral Science Foundation [2020M671630] to Cai Chen, and the High-end Talent Support Program of Yangzhou University to Chengyi Song.

Author Contributions

Chengyi Song, Klaus Wimmer, and Kui Li designed the experiments, and Cai Chen, Enrico D'Alessandro, Yao Zheng, Naisu Yang performed most of the experiments and analysed most of the results. Eduard Murani, and Domenico Giosa, Xiaoyan Wang and Bo Gao contributed to the experiments. Cai Chen, Enrico D'Alessandro, Klaus Wimmer, and Chengyi Song wrote the manuscript.

Acknowledgements

Not applicable

Supplementary Information

Additional file 1:

Additional file 1: Fig. S1. Alignment of the sequences of SINEA1-A11 subfamilies.

Additional file 1: Fig. S2. Insertion ages of SINEB and SINEC families.

Additional file 1: Fig. S3. Distribution of the differential SINE RIP alleles between each pair of genomes.

Additional file 1: Table S1. Predicted polymorphic ratio of SINE insertions from different subfamilies located in intergenic and intragenic regions.

Additional file 1: Table S2. Polymorphic ratio of randomly selected polymorphic and non-polymorphic SINE insertions following PCR verification.

Additional file 1: Table S3. Summary of the number of SINE insertions in the protocol used for annotating SINE RIPs.

Additional file 1: Table S4. Positive ratios of SINE RIPs obtained by PCR verification for rare SINE RIPs.

Additional file 1: Table S5. Positive ratios of the 36,284 SINE RIPs obtained by PCR verification with limited samples.

Additional file 1: Table S6. Density of SINE RIPs in each chromosome.

Additional file 1: Table S7. Characterization of 16 SINE RIPs analysed in 23 pig populations.

Additional file 1: Table S8. The pig genomes used for the SINE RIP screen protocol.

Additional file 2.xlsx

Additional file 3.xlsx

References

1. Platt, R.N.; Vandewege, M.W.; Ray, D.A. Mammalian transposable elements and their impacts on genome evolution. *Chromosom. Res.* **2018**, *26*, 25–43.
2. Eickbush, T.H.; Malik, H.S. Origins and Evolution of Retrotransposons. In *Mobile DNA II*; American Society of Microbiology, 2014; pp. 1111–1144.
3. Denli, A.M.; Narvaiza, I.; Kerman, B.E.; Pena, M.; Benner, C.; Marchetto, M.C.N.; Diedrich, J.K.; Aslanian, A.; Ma, J.; Moresco, J.J.; et al. Primate-Specific ORF0 Contributes to Retrotransposon-Mediated Diversity. *Cell* **2015**, *163*, 583–893.
4. Dewannieux, M.; Esnault, C.; Heidmann, T. LINE-mediated retrotransposition of marked Alu sequences. *Nat. Genet.* **2003**, *35*, 41–48.
5. Platt, R.N.; Vandewege, M.W.; Ray, D.A. Mammalian transposable elements and their impacts on genome evolution. *Chromosom. Res.* **2018**, *26*, 25–43.
6. Ben-David, S.; Yaakov, B.; Kashkush, K. Genome-wide analysis of short interspersed nuclear elements SINES revealed high sequence conservation, gene association and retrotranspositional activity in wheat. *Plant J.* **2013**, *76*, 201–210.
7. Seibt, K.M.; Wenke, T.; Muders, K.; Truberg, B.; Schmidt, T. Short interspersed nuclear elements (SINES) are abundant in Solanaceae and have a family-specific impact on gene structure and genome organization. *Plant J.* **2016**, *86*, 268–285.
8. Almeida, L.M.; Silva, I.T.; Silva, W.A.; Castro, J.P.; Riggs, P.K.; Carareto, C.M.; Amaral, M.E.J. The contribution of transposable elements to *Bos taurus* gene structure. *Gene* **2007**, *390*, 180–189.
9. Zhang, Y.; Romanish, M.T.; Mager, D.L. Distributions of Transposable Elements Reveal Hazardous Zones in Mammalian Introns. *PLoS Comput. Biol.* **2011**, *7*, e1002046.

10. Chen, C.; Wang, W.; Wang, X.; Shen, D.; Wang, S.; Wang, Y.; Gao, B.; Wimmers, K.; Mao, J.; Li, K.; et al. Retrotransposons evolution and impact on lncRNA and protein coding genes in pigs. *Mob. DNA* **2019**, *10*, 19.
11. Schwichtenberg, K.; Wenke, T.; Zakrzewski, F.; Seibt, K.M.; Minoche, A.; Dohm, J.C.; Weisshaar, B.; Himmelbauer, H.; Schmidt, T. Diversification, evolution and methylation of short interspersed nuclear element families in sugar beet and related Amaranthaceae species. *Plant J.* **2016**, *85*, 229–244.
12. Jordan, I.K.; Rogozin, I.B.; Glazko, G. V.; Koonin, E. V. Origin of a substantial fraction of human regulatory sequences from transposable elements. *Trends Genet.* 2003, *19*, 68–72.
13. Göke, J.; Ng, H.H. CTRL + INSERT: retrotransposons and their contribution to regulation and innovation of the transcriptome. *EMBO Rep.* **2016**, *17*, 1131–1144.
14. Sundaram, V.; Cheng, Y.; Ma, Z.; Li, D.; Xing, X.; Edge, P.; Snyder, M.P.; Wang, T. Widespread contribution of transposable elements to the innovation of gene regulatory networks. *Genome Res.* **2014**, *24*, 1963–1976.
15. Ichiyanagi, K. Epigenetic regulation of transcription and possible functions of mammalian short interspersed elements, SINEs. *Genes Genet. Syst.* **2013**, *88*, 19–29.
16. Estécio, M.R.H.; Gallegos, J.; Dekmezian, M.; Lu, Y.; Liang, S.; Issa, J.-P.J. SINE Retrotransposons Cause Epigenetic Reprogramming of Adjacent Gene Promoters. *Mol. Cancer Res.* **2012**, *10*, 1332–1342.
17. Fukuda, K.; Inoguchi, Y.; Ichiyanagi, K.; Ichiyanagi, T.; Go, Y.; Nagano, M.; Yanagawa, Y.; Takaesu, N.; Ohkawa, Y.; Imai, H.; et al. Evolution of the sperm methylome of primates is associated with retrotransposon insertions and genome instability. *Hum. Mol. Genet.* **2017**, *26*, 3508–3519.
18. Lehnert, S.; Van Loo, P.; Thilakarathne, P.J.; Marynen, P.; Verbeke, G.; Schuit, F.C. Evidence for Co-Evolution between Human MicroRNAs and Alu-Repeats. *PLoS One* **2009**, *4*, e4456.
19. Piriyaongsa, J.; Mariño-Ramírez, L.; Jordan, I.K. Origin and Evolution of Human microRNAs From Transposable Elements. *Genetics* **2007**, *176*, 1323–1337.
20. Kramerov, D.A.; Vassetzky, N.S. Origin and evolution of SINEs in eukaryotic genomes. *Heredity (Edinb).* 2011, *107*, 487–495.
21. Mariner, P.D.; Walters, R.D.; Espinoza, C.A.; Drullinger, L.F.; Wagner, S.D.; Kugel, J.F.; Goodrich, J.A. Human Alu RNA Is a Modular Transacting Repressor of mRNA Transcription during Heat Shock. *Mol. Cell* **2008**, *29*, 499–509.
22. Allen, T.A.; Von Kaenel, S.; Goodrich, J.A.; Kugel, J.F. The SINE-encoded mouse B2 RNA represses mRNA transcription in response to heat shock. *Nat. Struct. Mol. Biol.* **2004**, *11*, 816–821.
23. Lucas, B.A.; Lavi, E.; Shiue, L.; Cho, H.; Katzman, S.; Miyoshi, K.; Siomi, M.C.; Carmel, L.; Ares, M.; Maquat, L.E. Evidence for convergent evolution of SINE-directed Staufen-mediated mRNA decay. *Proc. Natl. Acad. Sci.* **2018**, *115*, 968–973.
24. Carrieri, C.; Cimatti, L.; Biagioli, M.; Beugnet, A.; Zucchelli, S.; Fedele, S.; Pesce, E.; Ferrer, I.; Collavin, L.; Santoro, C.; et al. Long non-coding antisense RNA controls Uchl1 translation through an embedded SINEB2 repeat. *Nature* **2012**, *491*, 454–457.

25. Bon, C.; Luffarelli, R.; Russo, R.; Fortuni, S.; Pierattini, B.; Santulli, C.; Fimiani, C.; Persichetti, F.; Cotella, D.; Mallamaci, A.; et al. SINEUP non-coding RNAs rescue defective frataxin expression and activity in a cellular model of Friedreich's Ataxia. *Nucleic Acids Res.* **2019**, *47*, 10728–10743.
26. Flavell, A.J.; Knox, M.R.; Pearce, S.R.; Ellis, T.H.N. Retrotransposon-based insertion polymorphisms (RBIP) for high throughput marker analysis. *Plant J.* **1998**, *16*, 643–650.
27. Kalendar, R.; Schulman, A.H. IRAP and REMAP for retrotransposon-based genotyping and fingerprinting. *Nat. Protoc.* **2006**, *1*, 2478–2484.
28. Kalendar, R.; Flavell, A.J.; Ellis, T.H.; Sjakste, T.; Moisy, C.; Schulman, A.H. Analysis of plant diversity with retrotransposon-based molecular markers. *Hered.* **2011**, *106*, 520–530.
29. Lee, J.; Mun, S.; Kim, D.H.; Cho, C.S.; Oh, D.Y.; Han, K. Chicken (*Gallus gallus*) endogenous retrovirus generates genomic variations in the chicken genome. *Mob. DNA* **2017**, *8*, 2.
30. Mandoulakani, B.A.; Piri, Y.; Darvishzadeh, R.; Bernoosi, I.; Jafari, M. Retroelement Insertional Polymorphism and Genetic Diversity in *Medicago sativa* Populations Revealed by IRAP and REMAP Markers. *Plant Mol. Biol. Report.* **2012**, *30*, 286–296.
31. Tam, S.M.; Mhiri, C.; Vogelaar, A.; Kerkveld, M.; Pearce, S.R.; Grandbastien, M.A. Comparative analyses of genetic diversities within tomato and pepper collections detected by retrotransposon-based SSAP, AFLP and SSR. *Theor. Appl. Genet.* **2005**, *110*, 819–831.
32. Chessa, B.; Pereira, F.; Arnaud, F.; Amorim, A.; Goyache, F.; Mainland, I.; Kao, R.R.; Pemberton, J.M.; Beraldi, D.; Stear, M.J.; et al. Revealing the history of sheep domestication using retrovirus integrations. *Science (80-)*. **2009**, *324*, 532–536.
33. Elleder, D.; Kim, O.; Padhi, A.; Bankert, J.G.; Simeonov, I.; Schuster, S.C.; Wittekindt, N.E.; Motameny, S.; Poss, M. Polymorphic Integrations of an Endogenous Gammaretrovirus in the Mule Deer Genome. *J. Virol.* **2012**, *86*, 2787–2796.
34. Hron, T.; Fabryova, H.; Elleder, D. Insight into the epigenetic landscape of a currently endogenizing gammaretrovirus in mule deer (*Odocoileus hemionus*). *Genomics* **2020**, *112*, 886–896.
35. Liu, C.; Ran, X.; Niu, X.; Li, S.; Wang, J.; Zhang, Q. Insertion of 275-bp SINE into first intron of PDIA4 gene is associated with litter size in Xiang pigs. *Anim. Reprod. Sci.* **2018**, *195*, 16–23.
36. Zheng, Y.; Chen, C.; Chen, W.; Wang, X.; Wang, W.; GAO, B.; Wimmers, K.; Mao, J.; Song, C. Two new SINE insertion polymorphisms in pig Vertnin (VRTN) gene revealed by comparative genomic alignment. *J. Integr. Agric.* **2020**, *19*, 2514–2522.
37. Fontanesi, L.; Scotti, E.; Buttazzoni, L.; Dall'Olio, S.; Russo, V. Investigation of a Short Interspersed Nuclear Element Polymorphic Site in the Porcine Vertnin Gene: Allele Frequencies and Association Study With Meat Quality, Carcass and Production Traits in Italian Large White pigs. *Ital. J. Anim. Sci.* **2014**, *13*, 3090.
38. Rooney, M.F.; Hill, E.W.; Kelly, V.P.; Porter, R.K. The “speed gene” effect of myostatin arises in Thoroughbred horses due to a promoter proximal SINE insertion. *PLoS One* **2018**, *13*, e0205664.
39. Gray, M.M.; Sutter, N.B.; Ostrander, E.A.; Wayne, R.K. The IGF1 small dog haplotype is derived from Middle Eastern grey wolves. *BMC Biol.* **2010**, *8*, 16.

40. Chuong, E.B.; Elde, N.C.; Feschotte, C. Retrotransposon insertion in SILV is responsible for merle patterning of the domestic dog. *Nat. Rev. Genet.* **2017**, *18*, 71–86.
41. Fukai, E.; Soyano, T.; Umehara, Y.; Nakayama, S.; Hirakawa, H.; Tabata, S.; Sato, S.; Hayashi, M. Establishment of a Lotus japonicus gene tagging population using the exon-targeting endogenous retrotransposon LORE1. *Plant J.* **2012**, *69*, 720–730.
42. Jiang, S.Y.; Ramachandran, S. Genome-Wide Survey and Comparative Analysis of LTR Retrotransposons and Their Captured Genes in Rice and Sorghum. *PLoS One* **2013**, *8*.
43. Kumar, A.; Hirochika, H. Applications of retrotransposons as genetic tools in plant biology. *Trends Plant Sci.* **2001**, *6*, 127–134.
44. Cordaux, R.; Batzer, M.A. The impact of retrotransposons on human genome evolution. *Nat. Rev. Genet.* **2009**, *10*, 691–703.
45. Kazazian, H.H.; Moran, J. V. The impact of L1 retrotransposons on the human genome. *Nat. Genet.* **1998**, *19*, 19–24.
46. Kofler, R.; Betancourt, A.J.; Schlötterer, C. Sequencing of Pooled DNA Samples (Pool-Seq) Uncovers Complex Dynamics of Transposable Element Insertions in *Drosophila melanogaster*. *PLoS Genet.* **2012**, *8*, e1002487.
47. Laricchia, K.M.; Zdraljevic, S.; Cook, D.E.; Andersen, E.C. Natural Variation in the Distribution and Abundance of Transposable Elements Across the *Caenorhabditis elegans* Species. *Mol. Biol. Evol.* **2017**, *34*, 2187–2202.
48. Rishishwar, L.; Tellez Villa, C.E.; Jordan, I.K. Transposable element polymorphisms recapitulate human evolution. *Mob. DNA* **2015**.
49. Hancks, D.C.; Kazazian, H.H. Active human retrotransposons: Variation and disease. *Curr. Opin. Genet. Dev.* **2012**, *22*, 191–203.
50. Liu, Z.; Wang, T.; Wang, L.; Zhao, H.; Yue, E.; Yan, Y.; Irshad, F.; Zhou, L.; Duan, M.; Xu, J. RTRIP: a comprehensive profile of transposon insertion polymorphisms in rice. *Plant Biotechnol. J.* **2020**, *18*, 2379–2381.
51. Altshuler, D.L.; Durbin, R.M.; Abecasis, G.R.; Bentley, D.R.; Chakravarti, A.; Clark, A.G.; Collins, F.S.; De La Vega, F.M.; Donnelly, P.; Egholm, M.; et al. A map of human genome variation from population-scale sequencing. *Nature* **2010**, *467*, 1061–1073.
52. Altshuler, D.M.; Durbin, R.M.; Abecasis, G.R.; Bentley, D.R.; Chakravarti, A.; Clark, A.G.; Donnelly, P.; Eichler, E.E.; Flicek, P.; Gabriel, S.B.; et al. An integrated map of genetic variation from 1,092 human genomes. *Nature* **2012**, *491*, 56–65.
53. Carpentier, M.-C.; Manfroi, E.; Wei, F.-J.; Wu, H.-P.; Lasserre, E.; Llauro, C.; Debladis, E.; Akakpo, R.; Hsing, Y.-I.; Panaud, O. Retrotranspositional landscape of Asian rice revealed by 3000 genomes. *Nat. Commun.* **2019**, *10*, 24.
54. Zhang, Y.; Maksakova, I.A.; Gagnier, L.; van de Lagemaat, L.N.; Mager, D.L. Genome-wide assessments reveal extremely high levels of polymorphism of two active families of mouse endogenous retroviral elements. *PLoS Genet.* **2008**, *4*, e1000007.

55. Wang, W. Short interspersed elements (SINEs) are a major source of canine genomic diversity. *Genome Res.* **2005**, *15*, 1798–1808.
56. Kalla, S.E.; Moghadam, H.K.; Tomlinson, M.; Seebald, A.; Allen, J.J.; Whitney, J.; Choi, J.D.; Sutter, N.B. Polymorphic SINEC_Cf Retrotransposons in the Genome of the Dog (*Canis familiaris*). *bioRxiv* **2020**, 2020.10.27.358119.
57. Asia, S.E.; Asia, S.E.; Genome, S.; Consortium, S.; Information, S.; Tables, S.; Information, S. Analyses of pig genomes provide insight into porcine demography and evolution. *Nature* **2012**, *491*, 393–398.
58. Lindblad-Toh, K.; Wade, C.M.; Mikkelsen, T.S.; Karlsson, E.K.; Jaffe, D.B.; Kamal, M.; Clamp, M.; Chang, J.L.; Kulbokas, E.J.; Zody, M.C.; et al. Genome sequence, comparative analysis and haplotype structure of the domestic dog. *Nature* **2005**, *438*, 803–819.
59. Almeida, L.M.; Silva, I.T.; Silva Jr., W.A.; Castro, J.P.; Riggs, P.K.; Carareto, C.M.; Amaral, M.E.J. The contribution of transposable elements to *Bos taurus* gene structure. *Gene* **2007**, *390*, 180–189.
60. Burns, K.H.; Boeke, J.D. Human transposon tectonics. *Cell* **2012**, *149*, 740–752.
61. Levy, A.; Sela, N.; Ast, G. TranspoGene and microTranspoGene: Transposed elements influence on the transcriptome of seven vertebrates and invertebrates. *Nucleic Acids Res.* **2008**, *36*.
62. Hancks, D.C.; Kazazian, H.H. Roles for retrotransposon insertions in human disease. *Mob. DNA* **2016**, *7*, 9.
63. Dixit, M.; Poudel, S.B.; Yakar, S. Effects of GH/IGF axis on bone and cartilage. *Mol. Cell. Endocrinol.* **2021**, *519*, 111052.
64. Park, S.G.; Hannenhalli, S.; Choi, S.S. Conservation in first introns is positively associated with the number of exons within genes and the presence of regulatory epigenetic signals. *BMC Genomics* **2014**, *15*.
65. Mandal, A.K.; Pandey, R.; Jha, V.; Mukerji, M. Transcriptome-wide expansion of non-coding regulatory switches: Evidence from co-occurrence of Alu exonization, antisense and editing. *Nucleic Acids Res.* **2013**, *41*, 2121–2137.
66. Faulkner, G.J.; Kimura, Y.; Daub, C.O.; Wani, S.; Plessy, C.; Irvine, K.M.; Schroder, K.; Cloonan, N.; Steptoe, A.L.; Lassmann, T.; et al. The regulated retrotransposon transcriptome of mammalian cells. *Nat. Genet.* **2009**, *41*, 563–571.
67. Ba, N. V.; Arakawa, A.; Ishihara, S.; Nam, L.Q.; Thuy, T.T.T.; Dinh, N.C.; Ninh, P.H.; Cuc, N.T.K.; Kikuchi, K.; Pham, L.D.; et al. Evaluation of genetic richness among Vietnamese native pig breeds using microsatellite markers. *Anim. Sci. J.* **2020**, *91*.
68. Fabuel, E.; Barragán, C.; Silió, L.; Rodríguez, M.C.; Toro, M.A. Analysis of genetic diversity and conservation priorities in Iberian pigs based on microsatellite markers. *Heredity (Edinb).* **2004**, *93*, 104–113.
69. Fang, M.; Hu, X.; Jiang, T.; Braunschweig, M.; Hu, L.; Du, Z.; Feng, J.; Zhang, Q.; Wu, C.; Li, N. The phylogeny of Chinese indigenous pig breeds inferred from microsatellite markers. *Anim. Genet.* **2005**, *36*, 7–13.

70. Pham, L.D.; Do, D.N.; Nam, L.Q.; van Ba, N.; Minh, L.T.A.; Hoan, T.X.; Cuong, V.C.; Kadarmideen, H.N. Molecular genetic diversity and genetic structure of Vietnamese indigenous pig populations. *J. Anim. Breed. Genet.* **2014**, *131*, 379–386.
71. Fan, B.; Wang, Z.-G.; Li, Y.-J.; Zhao, X.-L.; Liu, B.; Zhao, S.-H.; Yu, M.; Li, M.-H.; Chen, S.-L.; Xiong, T.-A.; et al. Genetic variation analysis within and among Chinese indigenous swine populations using microsatellite markers. *Anim. Genet.* **2002**, *33*, 422–427.
72. Boitard, S.; Chevalet, C.; Mercat, M.-J.; Meriaux, J.C.; Sanchez, A.; Tibau, J.; Sancristobal, M. Genetic variability, structure and assignment of Spanish and French pig populations based on a large sampling. *Anim. Genet.* **2010**, *41*, 608–618.
73. Vicente, A.A.; Carolino, M.I.; Sousa, M.C.O.; Ginja, C.; Silva, F.S.; Martinez, A.M.; Vega-Pla, J.L.; Carolino, N.; Gama, L.T. Genetic diversity in native and commercial breeds of pigs in Portugal assessed by microsatellites. *J. Anim. Sci.* **2008**, *86*, 2496–2507.
74. Kruglyak, L. The use of a genetic map of biallelic markers in linkage studies. *Nat. Genet.* **1997**, *17*, 21–24.
75. Vignal, A.; Milan, D.; SanCristobal, M.; Eggen, A. A review on SNP and other types of molecular markers and their use in animal genetics. *Genet. Sel. Evol.* **2002**, *34*, 275–305.
76. Yang, W.; Kang, X.; Yang, Q.; Lin, Y.; Fang, M. Review on the development of genotyping methods for assessing farm animal diversity. *J. Anim. Sci. Biotechnol.* **2013**, *4*.
77. Burg, K. Molecular Markers for Genetic Diversity. In *Progress in Botany Vol. 79*, Cánovas, F.M., Lüttge, U., Matyssek, R., Eds.; Springer International Publishing: Cham, 2017; pp. 33–47 ISBN 978-3-319-71413-4.
78. Grover, A.; Sharma, P.C. Development and use of molecular markers: past and present. *Crit. Rev. Biotechnol.* **2016**, *36*, 290–302.
79. Huang, M.; Yang, B.; Chen, H.; Zhang, H.; Wu, Z.; Ai, H.; Ren, J.; Huang, L. The fine-scale genetic structure and selection signals of Chinese indigenous pigs. *Evol. Appl.* **2020**, *13*, 458–475.
80. Traspov, A.; Deng, W.; Kostyunina, O.; Ji, J.; Shatokhin, K.; Lugovoy, S.; Zinovieva, N.; Yang, B.; Huang, L. Population structure and genome characterization of local pig breeds in Russia, Belorussia, Kazakhstan and Ukraine. *Genet. Sel. Evol.* **2016**, *48*, 16.
81. Muñoz, M.; Bozzi, R.; García-Casco, J.; Núñez, Y.; Ribani, A.; Franci, O.; García, F.; Škrlep, M.; Schiavo, G.; Bovo, S.; et al. Genomic diversity, linkage disequilibrium and selection signatures in European local pig breeds assessed with a high density SNP chip. *Sci. Rep.* **2019**, *9*, 13546.
82. Van Ba, N.; Nam, L.Q.; Do, D.N.; Van Hau, N.; Pham, L.D. An assessment of genetic diversity and population structures of fifteen Vietnamese indigenous pig breeds for supporting the decision making on conservation strategies. *Trop. Anim. Health Prod.* **2020**, *52*, 1033–1041.
83. Diao, S.; Huang, S.; Xu, Z.; Ye, S.; Yuan, X.; Chen, Z.; Zhang, H.; Zhang, Z.; Li, J. Genetic Diversity of Indigenous Pigs from South China Area Revealed by SNP Array. *Animals* **2019**, *9*, 361.
84. Pallares, L.F. Searching for solutions to the missing heritability problem. *Elife* **2019**, *8*.

85. Gibson, G. Rare and common variants: twenty arguments. *Nat. Rev. Genet.* **2012**, *13*, 135–145.
86. Giuffra, E.; Kijas, J.M.H.; Amarger, V.; Carlborg, Ö.; Jeon, J.T.; Andersson, L. The origin of the domestic pig: Independent domestication and subsequent introgression. *Genetics* **2000**, *154*, 1785–1791.
87. Kijas, J.M.H.; Andersson, L. A Phylogenetic Study of the Origin of the Domestic Pig Estimated from the Near-Complete mtDNA Genome. *J. Mol. Evol.* **2001**, *52*, 302–308.
88. Choi, S.K.; Lee, J.E.; Kim, Y.J.; Min, M.S.; Voloshina, I.; Myslenkov, A.; Oh, J.G.; Kim, T.H.; Markov, N.; Seryodkin, I.; et al. Genetic structure of wild boar (*Sus scrofa*) populations from East Asia based on microsatellite loci analyses. *BMC Genet.* **2014**, *15*, 85.
89. Wu, G.-S.; Yao, Y.-G.; Qu, K.-X.; Ding, Z.-L.; Li, H.; Palanichamy, M.G.; Duan, Z.-Y.; Li, N.; Chen, Y.-S.; Zhang, Y.-P. Population phylogenomic analysis of mitochondrial DNA in wild boars and domestic pigs revealed multiple domestication events in East Asia. *Genome Biol.* **2007**, *8*, R245.
90. Jin, L.; Zhang, M.; Ma, J.; Zhang, J.; Zhou, C.; Liu, Y.; Wang, T.; Jiang, A. an; Chen, L.; Wang, J.; et al. Mitochondrial DNA Evidence Indicates the Local Origin of Domestic Pigs in the Upstream Region of the Yangtze River. *PLoS One* **2012**, *7*.
91. Yang, S.; Zhang, H.; Mao, H.; Yan, D.; Lu, S.; Lian, L.; Zhao, G.; Yan, Y.; Deng, W.; Shi, X.; et al. The local origin of the Tibetan pig and additional insights into the origin of Asian pigs. *PLoS One* **2011**, *6*.
92. Yuan, J.; Luo, Y.; Wang, Z.; Xiang, H.; Zhao, X. Exploring the origin of domesticated pigs in the Yellow River area using information from ancient DNA. *Chinese Sci. Bull.* **2012**, *57*, 1011–1018.
93. Xiang, H.; Gao, J.; Cai, D.; Luo, Y.; Yu, B.; Liu, L.; Liu, R.; Zhou, H.; Chen, X.; Dun, W.; et al. Origin and dispersal of early domestic pigs in northern China. *Sci. Rep.* **2017**, *7*, 5602.
94. Tarailo-Graovac, M.; Chen, N. Using RepeatMasker to Identify Repetitive Elements in Genomic Sequences. *Curr. Protoc. Bioinforma.* **2009**, *25*.
95. Kimura, M. A simple method for estimating evolutionary rates of base substitutions through comparative studies of nucleotide sequences. *J. Mol. Evol.* **1980**, *16*, 111–120.
96. Larkin, M.A.; Blackshields, G.; Brown, N.P.; Chenna, R.; McGettigan, P.A.; McWilliam, H.; Valentin, F.; Wallace, I.M.; Wilm, A.; Lopez, R.; et al. Clustal W and Clustal X version 2.0. *Bioinformatics* **2007**, *23*, 2947–2948.
97. Quinlan, A.R.; Hall, I.M. BEDTools: A flexible suite of utilities for comparing genomic features. *Bioinformatics* **2010**.
98. Kent, W.J. BLAT—The BLAST-Like Alignment Tool. *Genome Res.* **2002**, *12*, 656–664.
99. Altschul, S.F.; Gish, W.; Miller, W.; Myers, E.W.; Lipman, D.J. Basic local alignment search tool. *J. Mol. Biol.* **1990**, *215*, 403–410.
100. Kolde, R. Pheatmap: pretty heatmaps. *R Packag. version* **2012**, *1*.

Tables

Table 1. Summary of SINE RIPs distributed among pig genomes.

No. of genomes containing SINE RIPs	No. of insertion allele	No. of deletion allele	No. of total SINE RIP allele
1	11452	N	11452
2	4042	N	4042
3	2436	N	2436
4	1575	1730	3305
5	1116	1452	2568
6	783	1478	2261
7	600	1344	1944
8	452	1270	1722
9	265	1380	1645
10	136	1091	1227
11	106	957	1063
12	51	908	959
13	26	620	646
14	18	567	585
15	3	426	429
Total	23061	13223	36284

Table 2. Intersection of SINE RIPs with genic regions.

Type of biogenic regions	No. of SINE RIPs	Denisty #/Mb	No. of gene contain SINE RIPs	Percentage
LncRNA gene	4443	14.75	2504	17.30%
LncRNA gene-exon	157	9.43	154	17.30%
LncRNA gene-first-exon	68	10.15	67	1.06%
LncRNA gene-last-exon	89	9.44	87	0.46%
LncRNA gene-intron	4305	14.98	2389	0.60%
LncRNA gene-intron1	2900	15.40	1846	16.50%
LncRNA gene-intron2	1228	14.10	631	12.75%
LncRNA gene-intron3	513	13.61	284	4.36%
LncRNA gene-intron4	219	13.17	137	1.96%
LncRNA gene-intron5	145	15.88	66	0.95%
LncRNA gene- 5' flank (5kb)	1007	14.28	946	6.53%
LncRNA gene-3' flank (5kb)	1116	15.80	1059	7.32%
Protein coding gene	16787	14.67	6154	29.78%
Protein coding gene-exon	260	3.11	245	1.19%
Protein coding gene-CDS	8	0.22	8	0.04%
Protein coding gene-5'UTR	98	5.39	93	0.45%
Protein coding gene-3'UTR	151	4.74	147	0.71%
Protein coding gene-intron	16635	15.41	6070	29.37%
Protein coding gene-intron1	4954	15.03	2636	12.76%
Protein coding gene-intron2	3937	14.78	2102	10.17%
Protein coding gene-intron3	2726	14.24	1546	7.48%
Protein coding gene-intron4	2156	14.86	1262	6.11%
Protein coding gene-intron5	1739	15.04	1094	5.29%
Protein coding gene-5' flank (5kb)	1589	15.96	1509	7.30%
Protein coding gene-3' flank (5kb)	1435	14.77	1400	6.77%
Intergenic	14688	14.41	N	N
Random5kb (N=20000)	1380	14.07	N	N

Intragenic	21596	14.57	10347	23.09%
------------	-------	-------	-------	--------

Figures

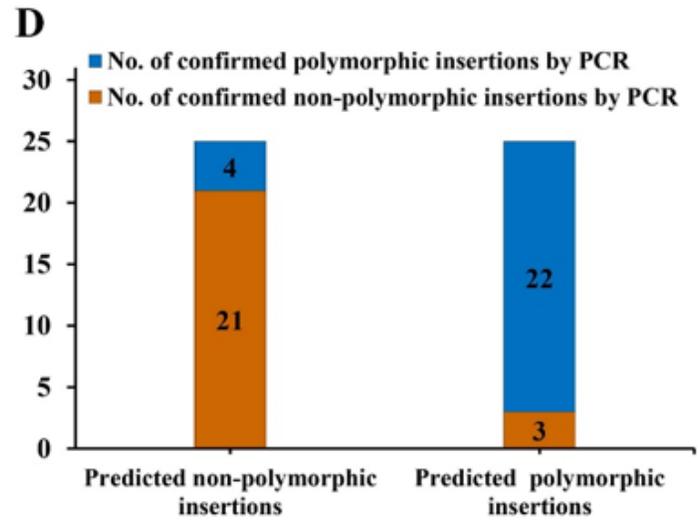
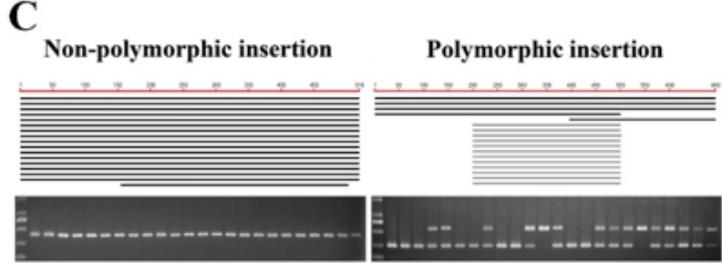
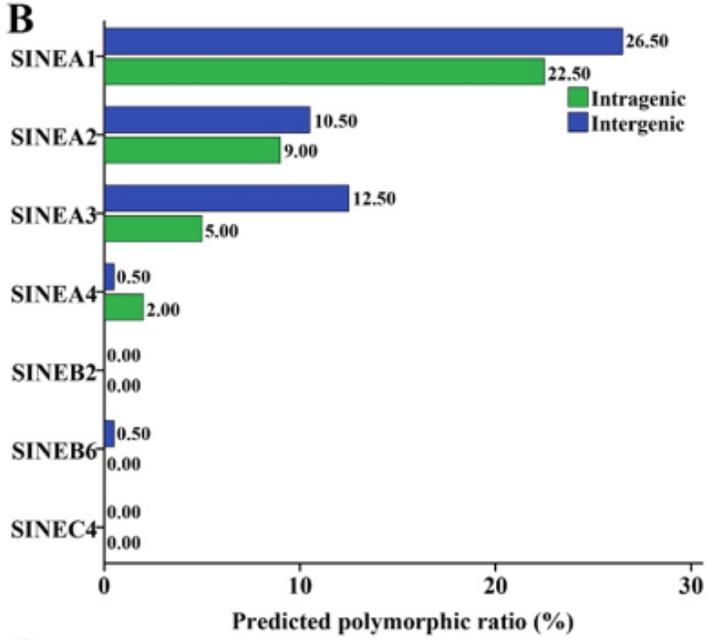
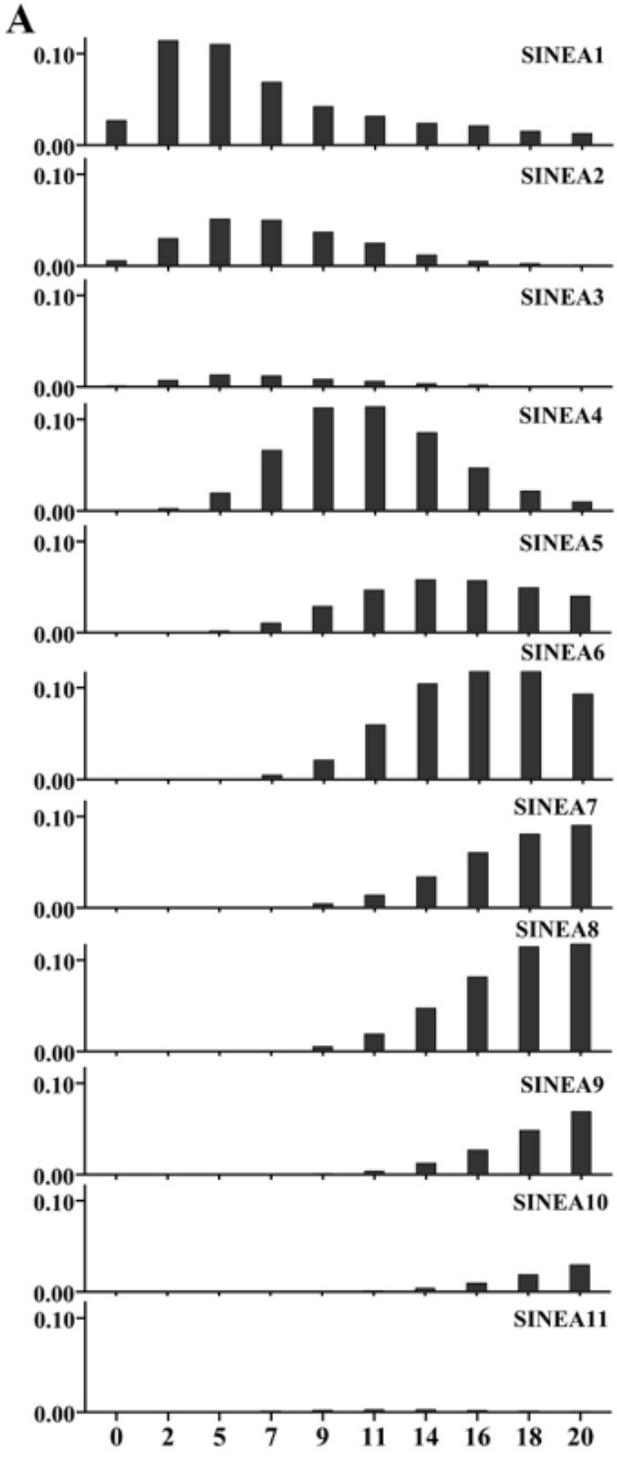


Figure 1

High polymorphisms of young SINE subfamilies. (A) Insertion ages of SINEA subfamilies. The X-axis represents insertion ages (million years ago, Mya), and the Y-axis represents the genome coverage (%) of SINE subfamilies. (B) Predicted polymorphic ratio of SINE insertions from seven SINE subfamilies representing different insertion ages. (C) Representative results of the BLAST prediction and PCR verification. (D) The results of PCR verification for 25 predicted polymorphic and 25 non-polymorphic insertions from different SINE subfamilies (primers listed in Supplementary File S1).

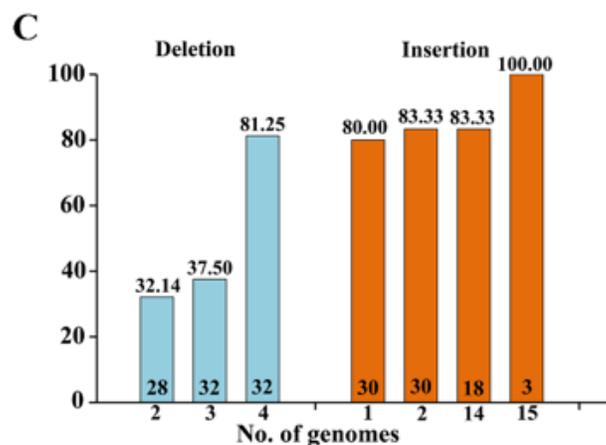
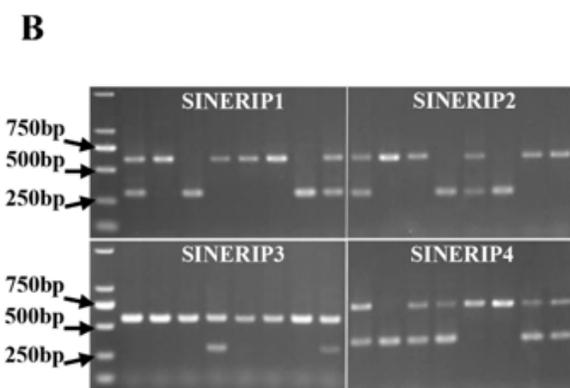
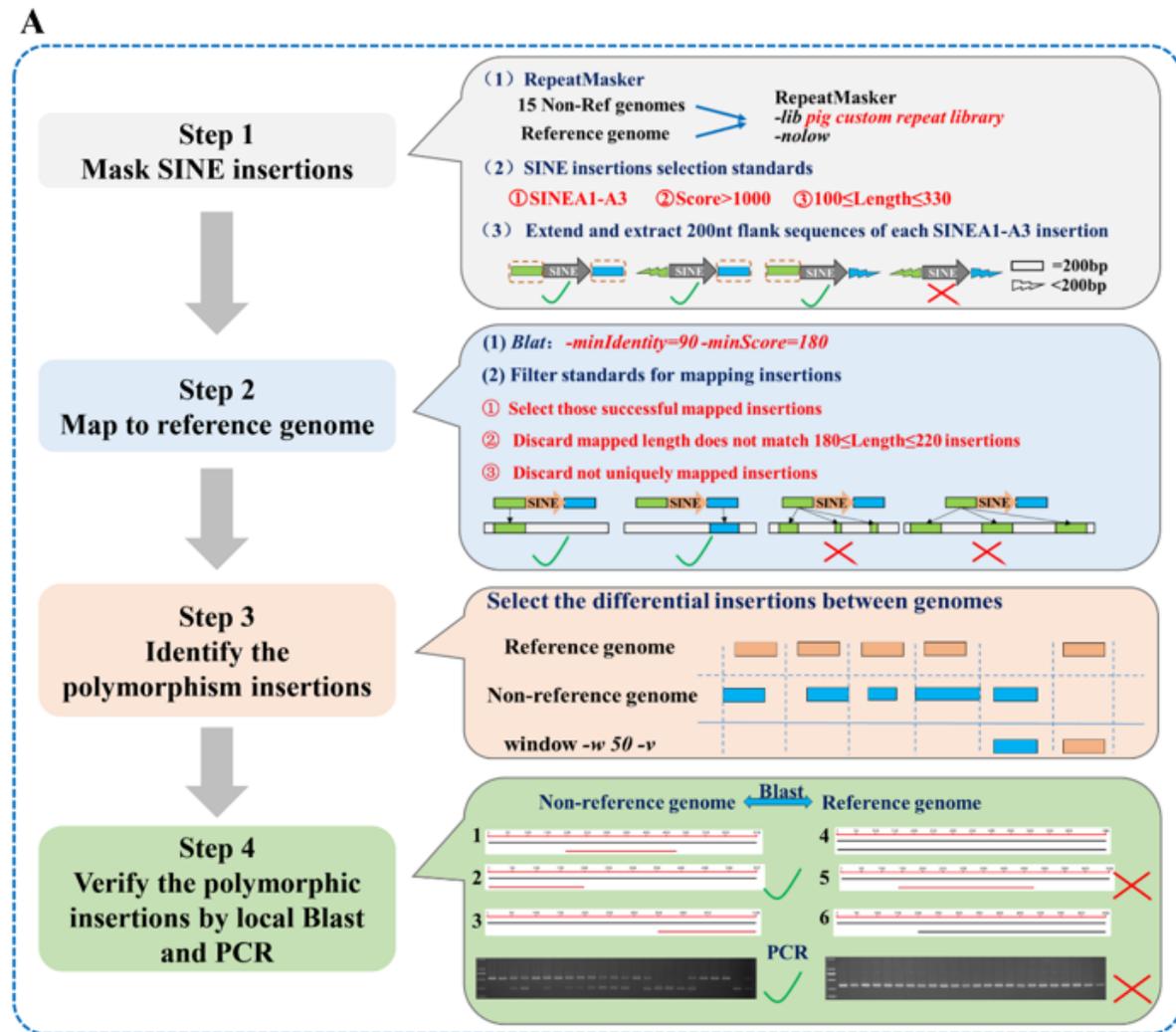


Figure 2

SINE RIP annotation and verification protocol. (A) Main steps and methods of SINE RIP annotation. (B) Representative electrophoresis results of PCR verification of SINE RIPs, and pooled DNA samples from each breed (Bama, Landrace, Large White, Duroc, Erhualian, Meishan, Sujiang, and Wuzhishan) were used for PCR amplification. (C) Results of PCR verification for rare SINE deletion and insertion alleles, (primers listed in Supplementary File S1). The X-axis represents the number of genomes presenting a certain allele (deletion or Insertion relative to the reference genome), The Y-axis represents the percentage positive rate of SINE RIPs confirmed by PCR.

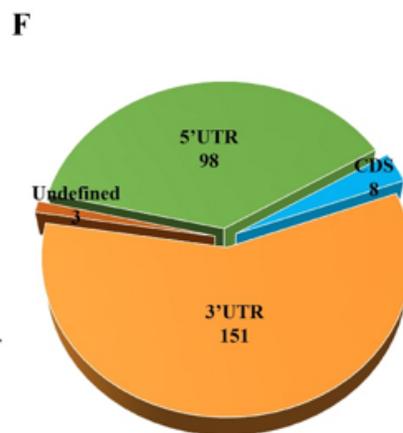
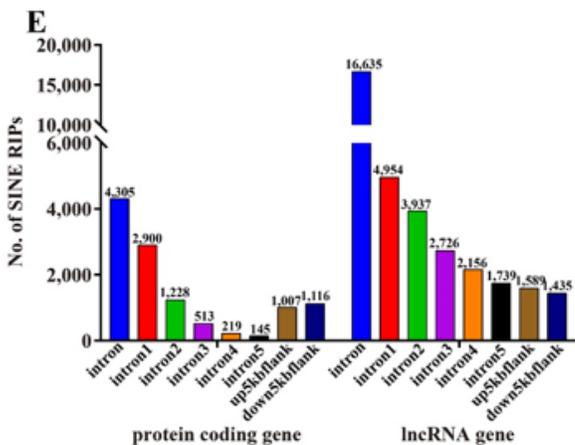
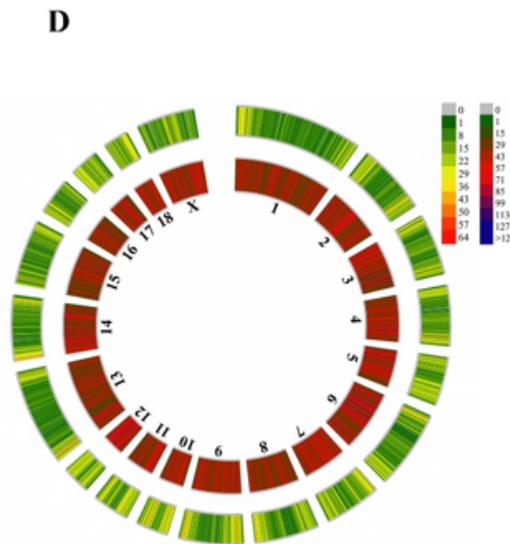
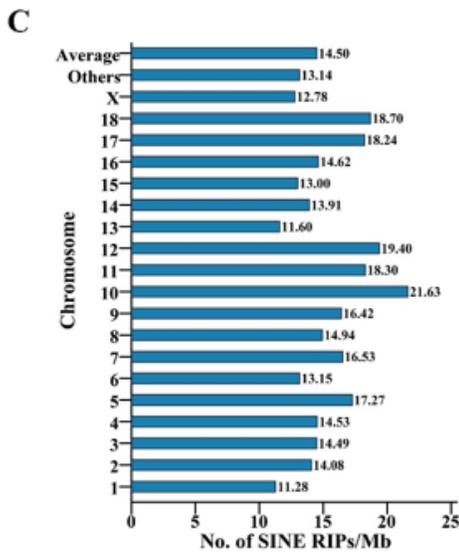
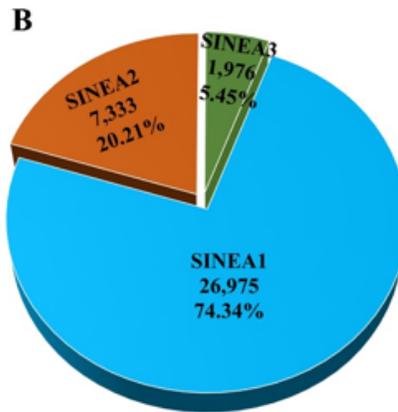
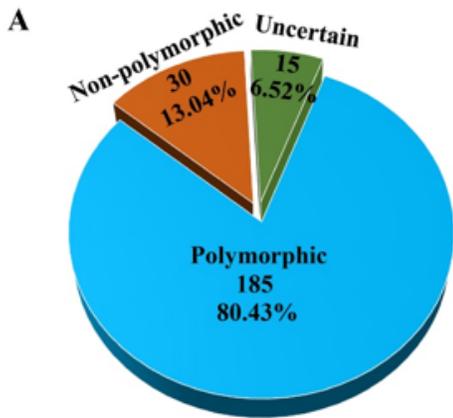


Figure 3

Distribution of SINE RIPs. (A) Summary results of PCR evaluation for 230 randomly selected SINE RIPs (primers listed in Supplementary File S1). (B) Distribution of SINE RIPs across SINEA1–SINEA3 subfamilies. (C) Density of SINE RIPs on each chromosome (RIPs/Mb). (D) Distribution of SINE RIPs (outer ring) and insertions (inner ring) on each chromosome. (E) Distribution of SINE RIPs in the introns and flanking regions of the lncRNA and protein-coding genes. (F) Distributions of SINE RIPs in the 5' UTR, 3' UTR, and CDS regions.

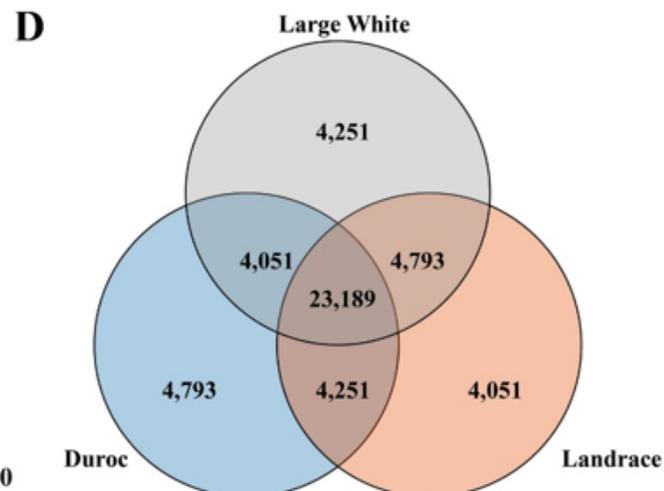
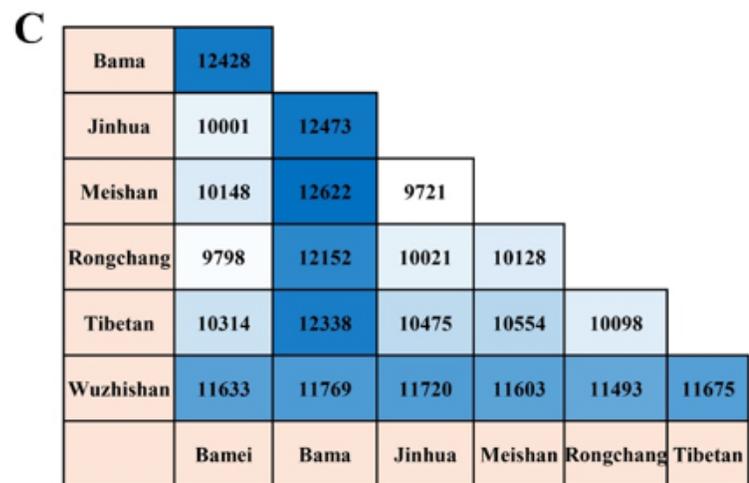
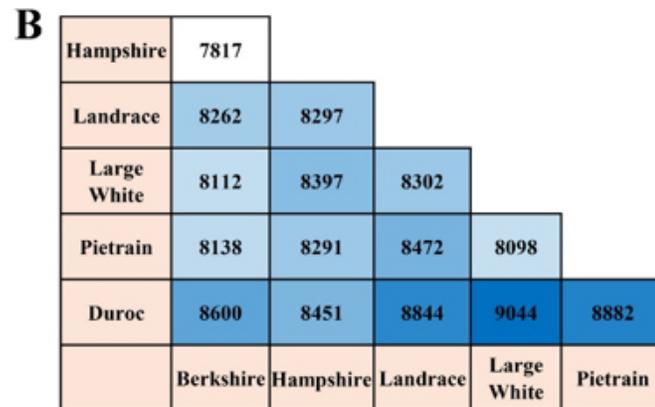
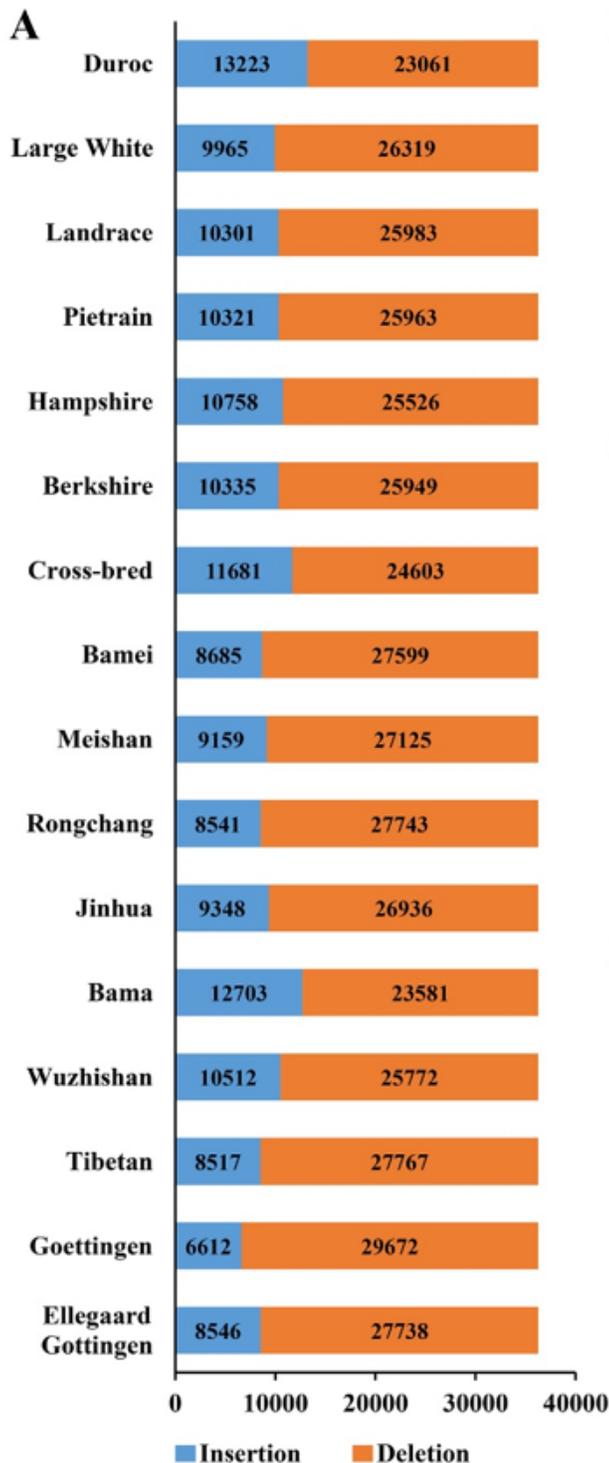


Figure 4

Distribution of the SINE RIP alleles in each genome and between pairs of commercial pig breeds. (A) The numbers of deletion and insertion alleles for 36,284 SINE RIPs in each genome. (B) Distribution of the different SINE RIP alleles between different commercial breeds. (C) Distribution of the different SINE RIP alleles between different Chinese native breeds, (D) Distribution of the differential SINE RIP alleles between the three most common commercial breeds.

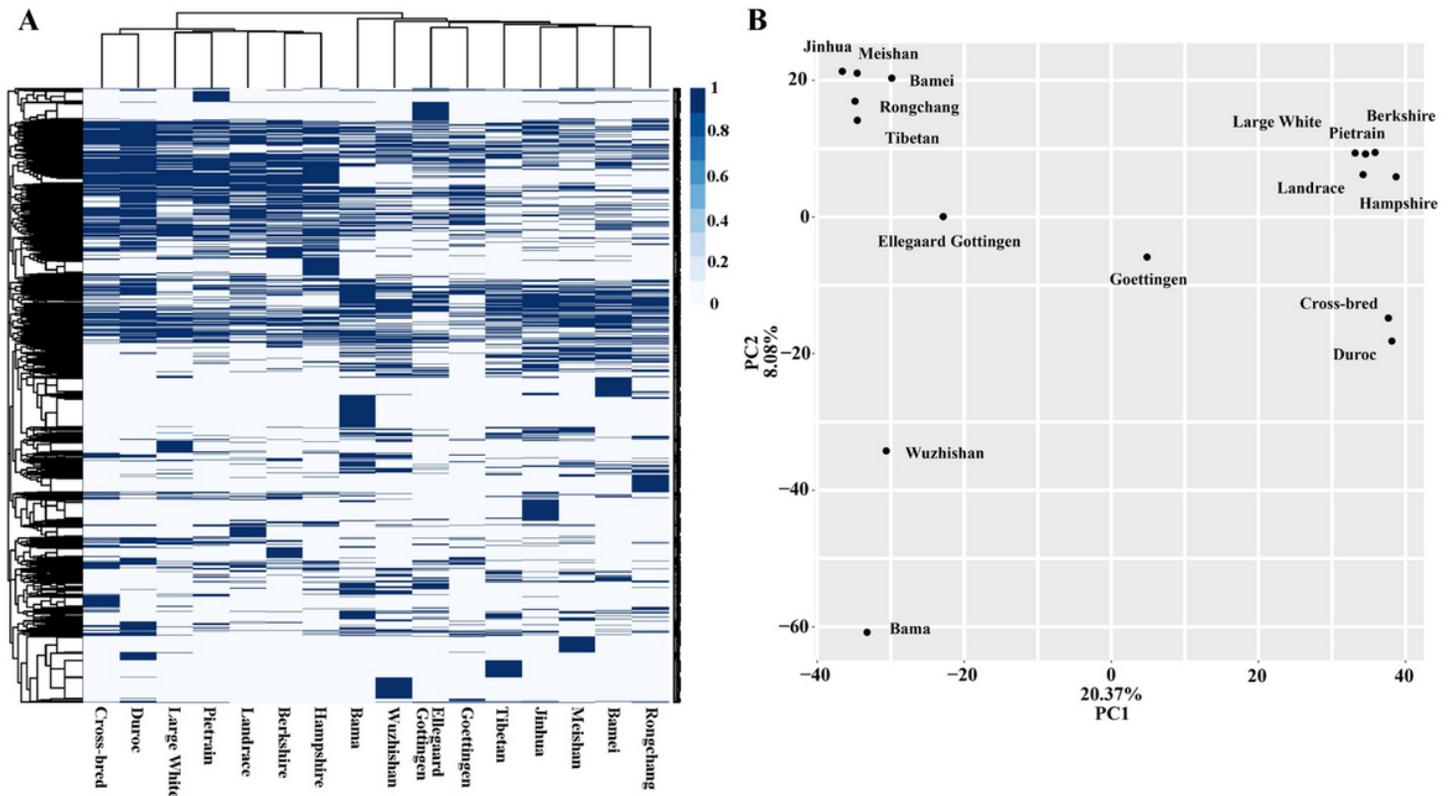


Figure 5

Genetic relationship analysis by heat mapping (A) and PCA (B) based on all SINE RIPs.

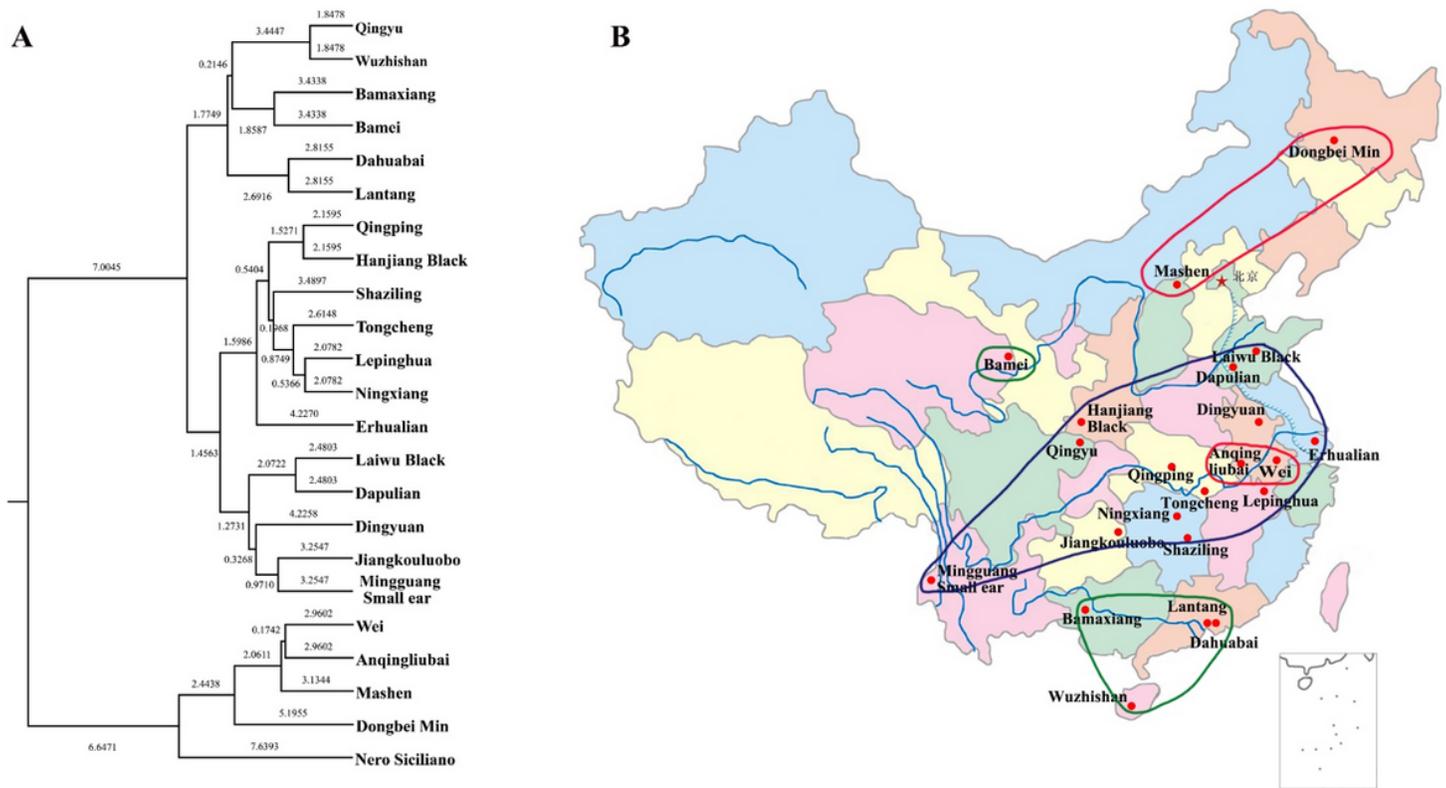


Figure 6

Population genetic analysis. (A) Cluster analysis for 23 populations with 16 SINE RIPs. (B) Distribution of Chinese native pigs used for analysis (primers listed in Supplementary File S1). Note: The designations employed and the presentation of the material on this map do not imply the expression of any opinion whatsoever on the part of Research Square concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. This map has been provided by the authors.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Additionalfile1.docx](#)
- [Additionalfile2.xlsx](#)
- [Additionalfile3.xlsx](#)