

Recovery of non-reference sequences missing from the human reference genome

Ran Li

Northwest Agriculture and Forestry University

Xiaomeng Tian

Northwest Agriculture and Forestry University

Peng Yang

Northwest Agriculture and Forestry University

Yingzhi Fan

Northwest Agriculture and Forestry University

Ming Li

Northwest Agriculture and Forestry University

Hongxiang Zheng

Fudan University

Xihong Wang

Northwest Agriculture and Forestry University

Yu Jiang (✉ yu.jiang@nwfau.edu.cn)

Northwest Agriculture and Forestry University <https://orcid.org/0000-0003-4821-3585>

Research article

Keywords: pan-genome, human genome, genomic variation, alternate alleles

Posted Date: July 19th, 2019

DOI: <https://doi.org/10.21203/rs.2.11742/v1>

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

[Read Full License](#)

Version of Record: A version of this preprint was published on October 16th, 2019. See the published version at <https://doi.org/10.1186/s12864-019-6107-1>.

Abstract

The non-reference sequences (NRS) represent structure variations in human genome with potential functional significance. However, besides the known insertions, it is currently unknown whether other types of structure variations with NRS exist. Here, we compared 31 human de novo assemblies with the current reference genome to identify the NRS and their location. We resolved the precise location of 6,113 NRS adding up to 12.8 Mb. Besides 1,571 insertions, we detected 3,041 alternate alleles, which were defined as having less than 90% (or none) identity with the reference alleles. These alternate alleles overlapped with 1,143 protein-coding genes including a putative novel MHC haplotype. Further, we demonstrated that the alternate alleles and their flanking regions had high content of tandem repeats, indicating that their origin was associated with tandem repeats. Our study enriched the spectrum of human genetic variations.

Background

The initial human reference genome was entirely linear (Lander et al., 2001). Despite introduction of alternate alleles for a graph-based representation, the current reference genome is largely derived from a single individual of African-European origin (Schneider et al., 2017), limiting its representation of diverse populations. Lines of evidence in recent years have revealed that individuals still carry sequences that are not represented in the reference genome. These sequences could be an important type of structural variation underlying disease associations or complex traits (Kehr et al., 2017). The discovery of non-reference sequences (NRS) will be a prerequisite for an more complete graph-based genome, thereby enabling improved genomic analyses and understanding of genomic architecture (Church et al., 2015).

Extensive efforts have been devoted in recent years to discover NRS. Based on a large amount of whole-genome sequencing data, two studies focused specifically on the discovery of NRS and identified as much as ~300 Mb novel sequences from a large number of re-sequencing data (Kehr et al., 2017; Sherman et al., 2018). However, the identified sequences were obtained by assembling unaligned short reads, posing a challenge to placing them in the reference genome. A recent study used long-read sequencing data from 15 samples to identify 32,838 insertions of presence in at least two samples but without exploring novel sequences within the insertions (Audano et al., 2019). Additionally, these studies have mainly focused on insertion events. In fact, some sequences belong to complex structural variants (e.g., two alleles with high divergence instead of simply introducing additional sequences) (Li et al., 2010; Sherman et al., 2018). De novo assembly is a promising approach for building the complete human pan-genome (Li et al., 2010). Using an assembly-versus-assembly approach, we discovered not only insertions but also sequences which are an alternate representation of a locus in the haploid genome. A well characterization of the insertions and alternate alleles in the human genome is necessary for a better understanding of their biological significance.

The identification of insertions and alternate alleles requires high-quality de novo assemblies. Fortunately, a number of human de novo assemblies have been generated via long-read sequencing (LRS) (Audano et

al., 2019; Cho et al., 2016; English et al., 2015; Jain et al., 2018; Pendleton et al., 2015; Shi et al., 2016), and these assemblies have covered major human ethnic groups. The unprecedented availability of these genomic resources enables us to depict the full spectrum of NRS, especially those representing alternate alleles. In this study, we compared 31 de novo assemblies (including 17 LRS assemblies) with the human reference genome to identify putative alternate alleles, most of which are newly reported in this study.

Methods

De novo assemblies used in this study

The human reference genome GRCh38.p12 (hg38) was used as the guiding genome for comparison. The hg38 consists of the primary GRCh38 assembly, the mitochondria genome, unlocalized/unplaced scaffolds and alternate contigs. We downloaded 31 human de novo assemblies from the NCBI, including 17 from PacBio sequencing, 1 from nanopore sequencing, 13 from next generation sequencing and one from Sanger sequencing (Additional table 1). For the assemblies using LRS technology (PacBio and nanopore sequencing), we focused on assemblies that were released since 2015 and of high quality. All of them had high continuity (contig N50 >1 Mb; 15 of 18 with an N50 > 5 Mb) and high sequencing coverage (16 of 17 with coverage >50 X). For the assemblies from SRS, we used the 13 haploid genomes of next generation sequencing since they were generated in one study with high continuity (Wong et al., 2018). The HuRef genome of Sanger sequencing was also included in our study due to its high continuity (Levy et al., 2007).

Recovery of candidate NRS from each assembly

Each assembly was aligned to hg38 using LAST (-m100 -E0.05) (Kielbasa et al., 2011). The unaligned or low-identity sequences (<90%) to hg38 with a length of at least 100 bp were extracted. The unaligned or low-identity sequences identified by LAST were then aligned back to hg38 using BLAST v2.2.31 (megablast) (Camacho et al., 2009) to further remove regions that share $\geq 90\%$ identity. Then, the simple repeats, low complexity regions and microsatellites were removed based on the repeat annotation file, which was downloaded from the NCBI (*_rm.out.gz). The remaining sequences were merged by adjacent regions (≤ 200 bp), and the resulting NRS, which were at least 400 bp, were retained for each assembly. Finally, the resulting sequences were aligned to hg38 using BLAST (megablast) again to remove the regions that share $\geq 90\%$ identity. The resulting sequences were then merged by adjacent regions (≤ 200 bp), and only those of at least 400 bp were kept. BEDTools v2.25.0 was used in the above processes when assessing genomic features (Quinlan, 2014). The NRS from all 31 assemblies were then merged to remove redundancy and to generate a non-redundant call set using CD-HIT (-c 0.95 -aS 0.8 -d 0 -sf 1 -M 10000) (Fu et al., 2012). The resulting unified and non-redundant call set was used for further analysis.

Removal of contamination

We did not expect sequencing from bacteria, viruses or other non-mammalian species to be present in our identified sequences since the NCBI has a stringent quality control process when assemblies are

submitted. We used BLAST (megablast) to align the non-reference call set to the NCBI nt database. Only a small number of the sequences exhibited significant alignment with the non-mammalian species using 90% identity and 90% query coverage filter thresholds and were removed from the final call set.

Presence of NRS in de novo assemblies of 31 humans and in four great apes

We examined the presence of each NRS in 31 humans and four great ape de novo assemblies using BLAST (megablast). The four great apes included chimpanzee (GCA_002880755.3), bonobo (GCF_000258655.2), gorilla (GCA_900006655.3) and orangutan (GCF_002880775.1). The presence was determined for the NRS when having identity $\geq 95\%$ and coverage $\geq 80\%$ with the assembly.

Anchoring NRS on human chromosomes

All the sequences were anchored to human chromosomes based on the LAST alignment of their flanking sequences. The anchored position was reported as 'precisely placed' when both of the flanking sequences were near perfectly aligned (no gap alignment) to the reference genome. If the sequences have flanking sequences of only one end aligned to the reference genome or have flanking sequences of two ends aligned but with gap alignment rendering the inference of exact breakpoints, it was reported as 'unlocalized'. The remaining sequences were reported as 'unplaced'. Based on the breakends coordinates (the genomic position of the two breakpoints for each NRS), the breakpoint resolved sequences could be further classified as insertions when simply introducing one sequence fragment to the reference genome. For alternate alleles, the NRS should share less than 90% (or 0%) identity with the reference allele, and the reference allele were required to be at least 400 bp. Furthermore, the NRS and the reference allele should have a comparable length, with the ratio of the length to be between 1/3 and 3. The remaining sequences that did not meet the above criteria for insertions and alternate alleles were classified as ambiguous sequences.

Aligning NRS to the human expressed sequence tag database

We downloaded the human dbEST from the NCBI FTP site (ftp://ftp.ncbi.nlm.nih.gov/blast/db/FASTA/est_human.gz). Then, the NRS were aligned to the dbEST using BLAST. Since entries in the dbEST are short and represent only the ends of expressed genes, alignments with $\geq 95\%$ identity and an e-value $< 1e-5$ were considered as a hit regardless of the query coverage.

Comparison with published datasets

We compared with our results with four previously published results each of which reported a list of non-reference sequences (Audano et al., 2019; Kehr et al., 2017; Sherman et al., 2018; Wong et al., 2018). The comparison for each of the datasets was performed using a reciprocal strategy as previously described (Wong et al., 2018). We first aligned all the NRS to each of the datasets using BLAST. Alignments with $\geq 95\%$ sequence identity and $\geq 80\%$ query coverage were considered as real hits. We then aligned each of the four datasets to our NRS also with BLAST, and the alignments were filtered using the same criteria as

described above. Results from the two alignments steps were merged and a non-redundant list of NRS was reported.

Aligning RNA-seq reads to the primary call set

We downloaded a total of 87 RNA-seq data from the Geuvadis project (<https://www.ebi.ac.uk/Tools/geuvadis-das/>) and another study (Fagerberg et al., 2014). The information of the samples was described in Additional table 6. Then the RNA-seq reads were mapped to the extended version of reference (hg38 plus the primary call set) using HISAT2 with default parameters (Kim et al., 2015). Only the reads with both of the mates mapped and in proper pair were considered as high-quality alignment before counting the mapped reads on each sequence using Sambamba depth (-F "mapping_quality >= 30 and proper_pair") (Tarasov et al., 2015). A sequence was regarded to be transcribed when ≥ 10 mapped reads were found in at least two samples.

Results

Detection of candidate NRS

An assembly-versus-assembly approach was used for each assembly to identify unaligned sequences using the human reference genome (GRCh38.p12, hg38) (Fig. 1, Methods). Apart from hg38, our study included 31 de novo assemblies: 17 from LRS (PacBio or nanopore sequencing technology), 13 generated with next generation sequencing and one from Sanger sequencing (Additional table 1).

The final unaligned sequences spanned, on average, 12.9 Mb for LRS assemblies, which were much larger than the other assemblies (2.3 Mb on average) ($P < 0.01$, t test) (Fig. 2). For each assembly, 70-80% of the preliminary unaligned sequence contents belonged to simple sequence repeats or low complexity sequences and were removed to obtain the final unaligned sequences for each assembly. After removing redundant sequences and sequences shorter than 400 bp, we obtained the unaligned sequences from each assembly, which were then merged into a unified non-reference call set of 15,055 sequences (hereafter referred to as NRS) adding up to 129.1 Mb with a median length of 2,848 bp.

We aligned all the NRS to genomes of four great apes and found that 1,211 sequences were present in at least one of the four great apes (identity $\geq 95\%$ and coverage $\geq 80\%$). The presence of each NRS was also examined in each of the human de novo assemblies and those which were present in at least two of the 31 de novo assemblies and great ape assemblies were determined as non-private sequences. Taken together, 28.2% (4248) of the NRS spanning 29.3 Mb were non-private sequences, indicating that they are of high confidence (Methods and Additional table 2).

Placing candidate NRS to the reference genome

We next determined the genomic locations of the NRS by aligning their flanking sequences to hg38 (Methods), by which we resolved the precise location of 6,112 NRS (40.5% of the total NRS) adding up to 12.8 Mb. Another 2,711 NRS were anchored to chromosomes without a precise location due to sequence gaps. The remaining sequences cannot be placed on hg38 due to a lack of flanking sequences or conflict anchoring information from the two sides.

For the precisely placed sequences, we can determine whether they belong to insertions or alternate alleles as described in next section. For the unplaced NRS and those without precise locations, 2,855 were non-private sequences adding up to 25.8 Mb, indicating that they are of high confidence (Additional table 3). Although we could not determine the precise locations of these sequences, their discovery will greatly expand the repertoire of sequence diversity in the human genome.

Insertions within NRS

We first determined the insertion events for the precisely placed NRS. A total of 1,571 (3.2 Mb) were found to be insertions including 769 non-private sequences (Fig. 3a). Furthermore, 246 were present in more than half of the assemblies, indicating that they could represent major alleles. Notably, 56.8% (881) of the insertions, including 158 non-private ones, were firstly described in our study. Principal component analysis (PCA) of all the insertions based on their presence in the 16 LRS de novo assemblies of Pacbio sequencing showed a population-specific pattern (Fig. 3e). PC1 clusters African samples away from other populations, while PC2 further separates the East Asians from the Europeans, Americans and South Asians.

Alternate alleles within NRS

We further found that many NRS represent an alternate allele to their counterparts in the reference instead of insertions. To identify alternate alleles, the NRS should share less than 90% identity (or none) and have comparable length with reference alleles (Methods). In this way, 3,041 were identified as candidate alternate alleles. The remaining 1,500 precisely placed NRS did not meet our criteria of insertions or alternate alleles and thus were classified as ambiguous sequences. Unlike insertions, the alternate alleles represent allelic sequences (Fig. 4a and 4d). Notably, a long alternate sequence of 24,676 bp was anchored to chr6: 29,955,749-29,986,299, which belongs to the class I major histocompatibility complex (MHC gene) (Fig. 4b) and potentially harbors a novel HLA-B gene (Additional fig. 1). This allele was present in two African assemblies (YRI, NA19240), in one American assembly (ASM311317v1) and in the gorilla genome whereas absent in other assemblies. Moreover, the reference allele was found in chimpanzees, indicating the presence of ancestral polymorphism at this locus. Furthermore, this alternate allelic sequence was not reported in the-NCBI- nr/nt database or in the human HLA database (IPD-IMGT/HLA database), suggesting that this alternate sequence represents a putative novel MHC allele.

Among the alternate alleles, 1,348 intersected with the genic region of 1,143 protein-coding genes. The genomic distribution of the alternate alleles was dispersed throughout the genome, and those belonging to non-private sequences are shown in Fig. 5. A total of 59 non-private alternate alleles intersected with

the genic region, and five of them further intersected with the CDS region of eight genes: HLA-W, MICD, HCG9, DDX39BP2, LOC107985837, ZNF880, PLIN4 and LOC728715.

Only 2.6% (80 out of 3,041) of the alternate alleles have been identified before but were miss-classified as insertions. Therefore, most of the identified alternate alleles described in our study are newly reported. The discovery frequency of alternate alleles in human assemblies appears to be lower than that of insertions (Fig. 3a and 3b). Most alternate alleles were present in less than half of the 31 assemblies, indicating that many of them represent minor alleles in the human genome. Nevertheless, 144 alternate alleles were non-private (Additional table 4), with the longest one found in 17 assemblies and spanning 19,330 bp (genomic placement position: chr7, 62408641-62451864). The ambiguous sequences also included a number of putative insertions and alternate alleles (Fig. 3c) and deserve further efforts for verification. The length distributions of the alternate alleles and insertions did not differ (Fig. 3d).

Similar to the insertions, PCA also showed that PC1 clusters African samples away from other populations, while PC2 separates the East Asians from the Europeans, South Asians and Americans (Fig. 3f). We also explored the potentially transcribed sequences that either have mapped RNA-seq reads (≥ 10 reads in at least two samples) or hits to the human expressed sequence tag database (dbEST) (e-value $< 1e-5$). We totally identified 74 transcribed alternate alleles from RNA-seq data and 238 with hits to the human dbEST, resulting in a total of 244 potentially transcribed sequences (Additional table 5). One alternate allele was found to be expressed in a tissue-specific manner (Fig. 4b and 4c), which is potentially a long non-coding gene since we couldn't annotate it to any known gene. The putative novel MHC allele was also found to be expressed (Fig. 4e and 4f).

To explore the origin of the alternate alleles, we analyzed associated repeats with NRS. Transposable elements (TEs) compose approximately 45% of the human genome (Lander *et al.* 2001), and a previous study showed that insertions were associated with TEs (Wong *et al.*, 2018). The percent of TEs within alternate alleles (10.0%) was much lower than that of insertions (55.1%) (Fig. 6a). The flanking sequences (5 kb on each side) of the alternate alleles also had less TE content (33.3%) than the insertions (48.0%) (Fig. 6b), suggesting that the alternate alleles are not associated with TEs. We then screened the tandem repeat content among the sequences. The alternate alleles possessed a much higher content of tandem repeats either within the sequences (Fig. 6c) or in the flanking sequences (5 kb on each side, Fig. 6d) compared with the insertions. Notably, the reference allele also be enriched in tandem repeat when the alternate allele have a large content of tandem repeat ($R^2=0.65$, Fig. 6e and Additional fig. 3), thereby implying that tandem repeats are associated with alternate alleles.

Discussion

A comprehensive characterization of structural variations is essential for studies attempting to identify causative variants that affect phenotypic variations and complex genetic diseases. In this study, we identified 129.1 Mb candidate NRS (4.2% of the genome). Although many of the NRS were singletons, a considerable number of reliable NRS (29.3 Mb) were identified by their presence in at least two

assemblies, and most of them were newly described. The discovery of these NRS will contribute to a final, comprehensive pan-genome capturing all of the DNA present in humans.

More importantly, we discovered a large number of alternate alleles. The majority of the alternate alleles that we found have not been previously reported, which could be due to several reasons: (1) Most previous work has designed their studies to focus on insertions, whereas other types of NRS were largely ignored (Kehr et al., 2017; Sherman et al., 2018; Wong et al., 2018); (2) Many studies have mainly relied on short-reads data to obtain NRS (Kehr et al., 2017; Sherman et al., 2018), which would be less efficient for the discovery of long structural variations compared with an assembly-versus-assembly approach, as applied in our study; and (3) Many alternate alleles were singletons, suggesting that they are either of very low frequency for detection or simply false positives due to assembly errors. Nevertheless, we still detected 144 non-private alternate alleles. The current human genome (GRCh38.p12) includes 261 alternate loci that capture a limited amount of population diversity and improve read mapping for some data sets (Schneider et al., 2017). Therefore, the sequences that we identified will greatly advance our knowledge of the alternate alleles in the human genome.

There is growing interest in using genetic variants to augment the reference genome into a graph genome (Crysnanto et al., 2019; Pritt et al., 2018; Rakocevic et al., 2019). To create a representative graph genome, the full spectrum of structural variations, including the alternate alleles, should be understood clearly. With the reduction in sequencing costs and advances in sequencing technology, increased numbers of de novo assemblies will be generated, which will eventually refine the full spectrum of sequence diversities in the human genome.

References

- Audano, P.A., Sulovari, A., Graves-Lindsay, T.A., Cantsilieris, S., Sorensen, M., Welch, A.E., Dougherty, M.L., Nelson, B.J., Shah, A., and Dutcher, S.K. (2019). Characterizing the major structural variant alleles of the human genome. *Cell* 176, 663-675. e619.
- Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., and Madden, T.L. (2009). BLAST+: architecture and applications. *BMC Bioinformatics* 10, 421.
- Cho, Y.S., Kim, H., Kim, H.-M., Jho, S., Jun, J., Lee, Y.J., Chae, K.S., Kim, C.G., Kim, S., and Eriksson, A. (2016). An ethnically relevant consensus Korean reference genome is a step towards personal reference genomes. *Nature Communications* 7, 13637.
- Church, D.M., Schneider, V.A., Steinberg, K.M., Schatz, M.C., Quinlan, A.R., Chin, C.S., Kitts, P.A., Aken, B., Marth, G.T., Hoffman, M.M., *et al.* (2015). Extending reference assembly models. *Genome Biology* 16, 13.
- Crysnanto, D., Wurmser, C., and Pausch, H. (2019). Accurate sequence variant genotyping in cattle using variation-aware genome graphs. *BioRxiv*, 460345.

English, A.C., Salerno, W.J., Hampton, O.A., Gonzaga-Jauregui, C., Ambreth, S., Ritter, D.I., Beck, C.R., Davis, C.F., Dahdouli, M., and Ma, S. (2015). Assessing structural variation in a personal genome-towards a human reference diploid genome. *BMC Genomics* 16, 286.

Fagerberg, L., Hallstrom, B.M., Oksvold, P., Kampf, C., Djureinovic, D., Odeberg, J., Habuka, M., Tahmasebpour, S., Danielsson, A., Edlund, K., *et al.* (2014). Analysis of the human tissue-specific expression by genome-wide integration of transcriptomics and antibody-based proteomics. *Molecular & Cellular Proteomics* 13, 397-406.

Fu, L., Niu, B., Zhu, Z., Wu, S., and Li, W. (2012). CD-HIT: accelerated for clustering the next-generation sequencing data. *Bioinformatics* 28, 3150-3152.

Jain, M., Koren, S., Miga, K.H., Quick, J., Rand, A.C., Sasani, T.A., Tyson, J.R., Beggs, A.D., Diltthey, A.T., and Fiddes, I.T. (2018). Nanopore sequencing and assembly of a human genome with ultra-long reads. *Nature biotechnology* 36, 338.

Kehr, B., Helgadottir, A., Melsted, P., Jonsson, H., Helgason, H., Jonasdottir, A., Jonasdottir, A., Sigurdsson, A., Gylfason, A., Halldorsson, G.H., *et al.* (2017). Diversity in non-repetitive human sequences not found in the reference genome. *Nature Genetics* 49, 588.

Kielbasa, S.M., Wan, R., Sato, K., Horton, P., and Frith, M. (2011). Adaptive seeds tame genomic sequence comparison. *Genome Research*, gr. 113985.113110.

Kim, D., Langmead, B., and Salzberg, S.L. (2015). HISAT: a fast spliced aligner with low memory requirements. *Nature Methods* 12, 357-360.

Lander, E.S., Linton, L.M., Birren, B., Nusbaum, C., Zody, M.C., Baldwin, J., Devon, K., Dewar, K., Doyle, M., FitzHugh, W., *et al.* (2001). Initial sequencing and analysis of the human genome. *Nature* 409, 860.

Levy, S., Sutton, G., Ng, P.C., Feuk, L., Halpern, A.L., Walenz, B.P., Axelrod, N., Huang, J., Kirkness, E.F., Denisov, G., *et al.* (2007). The Diploid Genome Sequence of an Individual Human. *PLOS Biology* 5, e254.

Li, R., Li, Y., Zheng, H., Luo, R., Zhu, H., Li, Q., Qian, W., Ren, Y., Tian, G., and Li, J. (2010). Building the sequence map of the human pan-genome. *Nature biotechnology* 28, 57-63.

Pendleton, M., Sebra, R., Pang, A.W.C., Ummat, A., Franzen, O., Rausch, T., Stütz, A.M., Stedman, W., Anantharaman, T., and Hastie, A. (2015). Assembly and diploid architecture of an individual human genome via single-molecule technologies. *Nature methods* 12, 780.

Pritt, J., Chen, N.-C., and Langmead, B. (2018). FORGe: prioritizing variants for graph genomes. *Genome Biology* 19, 220.

Quinlan, A.R. (2014). BEDTools: the Swiss-army tool for genome feature analysis. *Current protocols in bioinformatics* 47, 11.12. 11-11.12. 34.

Rakocevic, G., Semenyuk, V., Lee, W.-P., Spencer, J., Browning, J., Johnson, I.J., Arsenijevic, V., Nadj, J., Ghose, K., Suci, M.C., *et al.* (2019). Fast and accurate genomic analyses using genome graphs. *Nature Genetics* 51, 354-362.

Schneider, V.A., Graves-Lindsay, T., Howe, K., Bouk, N., Chen, H.-C., Kitts, P.A., Murphy, T.D., Pruitt, K.D., Thibaud-Nissen, F., Albracht, D., *et al.* (2017). Evaluation of GRCh38 and de novo haploid genome assemblies demonstrates the enduring quality of the reference assembly. *Genome Research* 27, 849-864.

Sherman, R.M., Forman, J., Antonescu, V., Puiu, D., Daya, M., Rafaels, N., Boorgula, M.P., Chavan, S., Vergara, C., Ortega, V.E., *et al.* (2018). Assembly of a pan-genome from deep sequencing of 910 humans of African descent. *Nature Genetics* 51.

Shi, L., Guo, Y., Dong, C., Huddleston, J., Yang, H., Han, X., Fu, A., Li, Q., Li, N., and Gong, S. (2016). Long-read sequencing and de novo assembly of a Chinese genome. *Nature Communications* 7, 12065.

Tarasov, A., Vilella, A.J., Cuppen, E., Nijman, I.J., and Prins, P. (2015). Sambamba: fast processing of NGS alignment formats. *Bioinformatics* 31, 2032-2034.

Wong, K.H.Y., Levy-Sakin, M., and Kwok, P.Y. (2018). De novo human genome assemblies reveal spectrum of alternative haplotypes in diverse populations. *Nature Communications* 9, 9.

Declarations

Compliance and ethics

The authors declare that they have no competing interests.

Acknowledgements

This study was supported by research grants from the National Natural Science Foundation of China (31822052) to Y.J., National Natural Science Foundation of China (31802027), Doctoral Fund of Ministry of Education of China (No. 2018M631209), and “the Fundamental Research Funds for the Central Universities” (2452018127) to R.L. We also thank the High Performance Computing platform of Northwest A&F University.

Data availability

The datasets supporting the conclusions of this article are included within the article and its additional files. The identified NRS is provided as a Additional file (Additional table 7) and all other data supporting the findings of this study are available in its additional files.

Figures

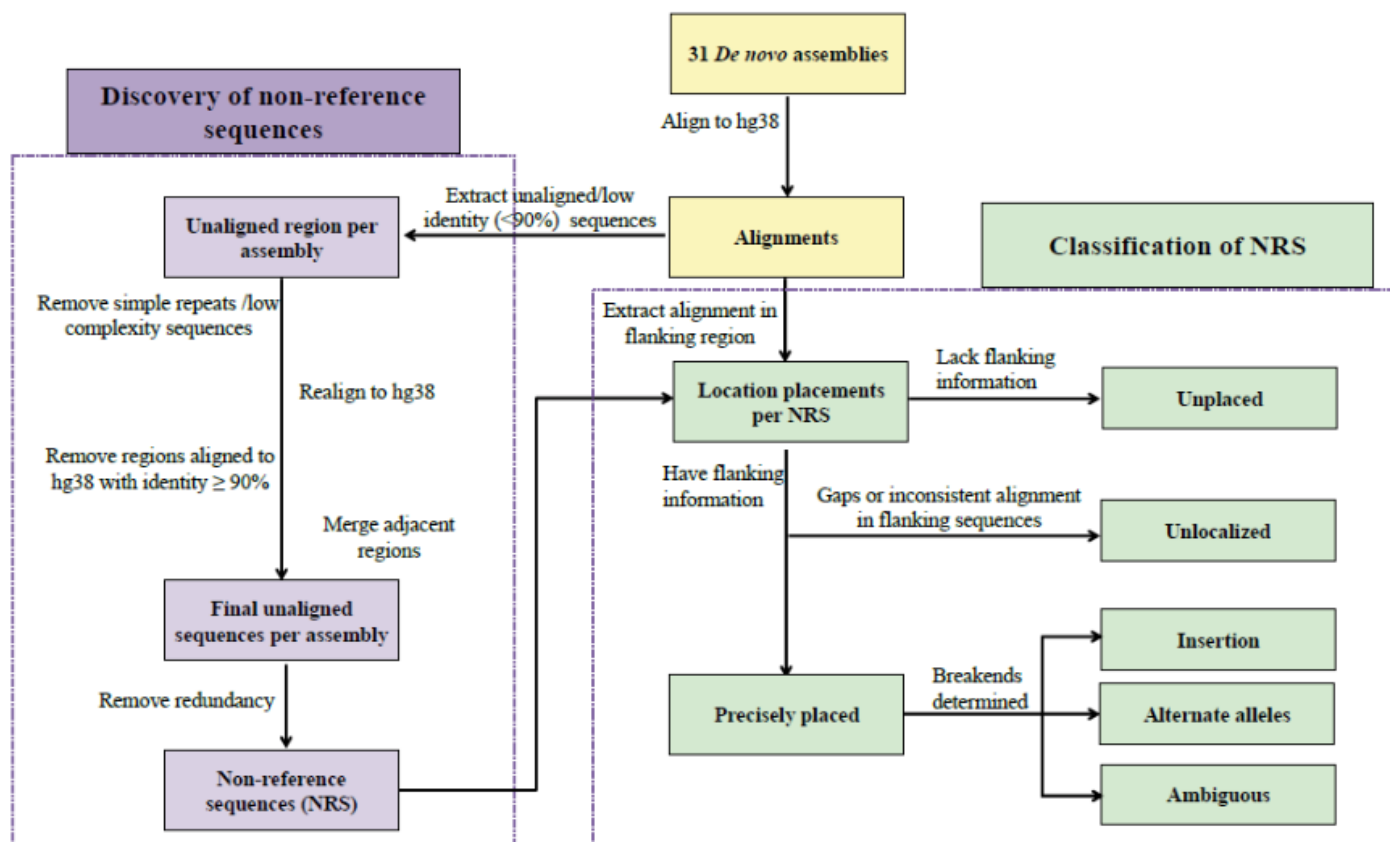


Figure 1

Overview of the workflow to identify non-reference sequences (NRS).

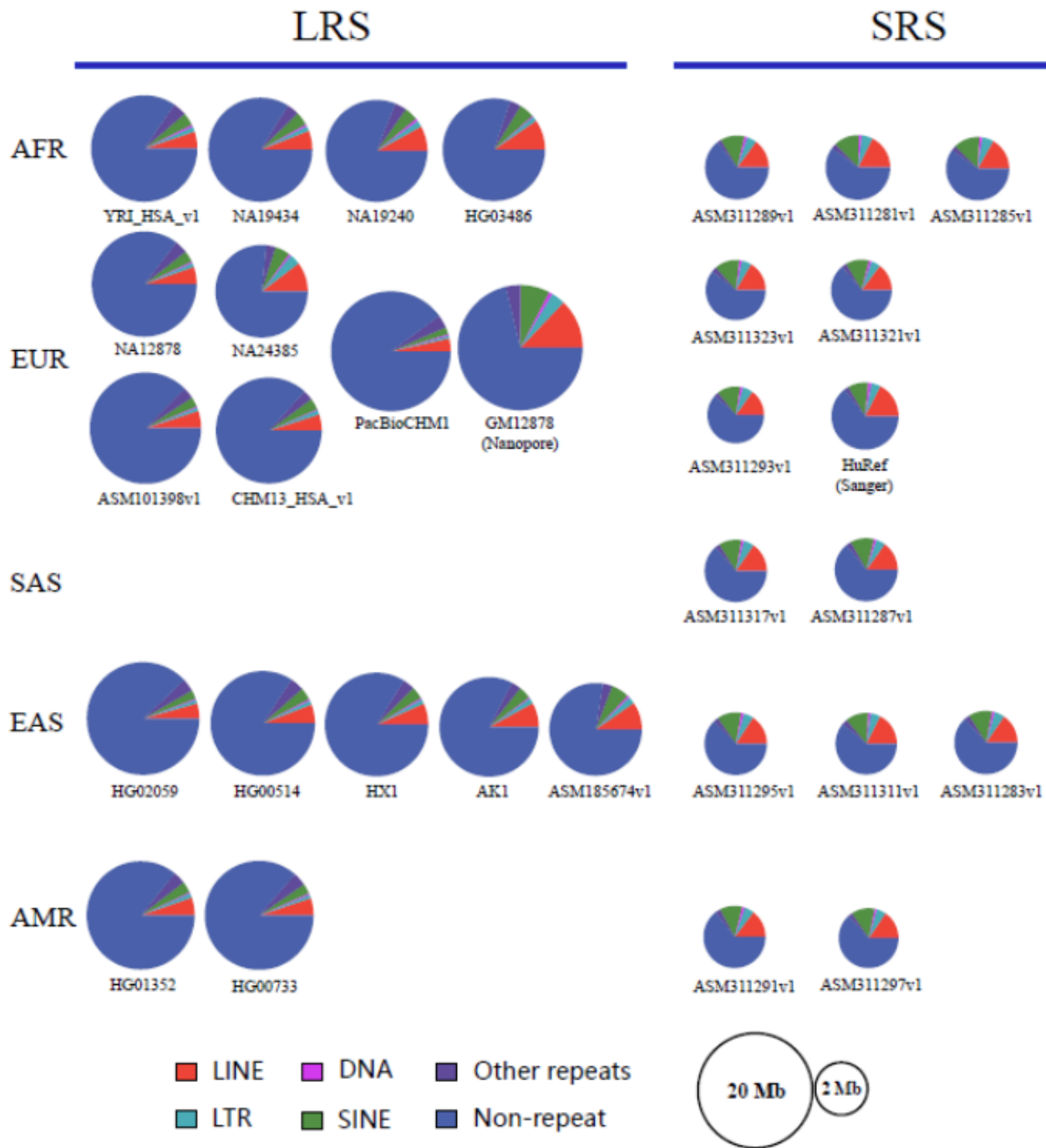


Figure 2

Non-reference sequences identified from each of the 31 human de novo assemblies. The repeat information was summarized from the repeat annotation files (*_rm.out.gz), which were generated with RepeatMasker and downloaded from the NCBI. The radius of each pie chart was log2 transformed.

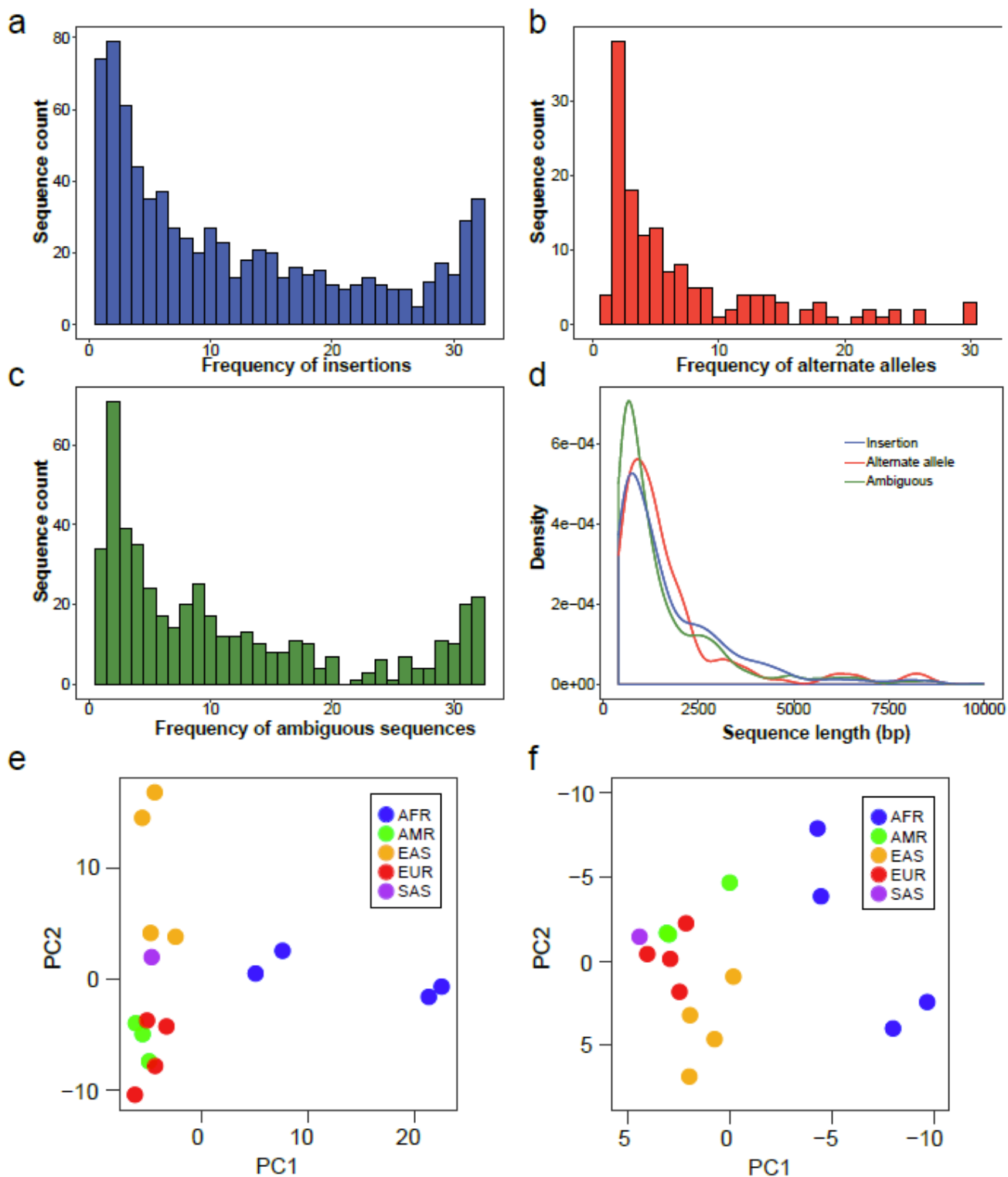


Figure 3

Characteristics of the insertions and alternate alleles. (a) Frequency of insertions within the 31 de novo assemblies; (b) Frequency of alternate alleles within the 31 de novo assemblies; (c) Frequency of ambiguous sequences within the 31 de novo assemblies; (d) Length distributions of the insertions, alternate alleles and ambiguous sequences; (e) The first two principal components based on the occurrence matrix of the insertions among the 16 de novo assemblies of Pacbio sequencing; (f) The first

two principal components based on the occurrence matrix of the alternate alleles among the 16 de novo assemblies of Pacbio sequencing.

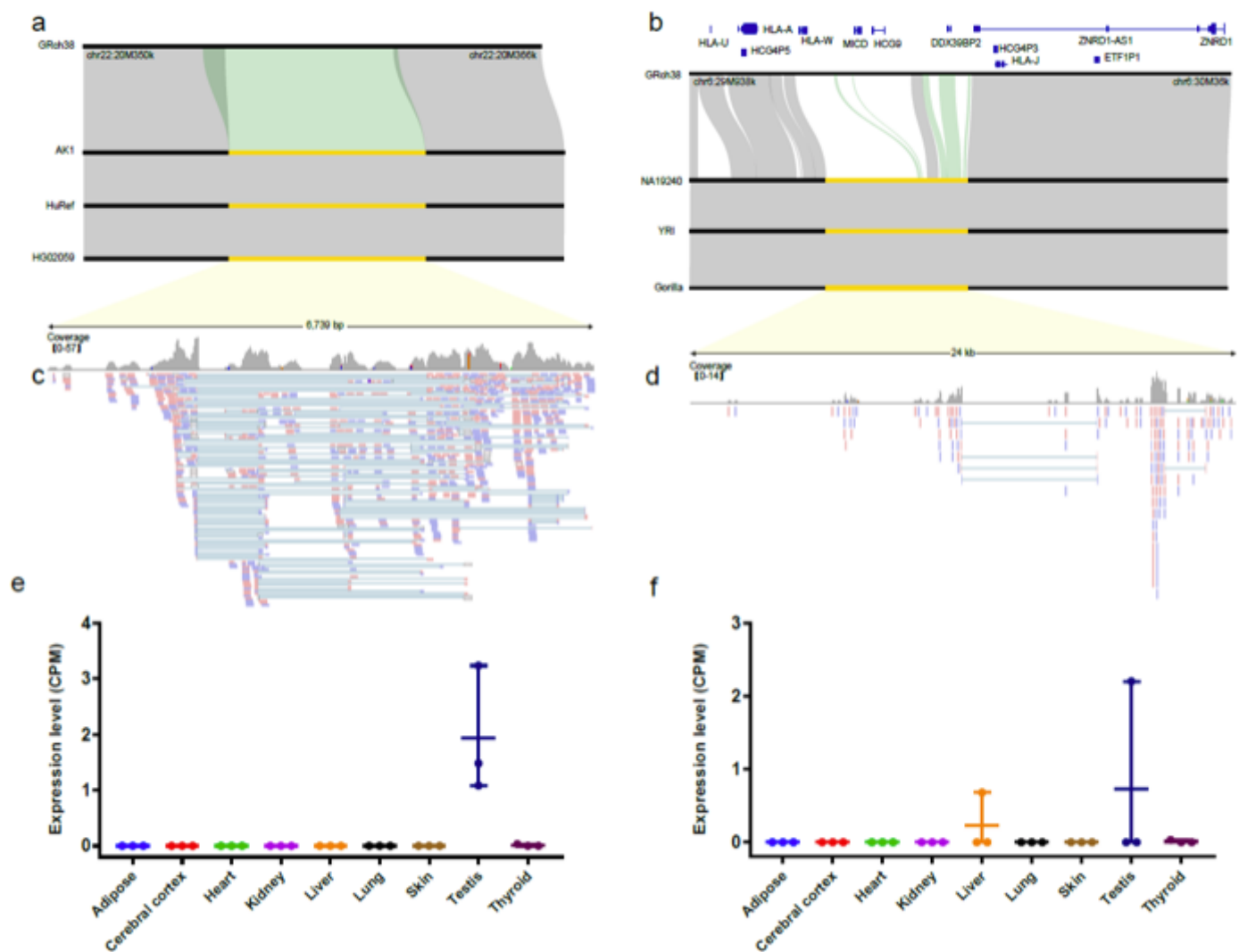


Figure 4

Examples of alternate alleles. (a-b) Alignment of the alternate alleles and their flanking sequences with hg38 and other assemblies. The blue lines in the top represent the gene annotations in hg38. The yellow segments represent the NRS. The gray block represents alignments that share $\geq 95\%$ identity, while the green block represents alignments that share $< 90\%$ identity. For each of the alignments, the reference sequence from hg38 is shown at the top following by the sequence where the NRS is derived from and sequences from two additional genomes. (c-d) RNA-seq reads mapping of the NRS shows expression potential. (e-f) Expression of the alternate alleles in nine tissues. Three replicates were used for each tissue. The expression level was measured using CPM (reads count per million of total mapped reads).

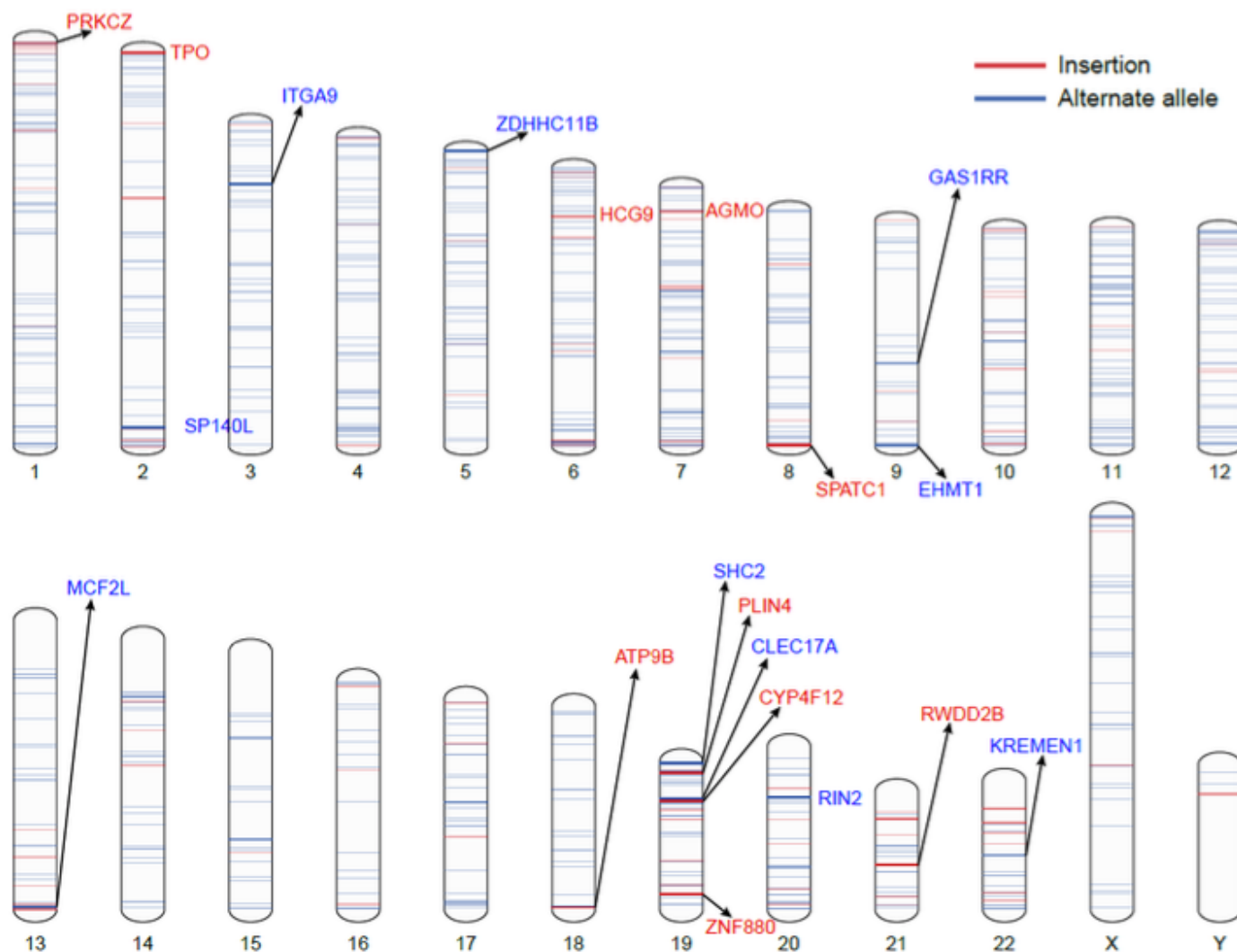


Figure 5

Locations of the non-private insertions and alternate alleles on the human reference genome (hg38). Red lines represent insertions, while blue lines represent alternate alleles. The gene symbol is shown for each of the ten longest insertions and ten longest alternate alleles overlapping genic regions. Line width is not drawn to scale.

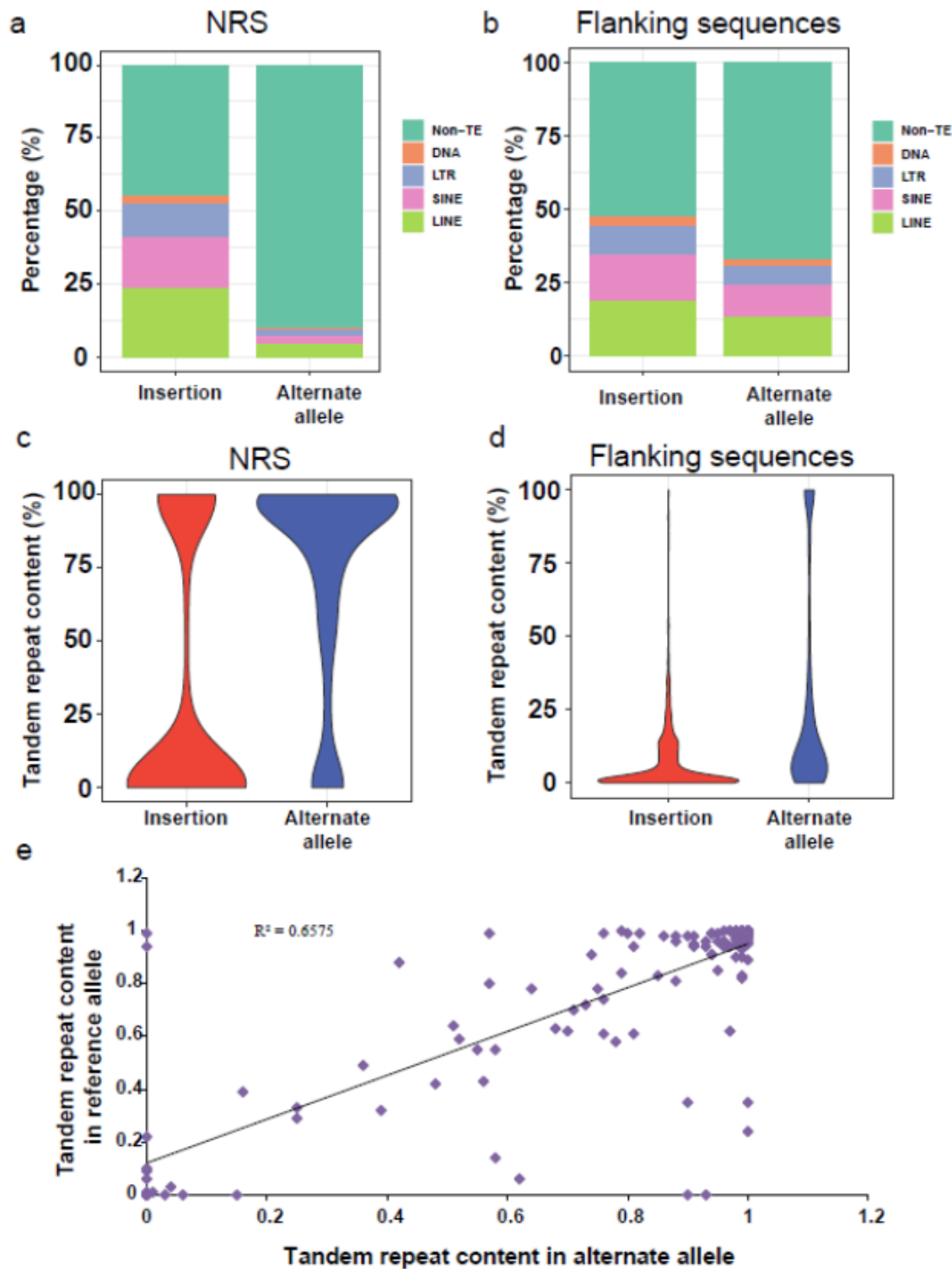


Figure 6

The repeat contents associated with the insertions and alternate alleles. (a) TE content within the insertions and alternate alleles; (b) TE content in the flanking region (5 kb on each side); (c) The tandem repeat content within the insertions and alternate alleles; (d) The tandem repeat content in the flanking region (5 kb on each side). The non-private sequences were included for statistics; (e) Dot plot showing the tandem repeat content in alternate allele (x axis) and in the corresponding reference allele.

Supplementary Files

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

- [supplement1.xls](#)
- [supplement2.xls](#)
- [supplement3.pdf](#)
- [supplement4.xls](#)
- [supplement5.pdf](#)
- [supplement6.xls](#)
- [supplement7.xls](#)
- [supplement8.xls](#)
- [supplement9.xls](#)