

Transposable element-mediated rearrangements are prevalent in human genomes

Parithi Balachandran, Isha A. Walawalkar, Jacob I. Flores, Jacob N. Dayton, Peter A. Audano & Christine R. Beck

Author Information

These authors contributed equally, Isha A. Walawalkar and Jacob I. Flores

Affiliations

The Jackson Laboratory for Genomic Medicine, Farmington CT 06032

Parithi Balachandran, Isha A. Walawalkar, Jacob I. Flores, Jacob N. Dayton, Peter A. Audano and Christine R. Beck

Department of Genetics and Genome Sciences, University of Connecticut Health Center, Farmington, CT 06030

Christine R. Beck

Institute for Systems Genomics, University of Connecticut, Storrs, CT 06269

Christine R. Beck

Contributions

C.R.B conceived of and supervised the study. P.B performed data preprocessing, SV calling, quality control, method development, data analysis and data interpretation. P.B, I.A.W, J.I.F and J.N.D designed primers for validation work. I.A.W, J.I.F and J.N.D performed PCR and Sanger sequencing. P.B and I.A.W performed junction reconstruction work. J.I.F designed and performed ddPCR experiment. P.A.A provided homology data for method development. P.B and C.R.B wrote the manuscript with input from all coauthors. All contributing authors have read, edited, and approved the manuscript.

Acknowledgements

We thank the members of the Beck lab for reading and editing the manuscript, in particular Alex V. Nesta and Ardian Ferraj. We also thank members of the Charles Lee lab and Miriam Konkel for their critical review of the manuscript. We thank Guruprasad Ananda for his advice on SV filtering and merging. We would also like to thank the Cellular Engineering group at The Jackson Laboratory, especially Mallory Ryan and Sofia Giansiracusa, for their help with ddPCR experiments. This work was supported in part by the National Institutes of Health grants R00GM120453 and R35GM133600 from the NIGMS and P30CA034196 from the NCI and startup funds from the University of Connecticut Health and The Jackson Laboratory to Christine R. Beck.

Corresponding author

Christine R. Beck

ABSTRACT

Transposable elements constitute about half of human genomes, and their role in generating human variation through retrotransposition is broadly studied and appreciated. Structural variants mediated by transposons, which we call transposable element-mediated rearrangements (TEMRs), are less well studied, and the mechanisms leading to their formation as well as their broader impact on human diversity are poorly understood. Here, we identify 493 unique TEMRs across the genomes of three individuals. While homology directed repair is the dominant driver of TEMRs, our sequence-resolved TEMR resource allows us to identify complex inversion breakpoints, triplications or other high copy number polymorphisms, and additional complexities. TEMRs are enriched in genic loci and can create potentially important risk alleles such as a deletion in *TRIM65*, a known cancer biomarker and therapeutic target. These findings expand our understanding of this important class of structural variation, the mechanisms responsible for their formation, and establish them as an important driver of human diversity.

INTRODUCTION

The development of a genome reference sequence over 20 years ago has driven greater knowledge of the number, effects, and variation of transposable elements (TEs) across human genomes¹. The current haploid human reference contains over 4 million annotated TEs^{2,3}, and by computational estimates TEs occupy up to two-thirds of human genome^{4,5}. TEs can alter the structure of genomes through transposition (*de novo* insertion)⁶, polymorphism⁷⁻¹⁰, transduction¹¹⁻¹⁴, and transposition-associated rearrangements¹⁵⁻¹⁸. In addition to retrotransposition, TE copies can play major roles in shaping human genomes, promoting polymorphism, and contributing to genomic instability.

Two homologous TEs can act as substrates for ectopic DNA repair resulting in structural variant (SV) formation we collectively call TE-mediated rearrangements (TEMRs). TEMRs are responsible for deletion of ~850 kbp of the human reference genome when compared to the chimpanzee draft genome^{1,19,20}, and were recently found to account for a loss of ~80 kbp in a Korean genome²¹. TEMRs are implicated in the expansion of segmental duplications in human genomes²². TEMRs are also associated with deletions of exonic regions, leading to cancer predisposition syndromes^{23,24}, and Mendelian diseases²⁵⁻²⁷. Furthermore, complex genomic rearrangements often harbor TEs at breakpoint junctions, and these TEMRs carry similar characteristics as simple deletions^{28,29}. Studying the genes and the genomic context of TEMR-associated instability can define the role these SVs play in mediating human disease and identify loci prone to this rearrangement mechanism^{30,31}. Although these analyses were important to establish the extent of TEMRs in comparative genomics and disease, the prevalence of TEMRs in SV callsets, the mechanisms that cause these rearrangements, the allele frequency of disease associated variants, and the full spectrum of how TEs contribute to SV formation are still poorly understood.

The vast majority of current studies investigating TEMR mechanisms point to non-allelic homologous recombination (NAHR) generating deletions in human genomes^{30,32}. Although NAHR dominates existing genome-wide TEMR studies based on reference genomes^{19,20}, a few cell-culture based systems have been derived to address mechanisms of TEMRs; these studies have indicated NAHR as well as non-homologous end joining (NHEJ)³³, single-strand annealing (SSA)^{34,35}, and microhomology-mediated end joining (MMEJ)³³ in the formation of these classes of SVs. Importantly, many of these studies indicate that more diverged repeats are unlikely to generate junctions in homologous regions of the two TEs, and these data are not reflected in human rearrangements³¹. Finally, a number of recent studies have detailed rearrangements with TEs at the junctions of inversions, duplications, and complex genomic rearrangements; these studies have indicated that repair by microhomology-mediated break

induced replication (MMBIR)^{28,31,36,37} can mediate diverse classes of TEMR. These mechanistic interrogations are impacted by the fact that most SV studies to date either lack precision or resort to hand curation of breakpoint junctions to determine potential mechanisms of formation^{19,20,31}. Thus, distinguishing the mechanisms of TEMR across large numbers of events would be prohibitive, and determining the scope of TEMR in generating human diversity has been difficult to assess. The widespread investigation of TEMRs across genomes allows an unbiased view of the diverse types of resultant SVs and the different mechanisms that drive TEMR; this type of approach has not been applied to genome wide analyses and will help interrogate how TEs lead to genomic instability and variation.

Identification of TEMRs and accurate characterization of their breakpoint junctions is challenging. Short-read sequencing (SRS) is able to identify most deletions, including TEMRs, outside of simple repeats and segmental duplications³⁸⁻⁴⁰ even with its inherent limitations⁴¹. Long-read sequencing (LRS) has overcome many of the challenges faced by SRS in identifying SVs across genomes^{38,42}. Recent advances in LRS have enabled long and contiguous *de novo* assembly of genomes and detection of SVs with precise junctions within repetitive regions of the genome^{16,43}. Although a majority of publicly available datasets are still SRS data⁴⁴⁻⁴⁶, LRS data is now in production for thousands of human genomes ([Human Genome Structural Variation Consortium \(HGSVC\)](#), [Human Pangenome Reference Consortium \(HPRC\)](#), [Solving the Unsolved Rare Disease \(Solve-RD\)](#) and [AllofUs](#)). Because LRS still has significant limitations including cost and input DNA quantities, we sought to develop a comprehensive TEMR identification method that using either short-read (Illumina) or long-read (PacBio CLR) sequencing data for discovery and characterization of these events. This approach is designed to scale across many samples and will enable high-throughput mechanistic inference for a growing number of sequenced genomes.

By applying our approach to three diverse genomes³⁸, we identified 493 nonredundant TEMRs using SRS and LRS. We ascertained the precise junctions of 70 randomly selected TEMRs with PCR and Sanger sequencing and verified the accuracy of breakpoint junctions from matching phased HiFi genome assemblies¹⁶. This allowed us to discern the precise junctions for all 493 TEMR events and infer mechanisms involved in the formation of TE-driven deletions, duplications, and inversions. Previously, the role of TEs in generating duplications and inversions was not appreciated. Additionally, we identified TEs mediating higher order amplifications and complex SVs indicating the range of rearrangements driven by TEs. *Alu* elements are a major (80.5%) contributor to TEMRs, primarily via homologous recombination mechanisms; yet the length of the homology at the junction of these *Alu* TEMRs are shorter (median of 15 bp) than what we expect from a traditional non-allelic homologous

recombination event (>100 bp). We show that TEs not only affect the genome through retrotransposition but are also a substrate for widespread rearrangements creating 600 kbp of structural alterations per human genome. Because TEMRs disproportionately affect genes, they are an important source to study phenotypic variation, disease, and human evolution.

RESULTS

IDENTIFYING TRANSPOSABLE ELEMENT-MEDIATED REARRANGEMENTS

To call SVs genome-wide, we analyzed SRS using Manta⁴⁷, LUMPY⁴⁸, and DELLY⁴⁹ and LRS using pbsv (<https://github.com/PacificBiosciences/pbsv>), Sniffles⁵⁰ and SVIM⁵¹. We generated a consensus callset using these individual SRS and LRS SV callers and ensemble heuristics to maximize accuracy [Methods]. We implemented our pipeline on Illumina and PacBio CLR data across three well-characterized diverse individuals: Puerto Rican HG00733 (PUR), Southern Han Chinese HG00514 (CHS), and Yoruban NA19240 (YRI)³⁸. We filtered these calls by setting thresholds for alignment support and excluding SVs in simple repeats. Additionally, we obtained phased HiFi assemblies for these three samples¹⁶, and used a new version of PAV¹⁶ for breakpoint homology. We merged the calls from all these methods and across all three individuals into a single nonredundant high-confidence callset of 5,297 SVs comprised of 4,997 deletions, 239 duplications and 61 inversions, with an average of 3,111 SVs per individual [Methods].

SVs with both breakpoints in different TE elements of the same class were categorized as TEMRs [Methods]. From our high-confidence callset of 5,297 SVs, we identified 493 nonredundant TEMRs (9.3%) across all three individuals [Fig. 1a]. We identified an average of 250 TEMRs per sample (225 from PUR, 229 from CHS, and 296 from YRI) and they collectively affected an average of 591 kbp per sample. The 493 TEMRs consisted of 397 (80.5%) *Alu* TEMRs and 96 (19.5%) LINE-1 TEMRs [Fig. 1a]. Interestingly, although the 90.3% (445) TEMRs were deletions, we also identified 33 duplications and 15 inversions, classes of TEMRs that were not previously surveyed in normal human genomes.

The size of TEMRs varies greatly from 93 bp to 25,425 bp with a median length of 1,342 bp [Supplementary Fig. 1], and these events are longer than the non-TEMRs (median 322 bp). Polymorphic MEIs account for a large proportion of SV calls (median 317 bp), upon excluding them we still observed a significant difference in median lengths between TEMR deletions and non-TEMR deletions (1,345 bp vs 537 bp; $p < 0.01$, Welch's *t*-test). TEMR duplications were also longer than non-TEMR duplications (1,085 bp vs 288 bp) but did not reach statistical significance. Conversely, TEMR inversions were significantly smaller than non-TEMR inversions (2,454 bp vs 8,870 bp, $p < 0.001$,

Welch's *t*-test), which are generally mediated by larger segmental duplications⁵² [Supplementary Fig. 2].

We examined the overall genomic architecture of TEMRs and found that 91% of deletions and duplications had junctions in TEs in the same/direct orientation and all inversions had junctions in TEs in the opposite/inverted orientation [Fig. 1b]. We found that *Alu* TEMRs (median length of 1,163 bp) are typically shorter than LINE-1 TEMRs (median length of 4,469 bp) ($p \leq 1e^{-5}$, Welch's *t*-test) [Fig. 1c], as expected³². Full-length LINE-1 elements (6 kbp)⁵³ are almost 20-fold longer than full-length *Alu* elements (300 bp)⁵⁴ and they therefore provide a longer substrate for recombination.

CHARACTERIZATION OF TEMR BREAKPOINT JUNCTIONS

To identify the breakpoint junctions of TEMRs with nucleotide accuracy, we randomly selected 70 TEMRs (66 deletions and 4 duplications) for PCR and Sanger sequencing [Supplementary Table S1]. By manually reconstructing TEMR breakpoint junctions with Sanger sequences [Fig 1d], we found that 55.7% (39) of TEMRs had breakpoint homology of at least 5 bp, 27.1% (19) had 1-4 bp breakpoint homology, and 17.1% (12) showed no breakpoint homology. A higher percentage and length of breakpoint homology is expected since TEMRs are mediated by homologous sequences.

We then compared the breakpoint junctions identified by all the SV callers used in this study with the corresponding manually curated junctions. Of the 58 out of 70 events with breakpoint homology, PAV calls from HiFi phased assemblies supported 57 (98.3%) events and the imprecise junction for the one remaining event was due to a nearby SNV. For the 12 TEMRs without breakpoint homology, PAV called identical breakpoints for six TEMRs, and the remaining six TEMRs had inaccurate breakpoint junctions due to the presence of indels at the junction. We observed that Manta was able to identify the indels present at the junctions for these six TEMRs and accurately call the breakpoints. Among the read-based callers, Manta had the highest breakpoint precision of 88% (51 out of 58) for TEMRs with breakpoint homology and 91% (11 out of 12) for TEMRs without breakpoint homology, although Manta failed to discover three TEMRs. Interestingly, when Manta identifies a homology or an indel at the breakpoint junction of an event, the junction was 100% precise. The breakpoint precision statistics for all SV callers used in this study can be found in Supplementary Table S2. Using these results as a guide, we annotated breakpoint homology for deletions and duplications in our nonredundant TEMR callset using Manta, when available, and HiFi assembly callset otherwise. Since inversions had imprecise breakpoint junctions from SV callers, we used the manually curated junctions with breakpoint homology for the 15 inversions in our callset.

No single SV caller was able to precisely identify breakpoint junctions and junction homologies for all 70 TEMRs. When we extracted and inspected the DNA sequence from HiFi assemblies around the seven TEMRs with inaccurate breakpoint junctions [Methods], we were able to successfully identify the SNVs and indels that were causing inaccuracies. Because of limitations imposed by SRS sensitivity and methods to accurately place breakpoints with LRS, no single approach is able to detect and accurately characterize breakpoints for all TEMRs.

TRANSPOSABLE ELEMENTS MEDIATE SVs BY DISTINCT MECHANISMS

To investigate the DNA repair mechanisms involved TEMR formation, we mapped TEs flanking TEMR events against their corresponding consensus sequences (*Alu*⁵⁴ and L1.3⁵³), identified the position of 5' and 3' breakpoints, and annotated breakpoint homology [Fig. 2a] [Methods]. We categorized TEMRs as products of homologous recombination (HR) if there was a significant overlap within homologous locations in the consensus repeat sequence, and the overlap is identical to the breakpoint homology [Fig. 2b]. Otherwise, they were categorized as products of non-homologous repair events (NHE; both end-joining and replication-based mechanisms) [Fig. 2c]. We were able to systematically categorize 90.5% of the callset. The remaining 9.5% (47) of TEMRs that required manual inspection were comprised of events with breakpoints in the TE poly-A tail (long sequence of adenine nucleotides in plus strand / poly-T on minus strand) and truncated TEs.

Of all 493 TEMRs, 390 (79.1%) were categorized as HR (TEMR-HR, 354 *Alu* and 36 LINE-1) and 103 (20.9%) were categorized as NHE (TEMR-NHE, 43 *Alu* and 60 LINE-1) [Fig. 2d]. *Alu* elements (~300 bp) are composed of two ~150 bp homologous monomers, and we identified 10 direct *Alu* TEMRs with breakpoints in homologous regions of different monomers resulting in one chimeric *Alu* element with a single monomer (~150 bp) and nine chimeric *Alu* elements with 3 monomers (~450 bp) [Supplementary Fig. 3]. While this imbalance appears to be significant ($p = 0.021$, two-tailed Binomial test), the mechanism responsible for a longer chimera preference is unclear.

We next examined the GC content of the microhomology greater than 5 bp at breakpoint junctions and found microhomologies of *Alu* TEMRs to be significantly enriched for GC content (median of 57.1%) compared to the *Alu* sequences in GRCh38 (median of 51.5%; $p < 0.001$, Welch's *t*-test) and the whole GRCh38 (median of 39.7%; $p < 1e^{-5}$, Welch's *t*-test) [Fig. 2d]. We found no significant enrichment for GC content at the microhomology of LINE-1 TEMRs (median of 37.1%) compared to the LINE-1 sequences in GRCh38 (median of 33.5%). Higher GC content indicates a stronger strand annealing between G-C (three hydrogen bonds) compared to A-T (two hydrogen bonds) leads to more stability for repair intermediates, as was previously proposed³¹. Overall, the variation in the size of *Alu* and

LINE-1 TEMRs and the divergence of G-C content at the breakpoints of HR events mediated by the two TEs may mean that the mechanisms governing TEMRs can be dependent on the repeat class and/or size of the homology tract leading to stability of annealing prior to repair.

CHARACTERISTICS OF TEMR BREAKPOINTS

Microhomologies were present at the junctions of 446 TEMRs varying from 1 bp to 307 bp with a median of 13 bp [Fig. 2e]. TEMR-HRs have a median microhomology length of 17 bp and 93.3% of them have microhomologies 5 bp or longer. In contrast, TEMR-NHEs have a median microhomology length of 1 bp and nearly 96% of them have microhomologies shorter than 5 bp. Additionally, we observed the presence of small insertions ranging from 1 bp to 23 bp in 35% of TEMR-NHE events (median size of 2 bp), a known signature of end-joining repair mechanisms including micro-homology mediated end joining (MMEJ)⁵⁵. We found *Alu* TEMRs have a shorter median microhomology length of 16 bp compared to 34 bp for LINE-1 TEMRs. We examined the *Alu*, and LINE-1 elements involved in TEMRs and found that 36 *Alu* subfamilies and 54 LINE-1 subfamilies were associated with TEMRs [Supplementary Table S3]. We further found that *AluS* and *AluY* subfamilies were associated with 66.8% and 24.6% of *Alu* TEMRs and L1PA subfamilies were associated with 50% of the LINE-1 TEMRs [Supplementary Table S4]. This observation is in concordance with previous studies showing that younger TEs (fewer acquired mutations) are more likely to be involved in TEMRs¹⁹, and *AluS* TEMRs are likely enriched due to their relative abundance (678,131 *AluS* compared to 139,234 *AluY* elements in GRCh38³).

We next investigated whether TEMR-HR events were more likely to be mediated by similar TEs; TEs involved in TEMR-HRs have a median similarity of 82.6% whereas TEMR-NHE events occurred between repeats that were significantly more diverged, with only 59% median similarity [Methods]. This enrichment for similarity is a signature of recombination-mediated repair as opposed to non-homology-mediated mechanisms such as NHEJ. We found a significant difference in similarity between HR-driven and NHE-driven LINE-1 TEMRs (93.5% vs 47.9%, $p < 1e^{-5}$, Welch's *t*-test), and *Alu* TEMRs (82.2% vs 79.9%, $p < 0.01$, Welch's *t*-test) [Fig. 3a]. Interestingly, we also observed the median difference between HR-driven and NHE-driven LINE-1 TEMRs was 20-fold higher than the difference between HR-driven and NHE-driven *Alu* TEMRs. These trends appeared to be consistent across variant types with no observable difference between deletions, duplications, and inversions for both LINE-1 and *Alu* events [Supplementary Table S5], although the number of duplications and inversions were too small for a test of significance.

Additionally, we found that deletion and duplication TEMR-HRs were mediated by TEs in direct orientation, and that inverted orientation TEs were occasionally found at the junctions of NHE events. All inversions were mediated by TEs in opposite orientation and largely contain breakpoints consistent with HR (73%). We then plotted the microhomologies present at breakpoint junctions of 354 *Alu* TEMRs driven by recombination with respect to their relative position in an *Alu* consensus sequence⁵⁴ [Fig. 3b]. Consistent with previous studies, we observe peaks near the 3' ends of the A-Box and A' box^{19,31}.

INVERSIONS ARE ACCOMPANIED BY COMPLEXITIES AT THE BREAKPOINTS

Since inversions are known to harbor complex breakpoints junctions^{56,57}, we inspected the DNA sequence from HiFi assemblies for each of the 15 inversions in our final callset. Four of the seven *Alu*-mediated inversions and none of the eight LINE-1-mediated inversions contained additional complexities at both breakpoints. Complexities at *Alu*-mediated inversions included both deletions (38 bp to 251 bp) and insertions (14 bp to 23 bp) [Fig. 3c, Supplementary Fig. 4-5]. Interestingly, we found that *Alu*-mediated inversions with complex breakpoints (median 1,097 bp) were smaller than those TEMR inversions without complex breakpoints (median 4,372 bp). Furthermore, *Alu*-mediated inversions with homologous, simple breakpoints were shorter than LINE-1-mediated inversions (median 1,646 bp vs 5,487 bp; $p < 1e^{-7}$, Welch's *t*-test). We also observed that the similarity between the two *Alu* elements (median 81.7%) involved in TEMR inversions were lower compared to similarity between the two LINE-1s (median 97.1%) involved. These breakpoint characteristics, percent similarity, and variant size for smaller *Alu*-mediated events reflect known repair mechanism signatures, such as MMBIR by break repair and serial replication slippage in cis by aberrant annealing within replication bubbles^{36,58}, which contrasts with larger inversions mediated by NAHR⁵².

COMPLEX REARRANGEMENTS

While curating duplications determined using our ensemble pipeline with insertions from the HiFi assembly callset¹⁶ [Methods], we found four TEMRs containing higher-order amplification (multiple copy number variant – mCNV) that were flanked by homologous TEs. One of these included a triplication within intron 1 of *DNASE1* [Supplementary Fig. 6]. Additionally, we found a 6 kbp mCNV with 1 bp breakpoint microhomology on chr17 that harbored a 2 kbp deletion (non-TEMR) with 59 bp breakpoint microhomology within a copy [Fig. 4a, 4b and 4c]. Due to high copy numbers involved, we manually curated these mCNV TEMR breakpoints from assembly data, estimated copy number using read-depth and performed ddPCR to orthogonally validate the copy number status for the region between TEs [Fig. 4d] [Methods]. We found that read-depth-based approaches were able to accurately estimate the copy number status of events where the intervening region contained at least 7% of unique DNA sequence. A single event was only amenable to ddPCR due to high repetitive content (54.4% *Alu*

and 45.6% ERVL-MaLR) [Supplementary Fig. 7]. We compared these mCNVs with all 64 haplotypes from the latest HGSVC publication¹⁶ and found two mCNVs were present in other haplotypes [Supplementary Fig. 8]. These findings are indicative of how TEMRs can facilitate additional copy gains after initial rearrangements.

TEMRs LEAD TO POLYMORPHIC INSERTIONS IN HUMAN GENOMES

In addition to simple deletions, duplications, inversions, and complex rearrangements, TEMRs can lead to insertions of DNA^{16,28}. These insertions are often representative of a polymorphic deletion relative to ancestral humans that became part of the reference by chance; therefore, the undeleted allele manifests as an insertion in SV callsets. Recently, a 4 kbp insertion in the first intron of the *LCT* gene was described where the deletion of the region (reference allele) could be responsible for adaptive evolution in the human population¹⁶. To better understand the role of TEMRs representing insertions of ancestral sequence, we compared the results of earlier studies^{19,20} with the HiFi assembly data¹⁶, and inferred the polymorphism rate of TEMR deletions between chimpanzee and human genomes [Supplementary Fig. 9a and 9b]. By comparing deletion in chimpanzees to the insertion in humans, we estimate that 20% of *Alu*- and 15% of LINE-1-mediated deletions were polymorphic in the human population [Supplementary Fig. 9c]. This survey of ancestral regions inserted into human genomes by TEMRs also points to the role of TEs in genomic evolution and is probably a significant contributor to the role of TEs throughout speciation that is yet understudied.

TEMRs CONTRIBUTE TO VARIATION IN FUNCTIONAL REGIONS OF HUMAN GENOMES

Alu elements, especially older subfamilies are known to cluster in genic, GC-rich regions of the genome^{1,59-61}, and may facilitate genic *Alu* TEMRs. To investigate this, we first inspected whether TEMRs occur more frequently within regions of high TE density, and then we inspected how frequently TEMRs occur within genic regions of the genome. We found that *Alu* TEMRs were enriched in *Alu* dense regions of the genome compared to other parts of the genome ($p \leq 1e^{-5}$, Welch's t-test) [Fig. 5a]. The highest *Alu* density is observed on chromosomes 19 (25%) and 17 (18%)²⁸ compared to the whole genome (9.9%) [Methods]. Further, we found the highest density of *Alu* TEMRs span callable regions for chromosomes 19 (40 events, 7.5 per 10 Mbp) and 17 (37 events, 5.1 per 10 Mbp) compared to the whole genome (397 events, 1.45 per 10 Mbp) [Methods]. We also observed a similar pattern with LINE-1 TEMRs (96 events, 0.35 per 10 Mbp) in LINE-1 dense regions of the human genome ($p \leq 1e^{-5}$, Welch's t-test) [Supplementary Fig. 10].

To determine the effects of the 493 TEMRs identified in this study on genic regions, we intersected our 493 TEMRs with the RefSeq Gene database⁶² using the Ensembl variant effect predictor (VEP), which

included curated and predicted set of genes⁶³. We found that 5.7% (28) were exonic, 48.9% (241) were intronic, and an additional 10.5% (52) were within 5 kbp of a gene [Fig. 5b]. Additionally, 95.2% (256 out of 269) of genic variants affect the curated set of protein coding genes [Supplementary Table S6]. From TEMR deletions, we identified seven in 3' UTR regions, nine in 5' UTR regions, and 11 in coding sequence including 3 stop-loss events [Fig. 5c]. TEMRs overlapped 126 cis-regulatory elements⁶⁴ and 70 transcription factor binding sites (TFBS)⁶³ enriched in gene-proximal sites [Fig. 5c]. We noted that TEMRs compared to non-TEMRs are enriched in and surrounding genes (321 vs 2,740, $p \leq 0.001$, two-tailed Fisher's exact test), likely due to the higher intragenic concentration of *Alu* elements^{65,66} [Supplementary Table S7]. One *Alu*-mediated deletion in the protein coding gene *TRIM65* eliminates a last exon, including a stop codon and the 3' untranslated region [Fig. 5d]. We verified that isoforms containing the deleted exon are expressed in a recent dataset comprised of 30 breast cancer samples⁶⁷. Additionally, this deletion spans six regulatory elements annotated in the Open Regulatory Annotation dataset⁶⁸. *TRIM65* is a well-known cancer biomarker and a potential therapeutic target for colorectal cancer and lung cancer treatment^{69,70}. Further, we found 22 intragenic TEMRs overlapping hotspot OMIM (Online Mendelian Inheritance in Man) genes that are highly susceptible to *Alu*-driven genomic instability³¹. Our findings identify TEMRs as a prevalent reservoir of genic variation due to their distribution in the genome and reveal polymorphic SVs that may impact autosomal recessive disease-associated genes.

DISCUSSION

With nearly 4.5 million copies spread across the human genome, TEs are ideal substrates for rearrangements¹. Due to the high sequence similarity between subfamilies of the same TE, it has been challenging to identify and characterize TEMRs across whole genomes, and overcoming both bioinformatic and sequencing hurdles was important to complete these analyses.

We identified 493 TEMRs in a genome-wide analysis of three genomes characterized with both long-read and short-read sequencing. Contrary to *in vitro* systems, which require varying the identity of the TEs in each assay, the characterization of rearrangements in human genomes provides a natural experimental framework for how TEMRs are formed. In general, we found that longer TEs mediate larger rearrangements, indicating that the length of homology might play a significant role in determining the mechanism of repair and the size of the resulting SV. Similarly, the orientation of flanking repeats, the similarity of those sequences, and the junction homology are all indicators of whether a TEMR occurred by HR or non-homologous repair (NHE). Similar TEs in direct orientation with longer homology at the junction are all indicative of HR. Conversely, TEs in inverted orientation with shorter junction homology points to NHE. Overall, 79.7% of TEMRs appear to be mediated by HR, yet the majority of

HR-driven TEMRs have homology shorter than 40 bp. This is significantly shorter than the minimum processing segment required for nonallelic homologous recombination (NAHR) in mammalian genomes, where >100 bp of perfect sequence identity is needed⁷¹. Therefore, most of these events likely proceed in RAD51-independent mechanisms that benefit from annealing homologous sequences³⁵.

Recent studies have proposed single-strand annealing (SSA), microhomology-mediated break-induced replication (MMBIR), and microhomology-mediated end joining (MMEJ) as alternative mechanisms capable of driving TEMRs^{28,31-35,72}. Due to the presence of homology at the junctions of TEMRs and the lower percent (~80%) of similarity between the TEs involved, it is likely that all these mechanisms play a role in their formation. The prevalence of shorter deletions in our study indicates a significant role for SSA or Alt-EJ (alternative-end joining) in the formation of deletion TEMRs. SSA is indicated by the homeology of the repeats involved in ~80% of the rearrangements. If Alt-EJ was more common, shorter stretches of perfect microhomology would potentially be more common, as would rearrangements between non-homologous regions of the *Alu* or LINE-1 elements. Alternatively, a hybrid approach that invokes both SSA and end joining has been proposed^{33,35} with this mechanism, longer deletions could be accompanied by longer stretches of homeology (imperfect homology⁷³), as seen with the LINE-1 TEMRs in our study. In contrast to *in vitro* studies where more identity appears to be required for homology directed repair, the average divergence of TEs giving rise to a chimeric *Alu* at the TEMR junction was 80%. The overall distribution of junctions in the TEMR-HRs contains a peak near the 3' ends of the A and A' boxes (Fig. 3b), which could reflect preferential synapsis between GC-rich regions and *Alu* elements (Fig. 3c) or reflect the evolutionary constraint on the A and B boxes that yields more extensive homology tracts within these loci. Finally, the secondary structures inherent in *Alu* elements could lead to some regions being preferential sites of DNA repair.

We found significant sequence divergence between LINE-1s involved in TEMR-HR and TEMR-NHE but not between *Alu* elements (Fig. 3a). Intriguingly, though we examined the three genomes for NHEJ events, they remained relatively rare (even between divergent repeats) comprising only ~20% of the TEMRs. Nearly 60% of LINE-1 TEMRs were driven by NHE, whereas only 10% of *Alu* TEMRs were driven by NHE. The investigation of hundreds to thousands more of these events will enable insights that may be hard to gain from directed experimentation, including the binning of lengths of rearrangements, mechanisms, and divergence of repeats. Further genome analysis could uncover additional TE classes involved in rearrangements, including HERVs (human endogenous retroviruses)^{74,75}. Finally, we can extend our analyses to investigate mechanistic signatures of repair

pathways, including enrichment of PRDM9 binding sites near the SV junctions and the formation of single nucleotide mutations and indels *in cis* with TEMRs⁵⁵.

Since approximately 0.3% of human genetic diseases have been associated with *Alu*-mediated SVs, we intersected our TEMRs with OMIM genes that had a risk score for *Alu* TEMRs that were greater than one³¹. We found 22 intragenic TEMRs overlapping these genes with a median allele frequency of 0.13. It is likely that TEMRs continue to be underestimated in human genetics due to the historic prevalence of array comparative genomic hybridization and short-read sequencing for genome-wide SV analyses. Population estimates of allele frequencies of TEMRs are needed to assess the true impact of this mechanism on phenotypically relevant genes. Further, polymorphic TEMRs can lead to evolutionary implications when they occur within genes or include regulatory sequences¹⁶; the impact of TEMRs can be relevant both between species and within species^{16,19}.

In three genomes, we annotated 493 TEMRs, including deletions, duplications, inversions, mCNVs, insertions (ancestral deletion), and complex rearrangements. The expansion of our analysis to further assembly-based SV callsets will broaden this pool of rearrangements and allow for the discovery of novel events, which is likely to uncover new biology. For example, in the current study, we found 10 TEMRs with breakpoint junctions in poly-A tails, suggesting that the poly-A tails of non-homologous repeats can potentially drive TEMRs. Additionally, we found an instance where the same *Alu* was involved in two different SVs (both an insertion and a duplication in humans). Interestingly, the junctions of both events bore the presence of the same 30 bp microhomology at the breakpoint, signifying hotspots within an element [Supplementary Fig. 11]. Furthermore, we found that 30X coverage short-read HTS is able to capture nearly 85% of TEMRs identified by 75X short-read HTS data in this study [Supplementary Fig. 12], indicating that publicly available datasets could be used to understand the role of TEMRs across disease and population cohorts^{16,45,46,72}.

Overall, our results show that TEMRs are an important source of variation in human genomes that can arise from ectopic repair by distinct pathways, and can lead to diverse consequences. TEMRs mediate rearrangements spanning more than 500,000 bp in a human genome, and are therefore important to individual variation and to evolutionary processes; the investigation of additional high-quality assemblies will increase our understanding of the impact of transposable elements in mediating genomic rearrangements.

FIGURES

Figure 1:

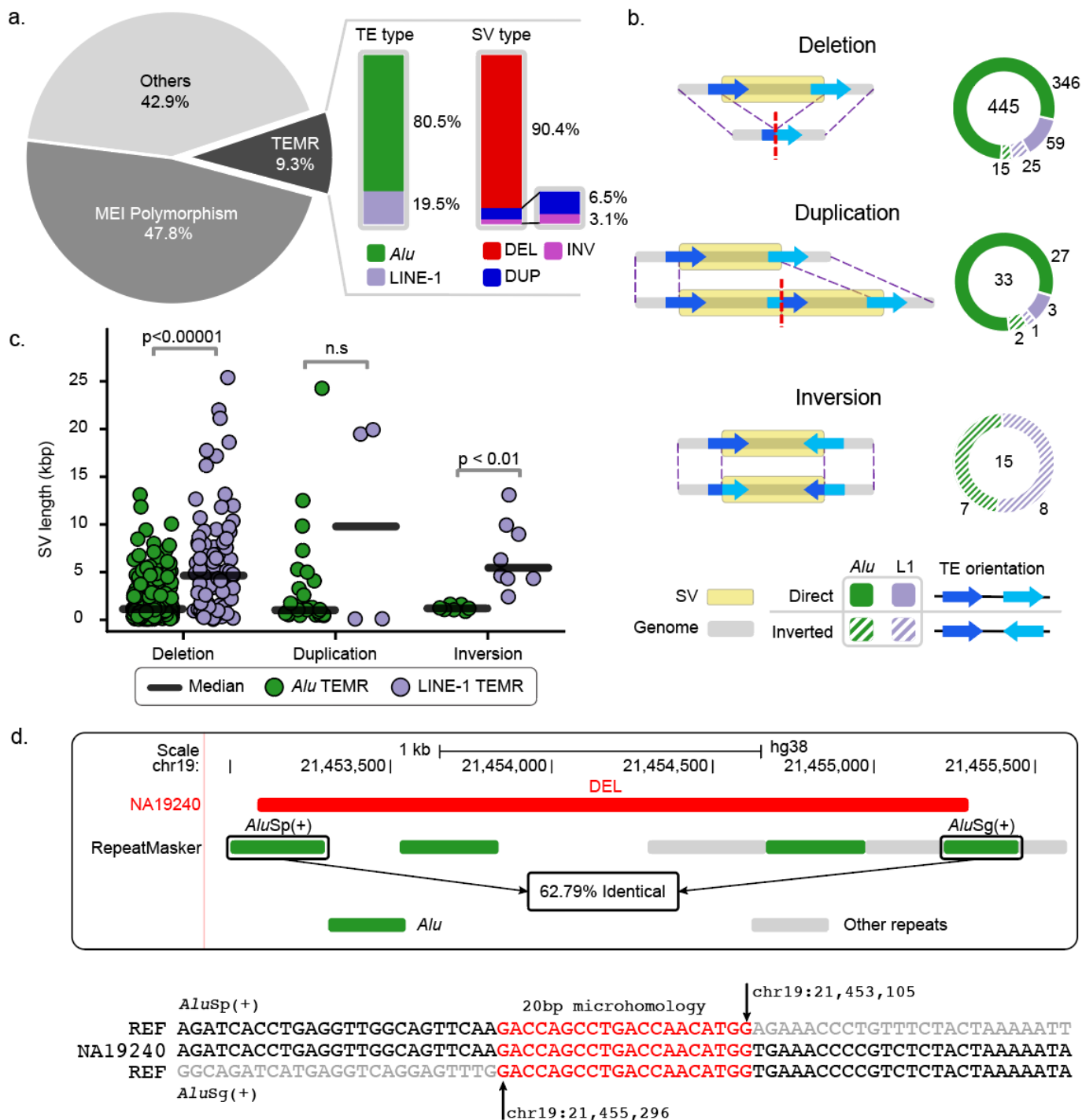


Figure 1: Generating confident TEMR callsets and breakpoint annotations across human genomes. **a**, TEMR callset summary (dark slice) by TE type, SV type and TE orientation. **b**, A diagram of TEMR structures showing distinct TEs (blue arrows). Deletion and duplication TEs are largely mediated by TEs in the same orientation (solid green and purple), and all inversions are mediated by TEs in opposite orientation (hatched green and purple). **c**, Median TEMR sizes differ by breakpoint TE types. Welch's *t*-test was used to calculate the p-value. **d**, A 2,210 bp deletion TEMR between two *Alu* repeats in direct orientation. Top panel: UCSC genome browser image showing the deletion with breakpoints in an *AluSp* and an *AluSg*. Bottom panel: Breakpoint reconstruction of the assembled deletion (middle, NA19240) against *Alu* consensus sequences (top and bottom) identifies a 20 bp breakpoint microhomology (red).

Figure 2

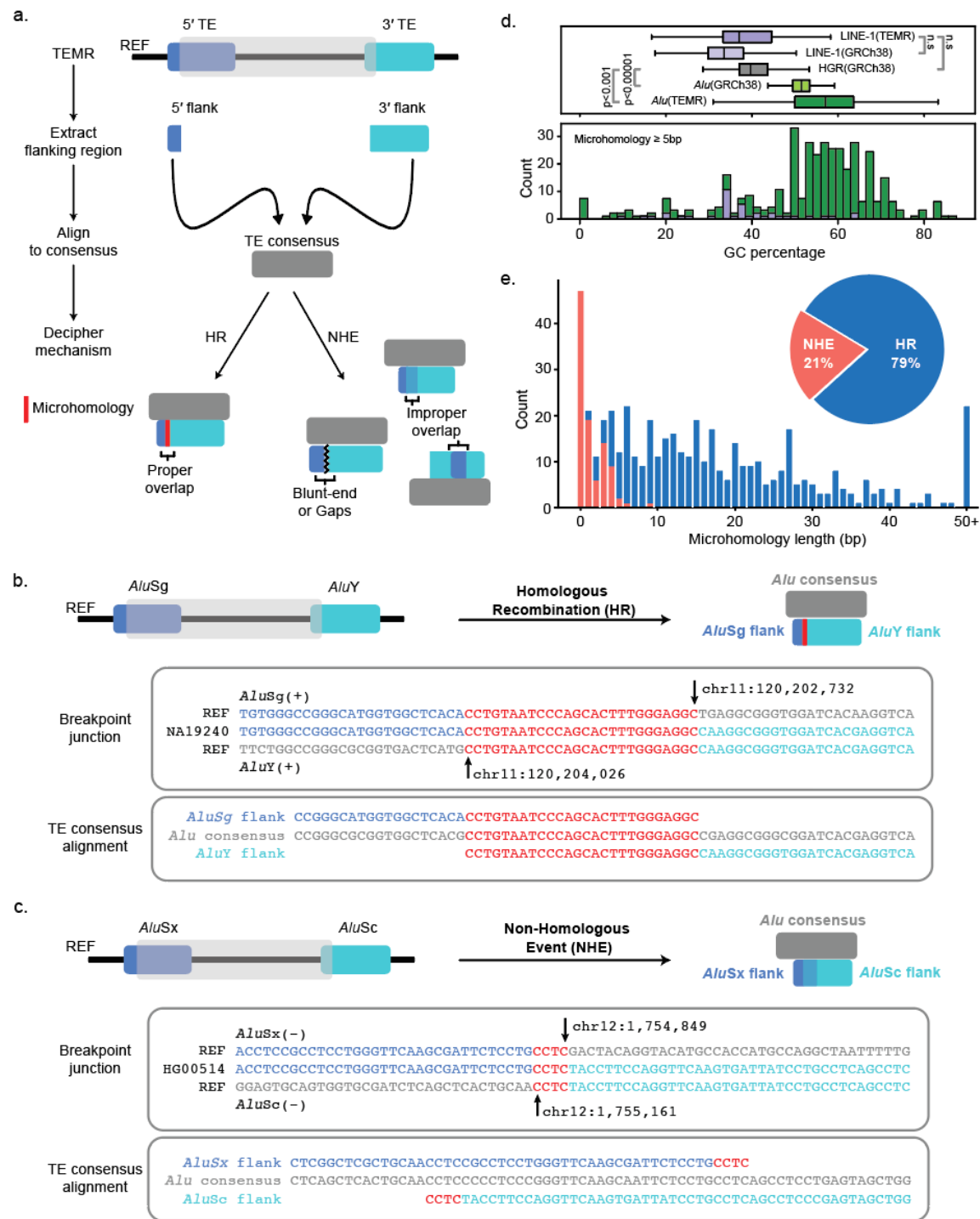


Figure 2. Identifying mechanistic signatures of TEMRs. **a**, TEMR events are classified by breakpoint characteristics guided by TE consensus sequences. For homologous breakpoints, the chimeric TE resulting from the TEMR event must reconstruct a full TE with microhomology at the breakpoint. HR, homologous recombination; NHE, non-homologous repair. **b**, An example of TEMR-HR (top). Breakpoint junction of 1,294 bp TEMR deletion in NA19240 (middle) and alignment between flanking *Alu* elements to a *Alu* consensus sequence (bottom). **c**, An example of TEMR-NHE (top). Breakpoint junction of 312 bp TEMR deletion in HG00514 (middle) and alignment between flanking *Alu* elements to a *Alu* consensus sequence (bottom). **d**, top: Microhomology GC content distribution for *Alu* TEMRs (dark green), reference *Alu* elements (light green), LINE-1 TEMRs (dark purple), reference LINE-1s (light purple) and the full human reference (gray). Bottom: Average GC content of TEMR breakpoint microhomologies for *Alu* (green) and LINE1 (purple) TEMRs. Microhomologies were restricted to 5+ bp for this analysis. Welch's t-test was used to calculate the p-value. **e**, The breakpoint microhomology distribution differs between NHE (orange) and HR (blue) TEMRs.

Figure 3

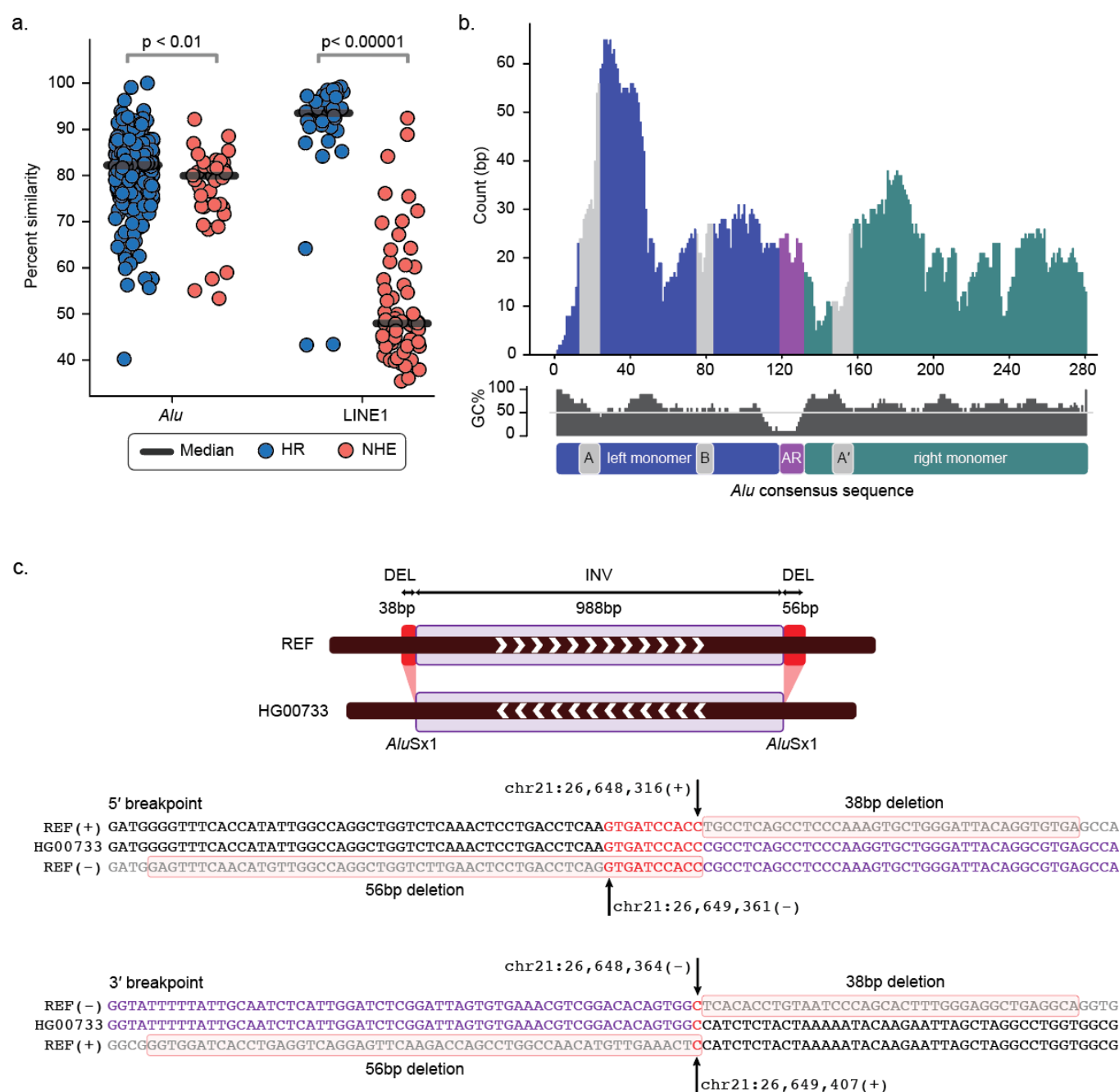


Figure 3: Features and complexities of TEMRs. **a**, Median similarity between TEs at TEMR breakpoints differs by TE type and mechanism. Welch's t-test was used to calculate the p-value. **b**, Distribution of the breakpoint microhomology along the *Alu* consensus sequence. *Alu* elements consist of a left monomer (indigo), right monomer (green), RNA polymerase III promoter regions (gray A-Box, B-Box and A'-box) and Adenosine Rich region (AR, purple). **c**, Example of a TEMR inversion with additional complex breakpoints (38 bp and 56 bp deletion indicated in pink shaded sections) implicating replication-based mechanisms of variant formation.

Figure 4

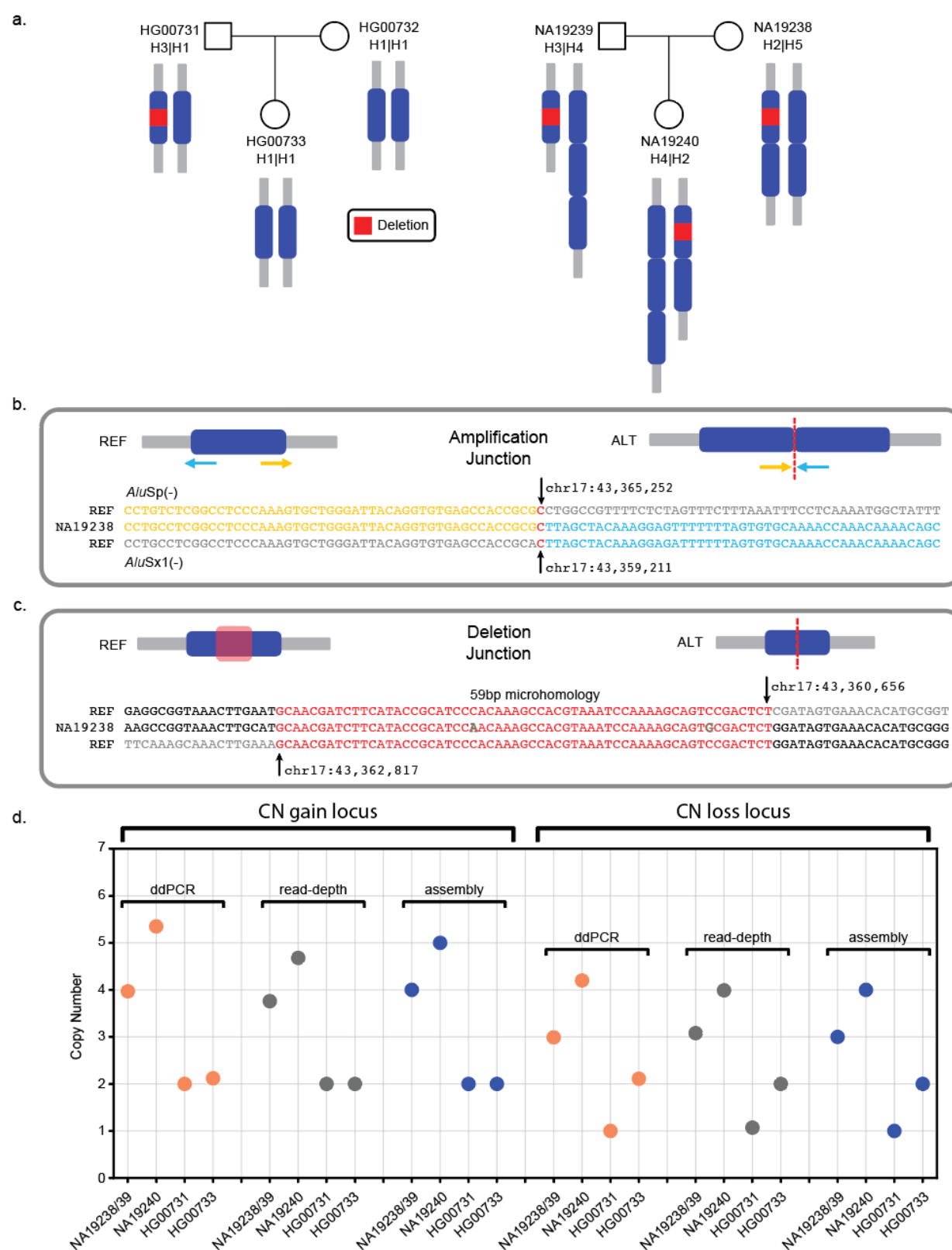


Figure 4. TEs mediate multi-copy CNVs (mCNVs). **a**, A 6 kbp multi-copy event (duplication and triplication) with a smaller 2.2 kbp deletion (red) in a subset of copies. **b**, Reconstruction of the *Alu*S-mediated mCNV TEMR breakpoint surrounding each 6 kbp copy shows 1 bp microhomology. **c**, Reconstruction of the 2 kbp inner-copy deletion shows 59 bp of near-perfect homology. **d**, Copy number status of the individuals containing the triplication found using ddPCR, read-depth analysis and assembly.

Figure 5

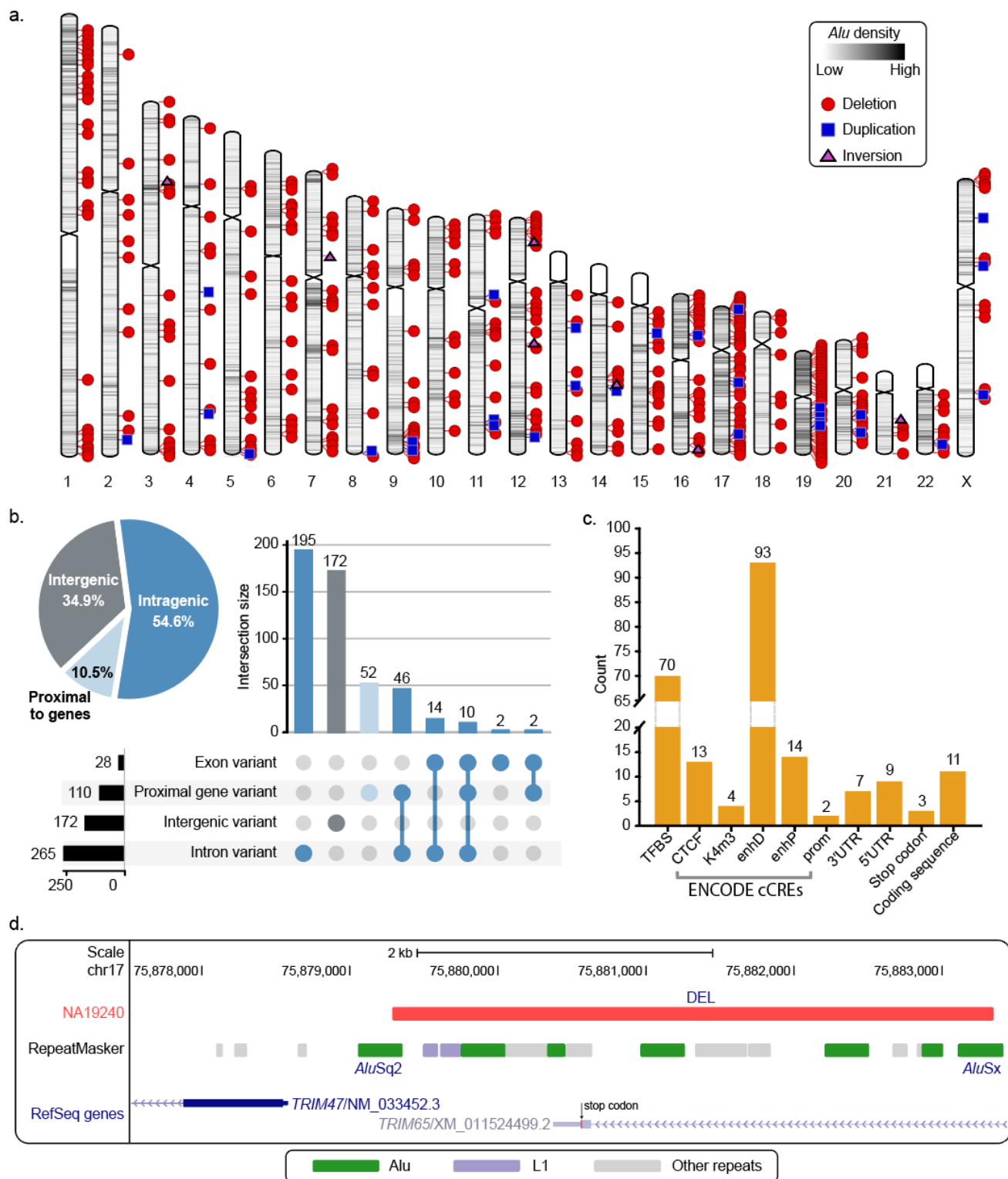


Figure 5: TEMRs disproportionately affect genes. **a**, All 390 *Alu* TEMRs by variant type (shape) and *Alu* density (shading). **b**, TEMRs disproportionately intersect genes (pie chart) affecting introns, exons, and gene-proximal regions that are often enriched with regulatory elements. **c**, Regulatory and coding regions affected by TEMRs. Regulatory annotations were retrieved from ENCODE cCREs, coding regions and TFBS predictions from Ensembl VEP. TFBS, transcription factor binding site; cCREs, Candidate Cis-Regulatory Elements; K4m3, DNase-H3K4me3; enhD, distal enhancer-like signature; enhP, proximal enhancer-like signature; prom, promoter-like signature; UTR, untranslated region. **d**, Examples of TEMRs deleting a stop codon in *TRIM65* (bottom).

METHODS

PIPELINE FOR IDENTIFYING, FILTERING, AND MERGING SVs

SV identification. PCR-free Illumina short-read Whole Genome Sequencing (WGS) data (75X coverage) and PacBio continuous long-read (CLR) data (40X coverage) for three samples were downloaded from the HGSV consortium (Puerto Rican (PUR) HG00733; Southern Han Chinese (CHS) HG00514; and Yoruban (YRI) NA19240)³⁸. Complete schematic for the short-read and long-read ensemble pipelines are shown in Supplementary Fig. 13. Raw paired-end sequencing data (FASTQ) were aligned against the human genome reference (GRCh38/hg38) using BWA MEM (v0.7.17)⁷⁶. Similarly, long-read sequences (FASTQ) were extracted from the native PacBio files (bax.h5) using pbh5tools (v0.8.0; <https://github.com/PacificBiosciences/pbh5tools>) and aligned to the HGR using NGMLR (v0.2.6)⁵⁰. The aligned files were converted (sam to bam), sorted, merged, and indexed using samtools (v1.7). SV calling was performed on the indexed short-read bam files using Manta (v1.3.2)⁴⁷, LUMPY (v0.2.13)⁴⁸, and DELLY (v0.7.8)⁴⁹ and indexed long-read bam files using Sniffles (v1.0.7)⁵⁰, SVIM (v1.4.0)⁵¹ and pbsv (v2.2.0; <https://github.com/PacificBiosciences/pbsv>). Read-depth information were annotated for the output generated by the three callers using Duphold (v0.2.1)⁷⁷. Settings with which each caller was run can be found in Supplementary Table S8.

SV filtering. After generating SV calls, resultant VCF files were mined for the coordinates (chrom, chromStart, chromEnd), svType, caller name, paired-read (PR), split-read (SR) and read-depth (RD) for SVs with “FILTER=PASS”, and this information was transferred to a tab separated file (TSV). For long-read SV calls, PR and SR were replaced with read-support (RS). We removed SVs within 500 bp of gaps and centromeres, and SVs that overlapped (50%) with simple repeats BEDTools *intersect* (v 2.29.2)⁷⁸. We retained deletions, duplications, and inversions of size ranging from 50 bp to 50 kbp.

SV merging. After testing multiple approaches to merge [Supplementary Fig. 14a and 14b] and filtering SVs, we applied read-based (SR ≥ 5 or PR ≥ 10 for short-reads and RS ≥ 5 for long-reads) and depth-based (RD < 0.7 for deletions and RD > 1.3 for duplications) filters to our callset [Supplementary Fig. 14c and 14d]. We merged and retained SVs from multiple callers using 80% reciprocal overlap (RO) with BEDTools *intersect* and a rank-based method.

Rank-based merging. We investigated the accuracy with which callers identified breakpoint junctions by comparing breakpoint junctions of deletion calls between our intersection callset (SVs identified by all 3 caller) and the truth set³⁸ using 80% RO. We calculated the deviation between the breakpoints at both 5' and 3' junctions. We ranked the three callers used in both short-read (1. Manta, 2. DELLY and

3. LUMPY) and long-read pipeline (1. pbsv, 2. Sniffles and 3. SVIM). When merging (80% RO) SVs in our ensemble callset, we retained the SV size and breakpoint data identified from the highest-ranking caller and removed SVs from other callers [Supplementary Fig. 15].

STATISTICS

Statistical analysis. All statistical testing were performed with scipy (v1.5.0)⁷⁹.

BREAKPOINT REFINEMENT

Revising breakpoint junctions. We overlapped our ensemble callset with HiFi assembly-based SV calls using 80% RO. Since duplications were also reported as insertions in the assembly callset, we followed a two-step approach to verify proper support for duplication from the assembly. First, we checked for insertions within our duplication (and 500 bp around the SV). Second, we compared the duplication size with the insertion size (90% match). We obtained the precise breakpoint junctions and microhomology information from the assembly callset¹⁶ and annotated our TSV files accordingly. Further, we assessed the accuracy of the microhomologies by validating 70 SVs using PCR and Sanger sequencing.

TE DENSITY

Computing TE density. We computed density of TEs corresponding to the TEMRs (for example, *Alu* density was calculated for each *Alu* TEMRs) by calculating the percentage of TE sequence present in the 50 kbp encompassing the SV (25 kbp on either side, from the center of an SV). As a control, one million random 50 kbp regions containing at least two TEs of the same family were selected across the reference genome using BEDTools *shuffle*. We used Welch's unequal variance *t*-test from scipy (v1.5.0)⁷⁹ to calculate the statistical significance of the TE density between TEMRs and our control group. We also inspected the density of TEs within the callable regions of each chromosome to avoid any sequencing discrepancies. We considered regions of the genome as callable regions by excluding gaps, centromeres, segmental duplication (identity $\geq 95\%$) and long (≥ 5 kbp) tandem repeats.

PERCENT SIMILARITY BETWEEN TWO TEs

Calculating percent similarity. We extracted the TE sequence in positive orientation (TEs on the negative strand were extracted and the reverse complement was generated to obtain the TE sequence) from the GRCh38 using RepeatMasker v3.0³ and BEDTools *getfasta*⁷⁸. The percent similarity was calculated by aligning the two sequences using swalign (v0.3.4) (<https://github.com/mbreese/swalign>) with the scoring matrix from EMBOSS Water⁸⁰.

EXTRACTING SEQUENCES FROM ASSEMBLY DATA

Sequence extraction. We used subseq (v 1.0) (<https://github.com/EichlerLab/seqtools>) to extract raw sequence (SV + 500 bp window) of both haplotypes from the assembled data of the corresponding individuals¹⁶. BLAT was used to confirm the presence of an SV. We used the sequence containing the SV to manually reconstruct the breakpoint junction.

DECIPHERING TEMR MECHANISM

Identifying mechanism. We obtained the TE sequences, location of the breakpoint junction, and the microhomology of all TEMRs. Then, we obtained the flanking regions outside the breakpoints and aligned it to a consensus TE (*Alu*⁵⁴ for *Alu* TEMRs and L1.3⁵³ for all LINE-1 TEMRs). We calculated the distance between the flanking regions (post alignment) along the consensus. We categorized TEMRs as HR-driven if the two flanking regions overlapped, and the overlap size is identical to the microhomology. For *Alu* TEMRs that failed in the above step and still had microhomology, we aligned them against each monomer in the consensus and repeated the process to identify TEMRs that recombined between different monomers (left monomer recombining with right). TEMRs that had no microhomology and flanking regions that failed to overlap along the consensus were categorized as NHE. TEMRs that were not characterized by the above two steps were manually curated, and a mechanism was assigned based on the formation of a chimera and length of homology at the junction.

COPY NUMBER ESTIMATION FOR mCNV USING READ-DEPTH

Calculating copy number. We used mosdepth (v0.3.2)⁸¹ to calculate the average per-base depth across the mCNV loci. To calculate the copy number at the mCNV loci we divided the average per-base depth at the mCNV loci by the average per-base depth across the whole genome (genome coverage).

ANALYZING ANCESTRAL DELETIONS FROM A PREVIOUSLY PUBLISHED DATASET

Ancestral TEMR deletion analysis. We download the *Alu* TEMRs¹⁹ and LINE-1 TEMRs²⁰, and lifted over the coordinates to GRCh38 (LiftOver from UCSC genome browser) along with the corresponding deletion size in the chimpanzee genome. We then intersected these coordinates with the insertions from the latest HGSVC publication¹⁶ using BEDTools *window*⁷⁸ (length of the consensus sequence was used as window size). We retained events if size of the deletion (from chimpanzee) and insertion (from human) had a 90% match.

BENCH EXPERIMENTS

Designing primers for PCR and Sanger sequencing. Primers were designed using the Primer3web interface (v4.1.0: <https://primer3.ut.ee/>). We obtained the DNA sequence of candidate SVs and their flanking regions (500 bp) using BEDTools *getfasta*⁷⁸ or UCSC Genome Browser (Get DNA in Window) and regions for primer pairs were identified within the flanking regions. To ensure high quality sequencing through the breakpoint, primers were designed to anneal to at least 75 bp from the predicted breakpoints. We required that at least one primer mapped (UCSC BLAT) to a unique region in the reference genome (GRCh38) and the primer pair uniquely amplified the locus of interest (UCSC in-silico PCR).

PCR reaction. PCR was conducted using the AccuPrime Taq DNA polymerase system (ThermoFisher Scientific, 12339024). The following conditions were used for 25 µl PCR reactions:

- 2.5 µl of AccuPrime 10× Buffer II
- 1.5 µl of 10 µM forward and reverse primers
- 3 µl of 5M Betaine
- 1 µl of DNA (20 ng/µl)
- 0.25 µl of AccuPrime Taq DNA polymerase
- 16.75 µl of nuclease free water

We utilized touchdown PCR cycling for optimal primer annealing and achieve specific amplification of the desired region. The thermocycler program was as follows:

- Step 1. 94°C for 2 min
- Step 2. 94°C for 30 sec
- Step 3. 63°C for 30 sec (with a 1°C ramp down per cycle)
- Step 4. 68°C for 1 min
- Step 5. Return to step 2 and repeat for 8 cycles
- Step 6. 94°C for 30 sec
- Step 7. 57°C for 30 sec
- Step 8. 68°C for 1 min
- Step 9. Return to step 6 and repeat for 25 cycles
- Step 10. 68°C for 1 min
- Step 11. 4°C infinite hold

PCR products combined with 5 µl of 10× OrangeG gel loading dye (0.4% w/v final concentration of OrangeG sodium salt, 40% w/v final concentration of sucrose, and sterile water) were run in 1% agarose gels with 0.1 µl ethidium bromide (10mg/ml Bio-Rad, 1610433) per milliliter of gel. 7 µl of a 1Kb Plus DNA ladder (ThermoFisher Scientific, 10787018) diluted in 10× BlueJuice Gel Loading Buffer (ThermoFisher Scientific, 10816015) was run alongside the reactions. Bands at the desired size were excised under a blue light and DNA was purified using the Zymoclean Gel DNA Recovery Kit (Zymo Research, D4008).

Sanger sequencing. Purified PCR products were subjected to Sanger dideoxy sequencing (Eton Biosciences) along with their respective PCR primer pairs. If sequencing using PCR primers was expected to cause premature polymerase slippage due to a homopolymeric region being present between the primer and the breakpoint, separate sequencing primers were designed to avoid such regions.

Droplet digital PCR (ddPCR) assay and probe design. Custom primers and probes were designed for each mCNV using Primer3Plus (<https://www.bioinformatics.nl/cgi-bin/primer3plus/primer3plus.cgi>). Primers and probe assays were designed using the parameters recommended by Bio-Rad's ddPCR application guide. We selected primer pairs that had at least one primer mapping (UCSC BLAT) to a unique region in the reference genome (GRCh38) and generated a shorter (<200 bp) product. The reference probe and assay used for the duplexed reaction was designed to diploid gene RPP30. A HEX fluorophore (Bio-Rad, 10031279) was used for the reference assay and FAM fluorophore for the target assay (Bio-Rad, 10031276). We used a 3.6:1 nM primer to probe concentration ratio. Restriction enzymes (Supplementary Table S7) from New England Biolabs (NEB) were used to fragment tandem duplications outside of predicted PCR product.

ddPCR Sample Preparation and Thermocycler Program: ddPCR reaction mix was assembled according ddPCR™ Copy Number Variation Assays (10000050421 Ver A) from Bio-Rad. The following concentrations were used for 22 µl of total volume prior to droplet formation:

- 11 µl of 2× ddPCR Supermix for Probes (No dUTP)
- 1.1 µl of 20× target assay
- 1.1 µl of 20× reference assay
- 1.5 µl of digested gDNA (20ng/ul)
- 7.3 µl of nuclease free water

Droplets were formed by Bio-Rad's QX200 AutoDG Droplet Digital PCR system. After sealing the plates, droplet samples were loaded in a thermocycler.

Thermocycler program for 40ul volume:

- Step 1. 95°C for 10 min
- Step 2. 94°C for 30 sec
- Step 3. 60°C for 1 min (with a 2°C/sec ramp per cycle)
- Step 4. Return to step 2 and repeat for 39 cycles
- Step 5. 98°C for 10 min
- Step 6. 4°C infinite hold

Plates were loaded into the QX200 Droplet Reader post PCR and copy number status was analyzed using QuantaSoft™ software from Bio-Rad.

References

1. Lander, E.S. *et al.* Initial sequencing and analysis of the human genome. *Nature* **409**, 860-921 (2001).
2. Jurka, J. Repbase update: a database and an electronic journal of repetitive elements. *Trends Genet* **16**, 418-20 (2000).
3. Smit, A.F.A., Hubley, R. & Green, P. RepeatMasker Open-3.0. (1996-2010).
4. de Koning, A.P., Gu, W., Castoe, T.A., Batzer, M.A. & Pollock, D.D. Repetitive elements may comprise over two-thirds of the human genome. *PLoS Genet* **7**, e1002384 (2011).
5. Wheeler, T.J. *et al.* Dfam: a database of repetitive DNA based on profile hidden Markov models. *Nucleic Acids Res* **41**, D70-82 (2013).
6. Kazazian, H.H., Jr. *et al.* Haemophilia A resulting from de novo insertion of L1 sequences represents a novel mechanism for mutation in man. *Nature* **332**, 164-6 (1988).
7. Stoneking, M. *et al.* Alu insertion polymorphisms and human evolution: evidence for a larger population size in Africa. *Genome Res* **7**, 1061-71 (1997).
8. Roy-Engel, A.M. *et al.* Alu insertion polymorphisms for the study of human genomic diversity. *Genetics* **159**, 279-90 (2001).
9. Brouha, B. *et al.* Hot L1s account for the bulk of retrotransposition in the human population. *Proc Natl Acad Sci U S A* **100**, 5280-5 (2003).
10. Beck, C.R. *et al.* LINE-1 retrotransposition activity in human genomes. *Cell* **141**, 1159-70 (2010).
11. Holmes, S.E., Dombroski, B.A., Krebs, C.M., Boehm, C.D. & Kazazian, H.H., Jr. A new retrotransposable human L1 element from the LRE2 locus on chromosome 1q produces a chimaeric insertion. *Nat Genet* **7**, 143-8 (1994).
12. Moran, J.V., DeBerardinis, R.J. & Kazazian, H.H., Jr. Exon shuffling by L1 retrotransposition. *Science* **283**, 1530-4 (1999).
13. Pickeral, O.K., Makalowski, W., Boguski, M.S. & Boeke, J.D. Frequent human genomic DNA transduction driven by LINE-1 retrotransposition. *Genome Res* **10**, 411-5 (2000).
14. Tubio, J.M.C. *et al.* Mobile DNA in cancer. Extensive transduction of nonrepetitive DNA mediated by L1 retrotransposition in cancer genomes. *Science* **345**, 1251343 (2014).
15. Han, K. *et al.* Genomic rearrangements by LINE-1 insertion-mediated deletion in the human and chimpanzee lineages. *Nucleic Acids Res* **33**, 4040-52 (2005).
16. Ebert, P. *et al.* Haplotype-resolved diverse human genomes and integrated analysis of structural variation. *Science* **372**(2021).
17. Gilbert, N., Lutz-Prigge, S. & Moran, J.V. Genomic deletions created upon LINE-1 retrotransposition. *Cell* **110**, 315-25 (2002).
18. Symer, D.E. *et al.* Human L1 retrotransposition is associated with genetic instability in vivo. *Cell* **110**, 327-38 (2002).
19. Sen, S.K. *et al.* Human genomic deletions mediated by recombination between Alu elements. *Am J Hum Genet* **79**, 41-53 (2006).
20. Han, K. *et al.* L1 recombination-associated deletions generate human genomic variation. *Proc Natl Acad Sci U S A* **105**, 19366-71 (2008).
21. Mun, S. *et al.* A study of transposable element-associated structural variations (TASVs) using a de novo-assembled Korean genome. *Exp Mol Med* **53**, 615-630 (2021).
22. Bailey, J.A., Liu, G. & Eichler, E.E. An Alu transposition model for the origin and expansion of human segmental duplications. *Am J Hum Genet* **73**, 823-34 (2003).
23. Franke, G. *et al.* Alu-Alu recombination underlies the vast majority of large VHL germline deletions: Molecular characterization and genotype-phenotype correlations in VHL patients. *Hum Mutat* **30**, 776-86 (2009).
24. Flynn, E.K. *et al.* Comprehensive analysis of pathogenic deletion variants in Fanconi anemia genes. *Hum Mutat* **35**, 1342-53 (2014).
25. Lehrman, M.A. *et al.* Mutation in LDL receptor: Alu-Alu recombination deletes exons encoding transmembrane and cytoplasmic domains. *Science* **227**, 140-6 (1985).

26. Temtamy, S.A. *et al.* Long interspersed nuclear element-1 (LINE1)-mediated deletion of EVC, EVC2, C4orf6, and STK32B in Ellis-van Creveld syndrome with borderline intelligence. *Hum Mutat* **29**, 931-8 (2008).
27. Boone, P.M. *et al.* Alu-specific microhomology-mediated deletion of the final exon of SPAST in three unrelated subjects with hereditary spastic paraplegia. *Genet Med* **13**, 582-92 (2011).
28. Gu, S. *et al.* Alu-mediated diverse and complex pathogenic copy-number variants within human chromosome 17 at p13.3. *Hum Mol Genet* **24**, 4061-77 (2015).
29. Pettersson, M. *et al.* Alu-Alu mediated intragenic duplications in IFT81 and MATN3 are associated with skeletal dysplasias. *Hum Mutat* **39**, 1456-1467 (2018).
30. Startek, M. *et al.* Genome-wide analyses of LINE-LINE-mediated nonallelic homologous recombination. *Nucleic Acids Res* **43**, 2188-98 (2015).
31. Song, X. *et al.* Predicting human genes susceptible to genomic instability associated with Alu/Alu-mediated rearrangements. *Genome Res* **28**, 1228-1242 (2018).
32. Hedges, D.J. & Deininger, P.L. Inviting instability: Transposable elements, double-strand breaks, and the maintenance of genome integrity. *Mutat Res* **616**, 46-59 (2007).
33. Morales, M.E. *et al.* The contribution of alu elements to mutagenic DNA double-strand break repair. *PLoS Genet* **11**, e1005016 (2015).
34. White, T.B., Morales, M.E. & Deininger, P.L. Alu elements and DNA double-strand break repair. *Mob Genet Elements* **5**, 81-85 (2015).
35. Bhargava, R., Onyango, D.O. & Stark, J.M. Regulation of Single-Strand Annealing and its Role in Genome Maintenance. *Trends Genet* **32**, 566-575 (2016).
36. Hastings, P.J., Ira, G. & Lupski, J.R. A microhomology-mediated break-induced replication model for the origin of human copy number variation. *PLoS Genet* **5**, e1000327 (2009).
37. Beck, C.R. *et al.* Megabase Length Hypermethylation Accompanies Human Structural Variation at 17p11.2. *Cell* **176**, 1310-1324 e10 (2019).
38. Chaisson, M.J.P. *et al.* Multi-platform discovery of haplotype-resolved structural variation in human genomes. *Nat Commun* **10**, 1784 (2019).
39. Li, Y. *et al.* Patterns of somatic structural variation in human cancer genomes. *Nature* **578**, 112-121 (2020).
40. Zhao, X. *et al.* Expectations and blind spots for structural variation detection from long-read assemblies and short-read genome sequencing technologies. *Am J Hum Genet* **108**, 919-928 (2021).
41. Mahmoud, M. *et al.* Structural variant calling: the long and the short of it. *Genome Biol* **20**, 246 (2019).
42. Chaisson, M.J. *et al.* Resolving the complexity of the human genome using single-molecule sequencing. *Nature* **517**, 608-11 (2015).
43. Wenger, A.M. *et al.* Accurate circular consensus long-read sequencing improves variant detection and assembly of a human genome. *Nat Biotechnol* **37**, 1155-1162 (2019).
44. 1000 Genomes Project Consortium *et al.* A global reference for human genetic variation. *Nature* **526**, 68-74 (2015).
45. ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium. Pan-cancer analysis of whole genomes. *Nature* **578**, 82-93 (2020).
46. Karczewski, K.J. *et al.* The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* **581**, 434-443 (2020).
47. Chen, X. *et al.* Manta: rapid detection of structural variants and indels for germline and cancer sequencing applications. *Bioinformatics* **32**, 1220-2 (2016).
48. Layer, R.M., Chiang, C., Quinlan, A.R. & Hall, I.M. LUMPY: a probabilistic framework for structural variant discovery. *Genome Biol* **15**, R84 (2014).
49. Rausch, T. *et al.* DELLY: structural variant discovery by integrated paired-end and split-read analysis. *Bioinformatics* **28**, i333-i339 (2012).
50. Sedlazeck, F.J. *et al.* Accurate detection of complex structural variations using single-molecule sequencing. *Nat Methods* **15**, 461-468 (2018).
51. Heller, D. & Vingron, M. SVIM: structural variant identification using mapped long reads. *Bioinformatics* **35**, 2907-2915 (2019).

52. Porubsky, D. *et al.* Haplotype-resolved inversion landscape reveals hotspots of mutational recurrence associated with genomic disorders. (bioRxiv, 2021).
53. Dombroski, B.A., Scott, A.F. & Kazazian, H.H., Jr. Two additional potential retrotransposons isolated from a human L1 subfamily that contains an active retrotransposable element. *Proc Natl Acad Sci U S A* **90**, 6513-7 (1993).
54. Stewart, C. *et al.* A comprehensive map of mobile element insertion polymorphisms in humans. *PLoS Genet* **7**, e1002236 (2011).
55. Carvalho, C.M. & Lupski, J.R. Mechanisms underlying structural variant formation in genomic disorders. *Nat Rev Genet* **17**, 224-38 (2016).
56. Collins, R.L. *et al.* A structural variation reference for medical and population genetics. *Nature* **581**, 444-451 (2020).
57. Collins, R.L. *et al.* Defining the diverse spectrum of inversions, complex structural variation, and chromothripsis in the morbid human genome. *Genome Biol* **18**, 36 (2017).
58. Chen, J.M., Chuzhanova, N., Stenson, P.D., Ferec, C. & Cooper, D.N. Intrachromosomal serial replication slippage in trans gives rise to diverse genomic rearrangements involving inversions. *Hum Mutat* **26**, 362-73 (2005).
59. Sellis, D., Provata, A. & Almirantis, Y. Alu and LINE1 distributions in the human chromosomes: evidence of global genomic organization expressed in the form of power laws. *Mol Biol Evol* **24**, 2385-99 (2007).
60. Jurka, J., Kohany, O., Pavlicek, A., Kapitonov, V.V. & Jurka, M.V. Duplication, coclustering, and selection of human Alu retrotransposons. *Proc Natl Acad Sci U S A* **101**, 1268-72 (2004).
61. Brookfield, J.F. Selection on Alu sequences? *Curr Biol* **11**, R900-1 (2001).
62. Pruitt, K.D., Tatusova, T. & Maglott, D.R. NCBI Reference Sequence (RefSeq): a curated non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Res* **33**, D501-4 (2005).
63. McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol* **17**, 122 (2016).
64. Consortium, E.P. *et al.* Expanded encyclopaedias of DNA elements in the human and mouse genomes. *Nature* **583**, 699-710 (2020).
65. Deininger, P.L. & Batzer, M.A. Alu repeats and human disease. *Mol Genet Metab* **67**, 183-93 (1999).
66. Grover, D., Mukerji, M., Bhatnagar, P., Kannan, K. & Brahmachari, S.K. Alu repeat analysis in the complete human genome: trends and variations with respect to genomic composition. *Bioinformatics* **20**, 813-7 (2004).
67. Veiga, D.F.T. *et al.* A comprehensive long-read isoform analysis platform and sequencing resource for breast cancer. *Sci Adv* **8**, eabg6711 (2022).
68. Lesurf, R. *et al.* ORegAnno 3.0: a community-driven resource for curated regulatory annotation. *Nucleic Acids Res* **44**, D126-32 (2016).
69. Chen, D. *et al.* Ubiquitin ligase TRIM65 promotes colorectal cancer metastasis by targeting ARHGAP35 for protein degradation. *Oncogene* **38**, 6429-6444 (2019).
70. Wang, X.L. *et al.* Knockdown of TRIM65 inhibits lung cancer cell proliferation, migration and invasion: A therapeutic target in human lung cancer. *Oncotarget* **7**, 81527-81540 (2016).
71. Liskay, R.M., Letsou, A. & Stachelek, J.L. Homology requirement for efficient gene conversion between duplicated chromosomal sequences in mammalian cells. *Genetics* **115**, 161-7 (1987).
72. Morales, M.E. *et al.* Altered DNA repair creates novel Alu/Alu repeat-mediated deletions. *Hum Mutat* **42**, 600-613 (2021).
73. Liu, P. *et al.* An Organismal CNV Mutator Phenotype Restricted to Early Human Development. *Cell* **168**, 830-842 e7 (2017).
74. Shuvarikov, A. *et al.* Recurrent HERV-H-mediated 3q13.2-q13.31 deletions cause a syndrome of hypotonia and motor, language, and cognitive delays. *Hum Mutat* **34**, 1415-23 (2013).
75. Campbell, I.M. *et al.* Human endogenous retroviral elements promote genome instability via non-allelic homologous recombination. *BMC Biol* **12**, 74 (2014).
76. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. in *arXiv e-prints* (2013).
77. Pedersen, B.S. & Quinlan, A.R. Duphold: scalable, depth-based annotation and curation of high-confidence structural variant calls. *Gigascience* **8**(2019).

78. Quinlan, A.R. & Hall, I.M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841-2 (2010).
79. Virtanen, P. *et al.* SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat Methods* **17**, 261-272 (2020).
80. Madeira, F. *et al.* The EMBL-EBI search and sequence analysis tools APIs in 2019. *Nucleic acids research* **47**, W636-W641 (2019).
81. Pedersen, B.S. & Quinlan, A.R. Mosdepth: quick coverage calculation for genomes and exomes. *Bioinformatics* **34**, 867-868 (2018).