

Discovery and characterization of *Alu* repeat sequences via precise local read assembly

Julia H. Wildschutte¹, Alayna Baron¹, Nicolette M. Diroff¹, and Jeffrey M. Kidd^{1,2,*}

¹Department of Human Genetics, University of Michigan Medical School, Ann Arbor, MI, 48109, USA

²Department of Computational Medicine and Bioinformatics, University of Michigan Medical School, Ann Arbor, MI, 48109, USA

* To whom correspondence should be addressed. Tel: (734) 763-7083; Fax: (734) 763-3784; Email: jmkidd@med.umich.edu

ABSTRACT

Alu insertions have contributed to >11% of the human genome and ~30-35 *Alu* subfamilies remain actively mobile, yet the characterization of polymorphic *Alu* insertions from short-read data remains a challenge. We build on existing computational methods to combine *Alu* detection and *de novo* assembly of WGS data as a means to reconstruct the full sequence of insertion events from Illumina paired end reads. Comparison with published calls obtained using PacBio long-reads indicates a false discovery rate below 5%, at the cost of reduced sensitivity due to the colocation of reference and non-reference repeats. We generate a highly accurate call set of 1,614 completely assembled *Alu* variants from 53 samples from the Human Genome Diversity Project panel. We utilize the reconstructed alternative insertion haplotypes to genotype 1,010 fully assembled insertions, obtaining >99% accuracy. In our assembled sequences, we find evidence of non-classical insertion mechanisms and observe 5' truncation in 16% of *AluYa5* and *AluYb8* insertions. The sites of truncation coincide with stem-loop structures and SRP9/14 binding sites in the *Alu* RNA, implicating L1 ORF2p pausing in the generation of 5' truncations.

INTRODUCTION

Mobile elements (MEs) are discrete fragments of nuclear DNA that are capable of copied movement to other chromosomal locations within the genome (1). In humans, the ~300 bp *Alu* retroelements are the most successful and ubiquitous MEs, collectively amounting to >1.1 million genome copies and accounting for >11% of the nuclear genome (2,3). The vast majority of *Alu* insertions represent germline events that occurred millions of years ago and now exist as non-functional elements that are highly mutated and no longer capable of mobilization (3). However, subsets of MEs, including *Alu* and its autonomous partner *L1Hs*, remain active and continue to contribute to new ME insertions (MEIs), resulting in genomic variation between individuals (4) and between somatic tissues within an individual (5,6).

The human *Alu* consists of the *AluY*, *AluS*, and *AluJ* lineages, which can be further stratified into more than ~35 subfamilies based on sequence diversity and diagnostic mutations (2,4,7). Most human *Alu* elements are from the youngest lineage, *AluY*, whose members have been most actively mobilized during primate evolution (4,8). Of these, the *AluYa5* and *AluYb8* subfamilies have

contributed to the bulk of insertions in humans (9-12), although polymorphic insertions from >20 other *AluY* and >6 *AluS* subfamilies have also been reported (4,13), implying polymorphic insertions of other lineages may still be segregating. In contemporary humans, the retrotransposition of active *Alu* copies results in *de novo* germline insertions at a frequency of ~1:20 live births (10,14). Over 60 novel *Alu* insertions have been shown to cause mutations leading to disease, either as a direct consequence of insertional mutagenesis, or by providing a template of highly repetitive sequence that has since facilitated chromosomal rearrangements and structural variation (2,15-18). Thus, *Alu* insertions continue to shape the genomic landscape and are recognized as profound mediators of genomic structural variation.

Active copies of *Alu* are non-autonomous but contain an internal RNA Pol III promoter (19). Mediated by L1 encoded enzymes, *Alu* transcripts are mobilized by a 'copy-and-paste' mechanism referred to as target primed reverse transcription (TPRT) (6,20). Classical TPRT involves the reverse transcription of a single stranded *Alu* RNA to a double stranded DNA copy, during which two staggered single-stranded breaks are introduced in the target DNA of ~5 to ~25 bp that are later filled by cellular machinery. The resulting structure consists of a new *Alu* flanked by characteristic target site duplications (TSDs) and a poly-A tail of variable length. Together these serve as hallmarks of retrotransposition. Integration of the new copy is permanent; although *Alu* can be removed by otherwise encompassing deletions, there is no known mechanism for precise excision. Classical TPRT is responsible for the majority of *Alu* insertions, however a minority of insertions have undergone movement by detectable non-classical TPRT mechanisms (21-23).

The primary difficulty in identifying novel *Alu* insertion loci stems from the highly repeated nature of the element itself. Various approaches for large-scale analyses of *Alu* and other ME types have been developed that utilize next generation sequencing platforms. Scaled sequencing of targeted *Alu* junction libraries has permitted genome-wide detection, as implemented in techniques such as Transposon-Seq (24) and ME-Scan (25,26). Such targeted methods offer high specificity and sensitivity, but are restricted by the primers used for detection and are generally subfamily-specific. A broader detection of *Alu* variant locations is possible by using computational methods to search Illumina whole genome sequence (WGS) paired reads by 'anchored' mapping. This method seeks to identify discordant read pairs where one read maps uniquely to the reference (*i.e.*, the 'anchor') and its mate maps to the element type in query (9,27,28). However, in considering read pair data alone, in their simplest form these methods are limited in the recovery of variant-genome junctions that might otherwise be captured from split read information. Specialized algorithms now offer improved breakpoint accuracy from WGS data by consideration of split-reads, soft-clipped reads, and unmapped reads during variant detection (28-31). Beyond these basic requirements, existing programs differ in the implementation of additional filters and read-support criteria to identify the subset of calls likely to be true. For most callers this includes removing MEI candidate calls that are located near reference MEs of the same class due to their higher likelihood of being false predictions ((28,30-33) but see (29)). Thus, as in all methods, a trade-off exists between sensitivity and specificity.

Having been mostly applied to high-coverage WGS, these methods require modification when applied to lower coverage data. Generally, these read-based detection methods have been developed in the context of ME variant loci discovery and reporting. Assembly approaches have received increasing application to next-generation sequencing data for breakpoint identification. Tools such as TIGRA (34), a modification of the SGA assembler used in HYDRA-MULTI (35,36), and the use of a *de bruijn* graph-based approach in SVMerge (37) have been developed to assemble structural variant breakpoints from population scale or heterogeneous tumor sequencing studies. Assembly-based approaches have also lead to increased sensitivity and specificity for the detection of SNPs and small indels (38-42) and mobile elements (30). Since sequence changes within elements can determine their activity (8), the assembly of insertion sequences can better inform our understanding of element proliferation. Additionally, insertion haplotype reconstruction offers a direct way to determine genotypes across samples based on an analysis of sequence data supporting each allele.

Here, we utilize a classic overlap-layout-consensus assembly strategy applied to ME-insertion supporting reads to completely reconstruct and characterize *Alu* insertions. We apply this approach to pooled WGS data, from 53 individuals in 7 geographically diverse populations from the Human Genome Diversity Project (HGDP) panel (43,44), to enable a comprehensive characterization of fully intact polymorphic *Alu* elements. We assess the limitations of this approach by cross-comparison with *Alu* insertions identified from PacBio long sequencing reads in a complete hydatidiform mole (CHM1) (45). Finally we also demonstrate the ability to obtain accurate genotypes based on explicit mapping to reconstructed reference and alternative alleles. We present the analysis of 1,614 fully reconstituted *Alu* insertions from these samples, including breakpoint refinement and genotyping of 1,010 insertions, with >99% accuracy. These results provide a basis for future study of such MEIs in human disease and population variation, and should facilitate similar analyses in relevant non-human models.

MATERIAL AND METHODS

Samples

We analysed whole genome, 2x101 bp Illumina read sequence data from a subset of 53 samples and 7 populations from the HGDP: Cambodia (HGDP00711, HGDP00712, HGDP00713, HGDP00715, HGDP00716, HGDP00719, HGDP00720, HGDP00721), Pathan (HGDP00213, HGDP00222, HGDP00232, HGDP00237, HGDP00239, HGDP00243, HGDP00247, HGDP00258), Yakut (HGDP00948, HGDP00950, HGDP00955, HGDP00959, HGDP00960, HGDP00963, HGDP00964, HGDP00967), Maya (HGDP00854, HGDP00855, HGDP00856, HGDP00857, HGDP00858, HGDP00860, HGDP00868, HGDP00877), Mbuti Pygmy (HGDP00449, HGDP00456, HGDP00462, HGDP00471, HGDP00474, HGDP00476, HGDP01081), Mozabite (HGDP01258, HGDP01259, HGDP01262, HGDP01264, HGDP01267, HGDP01274, HGDP01275, HGDP01277), and San (HGDP00987, HGDP00991, HGDP00992, HGDP01029, HGDP01032, HGDP01036). WGS data was processed using BWA, GATK (46) and Picard (<http://picard.sourceforge.net>) as described previously (44) and is available at the Sequence Read Archive under accession SRP036155. Final datasets are ~7x coverage per sample. For analysis of CHM1 we utilized Illumina data obtained under accession

SRX652547 and ERX009608 for analysis of NA18506. Reads were mapped to the hg19/GrCh37 or hg18/build36 genomes as appropriate using same procedures described above.

Non-reference *Alu* discovery

We performed anchored read pair mapping of our WGS data using RetroSeq (28) to identify non-reference *Alu* variants relative to the GRCh37/hg19 genome assembly. To identify candidate insertion loci, RetroSeq 'discover' was run on individual BAM files for each sample to identify discordant read pairs with one read mapping uniquely to the reference genome and its pair to an *Alu* consensus or to an annotated *Alu* present in the reference. A FASTA file of the available *Alu* consensus sequences was obtained from RepBase (47) and reference *Alu* elements were excluded using existing RepeatMasker (48) annotations. Next, candidate insertion loci were assessed using RetroSeq 'call'. For this analysis, we combined the supporting read information discovered in each individual and ran the 'call' phase on a combined BAM consisting of all samples. In the 'call' phase, we required a minimum read support of 2 supporting read pairs per call (-reads). A maximum read-depth of 1000 (-depth; default is 200) was utilized for regions surrounding each call in order to accommodate the increased coverage of the merged BAM. Finally, any output call within 500 bp of an annotated *Alu* insertion was excluded using the bedtools window command (49) and RepeatMasker hg19 reference annotations (48). Unless otherwise noted, any other RetroSeq options were run at the default settings. Final *Alu* calls having met the further criteria of a filter tag FL=6,7, or 8 were selected for subsequent analysis.

Assembly of non-reference *Alu* elements

De novo assembly of insertion-supporting reads for each candidate insertion was performed using the CAP3 assembler to utilize the overlap-layout-consensus algorithm (50). *Alu* insertion-supporting read pairs were extracted from the BAM file corresponding to each sample. For each candidate insertion site, we defined a window of 200 bp around the predicted breakpoint, and within that window extracted 1) read-pairs reported to support an insertion at that site based on RetroSeq outputs, and 2) read-pairs with a soft-clipped segment at least 20bp in length with a mean quality ≥ 20 . We then performed CAP3 assembly using the extracted reads per site, with parameters chosen to account for shorter matches that could be expected from 101 bp reads (-c 25 -j 31 -o 16 -s 251 -z 1 -c 10). CAP3 also utilizes read-pair information to report scaffolds of contigs that are linked together but without an assembled overlap. We merged such contigs together, separated by 300 'N' characters to represent sequence gaps in the assembly. The resulting contigs and scaffolds were analysed using RepeatMasker (48) to identify *Alu*-containing contigs. Of these 2,971 candidate assembled sites were identified that contain an *Alu* element (≥ 30 bp match) and at least 30bp of flanking non-gap sequence in our assemblies.

Breakpoint determination

Exact breakpoints for the assembled *Alu* variants were recovered utilizing a multiple alignment-based approach similar to an approach previously described (51,52). Orientation of candidate insertion sequences relative to the reference genome was determined using BLAT (53). Candidate breakpoints were then identified using *miropeats* (54) followed by a semi-automated parsing process. In turn, a global alignment was obtained for sequences from the two insertion breakpoints to the corresponding segment on the reference genome using *stretcher* (55) with default parameters, to generate pairwise alignments for two sequences aligned independently against the third (*i.e.*, reference) sequence. A 3-way alignment was then created from the two pair-wise alignments by inserting gaps into the alignment as appropriate. Alignment columns were scored as having either a match among all three sequences ('*'), a match between the left insertion breakpoint and the reference ('1'), a match between the right insertion breakpoint and the reference ('2'), or mismatch among all sequences ('N'). We then computed an alignment score across the left and right breakpoint sequence, with matches between the target sequence and the genome sequence ('1' or '*' for the left breakpoint and '2' or '*' for the right breakpoint) resulting in a score of +1, a sequence mismatch among all three sequences a score of -1, and a match among the reference and the other breakpoint a score of -3. The same procedure was applied to the right breakpoint, except the score was tabulated from right to left across the 3-way alignment. The breakpoint was then interpreted as the position where the maximum cumulative score was reached respective to the reference; overlapping sequence coordinates on the reference allele indicate the extent of TSDs. Of note, 1) TSDs are defined from the alignment itself, without regard to the insertion boundaries, and 2) in scoring, a small degree of divergence among putative TSDs is permitted, resulting in longer TSDs for some sites than when 100% identity is required. Visualizations of the resulting aligned sequences with breakpoint annotations were constructed and subjected to manual review. When necessary, the sequences extracted for breakpoint alignment were adjusted and the alignment and scoring scheme described above then repeated until a final curated set of 1,614 assembled insertions was obtained. Information pertaining to insertion site (locus siteID and determined breakpoint), assembled sequence, subfamily, and predicted TSD is summarized in Table S1 for each site.

Sub-family assignment and analysis

A multiple sequence alignment was constructed of *Alu* sequences extracted from the set of 1,614 assembled insertions (Table S1) as well as 43 *Alu* consensus sequences obtained from RepBase using MUSCLE v3.8.31 (56) run with default parameters. Poly-A tail regions were trimmed from the resulting multiple sequence alignment and the proportion of sequence differences between each element and each sub-family consensus were tabulated. Alignments for the 1,010 sites suitable for genotyping were utilized to assess the extent of the recovered element length relative to the subfamily consensus.

Validation

A subset of 35 assembled *Alu* insertions were validated by Sanger sequencing. Primer pairs were designed to ensure amplification across both predicted breakpoints, and the subsequent mapping of

those sequenced products uniquely to the hg19 reference, permitting comparison of amplified fragment sizes. Coordinates for each insertion were considered based on unique mapping of CAP3 assembled contigs and subsequent breakpoint analysis to the hg19 build. We extracted ~500bp in either direction of each insertion from the hg19 reference (UCSC Genome Browser; <http://genome.ucsc.edu/>). The sequence was masked using RepeatMasker, and primers designed to include ~150bp to ~200bp in either direction of the predicted insertion, avoiding masked sequence when possible. Each primer set was analysed by *in silico* PCR and BLAT to the hg19 reference to ensure site-specific target amplification predictions overlapping each breakpoint, and to infer product size predictions for either allele. All primers were designed using Primer3v.0.4.0 (57) (<http://bioinfo.ut.ee/primer3-0.4.0/>) and purchased from IDT. Loci examined, primers, and samples analysed for each site are summarized in Table S2.

All PCRs were performed with ~50ng of genomic DNA as template along with 1.5-2.5 μM Mg^{++} , 200 μM dNTPs, 0.2 μM each primer, and 2.5 U Platinum Taq Polymerase (Invitrogen). Reactions were run under conditions of 2 min denaturation at 95 °C; 35 cycles of [95 °C 30 sec, 55 °C to 59 °C 30 sec, 72 °C 2 min]; and a final extension at 72°C for 10 min. For each PCR reaction, 10 μL were analysed by electrophoresis in 1% agarose in 1 x TBE. Products from at least one positive reaction per locus were sequenced to confirm amplification of the desired product and its mapping to the hg19 reference. When possible, PCR products from a homozygous individual were sequenced; otherwise the insertion-supporting fragment was gel-extracted (Qiagen), and the products eluted in water and subjected to sequencing. Traces obtained for each insertion allele were aligned to the corresponding reference allele and CAP3 assembled contig in order to confirm the presence of the *Alu* insertion, TSDs, and agreement in nucleotide sequence between the validated and assembled insertion. Individual alignments for corresponding to each validated site are in Figures S2 and S3.

Comparison with previous studies

For comparison with Chaisson *et al* calls on CHM1 (45), we utilized insertion positions based on hg19 coordinates from <http://eichlerlab.gs.washington.edu/publications/chm1-structural-variation/>. These data consisted of 1,254 total calls classified as “AluYsimple”, “AluSsimple”, “AluSTR”, or “AluMosaic”. Overlapping calls were counted from intersection of any call located within 100 bp. The 1,727 *Alu* insertion calls for NA18506 were based on alu-detect analysis (29) of ERX009608 as obtained from the Sequence Read Archive. Consistent with that study, calls were relative to the hg18 genome assembly, and any call located within 100 bp was counted as an intersecting site.

Genotyping

We performed *in silico* genotyping by mapping relevant reads to a representation of the complete insertion and reference alleles for each site. The reference allele consisted of 600 bp of sequence upstream and downstream of the start and end of any inferred TSD extracted from the hg19 reference. Based on the aligned breakpoints, insertion alleles were created by replacing the appropriate portion of this sequence with insertion sequence, accounting for inferred TSDs or target site deletions. For

each site, these insertion and reference alleles constituted the target genome for mapping of reads. A BWA index was constructed from each (bwa version 0.5.9). Mapping and analysis was performed separately for each sample and each site. We extracted read-pairs with at least one read having an original mapping within the coordinates of the targeted reference allele with a MAPQ ≥ 20 . The extracted read-pairs were then aligned to the site reference and alternative sequences using bwa aln and bwa sampe (version 0.5.9). We then calculated genotype likelihoods based on the number of read pairs mapping to the insertion or reference alleles, considering the resulting MAPQ values as error probabilities as previously described (58). Read-pairs with equal mappings between reference and insertion sequences have a MAPQ of 0 and do not contribute.

Genotypes were obtained from the resulting raw genotype likelihoods using one of two approaches. For sites on the autosomes and the pseudoautosomal region of the X chromosome, genotype likelihoods for *Alu* insertions were processed, along with previously calculated SNP genotypes using LD-aware refinement using Beagle 3.3.2 (with options maxlr=5000, niteration=10, nsamples=30, maxwindow=2000) (59). For sites on the X chromosome, genotypes were obtained using a ploidy-aware expectation-maximization (EM) algorithm that utilized the genotype likelihoods and assumes Hardy-Weinberg Equilibrium across all 53 samples. Briefly, we follow (58) to estimate the allele frequency for each site via EM. Using the estimated allele frequency, we then determine genotype prior probabilities for X-linked alleles assuming Hardy-Weinberg equilibrium. These genotype priors are then combined with the already computed genotype likelihoods to identify the sample genotype with the highest posterior probability. Principal component analysis was performed on the resulting autosomal genotypes using the smartpca program from the EIGENSOFT package (60). Predicted genotypes for all 1,010 sites are provided in Table S5.

Genotype Validation

In order to validate *in silico* genotyping and permit estimation of genotyping accuracy, 11 random insertion loci were screened from a panel of ten individuals utilizing gel band assays, for a total of 110 predicted genotypes (Table S6). Locus-specific primer sets flanking each insertion locus were designed as above (see Validation). Primer pairs per locus were then used in PCR amplification of each sample in the panel, and the products were analysed for predicted shifting patterns following electrophoresis. All PCRs were performed with a template of 0.25ng genomic DNA, in cycling conditions of 2min at 95°C; 35x [95°C 30 sec, 55°C-59°C 30 sec, 72°C 1 min], and a final 72°C extension of 3 min. 10uL were analysed in 1.2% agarose in 1xTBE. Results were interpreted by banding patterns that supported either the unoccupied or insertion allele, as based on predicted band sizes from *in silico* PCR and size information for the assembled insertion at that site. Samples utilized for sequence and PCR genotyping validations are indicated respectively in Table S6 and Figure S2.

RESULTS

Precise assembly of full-length *Alu* variants using read data

To generate an accurate and highly specific collection of non-reference polymorphic *Alu* variants from population-scale WGS data we combined methods utilizing read-based discovery of all possible insertion sites with *de novo* local assembly of supporting reads (Figure 1). We utilized WGS data from a subset of the HGP collection, specifically consisting of 2x101 bp paired-end libraries from 53 individuals across seven populations, with a median coverage of ~7x per genome. All samples were aligned to the human GRCh37/hg19 reference and processed as previously described (44). Given the coverage levels, we anticipated insertions that were private to a single individual were likely to be missed. However, we reasoned that borrowing read information across samples would increase our ability to detect rarer insertions that were nonetheless present in multiple samples, and pooled the data into a single merged BAM from all individuals for an effective coverage of ~429x. Candidate *Alu* insertions were then identified by applying RetroSeq (28) to the merged BAM. This particular program implements standard approaches to identify MEI-supporting read signatures, with performance characteristics that are comparable to other existing callers (for example in (28,31,32) and directly reports supporting reads associated with all candidate calls (Figure 1, Methods). To minimize false calls associated with reference elements, we removed any candidate call that mapped within 500 bp of an annotated *Alu* in the human reference. After filtering, this resulted in 41,365 putative *Alu* insertions with an assigned quality score 6 or higher.

We then attempted to reconstruct as many individual insertion variants as possible, including the complete *Alu* sequence, its breakpoints, and contiguous flanking sequence for each site. While recent efforts in short read assembly have focused on a *de bruijn* graph approach (61-63), we reasoned a local assembly using an overlap-layout-consensus approach would take full advantage of our data. For these purposes we utilized the program CAP3 (50) that was originally developed for the assembly of large-insert clones sequenced using capillary sequencing, but has also been applied to *de novo* assembly of short read RNA-seq (64) and metagenomic sequence data (65). For each putative site called by RetroSeq, we retrieved read pairs reported by RetroSeq and soft-slipped reads that mapped within 200 bp, requiring that the clipped portion was ≥20bp in length and had a mean quality of ≥20. Using these read sets for all putative sites, we then performed *de novo* assemblies with CAP3 (Figure 1) run with parameters adjusted for joining smaller overlaps present in shorter reads (see Methods). The resulting scaffolds were subjected to additional analyses as follows. We first identified assembled sequences that contained an *Alu* sequence using RepeatMasker (48). We further required the presence of ≥30 bp corresponding to *Alu* sequence (requiring ≥90% nucleotide identity) and recovery of ≥30 bp of flanking non-gap sequence at one end, resulting in 2,971 candidate assemblies. Because we excluded any predicted insertion from our call set that was within 500 bp of any annotated *Alu* element, each assembled scaffold was interpreted to represent the presence of a *bona-fide* non-reference insertion.

The resulting assemblies were aligned to the reference genome and breakpoints identified using a semi-automated procedure supplemented by manual curating (see methods). For these purposes, we adapted a procedure developed for the analysis of structural variant breakpoints represented in finished fosmid clone sequences (51) to utilize our locally assembled scaffolds (also see Methods). A

total of 1,614 *Alu*-containing contigs were reconstructed, each having the complete associated insertion and at least one breakpoint with ≥ 30 bp of mapped flanking sequence (Table S1, Figure S1). Five loci were predicted to have *Alu* insertions associated with deleted sequence relative to the hg19 assembly, indicating potential non-classical insertion mechanisms (21,66). From this set, 1,010 *Alu* insertions had both breakpoints flanked by at least 100 bp of non-gap assembled sequences. These 1,010 insertions were deemed to be of the highest quality and suitable for subsequent *in silico* genotyping (see below), and had sizes ranging from 77bp-495bp (median 315 bp) with predicted TSDs up to 98 bp (median 14 bp).

Sensitivity and specificity of insertion discovery using short-read assembly

A comprehensive assessment of the performance of MEI callers is hindered by the lack of an orthogonal “gold standard” call set for formal comparison. To better assess the potential limitations of *Alu* discovery using Illumina (2x100 bp) paired-end short reads, we applied our approach to two additional samples, which have been extensively characterized. The first comparison utilized Illumina WGS data generated from a complete hydatidiform mole (CHM1) (45). This sample offers particular advantage for these purposes, as it is essentially haploid throughout its genome and has been subjected to extensive characterization using long PacBio sequencing reads (mean 5.8 kbp). In their analysis, Chaisson *et al* reported 1,254 *Alu* insertions from this sample; of which 911 intersected within 100 bp of a candidate call based on the Illumina sequence data (Table 1). However, we note these raw RetroSeq calls include 18,501 predictions, implying an extremely high false discovery rate (FDR). Considering only the highest level of RetroSeq support still results in an implied FDR greater than 80%.

Mapping into highly repetitive regions using shorter sequencing reads results in an unacceptably high degree of false calls and necessitates filtering out regions near reference elements; this step is common to most mobile element callers (28,30,31,33). The longer reads available from the CHM1 sample (mean 5.8 kb) permitted interrogation of genome intervals harbouring a high number of repetitive elements. As a result, 54% of Chaisson *et al* *Alu* insertion calls are located within 500bp of an *Alu* sequence present in the human reference genome. Limiting the analysis to only those calls ≥ 500 bp from any reference *Alu* results in an increase in both sensitivity and precision. Requiring successful element assembly further increases the precision: 446 of 468 of our assembled calls intersect (within 100 bp) with CHM1 reported calls. Assuming the remaining calls are all errors implies a FDR below 5%, however we note additional analysis suggests this may be an over estimate (see Discussion).

To further investigate the precision of our assembly approach, we applied this approach to 2x100 bp Illumina WGS data from Yoruba sample NA18506, an individual that has been analysed using multiple MEI callers (27,29,30). Again excluding calls within 500 bp of an annotated insertion, initial filtering of RetroSeq calls on the NA18506 sample resulted in 1,375 putative *Alu* insertions having a quality score 6 or higher. A total of 820 *Alu* insertions were fully assembled, of which 774 intersect with calls reported by alu-detect (out of 1,727 total calls) (29), again implying a FDR around 5%.

Validation of assembled HGDP insertions

From our assembled set of 1,614 non-reference insertions in the HGDP samples, we selected 35 sites for experimental validation by PCR and sequencing, biased in favour of sites with unusual breakpoint characteristics (*i.e.*, 0-3 bp TSDs, TSDs larger than 25bp, sites with corresponding target site deletions, and insertions with predicted 5' truncations) (Table S2). In validation, our specific goals were to demonstrate 1) the presence of the *Alu* at that chromosomal location, 2) agreement between the assembled sequence with the cognate validated sequence for the insertion, and 3) contiguous sequence of each insertion with its mapped flanking regions. We obtained Sanger sequence of the insertion allele for each of the 30 assembled sites in up to two individuals predicted to have the insertion, utilizing samples that were homozygous for the insertion when possible. Sequencing was performed with primers situated both upstream and downstream of the insertion in order to account for uncertainty introduced from polymerase slippage at the poly-A tails.

For all 35/35 tested sites, we confirmed the presence of an *Alu* insertion in the tested sample. Subsequent analysis of the corresponding nucleotide alignments verified that the nucleotide sequence of each *Alu* recovered in CAP3 assembly was in complete agreement with the corresponding Sanger traces. Examples for three representative insertions are highlighted in Figure 2, illustrating the recovery of mapped *Alu*-containing contigs, breakpoint estimations at those sites, and alignment of the deduced nucleotide sequence to the CAP3 assembly. Detailed alignments corresponding to individual insertions, including visualized trace information, are provided in Figure S2.

We assessed *in silico* breakpoint estimations for each validated CAP3 assembly in comparison to the *Alu*-genome junctions as obtained from Sanger reads. Six of the 35 validated sites were predicted to have breakpoints within 100 bp of a gap in the CAP3 assembly; sequence comparison to the CAP3 assembly revealed correctly predicted breakpoints for just 1 of these 6 insertions, justifying their exclusion from subsequent *in silico* genotyping (described further below). Exact breakpoint and TSD sequences for the remaining 5 insertions were determined from Sanger sequencing (Figure S2).

Overall, 26/35 insertions had target sites that precisely agreed with the corresponding assembly. Representative examples are shown in Figure 2; properties for all validated sites are summarized in Table S3. Of the 28 validated sites that were later utilized for genotyping, just 3 were found to have incorrect breakpoints, in each case due to the absence of target site sequence adjacent to the *Alu*-poly-A tail in the CAP3 assembly. For an additional site (insertion at chr2:123330649), comparison to the Sanger traces revealed a longer inferred TSD than predicted (16 bp vs. 13 bp), resulting from a nucleotide change within the assembled poly-A stretch, thus altering the inferred target duplication by 3 bp. The remaining 24 sites (85.7%) correctly matched the breakpoints recovered from the CAP3 assembly. These included insertions with unusual breakpoint characteristics, for example, a full-length insertion having a validated TSD of 46bp (chr18:74638702). We also correctly recovered insertions with evidence of non-classical insertion mechanisms, including three sites that were correctly predicted to have short target site deletions relative to the pre-insertion allele ranging in size from 1 to

6 bp (at positions chr6:164161904, chr11:26601646, and chr12:73056650) (see also Figure 2B), and 5 elements with evidence of 5' truncation (described further below; Figure S3). We were able to completely reconstruct, with correct breakpoints, insertions that were within other repetitive sequence classes (Figure 2C). Finally, we note one insertion, located at chr11:35425392 for which the recovered CAP3 contig was found to be in complete agreement with corresponding Sanger reads, however our automated identification of the TSD was 'miscalled' due to the presence of concomitant variation at this site relative to the hg19 assembly (Figure S4).

Characteristics of assembled insertions

Given the accuracy of our assemblies, we sought to more comprehensively characterize our set of reconstructed *Alu* insertions. Previous studies of full-length polymorphic elements have been mostly limited to insertions taken from an assembled reference genome (3,12), examination of trace archive data (11), or from insertions having been captured in relatively long read data (9). By making use of contemporary WGS data in *de novo* assembly, the insertion sequence itself is accurately reconstructed for analysis. Thus, utilizing our assembled contigs, we readily extracted the corresponding 1,614 *Alu* nucleotide sequences and characterized each in terms of subfamily distributions and properties.

Based on sequence divergence from *Alu* subfamily consensus sequences obtained from the most recent RepBase update (47), we were able to assign 1,452 (90%) of our insertions to one of 30 subfamilies (Table 2). We found 162 elements that were equally diverged from more than one subfamily consensus and could not be conclusively classified. Insertions from *AluY* subfamilies made up >99% of all assigned calls, with *AluYa5* and *Yb8* collectively representing more than half (62.7%) of the set. This observation was expected, given that *AluY* insertions have contributed to nearly all *Alu* genomic variation in humans, with *AluYa5* and *Yb8* being the most active subfamilies (4,8). Also as expected, insertions derived from non-*AluY* lineages were a minority, together representing less than 1% of calls that could be assigned to a subfamily (also see Discussion). These data are generally similar to previous analyses of representative intact polymorphic *Alu* in humans (9-12).

To assess the length distribution of non-reference *Alu* variants from our call set, we focused on insertions assembled from the *AluYa5* and *AluYb8* subfamilies. We reasoned that analysis of these particular subfamilies should provide the most informative resource for comparison given their representation as the majority of identified variants. We further limited analysis to those *Alu* that were suitable for genotyping, as insertions that do not meet our criteria for genotyping may erroneously appear to be truncated due to an incomplete breakpoint assembly. This resulted in an analysis set of 351 *AluYa5* and 215 *AluYb8* insertions. Based on nucleotide alignments of the assembled insertions against their respective consensus, we examined the collective coverage of assembled elements, per subfamily, in comparison to the nucleotide positions relative to their respective consensus (Figure 3A and B).

We observed that 84.9% of *AluYa5* (298/351) and 81.4% *AluYb8* (175/215) variants were full-length, or within at least 5 bp of being full-length, consistent with previous reports of the genome-wide distribution of full-length *Alu* (23,66). Comparing the length distribution of all insertions revealed a detectable minority of 5' truncations that were present in both subfamilies and exhibited a similar distribution of the apparent truncation point (Figure 3). More specifically, a subset of insertions from either subfamily was truncated ~8-45 bp from the consensus start (9.9% or 35/351 *AluYa5*, and 13.4% or 29/215 *AluYb8* insertions), and a second subset was truncated ~55-171bp from the consensus start (5.1% or 18/351 *AluYa5*, and 5% or 11/215 *AluYb8*) (Table 3). Besides having apparent 5' truncations, all but two of these assembled insertions displayed characteristics of 'standard' *Alu*, including flanking TSDs and a poly-A tail of variable length (insertions at chr13:86166445 and chr11:26601646; also see Table S1). We note the observed distribution is similar to that from two previous analyses of 10,062 reference human *Alu* (as extracted from NCBI build33) (66), and of 1,402 intact polymorphic *Alu* from the then-current dbRIP (23); aspects of both are addressed further in the Discussion.

L1 and *Alu* insertions that are truncated but otherwise standard are thought to arise from non-classical TPRT (21,23,66,67). For example, one mechanism thought to contribute to 5' truncations is a microhomology-mediated pairing of nucleotides at the genomic target 5' end with the nascent *Alu* mRNA, resulting in premature completion of TPRT (23), and in turn leaving a detectable signature. We manually examined each three way alignment of the 53 *AluYa5* and 40 *AluYb8* assembled 5' truncation events for such evidence, specifically searching for nucleotides at the 5' break that were shared with the respective *Alu* consensus at that position (23). We observed a subset of insertions with detected microhomology, with 40.9% of truncation having 1 bp of matching sequence, and 15.1% of all truncations with ≥ 2 bp shared at the 5' break (details are summarized in Table S4), though we note limitations of interpreting a single shared nucleotide as a 'true' instance of microhomology. Given this observation, the data indicate premature TPRT may account for a subset of the truncated insertions.

Insertion breakpoint distribution

We analysed the distribution of assembled insertions relative to genes based on Gencode v19 annotations (68). Of the 1,614 assembled insertions, 865 (~53.5%) were found within genes, of which 643 (~39.8%) were located within protein coding genes. Although these values are slightly higher than has been reported in previous analysis (9,10), these values are lower than expected based on random permutations of our data (924, or ~57.2% expected within all gene regions and 688, or ~42.6% in protein coding genes). Just 10 insertions (~0.61% of all calls) were found within exons, all of which were located in untranslated regions and therefore would not be predicted to disrupt coding sequence. This value is much lower than expected based on random simulations (50 sites, or ~3.1 %; $p < 0.02$), indicating potential selection against retrotransposition into exons and other coding sequence, and consistent with previous studies indicating exonic depletion of *Alu* (9,10,26,27).

A total of 708 (~43.8%) of our assembled insertions were located within repetitive sequence. The majority these insertions were found within other retrotransposon-derived elements (459, or ~28.4%, were in LINEs and 124, or ~7.6% in LTRs), and in DNA transposons (69, or ~4.2%); 22 insertions were found in minor or unknown repetitive classes. This distribution is also consistent with that observed in previous survey of non-reference *Alu* insertions (9). Since we excluded any candidate call that was near an annotated *Alu* prior to assembly, no insertion from our callset was recovered within any existing *Alu*, though a handful of insertions were observed within non-*Alu* SINE classes (e.g., from the *Mir*, FLAM, or FRAM groups). We compared these data to randomized values based on simulated uniform placement of insertions, excluding regions in the hg19 reference assembly that mapped to gaps or are annotated as *Alu*, observing no significant difference relative to random uniform placement. Based on separate simulations permitting placement within annotated *Alu* elements, we estimate 10.5% of random insertions would be near annotated *Alus*, and hence excluded. However, the Chaisson *et al* data demonstrate 54% of insertions are within 500bp of existing *Alu* elements, clearly demonstrating the non-uniform patterns of *Alu* insertions (2), and highlighting the limitations of repeat discovery using existing short-read methodologies.

Genotyping

We identified a subset of 1,010 insertions that had sites with both breakpoints at least 100 bp away from an assembly gap that were suitable for genotyping using Illumina sequencing reads. For each site, we recreated the reference and insertion haplotypes based on 600 bp flanking the inferred insertion site. For each sample, we then remapped Illumina read-pairs that mapped to each reference location against both reconstructed sequences using bwa (see Materials and Methods). We then determined genotype likelihoods based on the mapping of reads to each alternative allele, with error probabilities as indicated by the read mapping quality (58). Of note, read pairs that map equally well to the reference and insertion sequences will have a MAPQ of 0 and are uninformative for establishing the genotypes. Final genotypes were determined for each sample based on the resulting genotype likelihoods (Table S5). For sites on the autosomes and the pseudo-autosomal regions of chrX, genotypes were obtained after LD-aware refinement using BEAGLE v3 (69), in a procedure that also included SNP genotypes as previously described (70). We compared the inferred genotypes for 11 autosomal sites with PCR-based genotyping across 10 samples, and found a total concordance rate of 99% (109/110) (Figure 4A and B; predicted genotypes are in Table S6 for direct comparison). The only error among the tested calls occurred when the inferred genotype was homozygous for the insertion allele, while PCR genotyping indicates that the site is heterozygous (chr10:19550721; HGDP00476). Finally, we performed a Principle Component Analysis (PCA) of the autosomal genotypes across all 53 samples (60). As expected, individual samples largely cluster together by population with the first PC separating African from non-African samples (Figure 4C). This result further confirms the high accuracy of the inferred genotypes.

DISCUSSION

We utilized Illumina WGS paired reads to fully reconstitute a high-specificity set of 1,614 non-reference *Alu* insertions from a subset of 53 genetically diverse individuals in 7 global populations from the HGDP (43,44). Experimental interrogation of a total of 35 sites confirmed the presence of a non-reference *Alu* insertion at that site. We confirmed the presence of several insertions with aberrant assembled breakpoint characteristics, including insertions for which the TSD was absent (for example at chr17:46617220) or of extreme length (chr18:74638702, 46 bp), as well as insertions with deleted sequence relative to the hg19 reference (chr6:164161904, chr11:26601646, and chr12:73056650), and insertions with 5' truncations (also see Figures S2 and S3). Validation of one of two assembled *AluJ* insertions (at chr12:73056650; predicted allele frequency of 0.056) correctly confirmed its bimorphic presence and consensus divergence of 14.2% (Figures S2 and S5). We suggest that such insertions of a now-inactive lineage represent ancestral insertions not yet fixed or subsequently lost from the population. For total of 1,010 insertions that had at least 100 bp of assembled sequence flanking both sides, we obtained a high level of breakpoint accuracy, having perfect agreement at 25/28 sites tested (89%), including those with aberrant breakpoints and/or truncated elements. Analysis of SNPs has demonstrated that improvements in accuracy can be obtained by separating the “discovery” and “genotyping” phases of analysis (71). We therefore performed genotyping by determining genotype-likelihoods based on remapping Illumina read-pairs to the reconstructed reference and alternative haplotypes, achieving an estimated 99% genotype concordance (109 of 110 genotypes analysed).

For each of the 35 validated insertions, comparisons with Sanger sequences of those sites revealed the correct nucleotide sequence of the *Alu* insertion itself was obtained in assembly. However, a closer comparison of the sequenced TSDs at individual sites indicated that elements located near edges of assembled contigs (e.g., excluding the complete TSD length) were more likely to have incompletely assembled breakpoints. Further examination of the individual reads supporting the assembled contig indicated that this was due to aberrant joining or incomplete TSD capture of reads that covered the poly-A tract (also refer to trace data from insertions at chr12:99227704, chr22:26997608, and chrX:5781742 in the supplement). One example of this comes from our assembly of the Y *Alu* Polymorphic element (YAP) (72) located at chrY:21611993, which contained an incomplete 3' TSD. Capillary sequencing in sample HGDP00213 revealed the correct 11 bp TSD (5' AAAGAAATATA), and confirmed the presence of YAP-specific nucleotide markers (at bases 64, 207, 243, and 268 relative to the *AluY8b* consensus), as recovered by our CAP3 assembly and consistent with previous reports (Figure S2) (72). We additionally note that even when alleles are fully (and correctly) reconstructed, interpretation of the variant may not be clear. An insertion at chr11:35425392 is illustrative of this complexity. At this site, our identification of the variant breakpoints was inaccurate due to the presence of concomitant variation at this site relative to the hg19 reference, as revealed by sequencing in other individuals without the insertion to better reconstruct the structure of the pre-insertion allele (Figure S4). Notably, the CAP3 assembled insertion and proximal genomic sequence was found to be in complete agreement with corresponding Sanger reads, despite the presence of this surrounding structural variation relative to the reference build.

Having an assembled, high-specificity call set of non-reference *Alu* variants also permitted analysis of element properties. Performing this step was meant to take particular advantage of these data, as existing MEI callers are generally designed to catalogue events detected from read-based signatures within the data. Examination of individual assembled insertions showed the vast majority of elements exhibited properties consistent with classical retrotransposition, specifically being full length and the presence of a TSD and poly-A tail of variable length. However, our analysis of the length distribution of the reconstructed *AluYa5* and *Yb8* insertions also revealed that 93 (~16.4%) of this subset had evidence of having been 5' truncated, despite appearing otherwise standard, indicating insertion by potential non-classical TPRT mechanisms. We also observed evidence of at least two groups of this subset, respectively truncated ~30 to 50 bp and ~160 to 180 bp from the canonical 5' edge (Table 2 and Figure 3).

These data are consistent with a previous manual curation of 1,402 intact polymorphic *Alu* from dbRIP that characterized full-length elements available at the time (23). In that study the authors identified 115 elements (~8.2%) with apparent 5' truncations ~8-45 bp from the *Alu* start (~8.2%) and 89 elements had ~55-171 bp truncations (6.3%) (23). The authors proposed a model of microhomology-mediated nucleotide pairing of the 5' end of the genomic strand with the *Alu* RNA, having observed 41.2% events with nucleotides at the 5' break shared with the *Alu* consensus at that position. However a single shared base supported the majority of the truncations; considering ≥ 2 bp accounted for 16.7% of their observed events. We searched our own data corresponding to all 5' truncation events, and observed insertions with similar levels of putative microhomology: 15.1% had at least 2 shared bases at the 5' edge, and 40.9% of insertions shared 1 base; although tentatively considered to represent true cases of microhomology, this is greater than 25% of sites expected at random. One other study reported similar instances of *Alu* truncation events (1,005/10,062 or ~10.5%), but found little to no statistical support for base overlap at the 5' breaks (~29% 1 bp; ~13% ≥ 2 bp) (66). Given that the 5' *Alu* end is particularly GC rich, this suggests such a 'mis'-pairing during TPRT would account for a minority of observed truncations. In support, we examined the nick site for truncations with and without putative signatures of microhomology and found no difference in preference, further confirming that both classes contained the canonical L1 ORF2p protein (ORF2p) nick site, 5' T₄/A₂ (the '/' indicating the site of cleavage) (73,74).

We note that secondary structure of the *Alu* RNA itself may drive the non-random distribution of 5' truncation points. The bases associated with the points of truncation, near ~45bp and ~180 from the *Alu* start are also coincident with the predicted hairpin structure in the folded RNA (8). The *Alu* RNA is reverse transcribed by the L1 encoded ORF2p, which pauses at sites of RNA secondary structure such as poly-purine tracts and stem-loops (75). Additionally, both truncation regions are located directly 3' to predicted SRP9/14 binding locations (6,76). Although SRP9/14 binding is necessary for efficient retrotransposition, the younger *AluS* and *AluY* subfamilies contain nucleotide substitutions that reduce SRP9/14 binding affinity, suggesting that efficient displacement of bound SRP9/14 is important for the successful propagation of these elements (4,8). This suggests that the characteristic location of 5' truncations may be a consequence of ORF2p pausing and premature disengaging from

the *Alu* RNA during reverse transcription. Regardless, the data indicate non-classical and/or premature TPRT mechanisms may account for a subset of the truncated insertions, although alternative mechanisms cannot be ruled out (21,22,67).

Although our assembled calls are of high quality, our discovery process suffers from the same limitations that are common to other studies utilizing NGS. For example, because of the variability in coverage across samples, we are likely missing sites present in only one or a small number of the analysed samples. Additionally, by requiring successful element assembly, we focus on a highly reliable call set that will have reduced sensitivity. To further explore these issues, we compared our approach with *Alu* insertions identified in the CHM1 hydatidiform mole using PacBio reads (45). This analysis highlights trade-offs in sensitivity and specificity that are inherent in any discovery approach and clearly demonstrates the challenges for discovering *Alu* insertions that are coincident with existing *Alu* elements in the reference (Table 1). Considering insertions near other *Alu* results in thousands of false calls, necessitating subsequent filtering steps (for example in (28,30,31) but see (29)). When considering only insertions that are distant from existing reference *Alu* insertions, our assembly approach has a moderate sensitivity (77%), and an FDR less than 5%. The true false-call rate of the insertions assembled from Illumina data is likely to be lower. Of the 16 assembled insertions that do not intersect with Chaisson *et al* calls, three correspond to *Alu* insertions reported by Chaisson *et al* that are near existing *Alu* elements, but remained in our call set because the initial RetroSeq prediction was at least 500bp away from a reference element. An additional three calls correspond to more complex variants reported by Chaisson *et al* involving *Alu* and other repetitive sequence. Counting these six sites reduces the apparent FDR of the assembly based approach to 3.4%. Since the Chaisson *et al* call set is itself likely missing some calls due to variable coverage and mapping ambiguities, we consider these rates to be merely approximations. Despite the ability of local-assembly approaches to recover *Alu* insertions with high precision, it is clear that analysis of insertions, and other types of structural variation, within highly repetitive sequence using comparatively short reads remains a major challenge. Analysis of insertion site preferences, population diversity, and insertion rates across individuals and somatic tissues should be cognizant of the severe challenges posed for accurate variant detection in repetitive regions.

ACKNOWLEDGEMENT

We thank Sarah Emery for technical advice, Ryan Mills for meaningful input and critical reading of the manuscript, John Moran for advice on RNA secondary structure, and Amanda Pendleton for discussion and editorial comments.

DATA ACCESS

The HGDP *Alu* sequence data from this study is available in the NCBI Database of genomic structural variation (dbVar; <http://www.ncbi.nlm.nih.gov/dbvar/>) under accession nstd109 and is also available in Supplemental Table S1. The pipeline for *Alu* assembly and breakpoint analysis is available at <https://github.com/KiddLab/insertion-assembly>.

AUTHOR CONTRIBUTION

JHW and JMK designed the study. JHW, AB and NMB performed necessary PCR, sequencing, and sequence-based analysis. JHW and JMK were responsible for all other data analysis. JHW and JMK wrote the paper. All authors have read and approved the final manuscript.

FUNDING

This work was supported by the National Institutes of Health [1DP5OD009154 to J.M.K]. JHW was the recipient of a NRSA Fellowship from the National Institute of Health [F32GM112339].

REFERENCES

1. Kazazian, H.H., Jr. (2004) Mobile elements: drivers of genome evolution. *Science*, **303**, 1626-1632.
2. Batzer, M.A. and Deininger, P.L. (2002) Alu repeats and human genomic diversity. *Nature reviews. Genetics*, **3**, 370-379.
3. 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**, 860-921.
4. Mills, R.E., Bennett, E.A., Iskow, R.C. and Devine, S.E. (2007) Which transposable elements are active in the human genome? *Trends in genetics : TIG*, **23**, 183-191.
5. Richardson, S.R., Morell, S. and Faulkner, G.J. (2014) L1 retrotransposons and somatic mosaicism in the brain. *Annual review of genetics*, **48**, 1-27.
6. Deininger, P. (2011) Alu elements: know the SINEs. *Genome biology*, **12**, 236.
7. Jurka, J. and Smith, T. (1988) A fundamental division in the Alu family of repeated sequences. *Proceedings of the National Academy of Sciences of the United States of America*, **85**, 4775-4778.
8. Bennett, E.A., Keller, H., Mills, R.E., Schmidt, S., Moran, J.V., Weichenrieder, O. and Devine, S.E. (2008) Active Alu retrotransposons in the human genome. *Genome research*, **18**, 1875-1883.
9. Stewart, C., Kural, D., Stromberg, M.P., Walker, J.A., Konkel, M.K., Stutz, A.M., Urban, A.E., Grubert, F., Lam, H.Y., Lee, W.P. *et al.* (2011) A comprehensive map of mobile element insertion polymorphisms in humans. *PLoS genetics*, **7**, e1002236.
10. Xing, J., Zhang, Y., Han, K., Salem, A.H., Sen, S.K., Huff, C.D., Zhou, Q., Kirkness, E.F., Levy, S., Batzer, M.A. *et al.* (2009) Mobile elements create structural variation: analysis of a complete human genome. *Genome research*, **19**, 1516-1526.
11. Bennett, E.A., Coleman, L.E., Tsui, C., Pittard, W.S. and Devine, S.E. (2004) Natural genetic variation caused by transposable elements in humans. *Genetics*, **168**, 933-951.
12. Wang, J., Song, L., Gonder, M.K., Azrak, S., Ray, D.A., Batzer, M.A., Tishkoff, S.A. and Liang, P. (2006) Whole genome computational comparative genomics: A fruitful approach for ascertaining Alu insertion polymorphisms. *Gene*, **365**, 11-20.
13. Carroll, M.L., Roy-Engel, A.M., Nguyen, S.V., Salem, A.H., Vogel, E., Vincent, B., Myers, J., Ahmad, Z., Nguyen, L., Sammarco, M. *et al.* (2001) Large-scale analysis of the Alu Ya5 and Yb8 subfamilies and their contribution to human genomic diversity. *Journal of molecular biology*, **311**, 17-40.
14. Cordaux, R., Hedges, D.J., Herke, S.W. and Batzer, M.A. (2006) Estimating the retrotransposition rate of human Alu elements. *Gene*, **373**, 134-137.
15. Hancks, D.C. and Kazazian, H.H., Jr. (2012) Active human retrotransposons: variation and disease. *Current opinion in genetics & development*, **22**, 191-203.

16. Callinan, P.A. and Batzer, M.A. (2006) Retrotransposable elements and human disease. *Genome dynamics*, **1**, 104-115.
17. Sen, S.K., Han, K., Wang, J., Lee, J., Wang, H., Callinan, P.A., Dyer, M., Cordaux, R., Liang, P. and Batzer, M.A. (2006) Human genomic deletions mediated by recombination between Alu elements. *American journal of human genetics*, **79**, 41-53.
18. Solyom, S. and Kazazian, H.H., Jr. (2012) Mobile elements in the human genome: implications for disease. *Genome medicine*, **4**, 12.
19. Fuhrman, S.A., Deininger, P.L., LaPorte, P., Friedmann, T. and Geiduschek, E.P. (1981) Analysis of transcription of the human Alu family ubiquitous repeating element by eukaryotic RNA polymerase III. *Nucleic acids research*, **9**, 6439-6456.
20. Dewannieux, M., Esnault, C. and Heidmann, T. (2003) LINE-mediated retrotransposition of marked Alu sequences. *Nature genetics*, **35**, 41-48.
21. Srikanta, D., Sen, S.K., Conlin, E.M. and Batzer, M.A. (2009) Internal priming: an opportunistic pathway for L1 and Alu retrotransposition in hominins. *Gene*, **448**, 233-241.
22. Callinan, P.A., Wang, J., Herke, S.W., Garber, R.K., Liang, P. and Batzer, M.A. (2005) Alu retrotransposition-mediated deletion. *Journal of molecular biology*, **348**, 791-800.
23. Chen, J.M., Ferec, C. and Cooper, D.N. (2007) Mechanism of Alu integration into the human genome. *Genomic medicine*, **1**, 9-17.
24. Iskow, R.C., McCabe, M.T., Mills, R.E., Torene, S., Pittard, W.S., Neuwald, A.F., Van Meir, E.G., Vertino, P.M. and Devine, S.E. (2010) Natural mutagenesis of human genomes by endogenous retrotransposons. *Cell*, **141**, 1253-1261.
25. Witherspoon, D.J., Xing, J., Zhang, Y., Watkins, W.S., Batzer, M.A. and Jorde, L.B. (2010) Mobile element scanning (ME-Scan) by targeted high-throughput sequencing. *BMC genomics*, **11**, 410.
26. Witherspoon, D.J., Zhang, Y., Xing, J., Watkins, W.S., Ha, H., Batzer, M.A. and Jorde, L.B. (2013) Mobile element scanning (ME-Scan) identifies thousands of novel Alu insertions in diverse human populations. *Genome research*, **23**, 1170-1181.
27. Hormozdiari, F., Alkan, C., Ventura, M., Hajirasouliha, I., Malig, M., Hach, F., Yorukoglu, D., Dao, P., Bakhshi, M., Sahinalp, S.C. *et al.* (2011) Alu repeat discovery and characterization within human genomes. *Genome research*, **21**, 840-849.
28. Keane, T.M., Wong, K. and Adams, D.J. (2013) RetroSeq: transposable element discovery from next-generation sequencing data. *Bioinformatics*, **29**, 389-390.
29. David, M., Mustafa, H. and Brudno, M. (2013) Detecting Alu insertions from high-throughput sequencing data. *Nucleic acids research*, **41**, e169.
30. Lee, E., Iskow, R., Yang, L., Gokcumen, O., Haseley, P., Luquette, L.J., 3rd, Lohr, J.G., Harris, C.C., Ding, L., Wilson, R.K. *et al.* (2012) Landscape of somatic retrotransposition in human cancers. *Science*, **337**, 967-971.
31. Thung, D., de Ligt, J., Vissers, L., Steehouwer, M., Kroon, M., de Vries, P., Slagboom, E.P., Ye, K., Veltman, J.A. and Hehir-Kwa, J.Y. (2014) Mobster: accurate detection of mobile element insertions in next generation sequencing data. *Genome biology*, **15**, 488.
32. Wu, J., Lee, W.P., Ward, A., Walker, J.A., Konkelt, M.K., Batzer, M.A. and Marth, G.T. (2014) Tangram: a comprehensive toolbox for mobile element insertion detection. *BMC genomics*, **15**, 795.
33. Tubio, J.M., Li, Y., Ju, Y.S., Martincorena, I., Cooke, S.L., Tojo, M., Gundem, G., Pipinikas, C.P., Zamora, J., Raine, K. *et al.* (2014) Mobile DNA in cancer. Extensive transduction of nonrepetitive DNA mediated by L1 retrotransposition in cancer genomes. *Science*, **345**, 1251343.
34. Chen, K., Chen, L., Fan, X., Wallis, J., Ding, L. and Weinstock, G. (2014) TIGRA: a targeted iterative graph routing assembler for breakpoint assembly. *Genome research*, **24**, 310-317.
35. Malhotra, A., Lindberg, M., Faust, G.G., Leibowitz, M.L., Clark, R.A., Layer, R.M., Quinlan, A.R. and Hall, I.M. (2013) Breakpoint profiling of 64 cancer genomes reveals numerous complex

- rearrangements spawned by homology-independent mechanisms. *Genome research*, **23**, 762-776.
36. Simpson, J.T. and Durbin, R. (2012) Efficient de novo assembly of large genomes using compressed data structures. *Genome research*, **22**, 549-556.
37. Wong, K., Keane, T.M., Stalker, J. and Adams, D.J. (2010) Enhanced structural variant and breakpoint detection using SVMerge by integration of multiple detection methods and local assembly. *Genome biology*, **11**, R128.
38. Narzisi, G., O'Rawe, J.A., Iossifov, I., Fang, H., Lee, Y.H., Wang, Z., Wu, Y., Lyon, G.J., Wigler, M. and Schatz, M.C. (2014) Accurate de novo and transmitted indel detection in exome-capture data using microassembly. *Nature methods*, **11**, 1033-1036.
39. Li, S., Li, R., Li, H., Lu, J., Li, Y., Bolund, L., Schierup, M.H. and Wang, J. (2013) SOAPindel: efficient identification of indels from short paired reads. *Genome research*, **23**, 195-200.
40. Iqbal, Z., Caccamo, M., Turner, I., Flicek, P. and McVean, G. (2012) De novo assembly and genotyping of variants using colored de Bruijn graphs. *Nature genetics*, **44**, 226-232.
41. Li, H. (2012) Exploring single-sample SNP and INDEL calling with whole-genome de novo assembly. *Bioinformatics*, **28**, 1838-1844.
42. Rimmer, A., Phan, H., Mathieson, I., Iqbal, Z., Twigg, S.R., Consortium, W.G.S., Wilkie, A.O., McVean, G. and Lunter, G. (2014) Integrating mapping-, assembly- and haplotype-based approaches for calling variants in clinical sequencing applications. *Nature genetics*, **46**, 912-918.
43. Cann, H.M., de Toma, C., Cazes, L., Legrand, M.F., Morel, V., Piouffre, L., Bodmer, J., Bodmer, W.F., Bonne-Tamir, B., Cambon-Thomsen, A. *et al.* (2002) A human genome diversity cell line panel. *Science*, **296**, 261-262.
44. Martin, A.R., Costa, H.A., Lappalainen, T., Henn, B.M., Kidd, J.M., Yee, M.C., Grubert, F., Cann, H.M., Snyder, M., Montgomery, S.B. *et al.* (2014) Transcriptome sequencing from diverse human populations reveals differentiated regulatory architecture. *PLoS genetics*, **10**, e1004549.
45. Chaisson, M.J., Huddleston, J., Dennis, M.Y., Sudmant, P.H., Malig, M., Hormozdiari, F., Antonacci, F., Surti, U., Sandstrom, R., Boitano, M. *et al.* (2015) Resolving the complexity of the human genome using single-molecule sequencing. *Nature*, **517**, 608-611.
46. McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernysky, A., Garimella, K., Altshuler, D., Gabriel, S., Daly, M. *et al.* (2010) The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome research*, **20**, 1297-1303.
47. Jurka, J., Kapitonov, V.V., Pavlicek, A., Klonowski, P., Kohany, O. and Walichiewicz, J. (2005) Repbase Update, a database of eukaryotic repetitive elements. *Cytogenetic and genome research*, **110**, 462-467.
48. Smit, A.F., Hubley, R. and Green, P. (1996-2010). 3.0 ed.
49. Quinlan, A.R. (2014) BEDTools: The Swiss-Army Tool for Genome Feature Analysis. *Current protocols in bioinformatics / editorial board, Andreas D. Baxeavanis ... [et al.]*, **47**, 11 12 11-11 12 34.
50. Huang, X. and Madan, A. (1999) CAP3: A DNA sequence assembly program. *Genome research*, **9**, 868-877.
51. Kidd, J.M., Graves, T., Newman, T.L., Fulton, R., Hayden, H.S., Malig, M., Kallicki, J., Kaul, R., Wilson, R.K. and Eichler, E.E. (2010) A human genome structural variation sequencing resource reveals insights into mutational mechanisms. *Cell*, **143**, 837-847.
52. Kidd, J.M., Sampas, N., Antonacci, F., Graves, T., Fulton, R., Hayden, H.S., Alkan, C., Malig, M., Ventura, M., Giannuzzi, G. *et al.* (2010) Characterization of missing human genome sequences and copy-number polymorphic insertions. *Nature methods*, **7**, 365-371.
53. Kent, W.J. (2002) BLAT--the BLAST-like alignment tool. *Genome research*, **12**, 656-664.

54. Parsons, J. (1995) Miropeats: graphical DNA sequence comparisons. *Comput Appl Biosci*, **11**, 615-619.
55. Rice, P., Longden, I. and Bleasby, A. (2000) EMBOSS: the European Molecular Biology Open Software Suite. *Trends in genetics : TIG*, **16**, 276-277.
56. Edgar, R.C. (2004) MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic acids research*, **32**, 1792-1797.
57. Koressaar, T. and Remm, M. (2007) Enhancements and modifications of primer design program Primer3. *Bioinformatics*, **23**, 1289-1291.
58. Li, H. (2011) A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics*, **27**, 2987-2993.
59. Browning, S.R. and Browning, B.L. (2007) Rapid and accurate haplotype phasing and missing-data inference for whole-genome association studies by use of localized haplotype clustering. *American journal of human genetics*, **81**, 1084-1097.
60. Patterson, N., Price, A.L. and Reich, D. (2006) Population structure and eigenanalysis. *PLoS genetics*, **2**, e190.
61. Lee, H. and Tang, H. (2012) Next-generation sequencing technologies and fragment assembly algorithms. *Methods in molecular biology*, **855**, 155-174.
62. Li, H. and Homer, N. (2010) A survey of sequence alignment algorithms for next-generation sequencing. *Briefings in bioinformatics*, **11**, 473-483.
63. Miller, J.R., Koren, S. and Sutton, G. (2010) Assembly algorithms for next-generation sequencing data. *Genomics*, **95**, 315-327.
64. Yang, Y. and Smith, S.A. (2013) Optimizing de novo assembly of short-read RNA-seq data for phylogenomics. *BMC genomics*, **14**, 328.
65. Reddy, R.M., Mohammed, M.H. and Mande, S.S. (2014) MetaCAA: A clustering-aided methodology for efficient assembly of metagenomic datasets. *Genomics*, **103**, 161-168.
66. Zingler, N., Willhoeft, U., Brose, H.P., Schoder, V., Jahns, T., Hanschmann, K.M., Morrish, T.A., Lower, J. and Schumann, G.G. (2005) Analysis of 5' junctions of human LINE-1 and Alu retrotransposons suggests an alternative model for 5'-end attachment requiring microhomology-mediated end-joining. *Genome research*, **15**, 780-789.
67. Srikanta, D., Sen, S.K., Huang, C.T., Conlin, E.M., Rhodes, R.M. and Batzer, M.A. (2009) An alternative pathway for Alu retrotransposition suggests a role in DNA double-strand break repair. *Genomics*, **93**, 205-212.
68. Harrow, J., Frankish, A., Gonzalez, J.M., Tapanari, E., Diekhans, M., Kokocinski, F., Aken, B.L., Barrell, D., Zadissa, A., Searle, S. *et al.* (2012) GENCODE: the reference human genome annotation for The ENCODE Project. *Genome research*, **22**, 1760-1774.
69. Browning, B.L. and Browning, S.R. (2009) A unified approach to genotype imputation and haplotype-phase inference for large data sets of trios and unrelated individuals. *American journal of human genetics*, **84**, 210-223.
70. DePristo, M.A., Banks, E., Poplin, R., Garimella, K.V., Maguire, J.R., Hartl, C., Philippakis, A.A., del Angel, G., Rivas, M.A., Hanna, M. *et al.* (2011) A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nature genetics*, **43**, 491-498.
71. Nielsen, R., Paul, J.S., Albrechtsen, A. and Song, Y.S. (2011) Genotype and SNP calling from next-generation sequencing data. *Nature reviews. Genetics*, **12**, 443-451.
72. Hammer, M.F. (1994) A recent insertion of an alu element on the Y chromosome is a useful marker for human population studies. *Molecular biology and evolution*, **11**, 749-761.
73. Feng, Q., Moran, J.V., Kazazian, H.H., Jr. and Boeke, J.D. (1996) Human L1 retrotransposon encodes a conserved endonuclease required for retrotransposition. *Cell*, **87**, 905-916.
74. Cost, G.J. and Boeke, J.D. (1998) Targeting of human retrotransposon integration is directed by the specificity of the L1 endonuclease for regions of unusual DNA structure. *Biochemistry*, **37**, 18081-18093.

75. Piskareva, O. and Schmatchenko, V. (2006) DNA polymerization by the reverse transcriptase of the human L1 retrotransposon on its own template in vitro. *FEBS letters*, **580**, 661-668.
76. Weichenrieder, O., Wild, K., Strub, K. and Cusack, S. (2000) Structure and assembly of the Alu domain of the mammalian signal recognition particle. *Nature*, **408**, 167-173.

TABLE AND FIGURES LEGENDS

Table 1. Sensitivity and Specificity Analysis.

A. All Calls

Call Set	Predicted Insertions	Overlap with Chaisson et al.	Total Chaisson et al. Sites	Chaisson et al. Sites With Overlap	Sensitivity	Precision	False Discovery Rate
RetroSeq Calls	18,501	911	1,254	896	71.5%	4.9%	95.1%
Support Level ≥ 6	12,721	908	1,254	893	71.2%	7.1%	92.9%
Support Level 8	5,294	889	1,254	874	69.7%	16.8%	83.2%
Assembled	468	449	1,254	449	35.8%	95.9%	4.1%

B. Calls at least 500bp distant from reference Alus

Call Set	Predicted Insertions	Overlap with Chaisson et al.	Total Chaisson et al. Sites	Chaisson et al. Sites With Overlap	Sensitivity	Precision	False Discovery Rate
RetroSeq Calls	924	479	578	473	81.8%	51.8%	48.2%
Support Level ≥ 6	615	479	578	473	81.8%	77.9%	22.1%
Support Level 8	520	474	578	468	81.0%	91.2%	8.8%
Assembled	468	446	578	446	77.2%	95.3%	4.7%

Table 2. Classification of assembled *Alu* variants.

Subfamily	Count	% Total	% Assigned
<i>AluY</i>			
<i>AluY</i>	87	5.4%	6.0%
<i>AluYa1</i>	1	0.1%	0.1%
<i>AluYa4</i>	83	5.1%	5.7%
<i>AluYa5</i>	537	33.3%	37.0%
<i>AluYa8</i>	5	0.3%	0.3%
<i>AluYb3a1</i>	5	0.3%	0.3%
<i>AluYb8</i>	374	23.2%	25.8%
<i>AluYb9</i>	50	3.1%	3.4%
<i>AluYc1</i>	107	6.6%	7.4%
<i>AluYc2</i>	12	0.7%	0.8%
<i>AluYc3</i>	1	0.1%	0.1%
<i>AluYd8</i>	10	0.6%	0.7%
<i>AluYe5</i>	68	4.2%	4.7%
<i>AluYf1</i>	7	0.4%	0.5%
<i>AluYg6</i>	50	3.1%	3.4%
<i>AluYh3</i>	1	0.1%	0.1%
<i>AluYh7</i>	5	0.3%	0.3%
<i>AluYi6</i>	16	1.0%	1.1%
<i>AluYi6_4d</i>	7	0.4%	0.5%
<i>AluYj4</i>	2	0.1%	0.1%
<i>AluYk11</i>	4	0.2%	0.3%
<i>AluYk12</i>	5	0.3%	0.3%
<i>AluYk13</i>	5	0.3%	0.3%
<i>AluS</i>			
<i>AluSc</i>	1	0.1%	0.1%
<i>AluSc8</i>	1	0.1%	0.1%
<i>AluSp</i>	1	0.1%	0.1%
<i>AluSq</i>	2	0.1%	0.1%
<i>AluSq2</i>	1	0.1%	0.1%
<i>AluSx3</i>	2	0.1%	0.1%
<i>AluJ</i>			
<i>AluJb</i>	2	0.1%	0.1%
Unclassified	162	10.0%	
Total	1614	100%	

Table 3. Truncation analysis of *Alu* variants.

Subfamily	Start	Count	% Total
<i>AluYa5</i>			
	1 - 5 bp	298	84.9%
	8 - 45 bp	35	9.9%
	57 - 166 bp	18	5.1%
<i>AluYb8</i>			
	1 - 5 bp	175	81.4%
	8 - 45 bp	29	13.4%
	57 - 166 bp	11	5.0%

Figure 1. Strategy for detection and assembly of non-reference *Alu* insertions. Approach for reconstruction of non-reference *Alu* insertions from WGS data. 1). WGS in aligned BAM format from 53 samples were merged to a single BAM file, and clusters of *Alu*-supporting read pairs identified using the RetroSeq program by Keane *et al.* 2). *Alu*-supporting read pairs and intersecting split reads were extracted for each candidate site, and 3). Subjected to a *de novo* assembly using the CAP3 overlap-layout assembler 4). *Alu*-containing contigs were then mapped to the reference genome to verify chromosomal coordinates and uniqueness of the call. 5). Breakpoints and putative TSDs from each contig were computationally predicted by 3-way alignment to determine overlap of the assembled upstream and downstream flanks with the pre-insertion site from the hg19 reference.

Figure 2. Sequence analysis of assembled non-reference *Alu* insertions. Breakpoints were determined based on alignment the 5' and 3' edge of each insertion sequence with the corresponding sequence from the hg19 reference. *Miropeats* annotation, aligned breakpoints, and Sanger sequences are shown for three representative insertions. A. Insertion at chr7:46102164 of 297 bp *Alu* with breakpoint overlap of 14 bp. *Upper left*: Alignment and breakpoint prediction of the assembled contig to hg19. Aligned breaks are shown in blue or red (leftmost or rightmost aligned nucleotides, respectively); the bracket indicates the *Alu* location in the contig relative to hg19. Repetitive elements in the reference and assembled contig are colored as: LINEs, green; SINEs, purple; LTRs, orange; DNA elements, pink. *Lower left*: 3-way alignment of *Alu*-flanking assembled stretches to hg19. A '1' or '2' indicates nucleotides aligned between the assembled contig and hg19 reference upstream or downstream of the *Alu* junction. A '*' indicates positions with the same base. Terminal nucleotides of the left and right breaks are colored as above; the black bar shows contig overlap. *Right*: Alignment of the assembled contig with Sanger sequence data to the hg19 empty allele and subfamily consensus for that insertion. Blue and red bars indicate left and right breaks; shading shows assembled and validated base changes from the subfamily consensus. B. Insertion at chr8:120800779 of 221 bp with a 14 bp TSD. C. Insertion at chr11:26601646 of 310 bp with a target site deletion of 3 bp ('Δ' in the alignment). All breakpoint and alignment indications are as described in panel A.

Figure 3. Length distribution of assembled *AluYa5* and *AluYb8* insertions. A. Scaled representation of the *AluYa5* and *AluYb8* consensus and element properties. The *Alu* is comprised of two arms (*left*, blue; *right*, grey) joined by an A-rich region and having a 3' poly-A tail. The A and B boxes indicate promoter regions. A 31bp insertion distinguishes the arms; the *AluYb8* has an extreme 3' 7 bp insertion relative to Ya5; the sequences are otherwise structurally conserved. Bases involved SRP/14 binding sites are shown by the gold bar. B. The size distribution of 351 *AluYa5* and 215 *AluYb8* assembled insertions relative to the respective subfamily consensus. The number of assembled insertions containing an aligned nucleotide is shown against the corresponding position in the consensus.

Figure 4. Genotyping of a subset of non-reference *Alu* insertions. Genotype validation was performed for 11 sites across 10 individuals. A. Strategy for primer design and allele detection. A single primer set was used for genotyping each locus, designed to target within 250bp of the assembled insertion

coordinates relative to hg19. B. Genotyping from PCR screens and band scoring. Banding patterns supporting the unoccupied or *Alu*-containing allele were assessed following locus-specific PCR; predicted band sizes were estimated by *in silico* PCR analysis and mapped *Alu* coordinates per site. The chromosomal location of each *Alu* is indicated at left. A '+' or '-' shows the relative position of each allele. Sample information is provided for population (above) and for each individual (below). C. Principle Component Analysis PCA was performed on genotype matrix for 1,010 autosomal sites genotypes across 53 populations. A projection of the samples onto the first two Principal Components is shown.

Figure S1. Assembled *Alu* insertions. Summary of assembled insertions compared to the hg19 reference for all sites, depiction as in Figure 2. Information including the siteID, insertion coordinate, insertion size, predicted TSD, and genotyping are provided for each of 1,1614 insertions.

Figure S2. Trace alignments of validated insertions. Nucleotide alignments of Sanger sequencing results with the corresponding CAP3 assemblies for each validated site. Trace alignments are shown both upstream and downstream the insertion for all validated sites. Labelling is as in Figure 2.

Figure S3. Trace alignments of truncation validations. Alignments of Sanger sequencing traces with the CAP3 assembled sequence for each site as in Figure S2. Alignments are shown relative to the appropriate consensus element. A red line indicates the corresponding validated truncated regions; a black line indicates Flanking TSDs.

Figure S4. Comparison of recovered alleles at chr11:35425392. The sequence from the hg19 reference genome and the CAP3 assembled contig are compared for an insertion on chromosome 11. Extent of matching sequence between haplotypes is in yellow. The location of an indel coincident with the insertion site is depicted in blue; putative TSDs are underlined; the *Alu* insertion is in purple. Sequence present in the insertion haplotype but absent from the reference is in red. The green text corresponds to the extent of similar sequence (89.5% similarity) flanking the *Alu* insertion.

Figure S5. *AluJ* consensus comparison. Assembled *AluJ* insertions are shown in alignment with the corresponding *AluJb* subfamily and major *AluY* and *AluS* consensus elements. *AluJ* elements are boxed. Black asterisks are used to indicate major *AluJ* characteristic bases as originally reported by Jurka and Smith (7); grey asterisks indicate additional bases unique to the assembled and consensus *AluJ* elements. Trace validation for the insertion located at chr12:73056650 is in Figure S2.

Figure 1

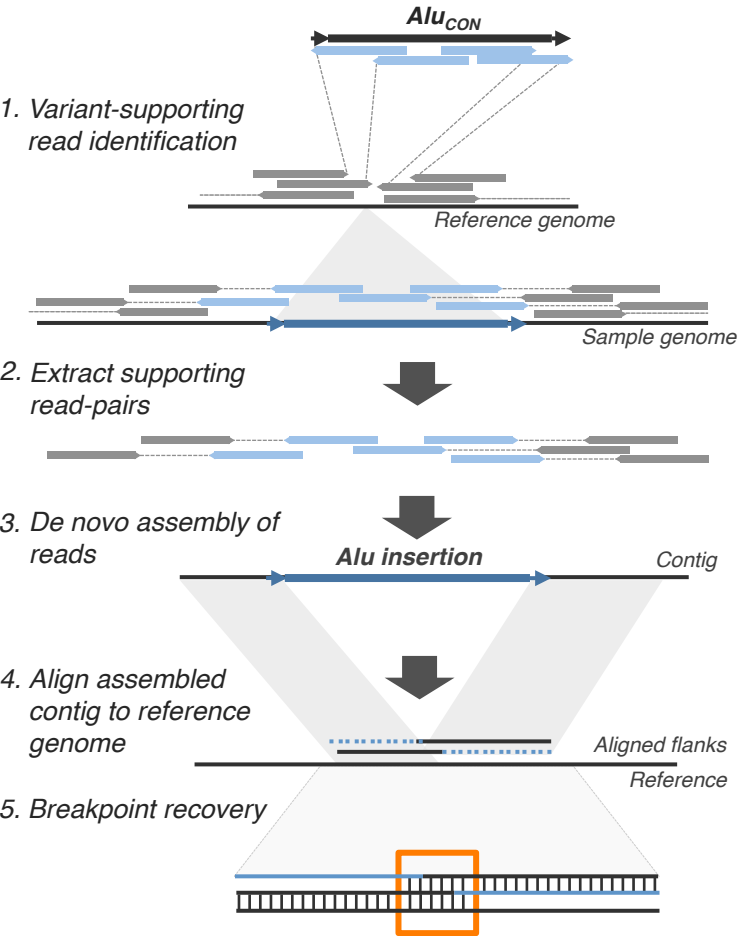
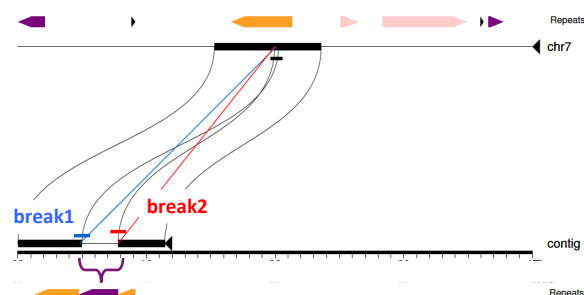


Figure 2

A.

Insertion site: chr7:46102164



siteID: chr7_46102169 Scafftig3 +
chr7:46102135-46102223
TSD size: 14 INS size: 297

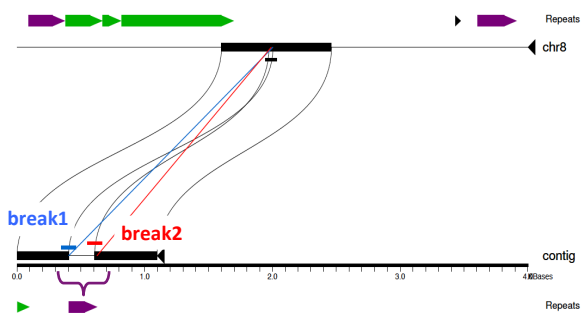
[illegible]

overlap

[illegible]

B.

Insertion site: chr8:120800779



chr8:120800733-120800820
TSD size: 14 INS size: 221

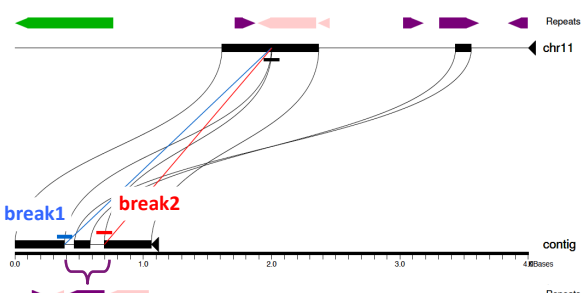
[illegible]

overlap

[illegible]

C.

Insertion site: chr11:26601646



chr11:26601587-26601709
TSD size: -3 INS size: 310

[illegible]

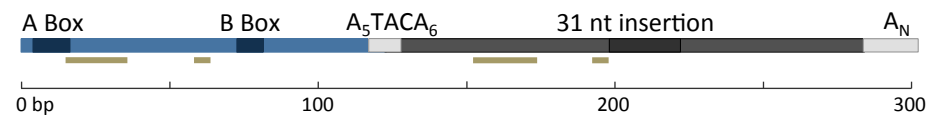
overlap

[illegible]

tg

Figure 3

A.



B.

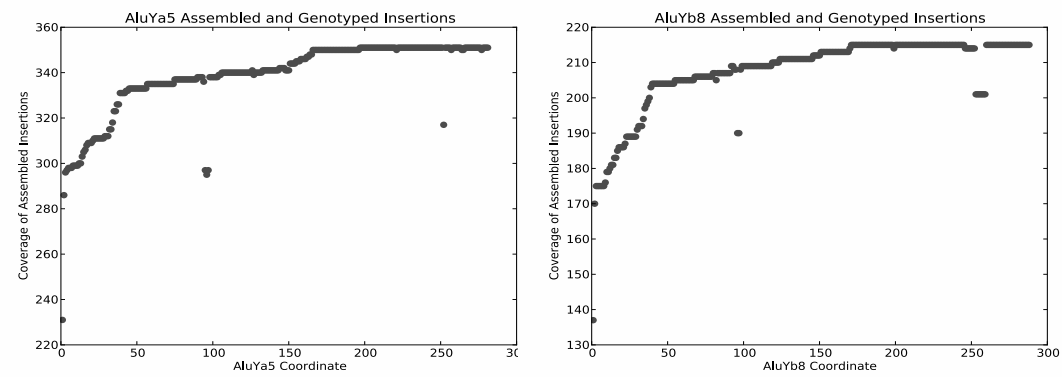


Figure 4

