1  # AluMine: alignment-free method for the discovery of

2  # polymorphic Alu element insertions

3  Tarmo Puurand, Viktoria Kukuškina, Fanny-Dhelia Pajuste and Maido Remm*

4  *Institute of Molecular and Cell Biology, University of Tartu, Tartu, Estonia*

5  * corresponding author

6

7  Tarmo Puurand (tarmo.puurand@ut.ee)

8  Viktoria Kukuškina (viktoria.tsernokozova@gmail.com)

9  Fanny-Dhelia Pajuste (fanny-dhelia.pajuste@ut.ee)

10  Maido Remm (maido.remm@ut.ee)

11    **ABSTRACT**

12    **Background**

13    Recently, alignment-free sequence analysis methods have gained popularity in the field of personal

14    genomics. These methods are based on counting frequencies of short $k$-mer sequences, thus

15    allowing faster and more robust analysis compared to traditional alignment-based methods.

16

17    **Results**

18    We have created a fast alignment-free method, AluMine, to analyze polymorphic insertions of Alu

19    elements in the human genome. We tested the method on 2,241 individuals from the Estonian

20    Genome Project and identified 28,962 potential polymorphic Alu element insertions. Each tested

21    individual had on average 1,574 Alu element insertions that were different from those in the

22    reference genome. In addition, we propose an alignment-free genotyping method that uses the

23    frequency of insertion/deletion-specific 32-mer pairs to call the genotype directly from raw

24    sequencing reads. Using this method, the concordance between the predicted and experimentally

25    observed genotypes was 98.7%. The running time of the discovery pipeline is approximately 2 hours

26    per individual. The genotyping of potential polymorphic insertions takes between 0.4 and 4 hours per

27    individual, depending on the hardware configuration.

28

29    **Conclusions**

30    AluMine provides tools that allow discovery of novel Alu element insertions and/or genotyping of

31    known Alu element insertions from personal genomes within few hours.

32

33    **KEYWORDS**

34    Alu repeat element, mobile element insertions, alignment-free sequence analysis

35 **INTRODUCTION**

36 Approximately 45% of the human genome contains repeated sequences. These repeated sequences

37 can be divided into tandem repeats and interspersed repeat elements (segmental duplications and

38 transposable elements). The most abundant transposable element in the human genome is the Alu

39 element. A typical Alu element is an approximately 300 bp long transposable nucleotide sequence.

40 [1–3]. The estimated number of full-length or partial Alu elements in the human genome is 1.1

41 million [4–7].

42

43 The presence or absence of some Alu elements is variable between individual genomes. Members of

44 the Alu AluY and AluS subfamilies actively retrotranspose themselves into new locations, thus

45 generating polymorphic Alu insertions [8–10]. A polymorphic Alu in this context refers to the

46 presence or absence of the entire element and not single nucleotide polymorphisms within the Alu

47 sequence. The insertion rate of Alu elements into new locations is approximately one insertion per

48 20 births [11,12]. Most of the variation in Alu elements is caused by the insertion of new elements.

49 Deletion of the entire Alu element is possible but occurs much less frequently than the insertion of

50 new elements. Polymorphic Alu insertions disturb the regulation of flanking genes and affect

51 phenotype. They cause changes in the genome that lead to disease [13–15]. Therefore,

52 computational methods that reliably detect polymorphic Alu element insertions from sequencing

53 data are needed.

54

55 Several methods for the identification of polymorphic Alu insertions have been developed that

56 include the following: VariationHunter [16,17], Hydra [18], TEA [19], RetroSeq [20], alu-detect [21]

57 and Tangram [22], MELT [23], T-lex2 [24], and STEAK [25]. All these methods are based on the

58 mapping of sequencing reads and the subsequent interpretation of mapping results. The discovery of

59 new insertions is typically based on split locations of a single read and/or the distance between

60 paired reads.

61

62     Several databases or datasets that describe polymorphic Alu insertions are available. The oldest

63     resource containing known polymorphic transposable elements is the dbRIP database [26]. It

64     contains insertions detected by comparison of Human Genome Project data with Celera genome

65     data. dbRIP also contains information about somatic Alu insertions that might be related to different

66     diseases. The most comprehensive Alu element dataset is available from the 1000 Genome Project

67     (1000G) [12,27]. A subset of these sequences has been validated by Sanger sequencing [9]. The

68     1000G dataset is currently the reference set for evaluating the accuracy of structural variant calls

69     generated by other methods. The dbRIP, 1000G, me-scan [28], TEA [19] and HGDP [29] datasets

70     together contain more than 10,000 polymorphic Alu insertions that were collected from hundreds of

71     individuals from different populations.

72

73     We have developed a set of novel, alignment-free methods for the rapid discovery of polymorphic

74     Alu insertions from fully sequenced individual genomes. In addition, we provide a method that calls

75     genotypes with previously known insertions directly from raw reads. Evaluation of these methods

76     was performed by computational simulations and PCR product size analysis.

77    **RESULTS**

78

79    **Rationale for the alignment-free discovery of Alu insertion sites**

80    We describe a novel method allowing both the discovery of new polymorphic Alu insertions and the

81    detection of known insertions directly from raw reads in next generation sequencing (NGS) data. Two

82    key steps within the discovery method are the a) identification of potential polymorphic Alu

83    insertions present in tested personal genomes but not in the reference genome (REF– discovery) and

84    the b) identification of potential polymorphic Alu elements present in the current reference genome

85    (REF+ discovery) that might be missing in the tested genomes.

86    All discovery pipelines use a 10 bp consensus sequence from the 5' end of the Alu (GGCCGGGCGC)

87    with one mismatch that we call Alu signatures. The REF– discovery pipeline identifies all occurrences

88    of Alu signatures in raw sequencing reads from an individual. A 25 bp flanking sequence from the 5'

89    region is recorded together with the discovered Alu signature sequence (Figure S1 in Additional file

90    1). Subsequently, the location of these 25 bp sequences in the reference genome is determined using

91    the custom-made software `gtester` (Kaplinski, unpublished). A new REF– element is reported if the

92    10 bp sequence in the raw reads is different from the 10 bp sequence in the reference genome.

93    The REF+ discovery pipeline uses the same Alu element signature to identify all locations in the

94    reference genome where the preceding 5 bp target site duplication motif (TSD) is present 270-350 bp

95    downstream from the signature sequence (see Figure S2 in Additional file 1 for details). Both

96    discovery pipelines generate a pair of 32-mers for each identified Alu element. These 32-mer pairs

97    are used for the subsequent genotyping of the Alu elements in other individuals. Two 32-mers in a

98    pair correspond to two possible alleles with or without the Alu element insertion. All candidate 32-

99    mer pairs are further filtered based on their genotypes in test individuals. The entire discovery

100    process is outlined in Figure 1.

101    The alignment-free genotyping of known Alu elements is based on counting the frequencies of 32-

102    mer pairs specific to Alu element breakpoints using the previously published FastGT software

103    package [30]. The principles of the generation of $k$-mer pairs specific to Alu insertion breakpoints are

104    shown in Figure 2. To detect polymorphic insertions, we use 25 bp from the reference genome

105    immediate to the 5' end of the potential Alu insertion point and then add either 7 bp from the Alu

106    element or 7 bp from the genomic sequence downstream of the second TSD motif (Figure 2A). The

107    names of two alleles are assigned based on their status in the reference genome; the allele that is

108    present in the reference genome is always called allele A, and the alternative allele is always called

109    allele B (Figure 2B). This allows us to use the same naming convention for alleles and genotypes used

110    by the FastGT package for single nucleotide variants.

111

112    **Compilation of the list of potential polymorphic Alu elements**

113    To test the applicability of the AluMine method to real data, we performed REF− element discovery

114    using 2,241 high-coverage genomes from the Estonian Genome Project and compiled a set of 32-mer

115    pairs for subsequent genotyping. REF− candidates consist of Alu elements that are present in the raw

116    reads from sequenced individuals but not in the reference genome. We searched the raw reads from

117    test individuals following the principles described above and detected 13,128 REF− Alu elements

118    overall.

119

120    REF+ discovery was performed using the human reference genome version 37. We searched for

121    potential REF+ candidates by using the following criteria: the element must have an intact Alu

122    signature sequence, have a TSD at least 5 bp long on both ends of the Alu element, have more than

123    100 bits similar to known Alu elements, and must not be present in the chimpanzee genome. Our

124    REF+ script detected 267,377 elements with an Alu signature sequence from the human reference

125    genome. However, only 15,834 (5.9%) of these passed all the abovementioned filtering criteria and

126    remained in the set of potential polymorphic elements. The proportion of different signature

127    sequences among the set of REF+ elements is shown in Table S1 in in Additional file 2. All the steps

128    involved in Alu element discovery are summarized in Table 1 together with the number of elements

129    that passed each step.

130

**Simulation tests of the discovery method**

132    We realize that although our discovery methods detected more than 13,000 REF– Alu element

133    insertions, some polymorphic Alu elements remain undiscovered in given individuals. There are two

134    obvious reasons why Alu variants are missed in the REF– discovery step: a) a low depth of coverage in

135    some individuals and b) difficulties with the unique localization of 25-mers in some genomic regions.

136

137    The effect of coverage on the discovery rate can be estimated from simulated data. We generated

138    data with 5× to 55× nucleotide-level coverage and analyzed how many REF– elements we would

139    discover from these with our method. The results are shown in Figure 3A. There is an association

140    between the depth of coverage and the discovery rate, which levels out at an approximately 40×

141    depth of coverage.

142

143    Another factor affecting the sensitivity of Alu element discovery is that the repeated structure of the

144    genome sequence prevents the unique localization of discovered Alu elements. The REF– discovery

145    method relies on the unique localization of the 25-mer in front of the Alu signature sequence. We

146    decided to perform a series of simulations with artificial Alu element insertions to determine what

147    fraction of them was discoverable by our REF– discovery method. For this, we inserted 1,000 typical

148    Alu elements into random locations of a diploid genome sequence and generated random

149    sequencing reads from this simulated genome using wgsim software [31]. The simulation was

150    repeated with 10 male and 10 female genomes using different mutation rates. Varying the mutation

151    rate helps to somewhat simulate older and younger Alu element insertions (older Alu elements have

152     accumulated more mutations) and estimate how their detection rate varies accordingly. We

153     observed that 20% to 23% of the elements remain undetected, depending on the mutation rate

154     (Figure 3B). The mutation rate has only a moderate effect on the sensitivity of detection; thus, we

155     assume that the age of the Alu element insertion does not significantly influence the number of

156     detected elements. Additionally, 7% of the inserted elements remained undiscovered because they

157     were inserted into inaccessible (N-rich) regions of the reference genome, and this number is

158     independent of mutation rate.

159

160     **Comparison with other Alu discovery methods**

161     When comparing the results of Alu discovery methods, we can compare two aspects. If the same

162     individuals are studied by many methods, we can estimate the overlap between identified elements.

163     Otherwise, we can compare the overall number of detected elements.

164

165     We were able to identify the overlap between Alu elements discovered from sample NA12878 within

166     the 1000G pilot project and the 1000G Phase3 project. AluMine discovered 60% (1204) of all

167     elements reported in the 1000G Pilot phase project plus an additional 443 elements (Figure 4). The

168     overlaps between methods are similar for REF+ and REF− elements.

169

170     To examine other methods, we were only able to compare the overall number of discovered REF−

171     elements. AluMine detected 1,116 and 1,127 REF− insertions in the CEPH individuals NA12877 and

172     NA12878 and 1,290 insertions in NA18506. alu-detect discovered on average 1,339 Alu insertions per

173     CEU individual [21]. Hormozdiari et al. detected 1,282 events in the CEU individual NA10851 with 22×

174     coverage and 1,720 events in the YRI individual NA18506 with 40× coverage [16]. TEA detected an

175     average of 791 Alu insertions in each individual genome derived from cancer samples [19]. In

176     genomes from Chinese individuals, Yu et al. discovered 1,111 Alu element insertions on average [32].

177     Thus, the overall number of detected REF– elements was similar for all methods.

178

179     The number of polymorphic REF+ elements (present in the reference genome) has been studied less

180     thoroughly. The number of human-specific REF+ insertions is at least 8,817 [33]. We identified 15,834

181     potential polymorphic REF+ elements, of which 1,762 were polymorphic in at least one individual in

182     the studied population.

183

184     **Frequency of non-reference Alu elements in tested individuals**

185     We scanned 2,241 Estonian individuals with the final filtered set of Alu elements to identify the

186     genotypes of all potential polymorphic Alu insertions in their genomes. All tested individuals had

187     some Alu elements that were different from those in the reference genome. The tested individuals

188     had 741 - 1,323 REF– elements (median 1,045) that were not present in the reference genome and

189     465 - 651 REF+ Alu elements (median 588) that were present in the reference genome but missing in

190     given individual (Figure 5).

191

192     One interesting question that can be addressed from the given data is the cumulative number of

193     REF– elements in a population. We discovered 14,455 REF– Alu elements from 2,241 tested

194     individuals. However, many of these were common within the population. Thus, saturation of the

195     total number of polymorphic elements is expected if sufficient number of individuals are sequenced.

196     The saturation rate of the REF– elements is shown in Figure 6. Obviously, the number of REF–

197     elements was still far from saturation. Each new individual genome sequence still contained 2-3

198     previously unseen REF– elements.

199

200     **Selection of 32-mers for genotyping**

201     In principle, we would like to call the genotypes with discovered Alu elements in other individuals

202    using pairs of specific 32-mers and FastGT genotyping software. Unfortunately, not all discovered Alu

203    elements are suitable for fast genotyping with a pair of short k-mers. Some of them tend to give

204    excessive counts from other regions of the genome, and some might be affected by common Single

205    Nucleotide Variants (SNVs). To select a set of Alu elements that gives reliable genotypes, we filtered

206    the Alu elements based on their genotyping results using data from the same 2,241 individuals that

207    were used for REF– element discovery. To this end, we merged 32-mers of REF– and REF+ Alu

208    elements with a set of SNV-specific 32-mers and determined the genotypes of these markers in test

209    individuals using the FastGT package. SNV-specific *k*-mers are required at this step because Alu

210    elements alone cannot provide reliable estimates of parameter values for the empirical Bayes

211    classifier used in FastGT. Additional filtering and removal of candidate elements was based on several

212    criteria. We removed elements that generated an excessive number of unexpected genotypes (a

213    diploid genotype is expected for autosomes, and a haploid genotype is expected for chrY), elements

214    that deviated from Hardy-Weinberg equilibrium and monomorphic REF– elements. The validation of

215    all tested markers together with their genotype counts is shown in Table S2 in Additional file 2. In the

216    final validated *k*-mer database, we included 9,712 polymorphic REF– elements that passed the

217    validation filters, including 1,762 polymorphic REF+ elements and 11,634 monomorphic REF+

218    elements. Although 87% of the candidate REF+ elements were monomorphic in the tested

219    individuals, the possibility exists that they are polymorphic in other populations; therefore, we did

220    not remove them from the *k*-mer database.

221

222    **Experimental validation of the genotyping method**

223    We decided to validate the alignment-free genotyping of polymorphic Alu elements with a subset of

224    newly discovered Alu elements. The validation was performed experimentally using PCR fragment

225    length polymorphism. We used four different Alu elements (1 REF– and 3 REF+ elements) and

226    determined their genotypes in 61 individuals. The individuals used in this validation did not belong to

227    the training set of 2,241 individuals and were sequenced independently. The electrophoretic gel

228    showing the PCR products of one REF– polymorphism is shown in Figure 7. The results for the three

229    REF+ individuals are shown in Figure 8. The computationally predicted genotypes and experimentally

230    determined genotypes conflicted in only 3 cases; thus, the concordance rate was 98.7%. The 32-mer

231    counts, predicted genotypes and experimental genotypes for each individual are shown in Table S3 in

232    Additional file 2.

233

234    **Performance**

235    The performance of the AluMine methods can be divided into three parts: the performance of the

236    REF– discovery pipeline, the performance of the REF+ discovery pipeline and the genotyping

237    performance. The REF+ pipeline was run on a server with a 2.27 GHz Intel Xeon CPU X7560 and 512

238    GB RAM. The REF– scripts and genotyping were run on cluster nodes with a 2.20 GHz Intel Xeon CPU

239    E5-2660 and 64 GB RAM.

240

241    The most time-consuming steps in the REF– discovery pipeline are a) searching for Alu signatures

242    from FASTQ files, which takes 2 hours per individual on a single CPU core, and b) finding their

243    locations in the reference genome using gtester software (2 hours for the first individual, 4

244    minutes for each subsequent individual). The increase in speed for subsequent individuals is due to

245    the large size of the gtester indices (approximately 60 GB). For the first individual, they are read

246    from a hard drive, and for subsequent individuals, the disk cache is used. None of the steps require

247    more than 8 GB of RAM.

248

249    The REF+ discovery pipeline contains the following three time-consuming steps: a) a search for 31

250    different Alu signatures from chromosomes of the reference genome (takes 14 minutes), b) a

251    homology search with all the candidates to confirm that they are Alu elements (2 minutes) and c) a

252    comparison with the chimpanzee genome to exclude fixed Alu elements (4 minutes, 28 GB RAM). All

253    these steps use a single processor. The REF+ discovery pipeline has to be run only once and should

254    not be repeated for each separate individual. Thus, in terms of performance, it occupies only a minor

255    part of the overall analysis.

256

257    The genotyping of individuals is performed with the previously published FastGT package [30]. The

258    performance of FastGT was analyzed in the original paper. In optimized conditions (>200 GB RAM

259    available, using FASTQ instead of BAM format, and using solid state drive), it can process one high

260    coverage individual within 30 minutes. However, we used FastGT on cluster nodes with a limited

261    amount of hard drive space and limited RAM. Therefore, in our settings, FastGT acquired sequence

262    data from BAM files through standard input, which limited its performance. In this way, we were able

263    to process one individual in 3-4 CPU hours.

264

265

266 **DISCUSSION**

267

268 **Parameter choice**

269 A common matter of discussion for alignment-free sequence analysis methods is the optimal length

270 of $k$-mers. In our case, the $k$-mers used for genotyping Alu elements had to be bipartite and contain

271 sufficient sequence from the genome and a couple of nucleotides from the Alu element (Figure 2).

272 The first part of the bipartite $k$-mer must guarantee the unique localization of the $k$-mer in the

273 human genome; the second part must allow distinguishing variants with and without the Alu element

274 at a given location. Both parts must fit into 32 nucleotides because we use the $k$-mer managing

275 software package GenomeTester4, which is able to handle $k$-mers with a maximum length of 32

276 nucleotides. In the current work, we chose to divide 32-mers into 25 + 7 nucleotides. Our previous

277 work demonstrated that all $k$-mers 22 to 32 nucleotides long should perform equally well to analyze

278 variations in the human genome (Figure 5 in [30]). Thus, we assume that we would obtain a rather

279 similar genotyping result with slightly different splits, such as 22 + 10 or 28 + 4 nucleotides. Using

280 fewer than 4 nucleotides from the Alu element would give too high of a chance to have an identical

281 sequence in the reference genome, and the program would not be able to distinguish variants with

282 and without Alu.

283

284 **Comparison with other software**

285 We compared the number of REF− elements discovered by different methods. However, the direct

286 comparison of these numbers to our data is complicated because different populations and

287 individuals were used in different reports. The number of discovered insertions was correlated with

288 the individual ancestry of the subjects: generally, fewer Alu insertions were discovered in CEU

289 individuals than in YRI individuals [12]. Additionally, the depth of coverage had a strong effect on the

290 results, as shown in Figure 3A. All methods, including AluMine, detected approximately 1000 REF-

291 elements per genome. The slight differences were likely due to differences in the depth of coverage

292    and the different origins of the samples used.

293

294    Different detection methods have different biases. The premature termination of target primed

295    reverse transcription during the replication of Alu elements can generate truncated Alu element

296    insertions that are missing the 5' end of the element. It has been estimated that 16.4% of Alu

297    elements are truncated insertions [29]. Furthermore, some Alu element polymorphisms appear

298    through the deletion of existing elements (2%) [9] or mechanisms that do not involve

299    retrotransposition (less than 1%) [29]. Our REF+ method relies on the presence of TSDs, and the REF–

300    method relies on the presence of intact 5' ends in the Alu. Thus, we would not be able to detect

301    those events, which would explain the majority of the differences between our results and the

302    elements detected in the 1000G pilot phase (Figure 4).

303

304    **Future directions**

305    In principle, our discovery method can be used to search for novel Alu elements in any whole-

306    genome sequencing data. Transposable elements are known to occur in genes that are commonly

307    mutated in cancer and to disrupt the expression of target genes [13,19]. Our method allows the

308    discovery of novel Alu elements from sequences from tumors and matched normal blood samples,

309    allowing the study of the somatic insertion of Alu elements in cancer cells and their role in

310    tumorigenesis. The precompiled set of 32-mer pairs allows the genotyping of known Alu element

311    insertions in high-coverage sequencing data. This facilitates the use of Alu elements in genome-wide

312    association studies along with SNVs.

313

314    The alignment-free discovery method could also be adapted for the detection of other transposable

315    elements, such as L1 or SVA elements. However, the discovery of these elements is more

316    complicated because SVA elements contain a variable number of $(CCCTCT)_n$ repeats in their 5' end,

317    and L1 elements contain variable number of Gs in front of the GAGGAGCCAA signature sequence.

318

319 **CONCLUSIONS**

320 We have created a fast, alignment-free method, AluMine, to analyze polymorphic insertions of Alu

321 elements in the human genome. It consists of two pipelines for the discovery of novel polymorphic

322 insertions directly from raw sequencing reads. One discovery pipeline searches for Alu elements that

323 are present in a given individual but missing from the reference genome (REF– elements), and the

324 other searches for potential polymorphic Alu elements present in the reference genome but missing

325 in some individuals (REF+ elements). We applied the REF– discovery method to 2,241 individuals

326 from the Estonian population and identified 13,128 polymorphic REF– elements overall. We also

327 analyzed the reference genome and identified 15,834 potential polymorphic REF+ elements. Each

328 tested individual had on average 1,574 Alu element insertions (1,045 REF– and 588 REF+ elements)

329 that were different from those in the reference genome.

330

331 In addition, we propose an alignment-free genotyping method that uses the frequency of

332 insertion/deletion-specific 32-mer pairs to call the genotype directly from raw sequencing reads. We

333 tested the accuracy of the genotyping method experimentally using a PCR fragment length

334 polymorphism assay. The concordance between the predicted and experimentally observed

335 genotypes was 98.7%.

336

337 The running time of the REF– discovery pipeline is approximately 2 hours per individual, and the

338 running time of the REF+ discovery pipeline is 20 minutes. The genotyping of potential polymorphic

339 insertions takes between 0.4 and 4 hours per individual, depending on the hardware configuration.

340

341

342    **METHODS AND DATA**

343    **Data**

344    The reference genome GRCh37.p13 was used for all analyses.

345

346    **Discovery of REF– and REF+ elements**

347    The exact details of all discovery pipelines are described in the corresponding scripts

348    (pipeline_ref_plus.sh, pipeline_ref_minus.sh and pipeline_merging_and_filtering.sh) available from

349    GitHub (https://github.com/bioinfo-ut/AluMine).

350

351    **PCR protocol**

352    To prepare a 20 µl PCR master mix, we mixed 0.2 µl FIREPol DNA polymerase (Solis BioDyne, Estonia),

353    0.6 µl of 10 mM DNTP, 0.8 µl of a 20 mM primer mix, 2 µl of 25 mM MgCl2, 2 µl polymerase buffer,

354    and 14.4 µl Milli-Q water. For PCR, Applied Biosystems thermocyclers were used. The PCR was run

355    for 30 cycles using a 1 minute denaturation step at 95°C, a 1 minute annealing step at 55°C and a 1.5

356    minutes elongation step at 72°C. For gel electrophoresis, a 1.5% agarose gel (0.5 mM TBE + agarose

357    tablets + EtBr) was used. The PCR primer pairs used for the amplification of potential polymorphic

358    regions are shown in Table S4 in Additional file 2.

359

360    **Simulated Alu insertions**

361    To simulate polymorphic Alu insertions, we inserted 1000 heterozygous Alu elements into random

362    locations of the diploid reference genome together with a 15 bp target site duplication sequence and

363    a random length polyA sequence (5-80 bp). A male genome (5.98 Gbp) and a female genome (6.07

364    Gbp) were generated by merging two copies of autosomal chromosomes and the appropriate

365    number of sex chromosomes into a single FASTA file. Simulated sequencing reads were generated

366    using wgSim (version 0.3.1-r13) software from the SAMtools package [31]. The following parameters

367    were used: haplotype_mode = 1, base_error_rate = 0.005, outer_distance_between_the_two_ends

368    = 500, length_of_ reads = 151, cutoff_for_ambiguous_nucleotides=1.0, and number_of_reads =

369    306,000,000.

370 **ABBREVIATIONS**

371 1000G: 1000 Genome Project

372 NGS: Next Generation Sequencing

373 REF– Alu element: polymorphic Alu element present in at least one personal genome but not in the

374 reference genome

375 REF+ Alu element: polymorphic Alu element present in the reference genome, but missing in at least

376 one personal genome

377 TSD: Target Site Duplication motif

378 SNV: Single Nucleotide Variant

379

380 **DECLARATIONS**

381

382 **Acknowledgements**

383 The authors thank Lauris Kaplinski for advice on improving performance of Alu element discovery

384 algorithms and for adapting the *k*-mer counting software FastGT and gtester for this project.

385

386 **Funding**

387 This work was funded by institutional grant IUT34-11 from the Estonian Research Council and the EU

388 ERDF grant No. 2014-2020.4.01.15-0012 (Estonian Center of Excellence in Genomics and

389 Translational Medicine). The cost of the sequencing of individuals from the Estonian Genome Center

390 was partly covered by the Broad Institute (MA, USA) and the PerMed I project from the TERVE

391 program. Computation was partly carried out in the High Performance Computing Center of the

392 University of Tartu.

393

394 **Authors' contributions**

395 TP conceived the idea and performed most of the large-scale genomic analyses. MR wrote the scripts

396    for post-processing of the data, performed simulations and wrote the manuscript. VK performed all

397    the PCR experiments and helped to develop genome analysis methods. FDP provided help with data

398    management and visualization. All authors read and approved the final manuscript.

399

400    **Ethics approval and consent to participate**

401    The genome data were collected and used with ethical approval (Nr. 206T4, obtained for the project

402    SP1GVARENG).

403

404    **Consent for publication**

405    Not applicable.

406

407    **Availability of data and materials**

408    All scripts (pipeline_ref_plus.sh, pipeline_ref_minus.sh and pipeline_merging_and_filtering.sh) and

409    software (gtester) created for this study are available from GitHub (https://github.com/bioinfo-

410    ut/AluMine). The FastGT package used for genotyping the Alu insertions is also available from GitHub

411    (https://github.com/bioinfo-ut/GenomeTester4/blob/master/README.FastGT.md). *K*-mer lists for

412    genotyping Alu elements using FastGT are available from University of Tartu webpage

413    (http://bioinfo.ut.ee/FastGT/). The whole genome sequencing data that support the findings of this

414    study are available on request from Estonian Genome Centre (https://www.geenivaramu.ee/en) but

415    restrictions apply to the availability of these data, and so are not publicly available.

416

417    **Competing interests**

418    The authors declare that they have no competing interests.

419

420    **Additional Files**

421    *Puurand_2019_AdditionalFile1.pdf*

422     Additional file 1. Figure S1 and Figure S2 explaining the REF- and REF+ discovery algorithms. (PDF,

423     139 kb)

424     *Puurand_2019_AdditionalFile2.xlsx*

425     Additional file 2. Supplementary tables Table S1, Table S2, Table S3 and Table S4. (XLSX, 2.1 Mb)

426    **REFERENCES**

427    1. Houck CM, Rinehart FP, Schmid CW. A ubiquitous family of repeated DNA sequences in the human
428    genome. J Mol Biol. 1979;132:289–306.

429    2. Rubin CM, Houck CM, Deininger PL, Friedmann T, Schmid CW. Partial nucleotide sequence of the
430    300-nucleotide interspersed repeated human DNA sequences. Nature. 1980;284:372–4.

431    3. Schmid CW, Jelinek WR. The Alu family of dispersed repetitive sequences. Science.
432    1982;216:1065–70.

433    4. Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, et al. Initial sequencing and
434    analysis of the human genome. Nature. 2001;409:860–921.

435    5. Bailey JA, Gu Z, Clark RA, Reinert K, Samonte R V, Schwartz S, et al. Recent segmental duplications
436    in the human genome. Science. 2002;297:1003–7.

437    6. Sebat J, Lakshmi B, Troge J, Alexander J, Young J, Lundin P, et al. Large-scale copy number
438    polymorphism in the human genome. Science. 2004;305:525–8.

439    7. Batzer MA, Deininger PL. Alu repeats and human genomic diversity. Nat Rev Genet. 2002;3:370–9.

440    8. Bennett EA, Keller H, Mills RE, Schmidt S, Moran J V., Weichenrieder O, et al. Active Alu
441    retrotransposons in the human genome. Genome Res. 2008;18:1875–83.

442    9. Konkel MK, Walker JA, Hotard AB, Ranck MC, Fontenot CC, Storer J, et al. Sequence Analysis and
443    Characterization of Active Human Alu Subfamilies Based on the 1000 Genomes Pilot Project. Genome
444    Biol Evol. 2015;7:2608–22.

445    10. Lee J, Kim Y-J, Mun S, Kim H-S, Han K. Identification of human-specific AluS elements through
446    comparative genomics. Gene. 2015;555:208–16.

447    11. Xing J, Zhang Y, Han K, Salem AH, Sen SK, Huff CD, et al. Mobile elements create structural
448    variation: analysis of a complete human genome. Genome Res. 2009;19:1516–26.

449    12. Stewart C, Kural D, Strömberg MP, Walker JA, Konkel MK, Stütz AM, et al. A comprehensive map
450    of mobile element insertion polymorphisms in humans. PLoS Genet. 2011;7:e1002236.

451    13. Solyom S, Kazazian HH. Mobile elements in the human genome: implications for disease. Genome
452    Med. 2012;4:12.

453    14. Kazazian HH, Moran J V. Mobile DNA in Health and Disease. N Engl J Med. 2017;377:361–70.

454    15. Payer LM, Steranka JP, Yang WR, Kryatova M, Medabalimi S, Ardeljan D, et al. Structural variants
455    caused by Alu insertions are associated with risks for many human diseases. Proc Natl Acad Sci U S A.
456    2017;114:E3984–92.

457    16. Hormozdiari F, Alkan C, Ventura M, Hajirasouliha I, Malig M, Hach F, et al. Alu repeat discovery
458    and characterization within human genomes. Genome Res. 2011;21:840–9.

459    17. Hormozdiari F, Hajirasouliha I, Dao P, Hach F, Yorukoglu D, Alkan C, et al. Next-generation
460    VariationHunter: combinatorial algorithms for transposon insertion discovery. Bioinformatics.
461    2010;26:i350-7.

462    18. Quinlan AR, Clark RA, Sokolova S, Leibowitz ML, Zhang Y, Hurles ME, et al. Genome-wide mapping
463    and assembly of structural variant breakpoints in the mouse genome. Genome Res. 2010;20:623–35.

464    19. Lee E, Iskow R, Yang L, Gokcumen O, Haseley P, Luquette LJ, et al. Landscape of somatic
465    retrotransposition in human cancers. Science. 2012;337:967–71.

466    20. Keane TM, Wong K, Adams DJ. RetroSeq: Transposable element discovery from next-generation
467    sequencing data. Bioinformatics. 2013;29:389–90.

468    21. David M, Mustafa H, Brudno M. Detecting Alu insertions from high-throughput sequencing data.
469    Nucleic Acids Res. 2013;41:e169.

470    22. Wu J, Lee W-P, Ward A, Walker JA, Konkel MK, Batzer MA, et al. Tangram: a comprehensive
471    toolbox for mobile element insertion detection. BMC Genomics. 2014;15:795.

472    23. Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, et al. An integrated
473    map of structural variation in 2,504 human genomes. Nature. 2015;526:75–81.

474    24. Fiston-Lavier A-S, Barrón MG, Petrov DA, González J. T-lex2: genotyping, frequency estimation
475    and re-annotation of transposable elements using single or pooled next-generation sequencing data.
476    Nucleic Acids Res. 2015;43:e22.

477    25. Santander CG, Gambron P, Marchi E, Karamitros T, Katzourakis A, Magiorkinis G. STEAK: A specific
478    tool for transposable elements and retrovirus detection in high-throughput sequencing data. Virus
479    Evol. 2017;3:vex023.

480    26. Wang J, Song L, Grover D, Azrak S, Batzer MA, Liang P. dbRIP: a highly integrated database of
481    retrotransposon insertion polymorphisms in humans. Hum Mutat. 2006;27:323–9.

482    27. Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, et al. An integrated
483    map of structural variation in 2,504 human genomes. Nature. 2015;526:75–81.

484    28. Witherspoon DJ, Zhang Y, Xing J, Watkins WS, Ha H, Batzer M a, et al. Mobile element scanning
485    (ME-Scan) identifies thousands of novel Alu insertions in diverse human populations. Genome Res.
486    2013;23:1170–81.

487    29. Wildschutte JH, Baron A, Diroff NM, Kidd JM. Discovery and characterization of Alu repeat
488    sequences via precise local read assembly. Nucleic Acids Res. 2015;43:10292–307.

489    30. Pajuste F-D, Kaplinski L, Möls M, Puurand T, Lepamets M, Remm M. FastGT: an alignment-free
490    method for calling common SNVs directly from raw sequencing reads. Sci Rep. 2017;7:2537.

491    31. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence Alignment/Map
492    format and SAMtools. Bioinformatics. 2009;25:2078–9.

493    32. Yu Q, Zhang W, Zhang X, Zeng Y, Wang Y, Wang Y, et al. Population-wide sampling of
494    retrotransposon insertion polymorphisms using deep sequencing and efficient detection.
495    Gigascience. 2017;6:1–11.

496    33. Tang W, Mun S, Joshi A, Han K, Liang P. Mobile elements contribute to the uniqueness of human
497    genome with 15,000 human-specific insertions and 14 Mbp sequence increase. DNA Res.
498    2018;25:521–33.

499    34. Hulsen T, de Vlieg J, Alkema W. BioVenn - a web application for the comparison and visualization
500    of biological lists using area-proportional Venn diagrams. BMC Genomics. 2008;9:488.
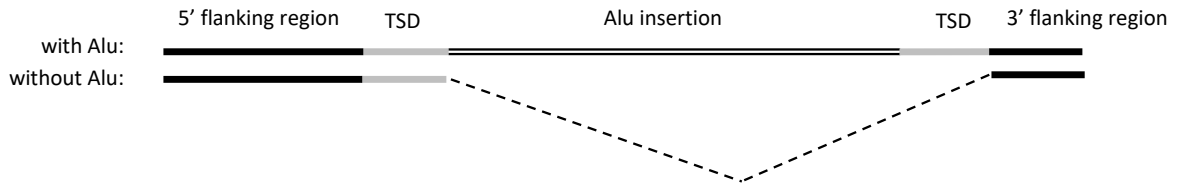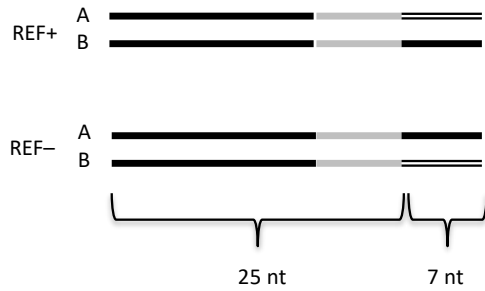
501

**FIGURES**



**Figure 1.** Overview of the discovery methods. Potential polymorphic Alu elements were identified from the raw reads of high-coverage WGS data (REF– Alu elements) and the reference genome (REF+ Alu elements). The candidate Alu elements were filtered using a subset of high-coverage individuals. A final set of 32-mers was used for the fast calling of polymorphic insertions from raw sequencing reads.

A)



B)



510

511 **Figure 2.** Principle of creating *k*-mer pairs for the calling (genotyping) of polymorphic Alu element
512 insertions. A) Genomic regions with or without an Alu element. B) A pair of 32-mers is created from
513 the insertion breakpoint region covering 25 nucleotides from the 5'-flanking region and 7 nucleotides
514 from either the Alu element or the 3'-flanking region. Allele A always represents the sequence from
515 the reference genome and allele B represents the alternative, non-reference allele.

516

517

518

519

520

521

522

523

**A**

**B**

524

**Figure 3.** (A) The number of discovered REF– Alu elements in individual NA12877 depending on the
depth of coverage. Various depth coverage levels were generated by randomly selecting a subset of
reads from the FASTQ file. (B) The frequency of false-negative and false-positive Alu elements found
in simulations. FN1 denotes false-negative findings that were undetectable because they are inserted
within unsequenced regions of the genome (N-rich regions). FN2 denotes false negatives that could
not be detected because they are inserted in nonunique regions of the genome. Error bars indicate
95% confidence intervals from 20 replicates.

**Figure 4.** Overlap between REF+ and REF– elements detected by AluMine, the 1000G pilot phase and 1000G Phase 3. The Venn diagram was created with BioVenn software [34].

537



538
539 **Figure 5.** Histogram showing the distribution of the number of non-reference REF– (light) and REF+
540 (dark) elements discovered per individual genome in 2,241 test individuals from the Estonian
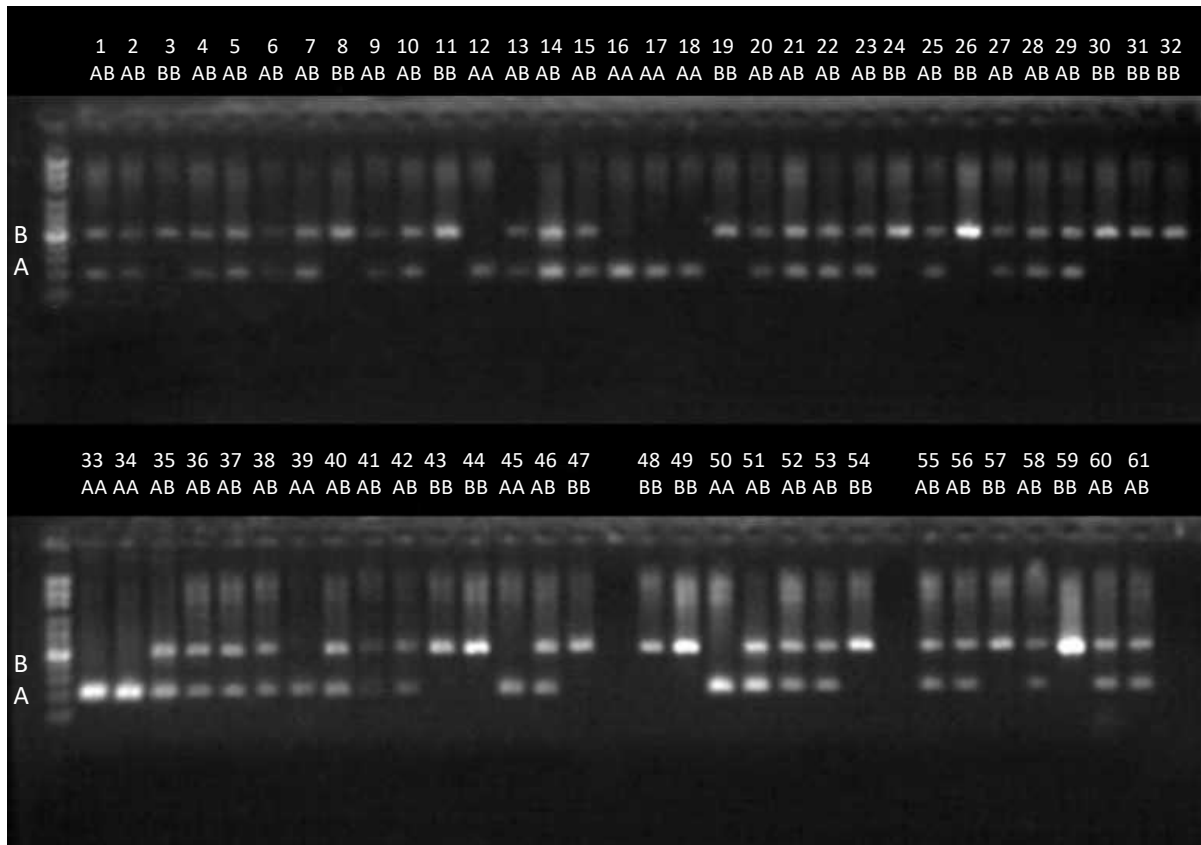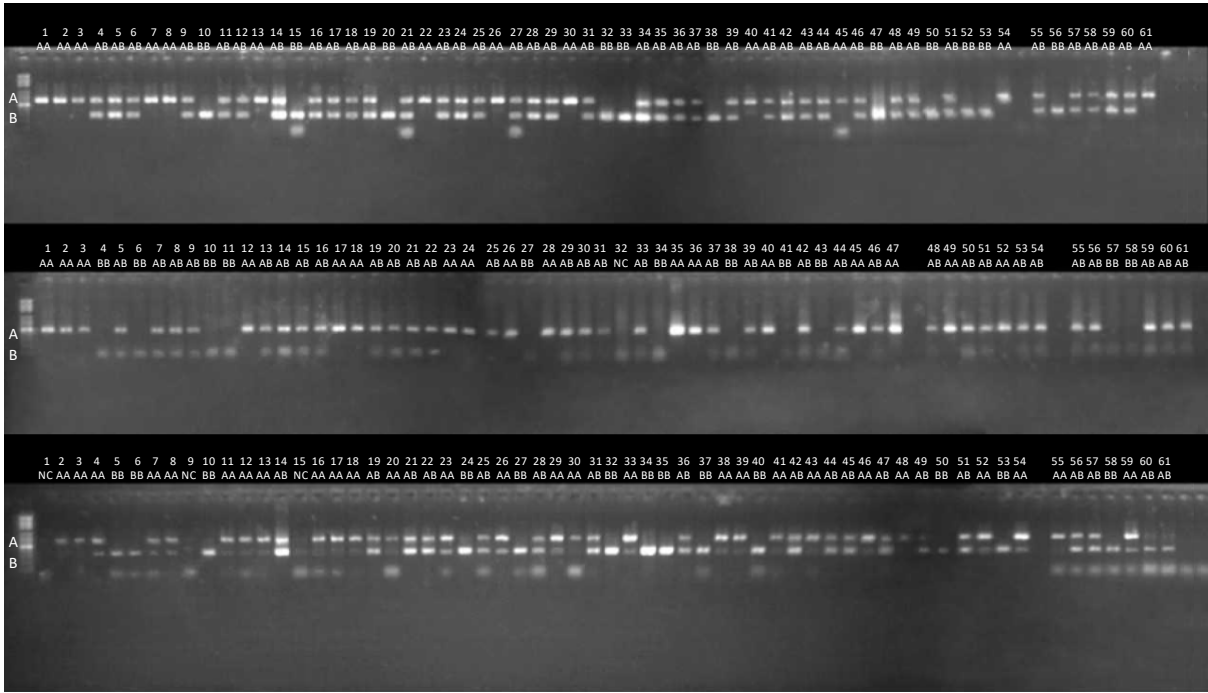541 Genome Project.
542

543



544
545 **Figure 6.** Cumulative frequency of REF– Alu elements discovered from studied individuals.

546

547



548

**Figure 7.** A gel electrophoretic image showing the experimental validation of polymorphic Alu
element insertion (REF– elements). One polymorphic Alu element from chr8:42039896 was tested by
PCR in DNA from 61 individuals. Lower bands show the absence of an Alu insertion (reference allele
A), and upper bands show its presence (alternative allele B).

553

**Figure 8.** A gel electrophoretic image showing the experimental validation of REF+ polymorphic Alu element insertions. Three locations from chr1:169160349, chr15:69049897 and chr3:95116523 were tested by PCR in DNA from 61 individuals. Upper bands show the presence of an Alu insertion (reference allele A), and lower bands show its absence (alternative allele B).

561     **TABLES**

562

563     **Table 1.** Number of REF– and REF+ candidates after different filtering steps

564

| **REF– filtering steps** | |
| --- | --- |
| REF– variations detected in 2,241 individuals | 572,081 |
| REF– candidates that can be located in the reference genome | 379,523 |
| REF– candidates that have unique location in the reference genome | 298,907 |
| REF– candidates after removal of duplicate, closely located and GC-rich k-mers | 13,128 |
| REF– elements that generate reliable genotypes | 9,712 |

| **REF+ filtering steps** | |
| --- | --- |
| Alu signature sequences detected in the reference genome | 267,377 |
| REF+ candidates with 5 bp TSD sequence within 270-350 bp | 110,938 |
| REF+ candidates with BLAST homology | 98,711 |
| REF+ candidates that are not present in chimpanzee genome | 16,434 |
| REF+ candidates after removal of duplicate k-mers | 15,834 |
| REF+ candidates that generate reliable genotypes | 13,396 |

565