

1 Targeted Genotyping of Variable Number Tandem Repeats with 2 adVNTR

3 Mehrdad Bakhtiari¹, Sharona Shleizer-Burko², Melissa Gymrek^{1,2}, Vikas Bansal³, and
4 Vineet Bafna*¹

5 ¹Department of Computer Science & Engineering, University of California, San Diego, La Jolla, CA 92093, USA

6 ²Department of Medicine, University of California, San Diego, La Jolla, CA 92093, USA

7 ³Department of Pediatrics, University of California, San Diego, La Jolla, CA 92093, USA

8 August 7, 2018

9 Abstract

10 Whole Genome Sequencing is increasingly used to identify Mendelian variants in clinical
11 pipelines. These pipelines focus on single nucleotide variants (SNVs) and also structural
12 variants, while ignoring more complex repeat sequence variants. We consider the problem
13 of genotyping *Variable Number Tandem Repeats* (VNTRs), composed of inexact tandem du-
14 plications of short (6-100bp) repeating units. VNTRs span 3% of the human genome, are
15 frequently present in coding regions, and have been implicated in multiple Mendelian disor-
16 ders. While existing tools recognize VNTR carrying sequence, genotyping VNTRs (determining
17 repeat unit count and sequence variation) from whole genome sequenced reads remains chal-
18 lenging. We describe a method, adVNTR, that uses Hidden Markov Models to model each
19 VNTR, count repeat units, and detect sequence variation. adVNTR models can be devel-
20 oped for short-read (Illumina) and single molecule (PacBio) whole genome and exome sequenc-
21 ing, and show good results on multiple simulated and real data sets. adVNTR is available at
22 <https://github.com/mehrdadbakhtiari/adVNTR>

23 **Keywords.** VNTR, Tandem Repeats, VNTR frameshift, Second-generation sequencing, Third-
24 generation sequencing

*Correspondence: vbafna@eng.ucsd.edu

25 1 Introduction

26 Next Generation Sequencing (NGS) is increasingly used to identify disease causing variants in clin-
27 ical and diagnostic settings, but variant detection pipelines focus primarily on single nucleotide
28 variants (SNVs) and small indels and to a lesser extent on structural variants. The human genome
29 contains repeated sequences such as segmental duplications, short tandem repeats, and minisatel-
30 lites which pose challenges for alignment and variant calling tools. Hence, these regions are typically
31 ignored during analysis of NGS data. In particular, *tandem repeats* correspond to locations where a
32 short DNA sequence or *Repeat Unit* (RU) is repeated in tandem multiple times. RUs of length less
33 than 6bp are classified as Short Tandem Repeats (STRs), while longer RUs spanning potentially
34 hundreds of nucleotides are denoted as *Variable Number Tandem Repeats* (VNTRs)(Shriver et al.,
35 1993; Wright, 1994).

36 VNTRs span 3% of the human genome and are often found in coding regions where the re-
37 peat unit length is a multiple of 3 resulting in tandem repeats in the amino acid sequence. More
38 than 1,200 VNTRs with a RU length of 10 or greater exist in the coding regions of the human
39 genome(Tyner et al., 2016). Compared to STRs, which have been extensively studied (Gymrek
40 et al., 2016; Ummat and Bashir, 2014; Liu et al., 2017; Willems et al., 2017; Dolzhenko et al.,
41 2017), VNTRs have not received as much attention. Nevertheless, multiple studies have linked
42 variation in VNTRs with Mendelian diseases (*e.g.*, Medullary cystic kidney disease(Kirby et al.,
43 2013), Myoclonus epilepsy(Lalioiti et al., 1997), and FSHD(Lemmers et al., 2002)) and complex
44 disorders such as bipolar disorder (Table 1). In some cases, the disease associated variants corre-
45 spond to point mutations in the VNTR sequence (Kirby et al., 2013; Ræder et al., 2006) while in
46 other cases, changes in the number of tandem repeats (RU count) show a statistical association
47 (or causal relationship) with disease risk. For example, the insulin gene (INS) VNTR has an RU
48 length of 14 bp with RU count varying from 26 to 200(Pugliese et al., 1997). Variation in this
49 VNTR has been associated with expression of the INS gene and risk for type 1 diabetes (OR =
50 2.2) (Durinovic-Belló et al., 2010). Notwithstanding these examples, the advent of genome-wide
51 SNP genotyping arrays led to VNTRs being largely ignored. They have been called ‘the forgotten
52 polymorphisms’(Brookes, 2013).

53 VNTRs were originally used as markers for linkage mapping since they are highly polymor-

54 phic with respect to the number of tandem repeats at a given VNTR locus(Gelfand et al., 2014).
 55 Traditionally, VNTR genotyping required labor intensive gel-based screens which limited the size
 56 of large population based studies of VNTRs (Orita et al., 1989). Whole genome sequencing has
 57 the potential to detect and genotype all types of genetic variation, including VNTRs. However,
 58 computational identification of variation in VNTRs from sequence remains challenging. Existing
 59 variant calling methods have been developed primarily to identify short sequence variants in unique
 60 DNA sequences that fall into a reference versus alternate allele framework, which is not well suited
 61 for detecting variation in VNTR sequences.

62 Genotyping VNTRs in a donor genome sequenced using short (Illumina) or longer single
 63 molecule reads, requires the following: (a) *recruitment of reads* containing the VNTR sequence;
 64 (b) *counting RUs* for each of the two haplotypes; (c) *identification of indels within VNTRs*; and
 65 (d) identification of mutations within the VNTR. Mapping tools such as BWA(Li and Durbin,
 66 2009) and Bowtie2(Langmead and Salzberg, 2012) can work for read recruitment for STRs, but
 67 are challenged by insertion/deletion of larger repeat units. Mapping issues also confound existing
 68 variant callers, including realignment tools such as GATK IndelRealigner(DePristo et al., 2011)
 69 if the total VNTR length is larger than the read length. This is because reads contained within
 70 the VNTR sequence have multiple equally likely mappings and therefore will be mapped randomly
 71 to different locations with low mapping quality(Kirby et al., 2013). Detection of point mutations
 72 in long VNTRs requires integrating information across the entire VNTR sequence. For VNTRs

Gene	Chr	Unit len	Number of units		Annotation	Inheritance	Disease
			Normal	Pathogenic			
<i>PER3</i>	1	54	4	5	coding	A	Bipolar disorder(Benedetti et al., 2008)
<i>MUC1</i>	1	60	11-12	single insertion	coding	M	MCKD1(Kirby et al., 2013)
<i>IL1RN</i>	2	86	3-6	2	intron	A	Stroke, CAD(Worrall et al., 2007)
<i>DUX4</i>	4	3.3kb	11-100	1-10		M	FSHD(Lemmers et al., 2002)
<i>DAT1</i>	5	44	7-11	10 (ADHD)	UTR	A	ADHD, Parkinson's(Franke et al., 2010; Kirchheiner et al., 2007)
<i>MUC21</i>	6	45	26-27	4 bp deletion	coding	A	Diffuse panbronchiolitis (DPB)(Hijikata et al., 2011)
<i>CEL</i>	9	33	11-21	single deletion	coding	M	Monogenic diabetes(Ræder et al., 2006)
<i>INS</i>	11	14-15	26-200	26-44 (T1D)	promoter	A	T1D;T2D;Obesity(Pugliese et al., 1997; Durinovic-Belló et al., 2010)
<i>DRD4</i>	11	48	2-11	7	coding	A	OCD, ADHD(LaHoste et al., 1996; Viswanath et al., 2013)
<i>ACAN</i>	15	57	27-33	13-25	coding	A	Osteochondritis dissecans(Eser et al., 2011)
<i>ZFX3</i>	16	12	4-5		coding	A	Kawasaki
<i>GP1BA</i>	17	39	1-4	2/3 genotype	coding	A	ATF in Stroke(Cervera et al., 2007)
<i>SERT</i>	17	16-17	9/10/12		intron	A	BPSD, Alzheimer's(Haddley et al., 2011; Pritchard et al., 2007)
<i>SERT</i>	17	22	14	16 (OCD)	promoter	A	OCD,Anxiety, Schizophrenia(Haddley et al., 2011)
<i>HIC1</i>	17	70	1-4	5+/5+	promoter	A	Metastatic Colorectal Cancer(Okazaki et al., 2017)
<i>MMP9</i>	20	12	5-6		coding	A	Kawasaki
<i>CSTB</i>	21	12	2-3	12+	5'UTR	M	Progressive myoclonic epilepsy 1A(Lalioti et al., 1997)
<i>MAOA</i>	X	30	2-5	4	promoter	A	Bipolar disorder(Byrd and Manuck, 2014)

Table 1: **Disease-linked VNTRs** are generally distinguished from STRs by a longer length (≥ 6) of the repeating unit. 'M' denotes Mendelian inheritance, while 'A' represents possibly complex inheritance captured via Association. As it is difficult to genotype VNTRs, most cases have been determined via association, but the inheritance mode could be high penetrance.

73 whose total sequence length (RU count times the RU length) is much longer than the read length,
74 detection of SNVs and indels is not feasible using existing variant callers. We focus mainly on
75 problems (a,b) relating to recruitment and RU counting. For problem (c), we focus on difficult case
76 of large (≥ 250 bp) VNTRs within coding regions where the indel shifts the translation frame. We
77 do not tackle problem (d) in this manuscript.

78 Other tools have addressed the problem of RU count estimation, focusing on the related problem
79 of STR genotyping. Some of these tools do not accept large repeating patterns as input (Willems
80 et al., 2017; Liu et al., 2017). Others require all repeat units to be near-identical (Dolzhenko et al.,
81 2017; Ummat and Bashir, 2014). In particular, ExpansionHunter (Dolzhenko et al., 2017) looks for
82 exact matches of short repeating sequence within flanking unique sequences, and works for STRs,
83 but not as well with the larger VNTRs with variations in RUs (Results). VNTRseek (Gelfand et al.,
84 2014) detects a VNTR-like pattern in reads and aligns it to tandem repeats, but uses a complex
85 alignment process making it difficult to run the tool. Alignment based tools need to align reads at
86 both unique ends, which may not be possible for short (Illumina) reads. Single molecule reads (*e.g.*,
87 PacBio (Eid et al., 2009), Nanopore (Clarke et al., 2009)) can span entire VNTR regions, but it is
88 difficult to estimate the RU count directly since the distance between the flanking regions varies
89 dramatically from read to read due to an excess of indel errors. For example, 14 reads spanning
90 the SERT VNTR in the in the PacBio sequencing data of NA12878 individual from Genome in
91 a Bottle (Zook et al., 2016) included fifteen distinct lengths between 292bp and 385bp, leading to
92 length-based RU count estimates 13, 14, 15, 16, and 18 for the diploid genome.

93 In contrast to methods like VNTRseek which seek to *discover/identify* VNTRs, we describe
94 a method, adVNTR, for *genotyping VNTRs* at targeted loci in a donor genome. For any target
95 VNTR in a donor, adVNTR reports an estimate of RU counts and point mutations within the
96 RUs. It trains Hidden Markov Models (HMMs) for each target VNTR locus, which provide the
97 following advantages: (i) it is sufficient to match any portions of the unique flanking regions for
98 read alignment; (ii) it is easier to separate homopolymer runs from other indels helping with
99 frameshift detection, and to estimate RU counts even in the presence of indels; (iii) each VNTR
100 can be modeled individually, and complex models can be constructed for VNTRs with complex
101 structure, along with VNTR specific confidence scores. For longer VNTRs not spanned by short
102 reads, adVNTR can still be used to detect indels, while providing lower bounds on RU counts.

103 Also, exact estimates for RU counts could be made for shorter VNTRs. Using simulated data as
104 well as whole-genome sequence data for a number of human individuals, we demonstrate the power
105 of adVNTR to genotype VNTR loci in the human genome.

106 2 Results

107 Our method, adVNTR, requires training of separate HMM models for each combination of target
108 VNTR and sequencing technologies. The detailed training procedure is described in Methods.
109 Given trained models, adVNTR genotypes the VNTRs in three stages: (i) Selection of reads that
110 contain VNTR locus (read recruitment); (ii) RU count estimation; and, (iii) variant detection. We
111 report results on performance of adVNTR in each of these stages using simulated and read datasets
112 based on short-read (Illumina) and single molecule (PacBio) technologies.

113 **HMM training.** Initial HMMs were trained using multiple alignments of RU sequences from
114 the reference assembly hg19(Lander et al., 2001), as described in methods. Similarly, HMMs
115 were trained for the left flanking and right flanking regions for each VNTR. The HMM models
116 were augmented using data from Genome in a Bottle (GIAB) project (NA12878 WGS). VNTR
117 models were trained for VNTRs in coding and promoter regions of the genome, for both Illumina
118 (1755 models) and PacBio (2944 models; Supplementary Material “Selecting Target VNTRs”).
119 Subsequently, we tested performance for (a) read-recruitment, (b) counting of Repeat Units, and
120 (c) detection of indels.

121 **Test Data.** To evaluate performance for *PacBio*, we simulated haplotypes for each of the 2944
122 VNTRs, revising the RU count to be ± 3 of the RU count in hg19, and setting 1 as the minimum
123 RU count. We simulated haplotype reads (15X coverage) using SimLoRD(Stöcker et al., 2016) and
124 aligned those reads to hg19 using Blasr(Chaisson and Tesler, 2012). For Illumina sequencing, we
125 used ART(Huang et al., 2011) to simulate haplotype WGS (shotgun 150bp) reads at 15X coverage
126 for each VNTR and simulated VNTR haplotype with changes in RU counts similar to PacBio. Pairs
127 of haplotypes were merged to get (30X coverage) diploid samples. The resulting data-sets were
128 called *PacBioSim* and *IlluminaSim*, respectively (Supplementary Material “Test Datasets”). To
129 evaluate performance of frameshift identification, we collected a set of 115 VNTRs (Supplementary
130 Material “Selecting Target VNTRs”). For each VNTR, we simulated haplotypes that contain a

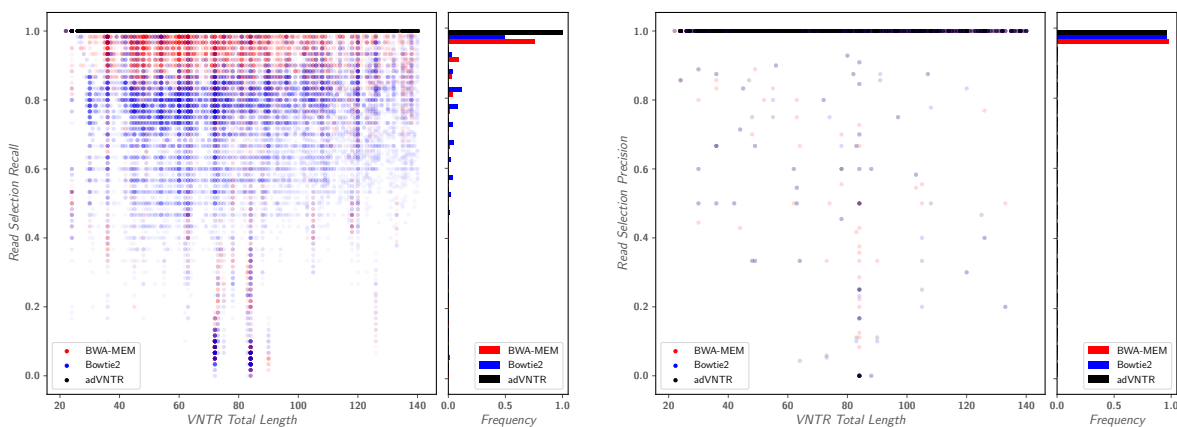


Figure 1: **Read recruitment quality on Illumina reads.** (A) Comparison of the recall ($\#$ true recruited reads/ $\#$ true reads) of adVNTR read recruitment against BWA-MEM and Bowtie2, as a function of VNTR length for 1775 VNTRs with different counts (31,788 tests). Each dot corresponds to a separate test. (B) Precision ($\#$ true recruited reads/ $\#$ recruited reads) of read recruitment.

131 deletion or an insertion in the VNTR (Supplementary Material “Test Datasets”). We simulated
132 reads from each of these haplotypes and merged pairs of haplotypes to obtain diploid samples. We
133 denote this data-set as *IlluminaFrameshift*.

134 **Read recruitment.** adVNTR takes a collection of VNTR models as input, and as a first step,
135 recruits reads that map to any of the VNTRs in the list. In testing recruitment for PacBio, we found
136 that alignment tools such as Blasr perform well in recruiting VNTR reads even in the presence of
137 deletions and insertions (data not shown) and used Blasr for all read recruitment. For *Illumina*
138 reads, we tested adVNTR read-recruitment for all 1775 VNTRs using IlluminaSim, and compared
139 against mapping tools BWA-MEM, Bowtie2, and BLAST. adVNTR achieves much greater recall
140 while maintaining or exceeding the precision of other tools (Fig. 1 and Fig. S3). Specifically,
141 adVNTR recall was 100% for 99.9% of the VNTRs, whereas the next best tool (BWA-MEM)
142 achieved this only for 68.2% of the VNTRs. The other mapping tools lose mapping sensitivity
143 when RU counts are increased or decreased (large indels), and perform best when the RU counts
144 are the same as reference (Fig. S2A-C), partially explaining their lower recall.

145 **VNTR genotyping (RU count estimation) with PacBio reads.** Recall that sequencing
146 (particularly homopolymer) errors can cause lengths to change, particularly for short RU lengths
147 and larger RU counts. To test adVNTR performance on PacBioSim, we compared against a naïve
148 method that estimates RU counts based on read length between the flanking regions from the

149 consensus of reads that cover VNTR. Detailed performance on three exemplars (INS, CSTB, and
150 HIC1) gene showed high genotype accuracy for adVNTR over a wide range of RU counts, and
151 coverage (Fig. 2A). Similar results were obtained for all 2944 VNTRs (Fig. 2B). Overall, 98.45%
152 of adVNTR estimates were correct while 26.45% of estimates made by naïve method were correct.
153 As it is difficult for the naïve method to call heterozygotes, we also compared on the subset of test
154 data with homozygous RU counts. 97.95% of adVNTR estimates were correct, while the consensus
155 method was correct in 66.16% of samples (Fig. S4). adVNTR estimates were uniformly good except
156 at low sequence coverage. To test for accuracy with changing RU counts, we simulated different
157 RU counts for individuals at 3 VNTRs (Table S4). adVNTR RU counts showed 100% accuracy in
158 each of the 52 different samples tested.

159 To test performance on real data where the true VNTR genotype was not known, we checked
160 for Mendelian inheritance consistency in the AJ trio from Genome in a Bottle (GIAB)(Zook et al.,
161 2016) and a Chinese Han trio from NCBI SRA (accession PRJEB12236). On four disease related
162 VNTRs, adVNTR predictions were consistent in each case (Fig. 2C). On the 2944 genic VNTRs,
163 the trio consistency of adVNTR calls was correlated with coverage. At a posterior probability
164 threshold of 0.99, 86.98% of the calls in the AJ trio, and 97.08% of the calls in the Chinese trio,
165 were consistent with Mendelian inheritance (Fig. 2E). Many of the discrepancies could be attributed
166 to low coverage and missing data. Increasing sequence coverage threshold from $5\times$ to $10\times$ increased
167 the average posterior probability from 0.91 to 0.98 and resulted in improved RU count accuracy
168 (Fig. S5). Also, many of these discrepancies in RU counts were off-by-one errors (Fig. S6). These
169 off-by-one discrepancies could be acceptable for Mendelian disease testing as the pathogenic cases
170 often have large changes in RU counts. Treating the off by one counts as correct, we found that
171 98.66% and 99.91% of the high confidence calls in AJ and Chinese trios, respectively, were consistent
172 (Fig. 2F). Finally, some of the off-by-one counts could be natural genetic variation.

173 We also performed a long range (LR)PCR experiment on the individual NA12878 to assess the
174 accuracy of the adVNTR genotypes using PacBio data (Table S2 and Table S3). The observed
175 PCR product lengths (black bands in Fig. 2D) were consistent with the adVNTR predictions (red
176 arrows), while being different from the hg19 reference RU count. adVNTR correctly predicted all
177 VNTRs to be heterozygous with the exception of SLC6A4, that was predicted to be homozygous.

178 While we could not get the VNTR discovery tool VNTRseek(Gelfand et al., 2014) to run on our

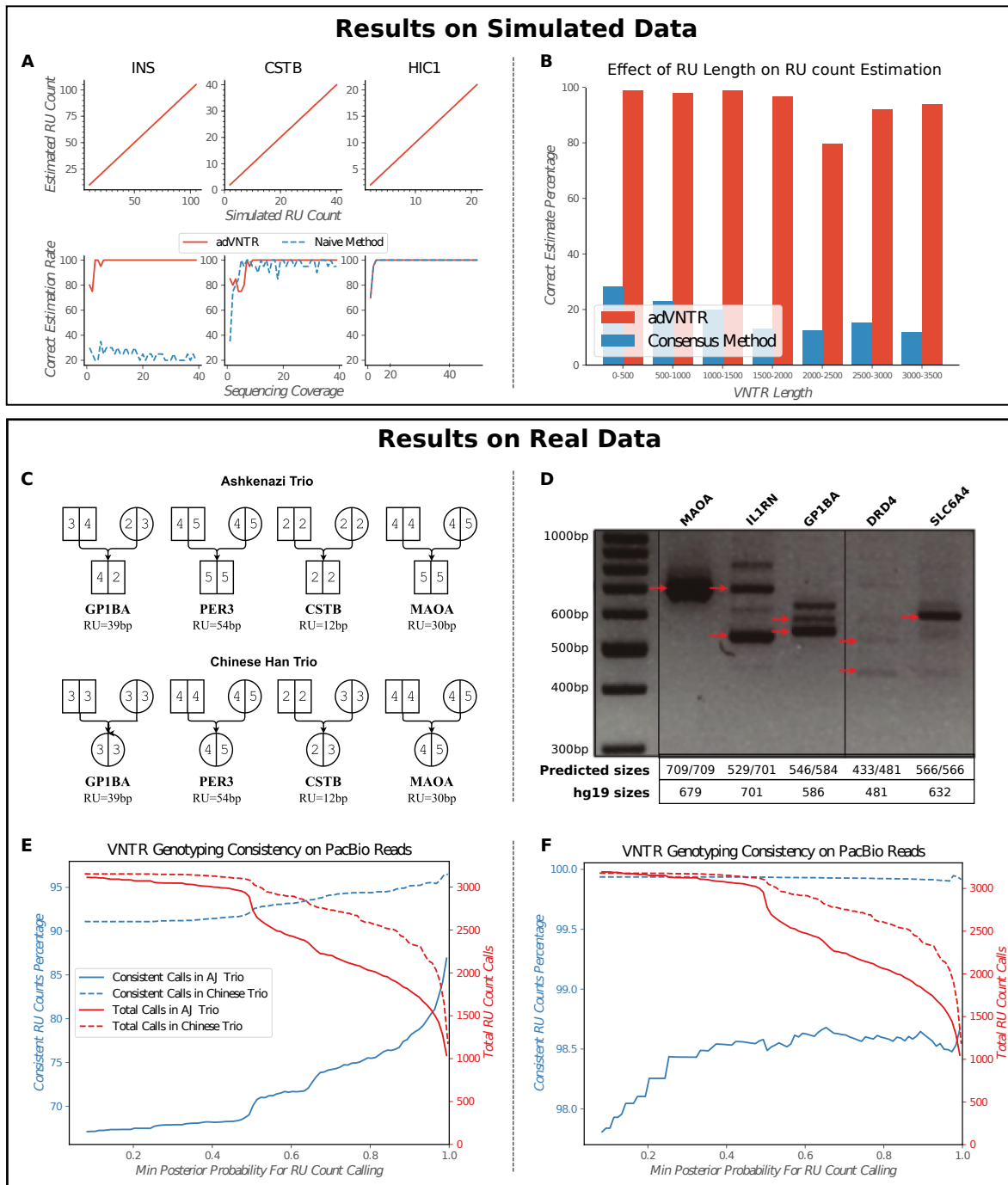


Figure 2: VNTR genotyping using PacBio data. (A) RU count estimation on simulated PacBio reads as a function of RU count and coverage for 3 medically relevant VNTRs: INS (RU length 14bp), CSTB (12 bp), and HIC1 (70bp). adVNTR performance is compared to a naïve method. (B) The effect of RU length on count accuracy over 2944 VNTRs (30418 tests). (C) Mendelian consistency of genotypes at 4 VNTR loci in the Chinese Han and Ashkenazi trios. Note that MAOA results are consistent with its location on Chr X. (D) LR-PCR based validation of genotypes at 5 disease-linked VNTRs in NA12878. Red arrow correspond to VNTR lengths estimated by multiplying predicted RU counts with RU lengths. (E) Fraction of consistent calls and number of calls across 2944 VNTRs in AJ and Chinese trios from GIAB and NCBI-SRA. (F) Fraction of consistent calls allowing for off-by-one errors.

179 machine (personal communication), we observed that the authors had predicted 125 VNTRs in the
180 Watson sequenced genome(Wheeler et al., 2008), and 75 VNTRs in two trios as being polymorphic.
181 In contrast, analysis of the PacBio sequencing data identified >500 examples of polymorphic VNTRs
182 that overlap with coding regions. The results suggest that variation in RU counts of VNTRs and
183 their role in influencing phenotypes might be greater than previously estimated.

184 **RU counting with Illumina.** The adVNTR estimate correctly matched both RU counts in 91.6%
185 of the cases in the IlluminaSim dataset (1775 VNTRs with up to 21 diploid RU counts each) and
186 matched at least one RU count in 97% of the cases (Fig. 3A,B). Most of the discrepancies occurred
187 in VNTRs with longer lengths not covered by Illumina reads (Fig. 3C,D). While there was a drop
188 in accuracy for increasing lengths, 84% of the genic VNTRs are shorter than 150bp, and could be
189 genotyped with 94.6% accuracy. Tools such as VNTRseek require at least 20bp flanking each side
190 of the VNTR and do not return a result for VNTRs with total length greater than 110bp, while
191 adVNTR could predict the genotype correctly in a majority of those cases (Supplementary Material
192 “VNTRseek”). ExpansionHunter, a tool designed primarily for STR genotyping (Dolzhenko et al.,
193 2017) provided incorrect estimates in over 90% cases from this data-set (Fig. S7). ExpansionHunter
194 makes the assumption that the different RUs are mostly identical in sequence which is valid for
195 STRs but not for most VNTRs, and we tested this through 52 samples on three VNTRs. adVNTR
196 predicted the correct genotype in all but 6 cases, with erroneous calls only in the case of high RU
197 counts where the read length did not span the VNTR perfectly, while ExpansionHunter did not
198 return the correct estimate in most cases (Table S4).

199 On the AJ trio from GIAB, 98.08% of the high confidence adVNTR calls were consistent
200 with Mendelian inheritance (Fig. 3E). Note that 95.93% of all calls were high confidence (posterior
201 probability ≥ 0.99). We validated adVNTR calls on 12 VNTRs using Gel electrophoresis (Table S3).
202 adVNTR predicted the correct RU counts in all cases, except in two cases where the PCR primers
203 failed to produce a band (Fig. 3F, S8). We also compared adVNTR against ExpansionHunter on
204 7 disease related short VNTRs in the AJ trio and obtained similar results (Table S5).

205 To test adVNTR for population-scale studies of VNTR genotypes using WGS data replacing
206 labor intensive gel electrophoresis(Byrd and Manuck, 2014; Cervera et al., 2007), we scanned the
207 PCR-free WGS data for 150 individuals (50 in each population) obtained from 1000 genomes

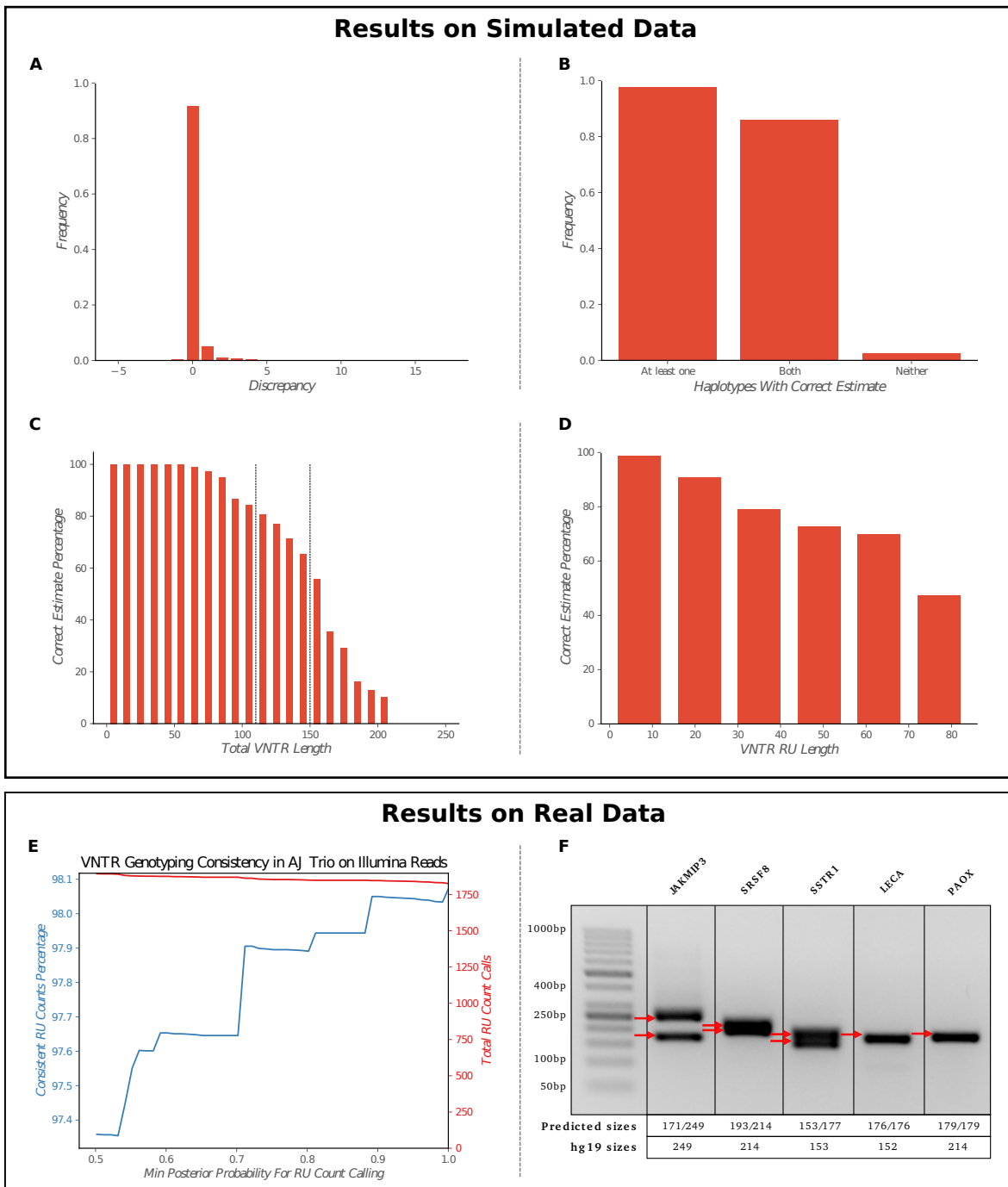


Figure 3: VNTR genotyping using Illumina data. (A-D) Correctness of RU count prediction for 1775 coding VNTRs in the IlluminaSim dataset, described by (A) RU count discrepancy, (B) haplotypes with correct estimates, (C) correctness as a function of VNTR length, and (D) RU length. (E) Consistency of adVNTR calls on the AJ trio WGS data from GIAB. Red line describes the cumulative number of calls made at specific posterior probability cut-offs. (F) Gel electrophoresis based validation of adVNTR calls on 5 short VNTRs using WGS of individual NA12878 from GIAB. Red arrows correspond to VNTR lengths estimated by multiplying the RU lengths with the estimated RU counts.

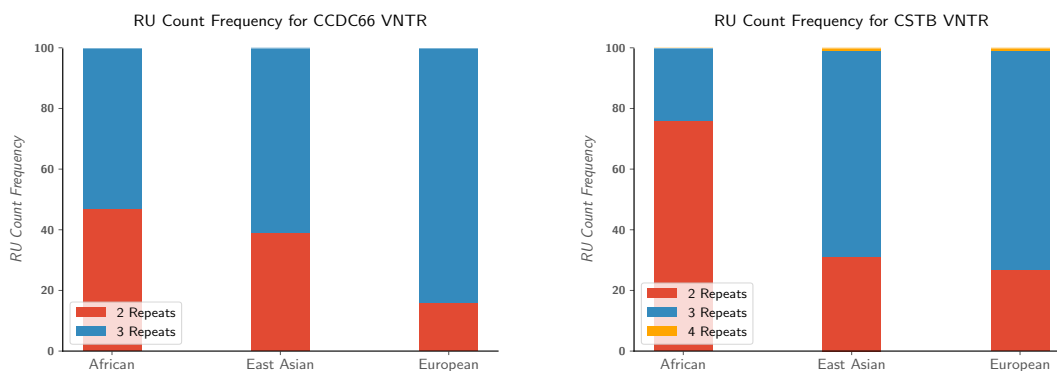


Figure 4: **Population-scale genotyping of VNTRs.** (A) RU count frequencies for the VNTR in CCDC66 gene, and (B) CSTB in African, Asian, and European population samples from 1000 genomes project. RU counts of 4 and higher in CSTB are associated with myoclonal epilepsy.

208 project(Consortium, 2015). We observed population specific RU counts (frequency difference >
209 10%) in 97 of 202 VNTRs tested (Table S7). Fig. 4 shows the RU count frequencies for a disease-
210 linked VNTR in the coding region of CSTB and a coding VNTR in CCDC66. The results suggest
211 an increase in VNTRs with higher RU counts with an increase in divergence time from Africa.
212 Thus RU3 is more prevalent in both VNTRs. We also observed RU4 in CSTB6 VNTR in the Asian
213 and European populations, where RU counts 4 and above have been associated with progressive
214 myoclonal epilepsy (Lalioti et al., 1997).

215 **VNTR mutation/indel detection.** As a proof of concept for other applications, we tested in-
216 del detection, focusing in particular on frameshifts in coding VNTRs. The CEL gene is known to
217 contain a VNTR where a deletion changes the coding frame. We simulated Illumina reads from
218 20 whole genomes after introducing a single insertion or deletion in the middle of the VNTR re-
219 gion in the CEL gene. As a negative control, we simulated 10 WGS experiments with a range of
220 sequence coverage values. We ran adVNTR, Samtools mpileup(Li, 2011), and GATK Haplotype-
221 Caller(DePristo et al., 2011) which uses GATK IndelRealigner, to identify frameshifts in each of the
222 simulated datasets, and the 10 control datasets. On the control data, none of the tools found any
223 variant. On the simulated indels, adVNTR made the correct prediction in each case (Suppl. Ta-
224 ble S6), while Samtools and GATK were unable to predict a single insertion or deletion. This result
225 is not surprising as the reads have poor alignment scores, and the indel can be mapped to multiple
226 locations (Suppl. Fig. S9)(Robinson et al., 2011). We note that mapping ambiguity in aligning
227 each read made it difficult to pinpoint the location of single indel. However, by integrating the
228 information across all reads, we could predict the occurrence of a frameshift in the VNTR. We next

229 tested adVNTR frameshift prediction on the 115 VNTRs in the IlluminaFrameshift dataset, simu-
230 lating 4090 total cases. Overall, the frameshifts in the VNTR regions were predicted with 51.7%
231 sensitivity and 86.8% specificity, in contrast with the 49.7, 43.5% sensitivity, specificity achieved
232 by GATK. Detailed performance of methods for each VNTR is available in Table S7. Note that
233 the performance is model specific and depends upon the similarity of different Repeat Units in a
234 VNTR. For 29 of the 115 VNTRs, adVNTR showed high sensitivity ($\geq 90\%$) and specificity (100%).

235 As frameshifts in the VNTR region of the CEL gene have been linked to a monogenic form
236 of diabetes (Ræder et al., 2006), we tested for frameshifts in CEL using whole Exome sequencing
237 (WES) data from 2,081 cases with Type 2 Diabetes (Fuchsberger et al., 2016) and compared the
238 numbers to 2,090 control individuals. WES data analysis is challenging as high GC-content makes
239 it difficult to PCR-amplify this VNTR. adVNTR found that while none of the controls had any
240 evidence of a frameshift, 8 of the 2,081 diabetes cases showed a frameshift in this VNTR region
241 (Suppl. Fig. S10).

242 **Compute requirements for genotyping.** adVNTR is multi-threaded. In genotyping mapped
243 PacBio reads at 30X coverage, adVNTR took 6 hours using Intel Xeon(R) 4-core CPUs (≤ 24
244 CPU-hours) to genotype all 2944 VNTRs, and 14:15 hours (≤ 57 CPU-hours) for 70X coverage.
245 For Illumina reads at 40X coverage, adVNTR took 87:30 cpu-hours on a single core to complete
246 read recruitment as well as genotyping of 1775 VNTRs.

247 **3 Discussion**

248 The problem of genotyping VNTRs (determining diploid RU counts and mutations) is increasingly
249 important for clinical pipelines seeking to find the genetic mechanisms of Mendelian disorders. As
250 VNTRs have not been extensively studied, existing research is often focused on their discovery.
251 One of the contributions of this paper is the separation of initial VNTR discovery from VNTR
252 genotyping, and a focus on the genotyping problem. adVNTR genotypes VNTRs using a hidden
253 markov model for each target VNTR, providing a uniform training framework, but still allowing us
254 to tailor the models for complex VNTRs on a case by case basis. The problem of mismapping due
255 to indels introduced by changing RU counts confounds most mapping based tools, but is solved here
256 by collapsing all RU copies and building HMMs that allow for variation in the RUs. adVNTR was

257 tested extensively on data from different sequencing technologies, including Illumina and PacBio.

258 Like other STR genotyping tools, adVNTR works best when reads span the VNTR. However,
259 even with this limitation, there are (a) close to 100,000 VNTRs in the genic regions of human
260 genome that can be spanned by Illumina reads; (b) indel detection is possible even when RU
261 counting is not, for long VNTRs; (c) lower bounds on RU counts can separate some pathogenic
262 cases from normal cases particularly when the normal VNTR length is shorter than the read
263 length, while the pathogenic case is much longer (e.g. *CSTB*). Finally, dropping costs for long read
264 sequencing (esp. PacBio, and Nanopore) will allow us to span and genotype over 158,000 genic
265 VNTRs.

266 The choice between short and long read technologies offers some trade-offs. Specifically, long
267 reads allow for the targeted genotyping of a larger set of VNTRs (559,804), and are becoming
268 increasingly cost-effective. However, the large numbers of indels in these technologies reduce the
269 accuracy somewhat, and they are best used when there is a big difference between normal and
270 pathogenic cases in terms of RU counts, or when the VNTRs are too long to be spanned by
271 Illumina.

272 In contrast, short-read Illumina sequencing is increasingly used for Mendelian pipelines, and
273 can be easily extended to include VNTR genotyping, with higher accuracy than PacBio. Also, the
274 large number of VNTRs (458,158) that can be spanned by Illumina reads makes it the technology
275 of choice for association testing and population based studies.

276 In this research, we also provided initial results on genotyping frameshift errors in coding VN-
277 TRs, focusing on the easier case when all RUs have the same length. Future work will focus on
278 extending the target VNTRs for RU counting and frameshift detection for VNTRs that are of
279 medical interest, population genetics of VNTRs, and algorithmic strategies for speeding up VNTR
280 discovery and genotyping.

281 4 Method

282 A VNTR sequence can be represented as $SR_1R_2 \dots R_uP$, where S and P are the unique flanking
283 regions, and $R_i(1 \leq i \leq u)$ correspond to the tandem repeats. For each i, j , R_i is similar in sequence
284 to R_j , and the number of occurrences, u , is denoted as the *RU count*. We do not impose a length

285 restriction on S and P , but assume that they are long enough to be unique in the genome. For
 286 genotyping a VNTR in a donor genome, we focus primarily on estimating the diploid RU counts
 287 (u_1, u_2) . However, many ($\sim 10^3$) VNTRs occur in coding regions, and mutations, particularly
 288 frameshift causing indels, are also relevant. Our method, adVNTR, models the problems of RU
 289 counting and mutation detection using HMMs trained for each target VNTR. adVNTR requires a
 290 one-time training of models for each combination of a VNTR and sequencing technology, although
 291 the user has the option to retrain models. Once models are trained, it has three stages for geno-
 292 typing: (i) Read recruitment; (ii) RU count estimation; and, (iii) variant (indel) detection. We
 293 describe the training procedure and the three modules below.

294 **HMM Training.** The goal of training is to estimate model parameters for each VNTR and each
 295 sequencing technology. Previous works have shown that an HMM with three groups of states could
 296 be used to find similarities between biological sequences (Eddy, 1996). In this model, a profile-
 297 HMMs can model a groups of sequences. Then, a new sequence can be aligned to a profile HMM to
 298 discover sequence family (Krogh et al., 1994). We use an HMM architecture with three parts, which
 have their own three groups of states (Fig. 5). The first part matches the 5' (left) flanking region

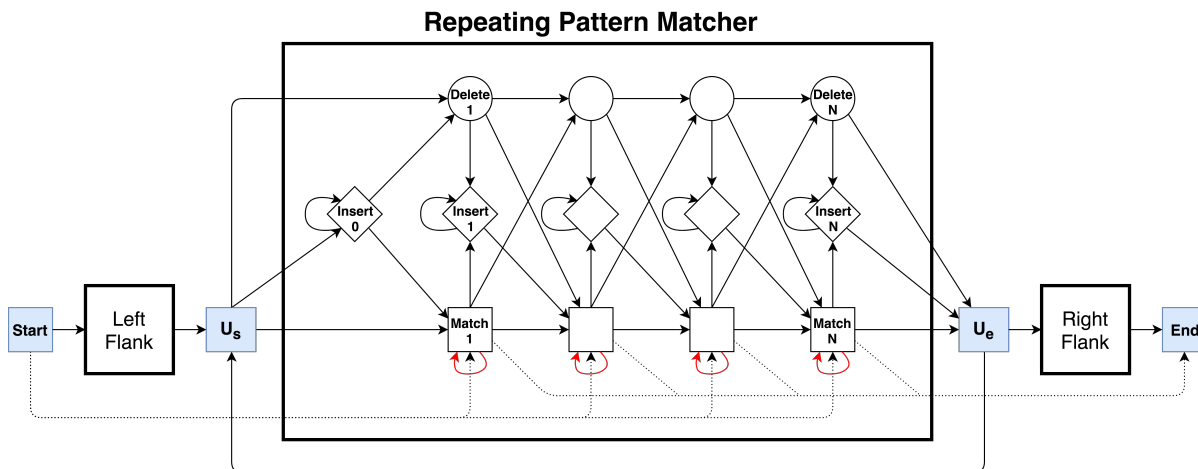


Figure 5: **The VNTR HMM.** The HMM is composed of 3 profile HMMs, one each for the left and right flanking unique regions, and one in the middle to match multiple and partial numbers of RUs. The special states U_s ('Unit-Start'), and U_e ('Unit-End') are used for RU counting. Dotted lines refer to special transitions for partial reads that do not span the entire region.

299

300 of the VNTR. The second part is an HMM which matches an arbitrary number of (approximately
 301 identical) repeating units. The last part matches the 3' (right) flanking region (Fig. S1). The
 302 RU pattern is matched with a profile HMM (RU HMM), with states for matches, deletions, and

303 insertions, and its model parameters are trained first. To train RU HMM for each VNTR, we
304 collected RU sequences from the reference assembly(Lander et al., 2001) and performed a multiple
305 sequence alignment(Eddy et al., 1995). Let $h(i, j)$ denote the number of observed transitions from
306 state i to state j in hidden path of each sequence in multiple alignment, and $h_i(\alpha)$ denote the number
307 of emissions of α in state i . We define permissible transition (arrows in Fig. 5) and match-state
308 emission probabilities as follows:

$$T(i, j) = \frac{h(i, j) + b_0}{\sum_{i \rightarrow l} (h(i, l) + b_0)}, \quad E_i(\alpha) = \frac{h_i(\alpha) + b_1}{\sum_{\alpha'} (h_i(\alpha') + b_1)} \quad \text{for } \alpha, \alpha' \in \{A, C, G, T\}.$$

309 Non-permissible transitions have probability 0, and $h_i(\alpha) = 1/4$ for insert state i and 0 for deletions.
310 The pseudocounts b_0 and b_1 were estimated by initially setting them to the error rate of the
311 sequencing technology, but they (along with other model parameters) were updated after aligning
312 Illumina or PacBio reads to the model. The RU HMM architecture was augmented by adding (a)
313 transitions from U_e to U_s to allow matching of variable number of RU; (b) adding the HMMs for
314 the matching of any portions of left and right flanking sequences; and (c) by adding transitions
315 to match reads that match either the left flanking or the right flanking region. In addition, reads
316 anchored to one of the unique regions can jump past the other HMM using dotted arrows.

317 While error correction tools for PacBio have been developed, most do not work for repetitive
318 regions,(Hackl et al., 2014; Salmela and Rivals, 2014; Au et al., 2012; Miclotte et al., 2016; Lee
319 et al., 2014) and others assume a single haplotype for error correction(Salmela et al., 2016; Berlin
320 et al., 2015). In contrast, the HMM allows us to model many of the common (homopolymer)
321 errors directly. Insertion deletion errors are common in single molecule sequencing particularly in
322 homopolymer runs of length ≥ 6 , and occur mostly as insertions in the homopolymer run(Chaisson
323 and Tesler, 2012). Consider a match state i with highest emission probability for nucleotide α . The
324 transition probability $T(i, i)$ from a match state i to itself was set based on the match probabilities
325 of α in previous $k = 6$ states. The model parameters were further updated using genome sequencing
326 data of NA12878 (Supplementary Material “Model Structure and Parameter Setting”).

327 **Read Recruitment.** The first step in adVNTR is to *recruit* all reads that match a portion of the
328 VNTR sequence. Alignment-based methods do not work well due to changes in RU counts (See
329 Results), but the adVNTR HMM allows for variable RU count. To speed up recruitment, we used

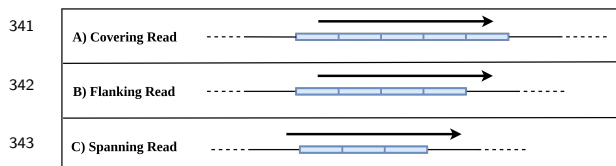
330 an Aho-Corasick keyword matching algorithm available as part of the Blast package(Altschul et al.,
 331 1990) to identify all reads that match a keyword from the VNTR patterns or the flanking regions.
 332 Note that the dictionary construction is a one-time process, and all reads must be scanned once
 333 for filtering. The keyword size and number of keywords were empirically chosen for each VNTR.
 334 Filtered reads were aligned to the HMM using the Viterbi algorithm. Only reads with matching
 335 probability higher than a specified threshold were retained. To compute the selection threshold
 336 for each VNTR, we aligned non-target genomic sequences that passed the keyword matching step
 337 to the HMM to form an empirical false distribution. Subsequently, we aligned VNTR encoding
 338 sequences to the HMM to form the score distribution of true reads. Then, we used a Naïve Bayes
 339 classifier to select a threshold.

Estimating VNTR RU Counts. All reads covering an RU element are aligned, or ‘matched’
 to the HMM using the Viterbi algorithm to create, in effect, a new multiple alignment. Recalling
 the Viterbi algorithm, let $V_{k,j}$ denote the highest (log) probability of emitting the first k letters
 of the sequence s_1, s_2, \dots, s_n and ending in state j of an HMM. Let, $\text{Prev}_{k,j}$ denote the state j'
 immediately prior to j in this optimum parse. Then,

$$V_{k,j} = \max_{j'} \{V_{k',j} + \log T(j', j) + \log E_j(s_k)\}, \quad (1)$$

$$\text{Prev}_{k,j} = \arg \max_{j'} \{V_{k',j} + \log T(j', j) + \log E_j(s_k)\}, \quad (2)$$

340 where, $k' = k - 1$ for match or insert states; $k' = k$ otherwise.



344 **Figure 6: Estimates of RU counts using re-**
 345 **recruited reads.** (A) $(k_1, k_2, k_3) = (1, 3, 1)$; RU count
 346 ≥ 5 . (B) $(k_1, k_2, k_3) = (0, 3, 1)$; RU count ≥ 4 (C)
 347 $(k_1, k_2, k_3) = (0, 3, 0)$; RU count = 3.
 348 example).

For each read, the Viterbi algorithm allows for
 the enumeration of the maximum likelihood (ML)
 path by going backwards from $\text{Prev}(\text{End}, n)$. Ignor-
 ing all but the U_s and U_e states in the Viterbi path,
 we get a pattern of the form $U_e^{k_1}(U_s U_e)^{k_2} U_s^{k_3}$ with
 $k_1, k_3 \in \{0, 1\}$, and $k_2 \geq 0$. We estimate the RU

349 One of the main reasons for erroneous RU counts is stutter during PCR amplification. The

350 PCR amplification process is similar to replication errors that result on genetic RU count variation
 351 during cell-division, except that there are multiple rounds of amplification. In each PCR round, the
 352 number of copies might change by 1 with some probability. Once a single event has occurred and
 353 an erroneous template is generated, the event of having another change is likely to be independent
 354 of the previous event(Gymrek, 2016). To model errors in read counts, we define parameter r_ϵ s.t.
 355 r_ϵ^Δ is the probability of RU counting error by $\pm\Delta$ in the estimation of the true count. Thus the
 356 probability of getting the correct count is $1 - r$, where

$$r = 2(r_\epsilon + r_\epsilon^2 + r_\epsilon^3 + \dots) = \frac{2r_\epsilon}{1 - r_\epsilon}$$

357 The analysis of reads at a VNTR gives us a multi-set of RU counts (or lower bounds) c_1, c_2, \dots, c_n .
 358 We assume that the donor genome is diploid but do not require any phasing information in the
 359 computation of the multi-set. Additionally, we allow the possibility that all reads are sampled from
 360 one haplotype with the RU count of the missing haplotype being X . We define $C = \{c_1, c_2, \dots, c_n\} \cup$
 361 $\{X\}$ and use C to get a list of possible genotypes (c_i, c_j) with $c_i \leq c_j$. Then, the conditional
 362 likelihood of a read with RU count c is given by:

$$\Pr(\text{RU} = c | (c_i, c_j)) = \begin{cases} 1 - r & c = c_i = c_j \\ \frac{1}{2}((1 - r) + r_\epsilon^{|c - c_j|}) & c = c_i \\ \frac{1}{2}((1 - r) + r_\epsilon^{|c - c_i|}) & c = c_j \\ \frac{1}{2}(r_\epsilon^{|c - c_j|} + r_\epsilon^{|c - c_i|}) & c \neq c_i, c \neq c_j \\ (\frac{1}{2})(1 - r) & c = c_i, c_j = X \end{cases}$$

363 Similarly, the likelihood of a read with a lower bound c on the RU count is given by:

$$\Pr(\text{RU} \geq c | (c_i, c_j)) = \begin{cases} (1 - r) & c \leq c_i \\ \frac{1}{2}(1 - r) & c_i < c \leq c_j \\ r & c > c_j \end{cases}$$

364 The likelihood of the data C is given by $\prod_{c_k \in C} \Pr(c_k | (c_i, c_j))$. The posterior genotype probabilities

365 can be computed using Bayes' theorem:

$$\Pr((c_i, c_j)|C) = \frac{\Pr(C|(c_i, c_j)) \Pr((c_i, c_j))}{\sum_{c_i', c_j' \in C} \Pr(C|(c_i', c_j')) \Pr((c_i', c_j'))} \quad (3)$$

366 We generally set equal priors. However, in the event that we only see reads with a single count c' ,
367 we choose $\Pr((c', c')) = \Pr((c', X)) = \frac{1}{2}$. The probability of "missing haplotype" event is modeled
368 as a Bernoulli process since in genome sequencing, sampling from either chromosome is done at
369 random and so, the probability of not observing a haplotype in each read (failure) is $1/2$. If we see
370 multiple counts, we set $\Pr((c', X)) = 0$ for all $c' \in C$, and give equal priors to all other genotypes.

371 **VNTR Mutation Detection.** It is not difficult to see that alignment based methods do not
372 work well in VNTRs. Changes in RU counts make it difficult to align reads even for mappers
373 that allow split-reads, as the gaps in different reads can be placed in different locations. A similar
374 problem appears with small indels, as there are multiple ways to align reads with an indel in a
375 Repeat Unit. The adVNTR HMM aligns all repeat units to the same HMM, and this has the effect
376 of aligning all mutations/indels in the same column. Consider the case where reads contain a total
377 of v nucleotides matching a VNTR RU of length ℓ , and RU count u . Moreover at a specific position
378 covered by d Repeats, suppose we observe ι indel transitions.

379 For a true indel mutation, we expect $\frac{u\ell}{v}$ fraction of transitions at a location to be an indel,
380 giving a likelihood of the observed data as $\text{Binom}(d, \iota, \frac{u\ell}{v})$. Alternatively, for a homopolymer run
381 of $i > 0$ nucleotides, let ε_i denote the per-nucleotide indel error rate. We modeled ε_1 empirically in
382 non-VNTR, non-polymorphic regions and confirmed prior results that ε_i increases with increasing
383 i (Margulies et al., 2005). Thus, the likelihood of seeing ι indel transitions due to sequencing error
384 in a homopolymer run of length i is $\text{Binom}(d, \iota, \varepsilon_i)$. We scored an indel in the VNTR using the
385 log-likelihood ratio

$$-2 \ln \left(\frac{\text{Binomial}(d, \iota, \frac{u\ell}{v})}{\text{Binomial}(d, \iota, \varepsilon_i)} \right), \quad (4)$$

386 which follows a χ^2 distribution. We select the indel if the nominal p -value is lower than 0.01.

387 Command line usage of adVNTR for RU count genotyping and frameshift identification is
388 available in Supplementary Material "Running adVNTR"

389 **Acknowledgements.** The analyses presented in this paper are based on the use of study data
390 downloaded from the dbGaP web site, under phs001095.v1.p1, phs001096.v1.p1 and phs001097.v1.p1.

391 References

- 392 Altschul, S. F., Gish, W., Miller, W., Myers, E. W., and Lipman, D. J., 1990. Basic local alignment search
393 tool. *Journal of molecular biology*, **215**(3):403–410.
- 394 Au, K. F., Underwood, J. G., Lee, L., and Wong, W. H., 2012. Improving PacBio long read accuracy by
395 short read alignment. *PloS one*, **7**(10):e46679.
- 396 Benedetti, F., Dallaspezia, S., Colombo, C., Pirovano, A., Marino, E., and Smeraldi, E., 2008. A length
397 polymorphism in the circadian clock gene *Per3* influences age at onset of bipolar disorder. *Neuroscience*
398 *letters*, **445**(2):184–187.
- 399 Benson, G., 1999. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic acids research*,
400 **27**(2):573.
- 401 Berlin, K., Koren, S., Chin, C.-S., Drake, J. P., Landolin, J. M., and Phillippy, A. M., 2015. Assembling large
402 genomes with single-molecule sequencing and locality-sensitive hashing. *Nature biotechnology*, **33**(6):623–
403 630.
- 404 Brookes, K., 2013. The VNTR in complex disorders: The forgotten polymorphisms? A functional way
405 forward? *Genomics*, **101**(5):273–281.
- 406 Byrd, A. L. and Manuck, S. B., 2014. MAOA, childhood maltreatment, and antisocial behavior: meta-
407 analysis of a gene-environment interaction. *Biological psychiatry*, **75**(1):9–17.
- 408 Cervera, A., Tassies, D., Obach, V., Amaro, S., Reverter, J., and Chamorro, A., 2007. The BC genotype of
409 the VNTR polymorphism of platelet glycoprotein *Iba* is overrepresented in patients with recurrent stroke
410 regardless of aspirin therapy. *Cerebrovascular Diseases*, **24**(2-3):242–246.
- 411 Chaisson, M. J. and Tesler, G., 2012. Mapping single molecule sequencing reads using basic local alignment
412 with successive refinement (BLASR): application and theory. *BMC bioinformatics*, **13**(1):238.
- 413 Clarke, J., Wu, H.-C., Jayasinghe, L., Patel, A., Reid, S., and Bayley, H., 2009. Continuous base identification
414 for single-molecule nanopore DNA sequencing. *Nature nanotechnology*, **4**(4):265–270.
- 415 Consortium, . G. P., 2015. A global reference for human genetic variation. *Nature*, **526**(7571):68–74.
- 416 DePristo, M. A., Banks, E., Poplin, R., Garimella, K. V., Maguire, J. R., Hartl, C., Philippakis, A. A.,
417 Del Angel, G., Rivas, M. A., Hanna, M., *et al.*, 2011. A framework for variation discovery and genotyping
418 using next-generation DNA sequencing data. *Nature genetics*, **43**(5):491–498.

- 419 Dolzhenko, E., van Vugt, J. J., Shaw, R. J., Bekritsky, M. A., van Blitterswijk, M., Narzisi, G., Ajay, S. S.,
420 Rajan, V., Lajoie, B., Johnson, N. H., *et al.*, 2017. Detection of long repeat expansions from PCR-free
421 whole-genome sequence data. *Genome Research*, :gr-225672.
- 422 Durinovic-Belló, I., Wu, R., Gersuk, V., Sanda, S., Shilling, H., and Nepom, G., 2010. Insulin gene VNTR
423 genotype associates with frequency and phenotype of the autoimmune response to proinsulin. *Genes and*
424 *immunity*, **11**(2):188–193.
- 425 Eddy, S. R., 1996. Hidden markov models. *Current opinion in structural biology*, **6**(3):361–365.
- 426 Eddy, S. R. et al., 1995. Multiple alignment using hidden markov models. In *Ismb*, volume 3, pages 114–120.
- 427 Eid, J., Fehr, A., Gray, J., Luong, K., Lyle, J., Otto, G., Peluso, P., Rank, D., Baybayan, P., Bettman, B.,
428 *et al.*, 2009. Real-time DNA sequencing from single polymerase molecules. *Science*, **323**(5910):133–138.
- 429 Eser, O., Eser, B., Cosar, M., Erdogan, M., Aslan, A., Yildiz, H., Solak, M., and Haktanir, A., 2011. Short
430 aggrecan gene repetitive alleles associated with lumbar degenerative disc disease in Turkish patients. *Genet*
431 *Mol Res*, **10**(3):1923–1930.
- 432 Franke, B., Vasquez, A. A., Johansson, S., Hoogman, M., Romanos, J., Boreatti-Hümmer, A., Heine, M.,
433 Jacob, C. P., Lesch, K.-P., Casas, M., *et al.*, 2010. Multicenter analysis of the SLC6A3/DAT1 VNTR
434 haplotype in persistent ADHD suggests differential involvement of the gene in childhood and persistent
435 ADHD. *Neuropsychopharmacology*, **35**(3):656.
- 436 Fuchsberger, C., Flannick, J., Teslovich, T. M., Mahajan, A., Agarwala, V., Gaulton, K. J., Ma, C., Fontanil-
437 las, P., Moutsianas, L., McCarthy, D. J., *et al.*, 2016. The genetic architecture of type 2 diabetes. *Nature*,
438 **536**(7614):41–47.
- 439 Galimberti, D., Scarpini, E., Venturelli, E., Strobel, A., Herterich, S., Fenoglio, C., Guidi, I., Scalabrini, D.,
440 Cortini, F., Bresolin, N., *et al.*, 2008. Association of a NOS1 promoter repeat with Alzheimer’s disease.
441 *Neurobiology of aging*, **29**(9):1359–1365.
- 442 Gelfand, Y., Hernandez, Y., Loving, J., and Benson, G., 2014. VNTRseek-a computational tool to detect
443 tandem repeat variants in high-throughput sequencing data. *Nucleic acids research*, **42**(14):8884–8894.
- 444 Gymrek, M., 2016. Pcr-free library preparation greatly reduces stutter noise at short tandem repeats.
445 *bioRxiv*, :043448.

- 446 Gymrek, M., Willems, T., Guilmatre, A., Zeng, H., Markus, B., Georgiev, S., Daly, M. J., Price, A. L.,
447 Pritchard, J. K., Sharp, A. J., *et al.*, 2016. Abundant contribution of short tandem repeats to gene
448 expression variation in humans. *Nature genetics*, **48**(1):22–29.
- 449 Hackl, T., Hedrich, R., Schultz, J., and Förster, F., 2014. proovread: large-scale high-accuracy PacBio
450 correction through iterative short read consensus. *Bioinformatics*, **30**(21):3004–3011.
- 451 Haddley, K., Bubb, V., Breen, G., Parades-Esquivel, U., and Quinn, J., 2011. Behavioural genetics of the
452 serotonin transporter. In *Behavioral Neurogenetics*, pages 503–535. Springer.
- 453 Hijikata, M., Matsushita, I., Tanaka, G., Tsuchiya, T., Ito, H., Tokunaga, K., Ohashi, J., Homma, S.,
454 Kobashi, Y., Taguchi, Y., *et al.*, 2011. Molecular cloning of two novel mucin-like genes in the disease-
455 susceptibility locus for diffuse panbronchiolitis. *Human genetics*, **129**(2):117–128.
- 456 Huang, W., Li, L., Myers, J. R., and Marth, G. T., 2011. ART: a next-generation sequencing read simulator.
457 *Bioinformatics*, **28**(4):593–594.
- 458 Kent, W. J., 2002. Blatthe blast-like alignment tool. *Genome research*, **12**(4):656–664.
- 459 Kirby, A., Gnirke, A., Jaffe, D. B., Barešová, V., Pochet, N., Blumenstiel, B., Ye, C., Aird, D., Stevens,
460 C., Robinson, J. T., *et al.*, 2013. Mutations causing medullary cystic kidney disease type 1 lie in a large
461 VNTR in MUC1 missed by massively parallel sequencing. *Nature genetics*, **45**(3):299–303.
- 462 Kirchheiner, J., Nickchen, K., Sasse, J., Bauer, M., Roots, I., and Brockmüller, J., 2007. A 40-basepair
463 VNTR polymorphism in the dopamine transporter (DAT1) gene and the rapid response to antidepressant
464 treatment. *The pharmacogenomics journal*, **7**(1):48.
- 465 Krogh, A., Brown, M., Mian, I. S., Sjölander, K., and Haussler, D., 1994. Hidden markov models in
466 computational biology: Applications to protein modeling. *Journal of molecular biology*, **235**(5):1501–
467 1531.
- 468 LaHoste, G., Swanson, J., Wigal, S., Glabe, C., Wigal, T., King, N., and Kennedy, J., 1996. Dopamine D4
469 receptor gene polymorphism is associated with attention deficit hyperactivity disorder. *Mol Psychiatry*,
470 **1**(2):121–124.
- 471 Lalioti, M. D., Scott, H. S., Buresi, C., Rossier, C., Bottani, A., Morris, M. A., Malafosse, A., and An-
472 tonarakis, S. E., 1997. Dodecamer repeat expansion in cystatin B gene in progressive myoclonus epilepsy.
473 *Nature*, **386**(6627):847.

- 474 Lander, E. S., Linton, L. M., Birren, B., Nusbaum, C., Zody, M. C., Baldwin, J., Devon, K., Dewar, K.,
475 Doyle, M., FitzHugh, W., *et al.*, 2001. Initial sequencing and analysis of the human genome. *Nature*,
476 **409**(6822):860–921.
- 477 Langmead, B. and Salzberg, S. L., 2012. Fast gapped-read alignment with Bowtie 2. *Nature methods*,
478 **9**(4):357–359.
- 479 Lee, H., Gurtowski, J., Yoo, S., Marcus, S., McCombie, W. R., and Schatz, M., 2014. Error correction and
480 assembly complexity of single molecule sequencing reads. *BioRxiv*, :006395.
- 481 Lemmers, R. J., de Kievit, P., Sandkuijl, L., Padberg, G. W., van Ommen, G.-J. B., Frants, R. R., and
482 van der Maarel, S. M., 2002. Facioscapulohumeral muscular dystrophy is uniquely associated with one of
483 the two variants of the 4q subtelomere. *Nature genetics*, **32**(2):235.
- 484 Li, H., 2011. A statistical framework for SNP calling, mutation discovery, association mapping and population
485 genetical parameter estimation from sequencing data. *Bioinformatics*, **27**(21):2987–2993.
- 486 Li, H. and Durbin, R., 2009. Fast and accurate short read alignment with Burrows–Wheeler transform.
487 *Bioinformatics*, **25**(14):1754–1760.
- 488 Liu, Q., Zhang, P., Wang, D., Gu, W., and Wang, K., 2017. Interrogating the unsequenceable genomic
489 trinucleotide repeat disorders by long-read sequencing. *Genome medicine*, **9**(1):65.
- 490 Margulies, M., Egholm, M., Altman, W. E., Attiya, S., Bader, J. S., Bemben, L. A., Berka, J., Braverman,
491 M. S., Chen, Y.-J., Chen, Z., *et al.*, 2005. Genome sequencing in microfabricated high-density picolitre
492 reactors. *Nature*, **437**(7057):376–380.
- 493 Miclotte, G., Heydari, M., Demeester, P., Rombauts, S., Van de Peer, Y., Audenaert, P., and Fostier,
494 J., 2016. Jabba: hybrid error correction for long sequencing reads. *Algorithms for Molecular Biology*,
495 **11**(1):10.
- 496 Okazaki, S., Schirripa, M., Loupakis, F., Cao, S., Zhang, W., Yang, D., Ning, Y., Berger, M. D., Miyamoto,
497 Y., Suenaga, M., *et al.*, 2017. Tandem repeat variation near the HIC1 (hypermethylated in cancer 1) pro-
498 moter predicts outcome of oxaliplatin-based chemotherapy in patients with metastatic colorectal cancer.
499 *Cancer*, .
- 500 Orita, M., Iwahana, H., Kanazawa, H., Hayashi, K., and Sekiya, T., 1989. Detection of polymorphisms
501 of human dna by gel electrophoresis as single-strand conformation polymorphisms. *Proceedings of the*
502 *National Academy of Sciences*, **86**(8):2766–2770.

- 503 Pritchard, A. L., Pritchard, C. W., Bentham, P., and Lendon, C. L., 2007. Role of serotonin transporter
504 polymorphisms in the behavioural and psychological symptoms in probable Alzheimer disease patients.
505 *Dementia and geriatric cognitive disorders*, **24**(3):201–206.
- 506 Pugliese, A., Zeller, M., Fernandez, A., Zalberg, L. J., Bartlett, R. J., Ricordi, C., Pietropaolo, M., Eisen-
507 barth, G. S., Bennett, S. T., and Patel, D. D., *et al.*, 1997. The insulin gene is transcribed in the human
508 thymus and transcription levels correlate with allelic variation at the INS VNTR-IDD3 susceptibility
509 locus for type 1 diabetes. *Nature genetics*, **15**(3):293–297.
- 510 Ræder, H., Johansson, S., Holm, P. I., Haldorsen, I. S., Mas, E., Sbarra, V., Nerøen, I., Eide, S. Å., Grevle,
511 L., Bjørkhaug, L., *et al.*, 2006. Mutations in the CEL VNTR cause a syndrome of diabetes and pancreatic
512 exocrine dysfunction. *Nature genetics*, **38**(1):54.
- 513 Robinson, J. T., Thorvaldsdóttir, H., Winckler, W., Guttman, M., Lander, E. S., Getz, G., and Mesirov,
514 J. P., 2011. Integrative genomics viewer. *Nature biotechnology*, **29**(1):24–26.
- 515 Salmela, L. and Rivals, E., 2014. LoRDEC: accurate and efficient long read error correction. *Bioinformatics*,
516 **30**(24):3506–3514.
- 517 Salmela, L., Walve, R., Rivals, E., and Ukkonen, E., 2016. Accurate self-correction of errors in long reads
518 using de Bruijn graphs. *Bioinformatics*, **33**(6):799–806.
- 519 Shriver, M. D., Jin, L., Chakraborty, R., and Boerwinkle, E., 1993. VNTR allele frequency distributions
520 under the stepwise mutation model: a computer simulation approach. *Genetics*, **134**(3):983–993.
- 521 Stöcker, B. K., Köster, J., and Rahmann, S., 2016. SimLoRD: Simulation of Long Read Data. *Bioinformatics*,
522 **32**(17):2704–2706.
- 523 Tyner, C., Barber, G. P., Casper, J., Clawson, H., Diekhans, M., Eisenhart, C., Fischer, C. M., Gibson, D.,
524 Gonzalez, J. N., Guruvadoo, L., *et al.*, 2016. The UCSC Genome Browser database: 2017 update. *Nucleic
525 acids research*, **45**(D1):D626–D634.
- 526 Ummat, A. and Bashir, A., 2014. Resolving complex tandem repeats with long reads. *Bioinformatics*,
527 **30**(24):3491–3498.
- 528 Viswanath, B., Purushottam, M., Kandavel, T., Reddy, Y. J., Jain, S., *et al.*, 2013. DRD4 gene and
529 obsessive compulsive disorder: do symptom dimensions have specific genetic correlates? *Progress in
530 Neuro-Psychopharmacology and Biological Psychiatry*, **41**:18–23.

- 531 Wheeler, D. A., Srinivasan, M., Egholm, M., Shen, Y., Chen, L., McGuire, A., He, W., Chen, Y.-J.,
532 Makhijani, V., Roth, G. T., *et al.*, 2008. The complete genome of an individual by massively parallel
533 DNA sequencing. *nature*, **452**(7189):872–876.
- 534 Willems, T., Zielinski, D., Yuan, J., Gordon, A., Gymrek, M., and Erlich, Y., 2017. Genome-wide profiling
535 of heritable and de novo STR variations. *Nature Methods*, .
- 536 Worrall, B. B., Brott, T. G., Brown, R. D., Brown, W. M., Rich, S. S., Arepalli, S., Wavrant-De Vrièze, F.,
537 Duckworth, J., Singleton, A. B., Hardy, J., *et al.*, 2007. IL1RN VNTR polymorphism in ischemic stroke.
538 *Stroke*, **38**(4):1189–1196.
- 539 Wright, J. M., 1994. Mutation at vntrs: Are minisatellites the evolutionary progeny of microsatellites?
540 *Genome*, **37**(2):345–347.
- 541 Zook, J. M., Catoe, D., McDaniel, J., Vang, L., Spies, N., Sidow, A., Weng, Z., Liu, Y., Mason, C. E.,
542 Alexander, N., *et al.*, 2016. Extensive sequencing of seven human genomes to characterize benchmark
543 reference materials. *Scientific data*, **3**.

Supplementary Material

544

545

546 A. Model Structure and Parameter Setting

547 Each VNTR is represented by three Hidden Markov Models. A detailed sketch of the Repeat Match
548 HMM is shown in Fig. 5. Here, we show the structure of two other parts in Fig. S1. We repeated
549 the blue silent states (*Start*, *U_S*, *U_e*, and *End*) to show how these three models are connected.

To set the transition and emission probabilities of repeat matcher, we used the parameter obtained by pair HMM of repeating units in reference genome. We set pseudocounts equal to error rate of sequencing technology in all three HMMs to allow for mutations and sequencing errors. After the initialization of each model, we updated them using sequencing data of NA12878 (Table S2). To update each model, we ran read recruitment on sequencing data of NA12878 and extracted repeating units as described in Methods. Then, we aligned the repeating units to the HMM, and used the new aligned reads to update HMM parameters. We measure fitness of model by the sum of log-likelihood of the recruited reads, as follows:

$$\text{fitness} = \sum_{r \in \text{reads}} \log(\text{likelihood}(r)),$$

550 where likelihood of read r is defined as the probability of most likely path in the HMM to emit
551 r . We continued to iterate the model alignment, and parameter update steps until convergence of
552 fitness values.

As described in Methods, we compute the likelihood using the Viterbi algorithm. Let $V_{k,j}$ denote the highest (log) probability of emitting the first k letters of the sequence s_1, s_2, \dots, s_n and ending in state j of an HMM. Let, $\text{Prev}_{k,j}$ denote the state j' immediately prior to j in this optimum parse. Then,

$$V_{k,j} = \max_{j'} \{V_{k',j} + \log T(j', j) + \log E_j(s_k)\},$$
$$\text{Prev}_{k,j} = \arg \max_{j'} \{V_{k',j} + \log T(j', j) + \log E_j(s_k)\},$$

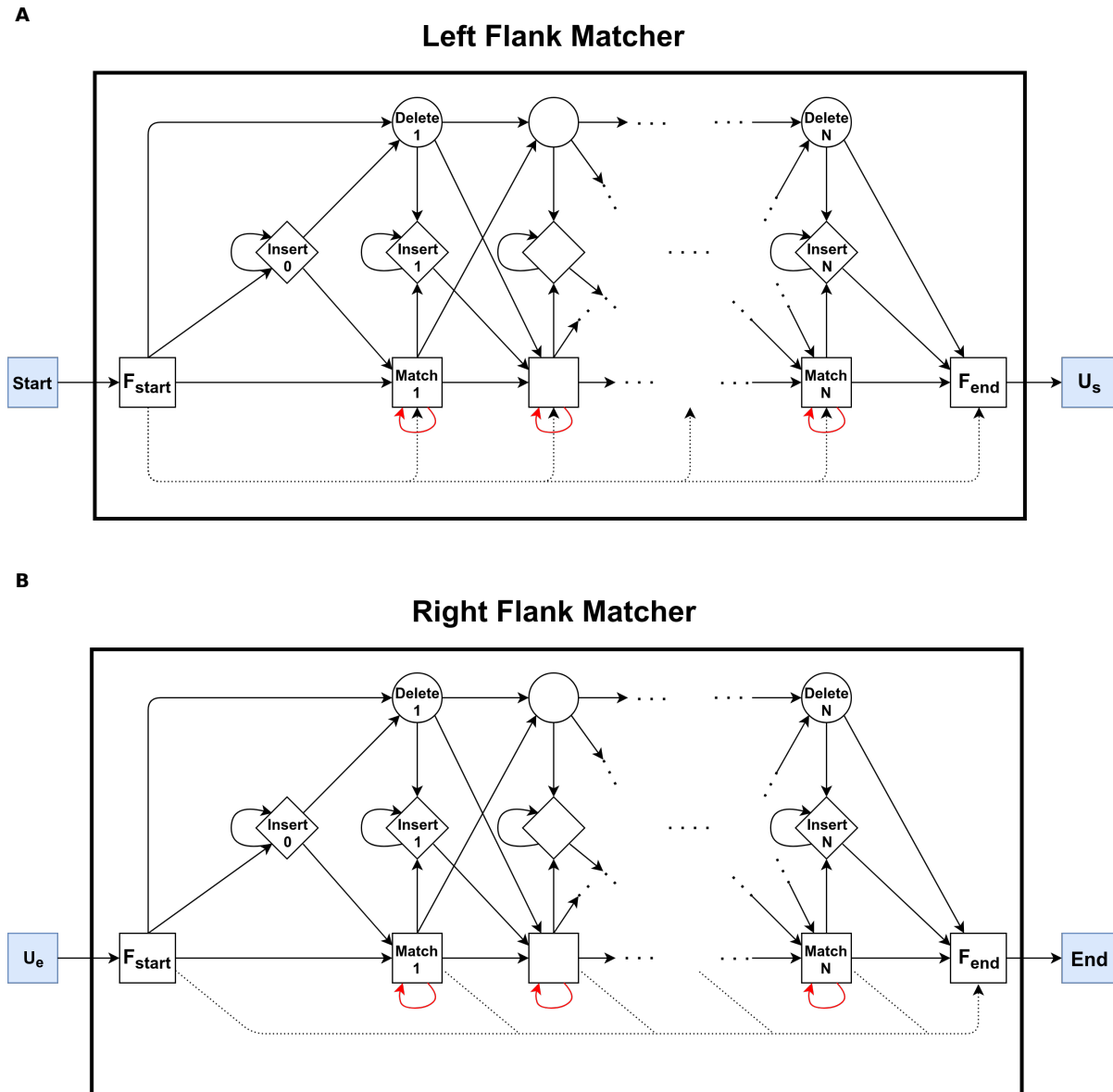


Figure S1: **Flanking region matchers.** (A) Shows the structure of Left Flank Matcher, which matches a suffix of left flanking region of the VNTR. In this part, the dotted edges allows skipping of adjacent base pairs at the beginning of the flanking region, and the rest of region (base pairs on the right) should be matched to the states and this is how matching of a suffix is insured. (B) Shows the structure of Right Flank Matcher, the model that matches a prefix of right flanking region of the VNTR. Here, dotted edges ensure the matching of a prefix of the flanking region sequence.

553 where, $k' = k - 1$ for match or insert states; $k' = k$ otherwise. Then, for a read sequence r with
 554 length n , $\max_j V_{n,j}$ over all states j in the HMM determines the maximum likelihood.

555 B. Selecting Target VNTRs

556 We selected sets of target models that could be analyzed based on their characteristics and the
557 sequencing technologies as follows: We started with the human VNTR list created by Tandem
558 Repeat Finder. To select the most important loci, we considered VNTRs that had an intersection
559 with coding regions of human genome. Next, we excluded cases where the flanking regions of
560 VNTR were not known (e.g. VNTR is close to telomere; the flanking region doesn't exist in
561 reference genome; and there is a sequence of 'N' adjacent to the VNTR.). Finally, we added
562 17 VNTRs that are in promoter or intron of the genes but are known to be linked to a disease
563 (Table 1). We removed VNTRs that appear multiple times in different loci of the genome with
564 identical patterns and flanking regions, but with different number of copies. To find such similar
565 VNTRs, we compared each pair of VNTRs by comparing the flanking regions and repeating unit
566 with BLAT (Kent, 2002) and eliminating the VNTRs if their similarity was higher than 75%.

567 This procedure resulted in 2944 'coding' VNTRs out of 3147 VNTRs that intersected with
568 coding regions of human genome. The 2944 VNTRs were used for PacBio analysis. For Illumia
569 analysis, we used a subset of 1775 VNTRs of the 2944, whose length was shorter than 140bp.
570 Finally to create a difficult test case for testing frame-shifts, we selected 115 of 2944 VNTRs for
571 which the total length was ≥ 250 bp, and all Repeat Units had the same length, and used those to
572 simulate indel (frameshift) data-sets.

573 C. Test Datasets

574 Multiple test cases were generated using the three lists containing 2944, 1775, and 115 VNTRs,
575 respectively as described in the previous section. We started by generating a distinct human
576 genomic sequence `VNTR_I_X_reference.fa` for each $I \in [1, 2944]$ and each value $X \in [-3, 3]$ (20,608
577 total sequences). Each sequence `VNTR_I_X_reference.fa` was identical to the human reference except
578 that it had X ' copies for I -th VNTR, where X ' takes the RU count in reference genome $\pm X$. To
579 increase the RU count of a VNTR, we added the repeating units from the first repeat to the last
580 unit, one at a time. We additionally generated ~ 4920 reference sequences `VNTR_I_Deletion_P.fa`
581 and `VNTR_I_Insertion_P.fa` for all $I \in [1, 115]$ VNTRs indexing the third list, and a single insertion
582 or deletion at the P th base pair of the I th VNTR. We set P to every position in the VNTR that

583 was a multiple of 10 and was at least 140bp apart from each side of the VNTR. These reference
584 templates were used for generating simulated datasets as follows:

585 **IlluminaSim Dataset.** We used the following command to simulate the reads from haplotypes
586 using ART:

```
587 art_illumina -ss HSXt -sam -i VNTR.I.X.reference.fa -l 150 -f 15 -o VNTR.I.X.set
```

588 Then, we merged every pair of haploid datasets with RU counts X and Y to get diploid
589 sequencing data with genotype (X,Y) for VNTR *I* by appending VNTR.I.X.set.fq to the end
590 of VNTR.I.Y.set.fq to get VNTR.I.XY.set.fq. Then, we aligned these diploid reads to the
591 reference genome using Bowtie2 as follows:

```
592 bowtie2 -x hg19.bowtie2_index -U VNTR.I.XY.set.fq -S VNTR.I.XY.aln.sam
```

593 **PacBioSim Dataset.** We used the following command to simulated the reads for *I*th VNTR
594 using SimLoRD:

```
595 simlord -rr VNTR.I.X.reference.fa -pi 0.12 -pd 0.02 -ps 0.02 -c 15 VNTR.I.X.pb.set.sam
```

596 Next, we merged each pair of reads (fastq files) to get the diploid set of reads at 30X coverage.

597 **PacBioLong Dataset.** The dataset is similar to PacBioSim but with higher RU counts for 3
598 VNTRs 120, 40, and 25 for VNTRs in *INS*, *CSTB*, and *HIC1* genes, which represent the
599 largest expansion known for these VNTRs. Again, we used SimLord to generate reads.

```
600 simlord -rr VNTR.I.X.reference.fa -pi 0.12 -pd 0.02 -ps 0.02 -c 30 VNTR.I.X.pb.set.sam
```

601 **PacBio Coverage Dataset.** We simulated different levels of coverage for the three VNTRs using:

```
602 simlord -rr VNTR.I.X.reference.fa -pi 0.12 -pd 0.02 -ps 0.02 -c C VNTR.I.X.C.set
```

603 Here, $1 \leq C \leq 40\times$.

604 **IlluminaFrameshift Dataset.** We simulated these datasets using following commands:

```
605 art_illumina -ss HSXt -sam -i VNTR.I.Deletion.P.fa -l 150 -f 15 -o VNTR.I.Deletion_p
```

```
606 art_illumina -ss HSXt -sam -i VNTR.I.Insertion.P.fa -l 150 -f 15 -o VNTR.I.Insertion_p
```

607 We also simulated reads from reference genome without the frameshift:

```
608 art_illumina -ss HSXt -sam -i hg19.fa -l 150 -f 15 -o normal.haplotype
```

609 Finally, we merged fastq read files of a haplotype with frameshift with that of normal hap-
 610 lotype to get the diploid sample at 30X coverage and aligned the reads with Bowtie2 similar
 611 to “IlluminaSim Dataset”.

Dataset Name	Profile	Depth	# of VNTRs
Illumina Genotyping Dataset	HiSeqX TruSeq	30X	1775
PacBio Genotyping Dataset	PacBio	30X	2944
PacBio Long Expansion Dataset	PacBio	30X	3
PacBio Coverage Dataset	PacBio	$1 \leq C \leq 40$	3
Frameshift Dataset	HiSeqX TruSeq	$C \in \{10, 20, 30, 40\}$	123

Table S1: **Simulated dataset summary.**

612 WGS data used for testing was taken from Genome in a Bottle, NCBI short read archive,
 Polaris, while exome data was obtained from GoT2D. See Table S2

Samples	Study	Profile	PCR free	Depth	Access
AJ Child (NA24385)	GIAB	PacBio	-	70X	http://jimb.stanford.edu/giab-resources
AJ Father (NA24149)	GIAB	PacBio	-	30X	http://jimb.stanford.edu/giab-resources
AJ Mother (NA24143)	GIAB	PacBio	-	30X	http://jimb.stanford.edu/giab-resources
Chinese Child (HG00514)	PRJEB12236	PacBio	-	70X	ncbi.nlm.nih.gov/sra/ERX1322863
Chinese Father (HG00512)	PRJEB12236	PacBio	-	35X	ncbi.nlm.nih.gov/sra/ERX1322861
Chinese Mother (HG00513)	PRJEB12236	PacBio	-	35X	ncbi.nlm.nih.gov/sra/ERX1322862
AJ Child (NA24385)	GIAB	HiSeq 2500	Y	40X	http://jimb.stanford.edu/giab-resources
AJ Father (NA24149)	GIAB	HiSeq 2500	Y	40X	http://jimb.stanford.edu/giab-resources
AJ Mother (NA24143)	GIAB	HiSeq 2500	Y	40X	http://jimb.stanford.edu/giab-resources
NA12878	GIAB	PacBio	-	70X	http://jimb.stanford.edu/giab-resources
NA12878	GIAB	HiSeq 2500	Y	30X	http://jimb.stanford.edu/giab-resources
150 Individuals of 1KGP	Polaris	HiSeq X	Y	30-40X	ebi.ac.uk/ena/data/view/PRJEB20654
Diabetes WES data	GoT2D	HiSeq 2000	N	82X	dbGaP: phs001095, phs001096, and phs001097

Table S2: **Real sequencing data used in tests.**

613

614 D. Running adVNTR

615 adVNTR is available at <https://github.com/mehrdadbakhtiari/adVNTR>. As stated in the repos-
 616 itory, the best way to install it is to use conda package manager and running `conda install`
 617 `advntr`. After installation, `advntr` command invokes the program with four possible commands
 618 `genotype`, `addmodel`, `viewmodel`, and `delmodel`. Detail of each command as well as complete tu-
 619 torial on installation and usage are available at <http://advntr.readthedocs.io/>. Also, passing
 620 `-h` argument to each command will show the correct command line usage of the command.

621 **E. VNTRseek**

622 In order to make a call on a VNTR, VNTRseek requires both ends to be anchored with a minimum
623 of 20bp on each side of VNTR. This limits the length of VNTRs that can be identified using VN-
624 TRseek is limited to 110bp using Illumina sequencing technology. Also, it compares each VNTR
625 in the sequencing reads to every VNTR in reference genome which makes the process computa-
626 tionally demanding, and inaccessible for large data-sets. For these reasons, extensive VNTRseek
627 comparisons were not conducted.

628 F. Supplementary Figures and Tables

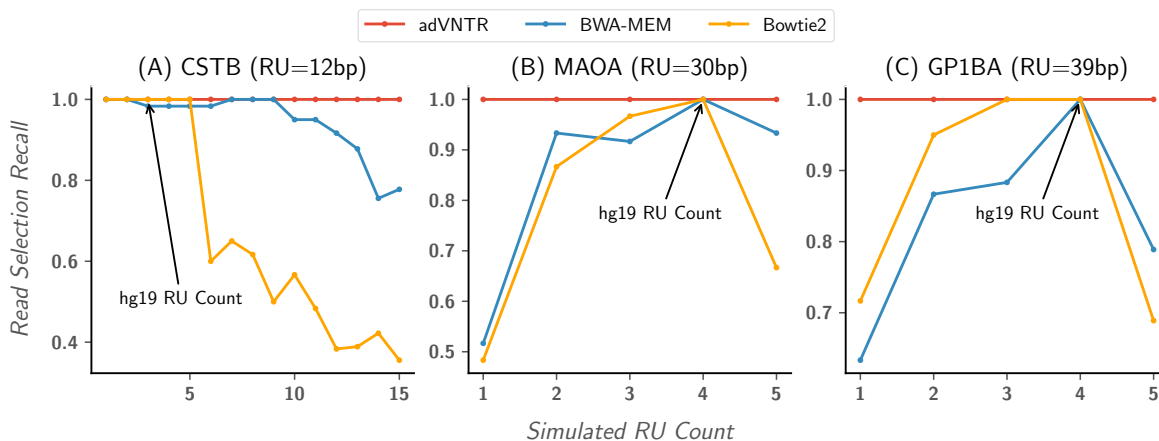


Figure S2: Sensitivity of Illumina read recruitment at specific VNTR loci. Comparison of adVNTR read selection with BWA-MEM and Bowtie2 mapping for Illumina reads (short VNTRs). Each plot shows the sensitivity of mapped/selected reads as a function of the number of repeats for different VNTRs. These plots show examples of alignment tools' behavior when RU count of VNTR deviates from the RU count in the reference genome. (A) Shows the comparison for the VNTR in *CSTB* gene, in which the pathogenic cases have more than 12 repeats and as it is shown alignment tools perform poorly in those cases. (B) Shows the comparison for the VNTR in *MAOA* gene, where the 4 repeats corresponds to both pathogenic case and number of repeats in reference genome. However, other tools perform poorly in normal cases. (C) Shows the comparison for the VNTR in *GP1BA* gene, and again, alignment tools only perform well when RU count is same as RU count in reference genome.

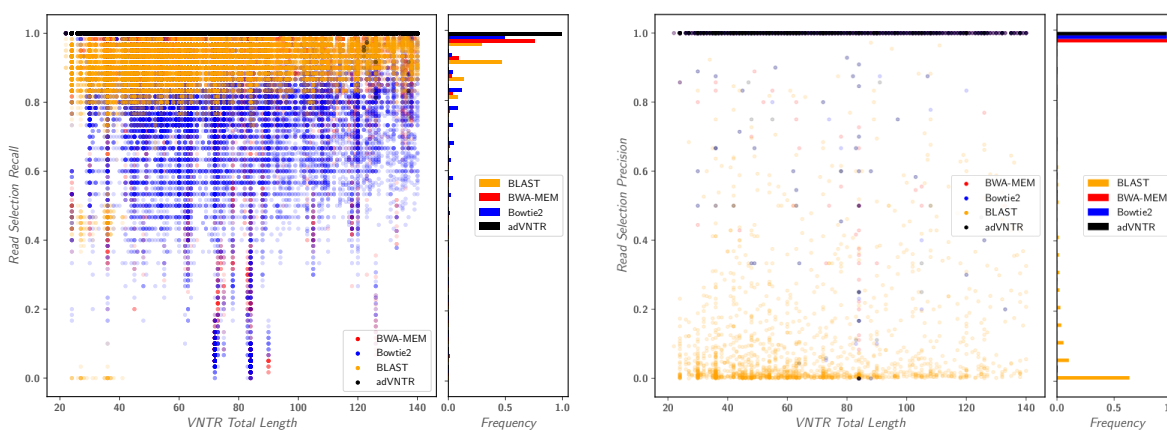


Figure S3: Read recruitment quality on Illumina reads. (A) Shows the comparison of the recall of adVNTR read recruitment with BWA-MEM, Bowtie2, and BLAST. (B) Shows the precision for read recruitment. These figures show that adVNTR has much higher recall compare to standard alignment tools without losing precision.

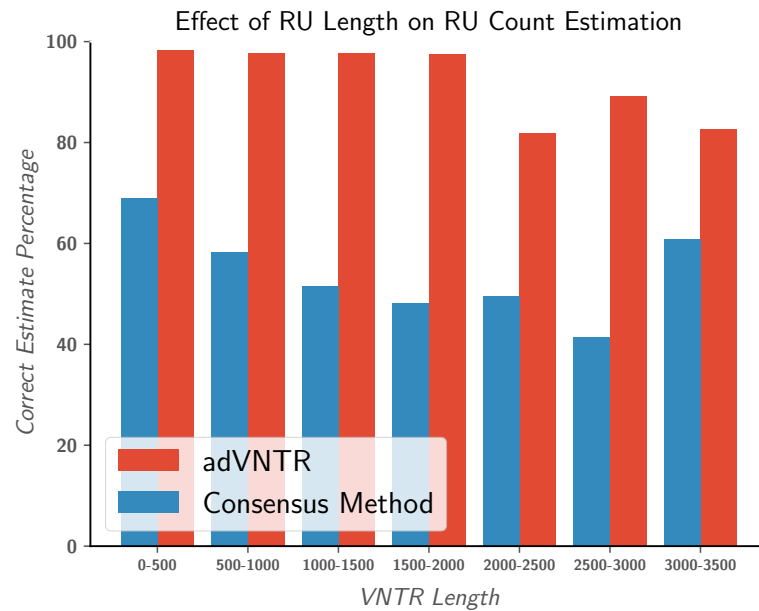


Figure S4: Comparison of adVNTR genotyping with consensus method on homozygous simulated data. adVNTR and consensus method comparison on homozygous testcases in *PacBioSim*.

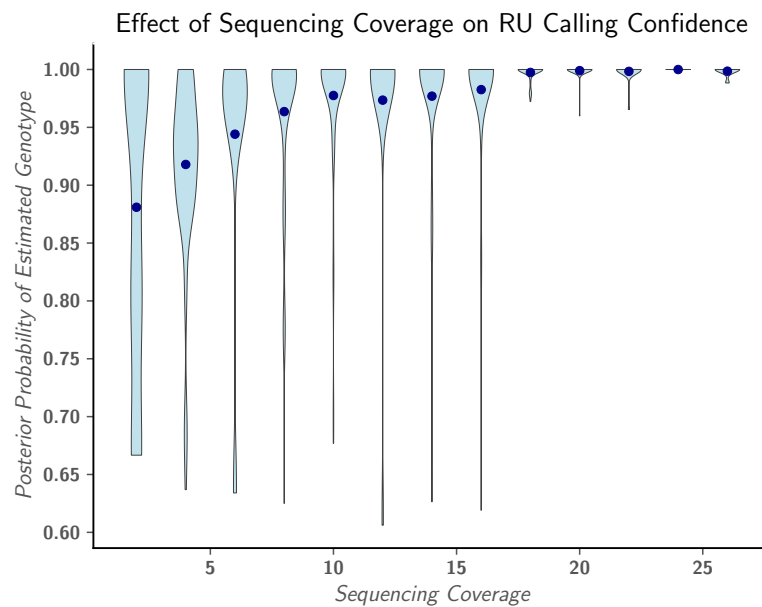


Figure S5: Association of PacBio sequencing coverage in VNTR region and posterior probability of RU count calling. The figure shows posterior probability of RU count estimation in AJ trio sequencing data from GIAB. Most of calls with low posterior probability (low confidence calls) result from low coverage in VNTR region. With at least 10 reads that span the VNTR, we will get 0.98 posterior probability for estimated genotype.

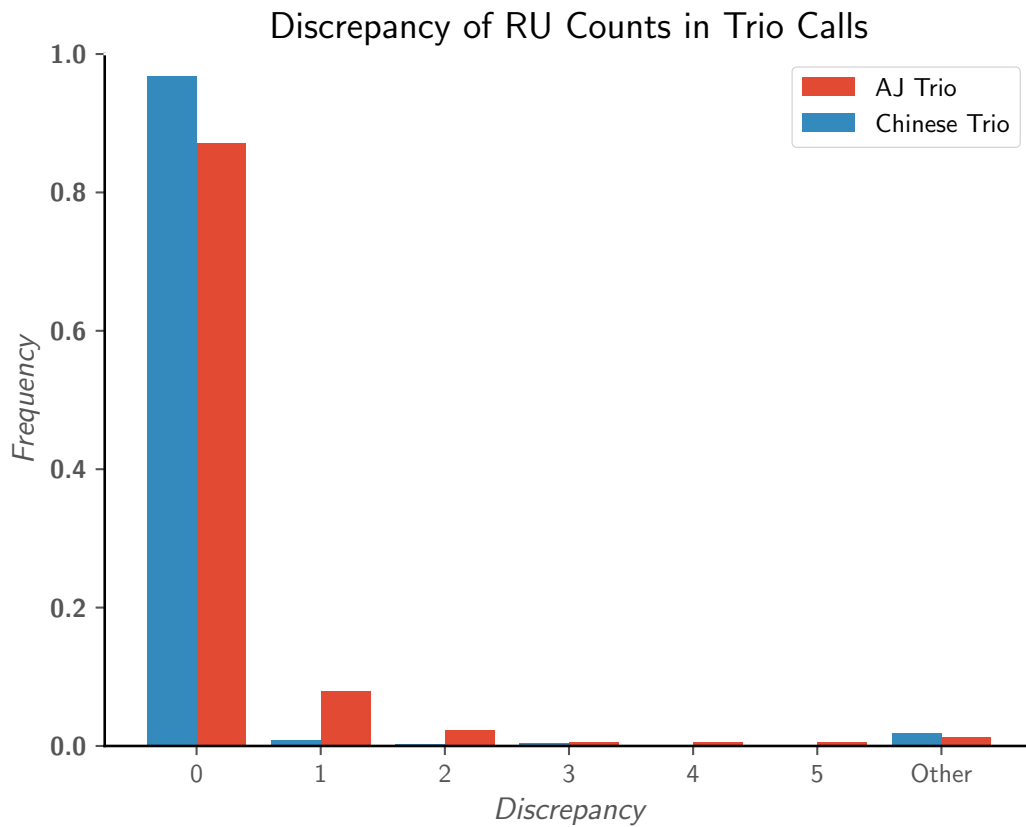


Figure S6: Distribution of discrepancies on trio calls using PacBio reads. This figure shows the distribution of discrepancies in adVNTR estimates on AJ and Chinese trios. As shown in the figure, most of non consistent calls in AJ trio have one discrepancy in estimated RU counts.

Gene	Locus (hg19)	Forward primer	Reverse primer	hg19 product length	Long reads
MAOA	chrX:43514348-43514468	GGCTACACCCACG TCTACTC	CACTCTTGGAGTC GGAGTCA	679	Y
IL1RN	chr2:113888105-113888449	ATTCCTGTCCTGG TAGTCTCC	AGAGGGGAGGGTC AGGTTAAT	701	Y
GP1BA	chr17:4837118-4837278	AGGACTGTGGTCA AGTTCCC	GCTTTGGTGGCTG ATCAAGT	586	Y
DRD4	chr11:639988-640180	CCGTGTGCTCCTT CTTCCTA	GACAGGAACCCAC CGACC	481	Y
SLC6A4	chr17:28564157-28564483	AGGGACTGAGCTG GACAAC	AGGCAGCAGACAA CTGTGTT	632	Y
JAKMIP3	chr10:133954073-133954190	CAAACAGACAGGA CGGACC	GTGCCCAGTCAG CTATCA	249	N
SRSF8	chr11:94800727-94800790	CAGGTGGCGCGCT ATG	GAGACCGGCTATA GCGAGAA	214	N
SSTR1	chr14:38679763-38679811	CGTCTTCCGTAAT GGCACCT	CCCTGGATAACCGT CCCTTT	153	N
C14orf180	chr14:105055118-105055145	CCTATACTGCGGC CGGG	CCTAGTTAGCCCT CAGGCAG	265	N
EIF3G	chr19:10229726-10229768	GGCAGAAGGGGAA AAACAGA	AGCTGACTCCTCC TTCCTAC	247	N
STK39	chr2:169103796-169103845	AACTGTTGAAGCC AGTAGGC	AGTTTCAAGTGGA AGGTCGT	408	N
BRWD1	chr21:40585353-40585415	TGCCCTATTTGTT CATTGGACT	TCCTTGCCAACAA GTCACTAC	249	N
CSTB	chr21:45196323-45196359	GAGGCACTTTGGC TTCGGA	GCGCCCGAAAGA CGATA	193	N
UBXN11	chr1:26608801-26608909	GCCTTTCCTACGT GCCTG	AGATCTTCAGCAC ATTCCCG	321	N
CLCA4	chr1:87045895-87045932	CTCAGAAGAAAAT GCAACCCAC	CACAGACAATACC AGCGTAGA	214	N
LCE4A	chr1:152681679-152681727	ATCCCCAAGTATC CCCCAAA	GACCTATGGTGTC TGTGGTG	152	N
PAOX	chr10:135202324-135202464	CAGTGGTTCCCTG CTGAGAA	GGCAATGAACCCA CAGAGAA	214	N

Table S3: **Primer for gel electrophoresis.** Last column shows whether we used the primers were used for a long range PCR. We used long range PCR to validate adVNTR calls on longer VNTRs (using PacBio reads).

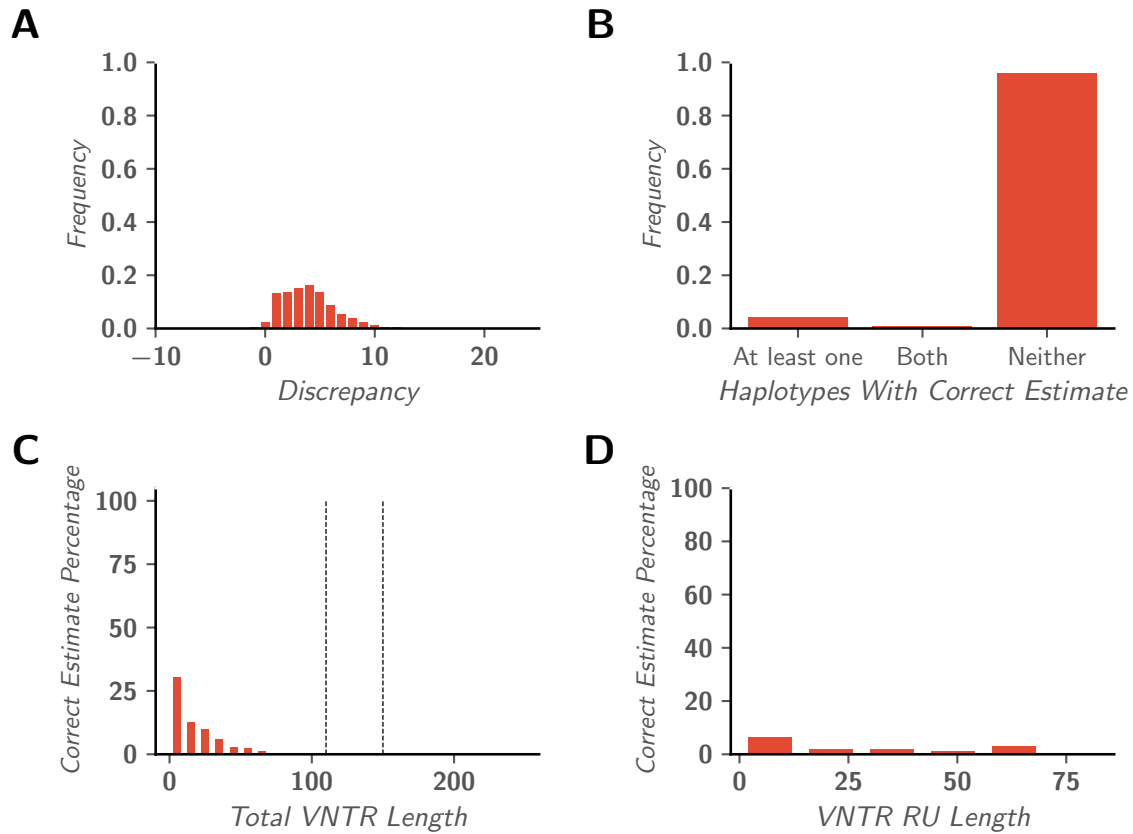


Figure S7: Expansion Hunter's performance on VNTR genotyping using Illumina reads. Expansion Hunter's performance on IlluminaSim dataset.

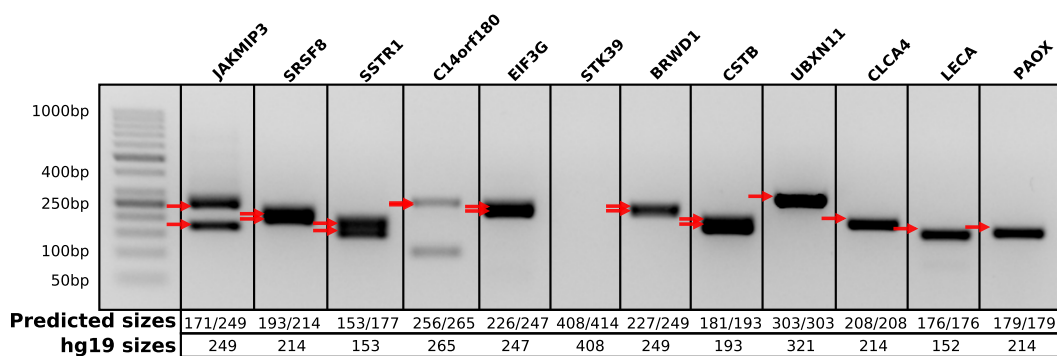


Figure S8: Validation of adVNTR genotyping on short VNTRs. In experiment for *C14orf180* the primers were repeated in another region of genome which resulted in having extra band. Even with zero copy of VNTR patterns, the distance of primers around VNTR is 238bp which means the extra band (~100bp) is resulted from another region of genome. Also, PCR amplification failed for *STK39* and no band is visible. Results of all other 10 experiments are consistent with adVNTR's estimates.

VNTR	Simulated Genotype	RU Count Discrepancy		
		PacBio Dataset		llumina Dataset
		adVNTR	Expansion Hunter	adVNTR
MAOA	1/1	0/0	-/-	0/0
MAOA	1/2	0/0	0/-1	0/0
MAOA	1/3	0/0	0/-2	0/0
MAOA	1/4	0/0	0/-3	0/0
MAOA	1/5	0/0	0/-4	0/0
MAOA	2/2	0/0	-1/-1	0/0
MAOA	2/3	0/0	0/-1	0/0
MAOA	2/4	0/0	-1/-3	0/0
MAOA	2/5	0/0	-1/-4	0/0
MAOA	3/3	0/0	-2/-2	0/0
MAOA	3/4	0/0	-2/-3	0/0
MAOA	3/5	0/0	-2/-4	0/0
MAOA	4/4	0/0	-3/-3	0/0
MAOA	4/5	0/0	-3/-4	0/0
MAOA	5/5	0/0	-4/-4	-1/-1
GP1BA	1/1	0/0	0/0	0/0
GP1BA	1/2	0/0	0/0	0/0
GP1BA	1/3	0/0	0/-1	0/0
GP1BA	1/4	0/0	1/-2	0/-1
GP1BA	2/2	0/0	0/0	0/0
GP1BA	2/3	0/0	0/-1	0/0
GP1BA	2/4	0/0	0/-2	0/-1
GP1BA	3/3	0/0	-1/-1	0/0
GP1BA	3/4	0/0	-1/-2	0/0
GP1BA	4/4	0/0	-2/-2	-1/0
CSTB	1/1	0/0	-/-	0/0
CSTB	1/2	0/0	1/0	0/0
CSTB	1/3	0/0	2/0	0/0
CSTB	1/4	0/0	3/0	0/0
CSTB	1/5	0/0	4/0	0/0
CSTB	1/6	0/0	4/-1	0/0
CSTB	1/7	0/0	3/-3	0/0
CSTB	1/8	0/0	4/-3	0/0
CSTB	1/9	0/0	3/-5	0/0
CSTB	1/10	0/0	4/-5	0/0
CSTB	1/11	0/0	4/-6	0/0
CSTB	1/12	0/0	4/-7	0/0
CSTB	1/13	0/0	4/-8	0/0
CSTB	1/14	0/0	3/-10	0/-1
CSTB	2/2	0/0	0/0	0/0
CSTB	2/3	0/0	1/0	0/0
CSTB	2/4	0/0	1/-1	0/0
CSTB	2/6	0/0	3/-1	0/0
CSTB	2/8	0/0	3/-3	0/0
CSTB	2/10	0/0	3/-5	0/0
CSTB	2/12	0/0	3/-7	0/0
CSTB	2/14	0/0	3/-9	0/-1
CSTB	3/3	0/0	-1/-1	0/0
CSTB	3/4	0/0	2/1	0/0
CSTB	3/6	0/0	2/-1	0/0
CSTB	3/8	0/0	2/-3	0/0
CSTB	3/10	0/0	2/-5	0/0

Table S4: **RU count genotyping results on simulated data.** For two cases, (MAOA 1/1 and CSTB 1/1) Expansion Hunter doesn't find any RU count.

VNTR	Estimated Genotype					
	adVNTR			ExpansionHunter		
AJ Child	AJ Mother	AJ Father	AJ Child	AJ Mother	AJ Father	
DRD4	4/5	4/5	4/4	-/-	-/-	-/-
ZFH3	4/4	4/4	4/4	3/3	-/-	3/3
GP1BA	2/5	2/3	3/4	2/2	1/1	2/2
SLC6A4	13/13	11/13	13/13	-/-	-/-	-/-
MMP9	3/3	3/3	3/3	-/-	-/-	-/-
CSTB	2/2	2/2	2/2	3/3	2/2	1/1
MAOA	5/5	4/5	4/4	-/-	-/-	-/-

Table S5: **Genotyping comparison on AJ trio using Illumina reads from GIAB.** Table shows the genotype found by adVNTR and ExpansionHunter in disease causing VNTRs that are shorter than Illumina reads. -/- denotes ExpansionHunter has not found any genotype for the VNTR. It worths mentioning the genotypes found by adVNTR for *MAOA* are not inconsistent as this VNTR is located on ChrX and the son has haploid RU counts inherited from mother.

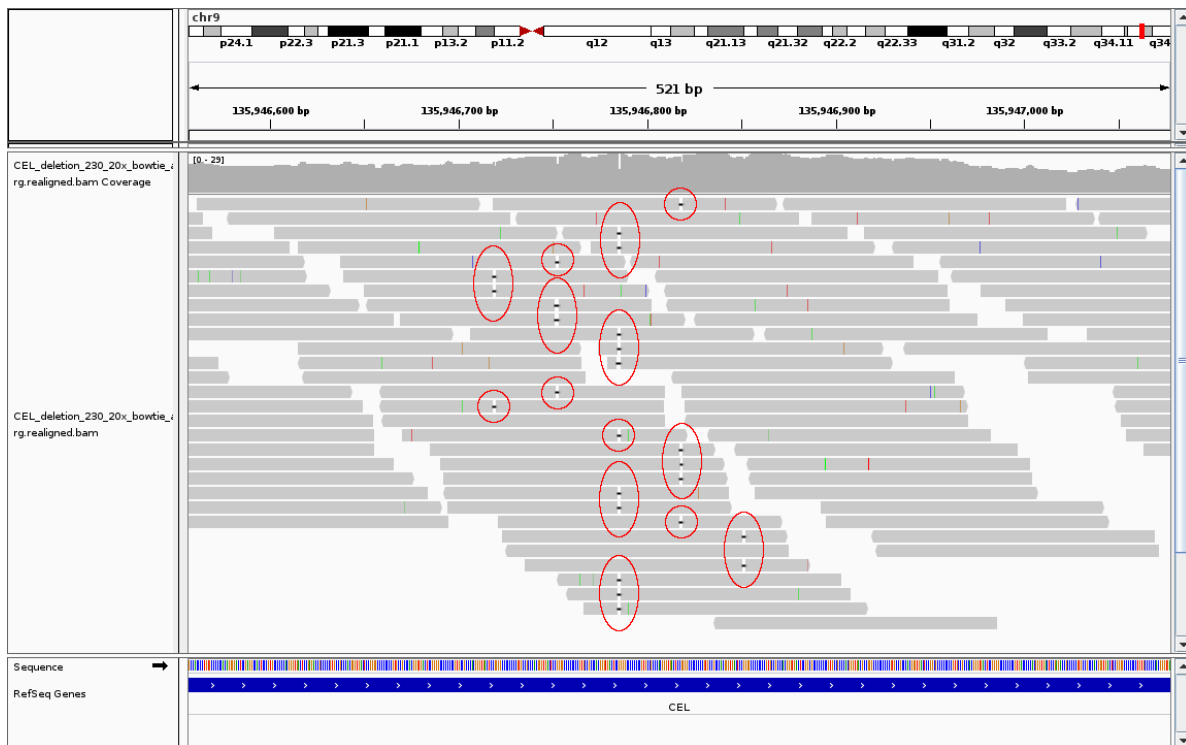


Figure S9: **Alignment stats with frameshift.** Alignment of a simulated data after running GATK IndelRealigner, when there is a deletion. With a sequencing mean of 30X, 25 reads contain the deletion but even after running realigner, deletions are mapped to five different repeating units.

Read1	GGCCACCCCTGTG-CCCCACAGGGGACTCCGA
Read2	GGCCACCCCTGTG-CCCCACAGGGGACTCCGA
Read3	GGCCACCCCTGTG-CCCCACAGGGGACTCCGA
Read4	GGCCACCCCTGTG-CCCCACAGGGGACTCCGA
Read5	GGCCACCCCTGTG-CCCCACAGGGGACTCCGA
Read6	GGCCACCCCTGTG-CCCCACAGGGGACTCCGA
ReferenceRepeatingUnit	GGCCACCCCTGTGCCCCACAGGGGACTCCGA *****

Figure S10: **Frameshift in CEL gene.** Multiple alignment of sequenced reads and reference repeating unit shows a deletion in diabetes patient genome. Due to low PCR amplification in GC rich VNTR region (84.8%), the coverage of VNTR region is 14X and 6 reads support the deletion.

		# of Samples	# of samples that frameshift has been identified		
			Samtools	Our Method	GATK
10X	Insertions	20	0	20	0
	Deletions	20	0	20	0
20X	Insertions	20	0	20	0
	Deletions	20	0	20	0
30X	Insertions	20	0	20	0
	Deletions	20	0	20	0
40X	Insertions	20	0	20	0
	Deletions	20	0	20	0

Table S6: Comparison of indel finding with Samtools and GATK