

Genotyping Variable Number Tandem Repeats in Human Genome

Mehrdad Bakhtiari

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

March 20, 2019

Abstract

Whole Genome Sequencing is increasingly used to identify Mendelian variants in clinical pipelines. These pipelines focus on single nucleotide variants (SNVs) and also structural variants, while ignoring more complex repeat sequence variants. We consider the problem of genotyping *Variable Number Tandem Repeats* (VNTRs), composed of inexact tandem duplications of short (6-100bp) repeating units. We show that VNTRs span 3% of the human genome, and are frequently present in coding regions. We present a survey of their implications in multiple Mendelian disorders. While existing tools recognize VNTR carrying sequence, genotyping VNTRs (determining repeat unit count and sequence variation) from high throughput sequenced reads remains challenging. We study this challenges and existing methods for identifying these loci in human reference genome, and describe a method, adVNTR, that uses Hidden Markov Models to model each VNTR, count repeat units, and detect sequence variation. We demonstrate the performance of the method on short-read (Illumina) and single molecule (PacBio) whole genome sequencing data using both simulations and real data.

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

1 Introduction

Next Generation Sequencing (NGS) is increasingly used to identify disease causing variants in clinical and diagnostic settings, but variant detection pipelines focus primarily on single nucleotide variants (SNVs) and small indels and to a lesser extent on structural variants. The human genome contains repeated sequences such as segmental duplications, short tandem repeats, and minisatellites (defined below) which pose challenges for alignment and variant calling tools. Hence, these regions are typically ignored during analysis of NGS data. In particular, *tandem repeats* correspond to locations where a short DNA sequence or *Repeat Unit* (RU) is repeated in tandem multiple times. RUs of 1-6bp are classified as Short Tandem Repeats (STRs), while longer RUs spanning potentially hundreds of nucleotides are denoted as *Variable Number Tandem Repeats* (VNTRs)(Shriver et al., 1993; Wright, 1994).

Disease relevance of VNTRs: VNTRs span 3% of the human genome and are often found in coding regions where the repeat unit length is a multiple of 3 resulting in tandem repeats in the amino acid sequence. More than 1,200 VNTRs with a RU length of 10 or greater exist in the coding regions of the human genome(Tyner et al., 2016). Compared to STRs, which have been extensively studied (Gymrek et al., 2016; Ummat and Bashir, 2014; Liu et al., 2017; Willems et al., 2017; Dolzhenko et al., 2017), VNTRs have not received as much attention. Nevertheless, multiple studies have linked variation in VNTRs with Mendelian diseases (*e.g.*, Medullary cystic kidney disease(Kirby et al., 2013), Myoclonus epilepsy(Lalioti et al., 1997), and FSHD(Lemmers et al., 2002)) and complex disorders such as bipolar disorder (Table 1). In some cases, the disease associated variants correspond to point mutations in the VNTR sequence (Kirby et al., 2013; Ræder et al., 2006) while in other cases, changes in the number of tandem repeats (RU count) show a statistical association (or causal relationship) with disease risk. For example, the insulin gene (INS) VNTR has an RU length of 14 bp with RU count varying from 26 to 200(Pugliese et al., 1997). Variation in this VNTR has been associated with expression of the INS gene and risk for type 1 diabetes (OR = 2.2) (Durinovic-Belló et al., 2010). Also, variation in the length of a VNTR in the third exon of DRD4 gene has been linked to numerous mental diseases, including ADHD and OCD(Viswanath et al., 2013). Notwithstanding these examples, the advent of genome-wide SNP genotyping arrays led to VNTRs being largely ignored. They have been called ‘the forgotten poly-

morphisms’(Brookes, 2013). Whole genome sequencing has the potential to detect and genotype all types of genetic variation, including VNTRs. However, computational identification of variation in VNTRs from next generation sequencing data remains challenging.

2 Challenges of VNTR genotyping

Lack of high throughput method for VNTR genotyping: VNTRs were originally used as markers for linkage mapping since they are highly polymorphic with respect to the number of tandem repeats at a given VNTR locus(Gelfand et al., 2014). Traditionally, VNTR genotyping required labor intensive and time consuming gel-based screens. In gel-based methods, a part of donor genome starting before the VNTR and finishing after that will be copied multiple times, and by reporting the average length of these copied fragments we will be able to estimate the number of repeats in the VNTR region. However, this approach limits the size of large population based studies of VNTRs.

Mapping difficulty in variant calling: Existing variant calling methods have been developed primarily to identify short sequence variants in unique DNA sequences that fall into a reference versus alternate allele framework, which is not well suited for detecting variation in VNTR sequences. As Fig.1 illustrates, for each read sequenced from any repeating unit in the VNTR region of the donor, there are multiple locations in the reference genome that result in the same mapping

Gene	Chr	Unit len	Number of units		Annotation	Disease
			Normal	Pathogenic		
<i>PER3</i>	1	54	4	5	coding	Bipolar disorderBenedetti et al. (2008)
<i>MUC1</i>	1	60	11-12	Insertion	coding	MCKD1Kirby et al. (2013)
<i>IL1RN</i>	2	86	3-6	2	intron	Stroke, CADWorrall et al. (2007)
<i>DUX4</i>	4	3.3kb	11-100	1-10		FSHDLemmers et al. (2002)
<i>DAT1</i>	5	44	7-11	10 (ADHD)	UTR	ADHD, Parkinson’s, BipolarFranke et al. (2010); Kirchheiner et al. (2007)
<i>MUC21</i>	6	45	26-27	4bp deletion	coding	Diffuse panbronchiolitis (DPB)Hijikata et al. (2011)
<i>CEL</i>	9	33	11-21	Deletion	coding	Monogenic diabetesRæder et al. (2006)
<i>INS</i>	11	14-15	26-200	26-44 (T1D)	promoter	T1D;T2DBrookes (2013); Pugliese et al. (1997); Durinovic-Belló et al. (2010)
<i>DRD4</i>	11	48	2-11	7	coding	OCD, ADHDLaHoste et al. (1996); Viswanath et al. (2013)
<i>ACAN</i>	15	57	27-33	13-25	coding	Osteochondritis dissecansEser et al. (2011)
<i>ZFHX3</i>	16	12	4-5		coding	Kawasaki
<i>GP1BA</i>	17	39	1-4	2/3 genotype	coding	ATF in StrokeCervera et al. (2007)
<i>SERT</i>	17	16-17	9/10/12		intron	BPSD, Alzheimer’sHaddley et al. (2011); Pritchard et al. (2007)
<i>SERT</i>	17	22	14	16 (OCD)	promoter	OCD,Anxiety, SchizophreniaHaddley et al. (2011)
<i>HIC1</i>	17	70	1-4	5+/5+	promoter	Metastatic Colorectal CancerOkazaki et al. (2017)
<i>MMP9</i>	20	12	5-6		coding	Kawasaki
<i>CSTB</i>	21	12	2-3	12+	5’UTR	Progressive myoclonic epilepsy 1ALalioti et al. (1997)
<i>MAOA</i>	X	30	2-5	4	promoter	Bipolar disorderByrd 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.

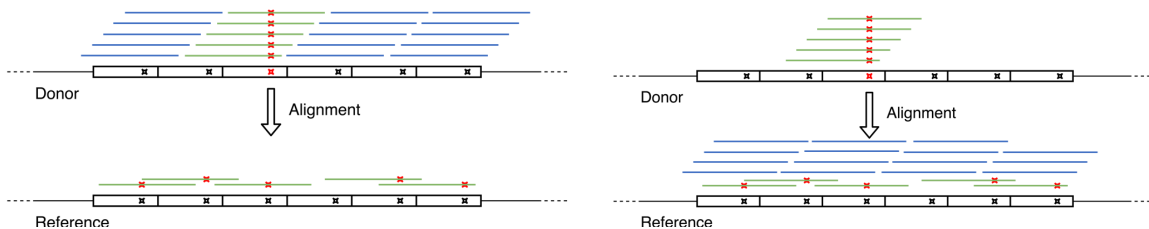


Figure 1: **Read alignment in VNTRs.** Green reads correspond to the repeating unit that contains a single variation (insertion/deletion) and blue reads correspond to other reads that are sequenced from VNTR region. After alignment, the green reads have multiple mapping location with only one mismatch which is the best alignment score. So, they will map to different location and after alignment of blue reads they will be treated as sequencing errors.

quality (*e.g.* alignment score). So, all the reads sequenced from the same repeating unit will go to different location. These mapping issues confound existing variant callers, including realignment tools such as GATK IndelRealigner(DePristo et al., 2011). Therefore, detection of point mutations in long VNTRs requires integrating information across the entire VNTR sequence. For VNTRs whose total sequence length (RU count times the RU length) is much longer than the read length, detection of SNVs and indels is not feasible using existing variant callers.

Mapping difficulty for varying RU count: Mapping tools such as BWA(Li and Durbin, 2009) and Bowtie2(Langmead and Salzberg, 2012) can work for read recruitment for STRs, but are challenged by insertion/deletion of larger repeat units. Similar to the case with presence of indel in the VNTR, variations in the number of RU counts will also affect the performance of alignment methods since they try to find the exact matching sequence within the reference genome. Fig. 4 demonstrate the performance of state-of-the-art mapping tools based on the number of RU counts in donor genome.

High error rate in long read technologies: Single molecule reads (*e.g.*, PacBio(Eid et al., 2009), Nanopore(Clarke et al., 2009)) can span entire VNTR regions, but it is difficult to estimate the RU count directly since the distance between the flanking regions varies dramatically from read to read due to an excess of indel errors. For example, the length-based RU count estimate from a VNTR in the SERT gene included five different values (13, 14, 15, 16, and 18) for a diploid genome. Fig. 5B shows the performance of a method that tries to infer copy number based on length of PacBio reads.

3 Existing methods

Finding a repetitive region of a string is a classical algorithmic problem, and the solutions to that problem are used on the human reference genome to find VNTRs (Benson, 1999). Other tools have addressed the problem of RU count estimation using NGS data, focusing on the related problem of STR genotyping. Some of these tools do not work with large repeating patterns (Willems et al., 2017; Liu et al., 2017). Others require all repeat units to be near-identical (Dolzhenko et al., 2017; Ummat and Bashir, 2014), while in a VNTR region the repeating units are approximate matches. In particular, ExpansionHunter (Dolzhenko et al., 2017) looks for exact matches of short repeating sequence within flanking unique sequences, and works for STRs, but not as well with the larger VNTRs with variations in RUs (Results). VNTRseek (Gelfand et al., 2014) detects a VNTR-like patterns in reads and reference genome and defines them as the vertices of a graph. To define an edge in the graph, it aligns the flanking regions of VNTRs in the reads to the reference genome to get similarity of flanking regions, and combines it with the similarity of repetitive pattern inside the flanking region based on alignment scores. Using a bipartite matching algorithm, it assigns every VNTR-like pattern in the NGS to a VNTR in the reference genome. In this approach, alignment based tools need to align reads at both unique ends, which may not be possible for short (Illumina) reads.

In contrast to methods like VNTRseek which seek to *discover/identify* VNTRs, we describe a method, adVNTR (Bakhtiari et al., 2018), for *genotyping VNTRs* at targeted loci in a donor genome. For any target VNTR in a donor, adVNTR reports an estimate of RU counts and point mutations within the RUs. It trains Hidden Markov Models (HMMs) for each target VNTR locus, which provide the following advantages: (i) it is sufficient to match any portions of the unique flanking regions for read alignment; (ii) it is easier to separate homopolymer runs from other indels helping with frameshift detection, and to estimate RU counts even in the presence of indels; (iii) each VNTR can be modeled individually, and complex models can be constructed for VNTRs with complex structure, along with VNTR specific confidence scores. For longer VNTRs not spanned by short reads, adVNTR can still be used to detect indels, while providing lower bounds on RU counts, or exact estimates for short VNTRs. adVNTR models can estimate RU counts for thousands of VNTRs on PacBio data. Using simulated data as well as whole-genome sequence data for a number

of human individuals, we demonstrate the power of adVNTR to genotype VNTR loci in the human genome.

4 Method

A VNTR sequence can be represented as $SR_1R_2\dots R_uP$, where S and P are the unique flanking regions, and $R_i(1 \leq i \leq u)$ correspond to the tandem repeats. For each i, j , R_i is similar in sequence to R_j , and the number of occurrences, u , is denoted as the *RU count*. We do not impose a length restriction on S and P , but assume that they are long enough to be unique in the genome. For genotyping a VNTR in a donor genome, we focus primarily on estimating the diploid RU counts (u_1, u_2) . However, many ($\sim 10^3$) VNTRs occur in coding regions, and mutations, particularly frameshift causing indels, are also relevant. Our method, adVNTR, models the problems of RU counting and mutation detection using HMMs trained for each target VNTR. adVNTR requires a one-time training of models for each combination of a VNTR and sequencing technology. Once models are trained, it has three stages for genotyping: (i) *recruitment of reads* containing the VNTR sequence; (ii) *counting RUs* for each of the two haplotypes; and, (iii) *identification of mutations*, specifically indels in coding regions. We describe the training procedure and the three modules below.

HMM Training. The goal of training is to estimate model parameters for each VNTR and each sequencing technology. We use an HMM architecture with three parts (Fig. 2). The first part matches the 5' (left) flanking region of the VNTR. The second part is an HMM which matches an arbitrary number of (approximately identical) repeating units. The last part matches the 3' (right) flanking region. The RU pattern is matched with a profile HMM (*RU HMM*), with states for matches, deletions, and insertions, and its model parameters are trained first. To train RU HMM for each VNTR, we collected RU sequences from the reference assembly Lander et al. (2001) and performed a multiple sequence alignment. Let $h(i, j)$ denote the number of observed transitions from state i to state j in hidden path of each sequence in multiple alignment, and $h_i(\alpha)$ denote the number of emissions of α in state i . We define permissible transition (arrows in Fig. 2) and

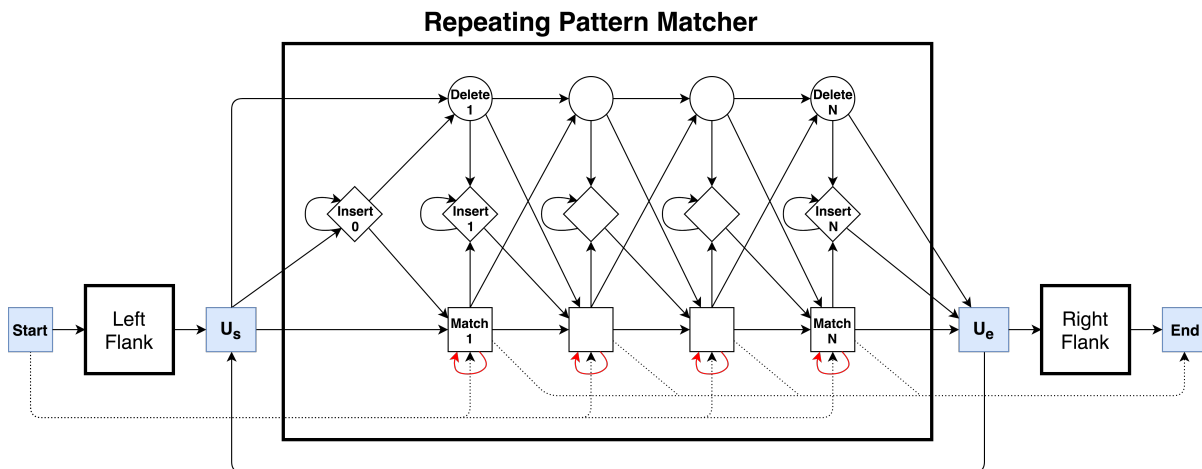


Figure 2: **The VNTR HMM.** The HMM is composed of 3 profile HMMs, one each for the left and right flanking unique regions, and one in the middle to match multiple and partial numbers of RUs. The special states U_s ('Unit-Start'), and U_e ('Unit-End') are used for RU counting. Dotted lines refer to special transitions for partial reads that do not span the entire region.

match-state emission probabilities as follows:

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

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

While error correction tools for PacBio have been developed, most do not work for repetitive regions,(Hackl et al., 2014; Salmela and Rivals, 2014; Au et al., 2012; Micolte et al., 2016; Lee et al., 2014) and others assume a single haplotype for error correction(Salmela et al., 2016; Berlin et al., 2015). In contrast, the HMM allows us to model many of the common (homopolymer) errors directly. Insertion deletion errors are common in single molecule sequencing particularly in homopolymer runs of length ≥ 6 , and occur mostly as insertions in the homopolymer run(Chaisson and Tesler, 2012). Consider a match state i with highest emission probability for nucleotide α . The

transition probability $T(i, i)$ from a match state i to itself was set based on the match probabilities of α in previous $k = 6$ states. The PacBio model parameters were updated using both simulated and real PacBio data.

Read Recruitment. The first step in adVNTR is to *recruit* all reads that match a portion of the VNTR sequence. Alignment-based methods do not work well due to changes in RU counts (See Results), but the adVNTR HMM allows for variable RU count. To speed up recruitment, we used an Aho-Corasick keyword matching algorithm(Altschul et al., 1990) to identify all reads that match a keyword from the VNTR patterns or the flanking regions. Note that the dictionary construction is a one-time process, and all reads must be scanned once for filtering. The keyword size and number of keywords were empirically chosen for each VNTR. Filtered reads had high sensitivity and were filtered by aligning to the HMM using the Viterbi algorithm. We aligned 10^7 non-target genomic sequences to the HMM to form an empirical null distribution, and used 10^{-5} as the p-value cut-off, correcting for multiple ($\sim 10^3$) target VNTRs.

Estimating VNTR RU Counts. Recall 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,

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

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

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

For each read, the Viterbi algorithm allows for the enumeration of the maximum likelihood (ML) path by going backwards from $\text{Prev}(\text{End}, n)$. Ignoring 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 count of the read as $k_1 + k_2 + k_3$, and mark it as a lower bound if $k_1 + k_3 > 0$ (see Fig. 3 for an example).

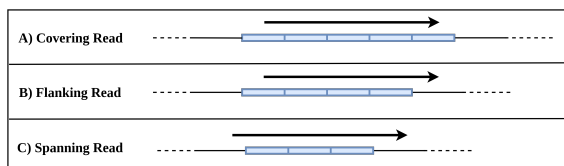


Figure 3: **Estimates of RU counts using recruited reads.** (A) $(k_1, k_2, k_3) = (1, 3, 1)$; RU count ≥ 5 . (B) $(k_1, k_2, k_3) = (0, 3, 1)$; RU count ≥ 4 (C) $(k_1, k_2, k_3) = (0, 3, 0)$; RU count = 3.

To model errors in read counts, we define parameter r_ϵ s.t. r_ϵ^Δ is the probability of RU counting error by $\pm\Delta$ in the estimation of the true count. Thus the

probability of getting the correct count is $1-r$, where

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

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 . Additionally, we allow the possibility that all reads are sampled from one haplotype with the RU count of the missing haplotype being X . We define $C = \{c_1, c_2, \dots, c_n\} \cup \{X\}$ and use C to get a list of possible genotypes (c_i, c_j) with $c_i \leq c_j$. Then, the conditional likelihood of a read with RU count c is given by:

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

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

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

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 can be computed using Bayes' theorem:

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

We generally set equal priors. However, in the event that we only see reads with a single count c' , we choose $\Pr((c', c')) = \Pr((c', X)) = \frac{1}{2}$. If we see multiple counts, we set $\Pr((c', X)) = 0$ for all $c' \in C$, and give equal priors to all other genotypes.

VNTR Mutation Detection. It is not difficult to see that alignment based methods do not work well in VNTRs. Changes in RU counts make it difficult to align reads even for mappers

that allow split-reads, as the gaps in different reads can be placed in different locations. A similar problem appears with small indels, as there are multiple ways to align reads with an indel in a Repeat Unit. The adVNTR HMM aligns all repeat units to the same HMM, and this has the effect of aligning all mutations/indels in the same column. Consider the case where reads contain a total of v nucleotides matching a VNTR RU of length ℓ , and RU count u . Moreover at a specific position covered by d Repeats, suppose we observe ι indel transitions.

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

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

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

5 Results

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

HMM training. Initial HMMs were trained using a multiple alignment of RU sequences from the reference assembly hg19 (Lander et al., 2001), as described in methods. Similarly, HMMs were trained for the left flanking and right flanking regions for each VNTR. To tailor HMMs for short reads, we used WGS data of a CEU trio from 1000 Genomes project and the AJ trio from Genome in a Bottle (GIAB) project. Correspondingly, to train models for PacBio reads, we re-estimated model parameters after aligning PacBio simulated reads using SimLoRD (Stöcker et al., 2016). A

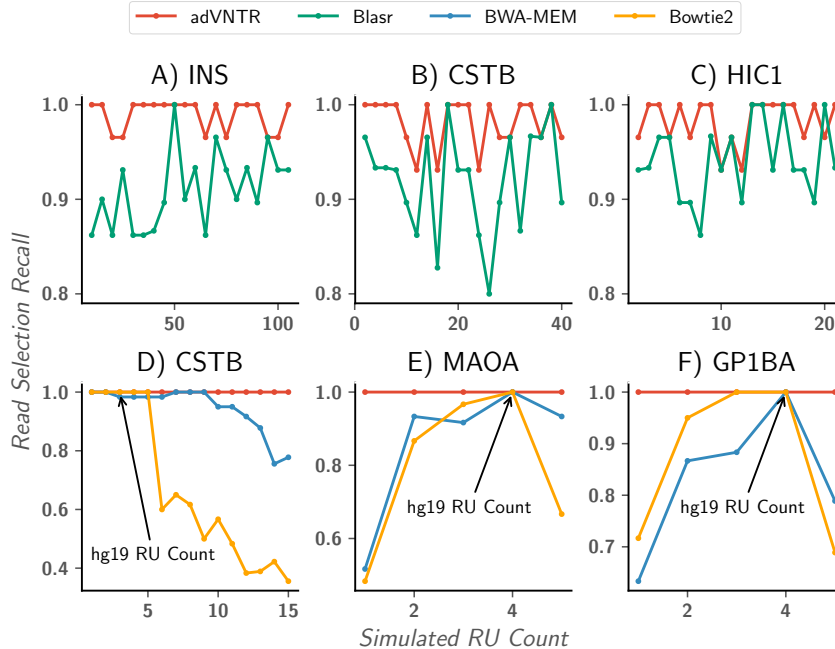


Figure 4: **Sensitivity of read recruitment at VNTR loci.** Comparison of adVNTR read selection with BWA-MEM and Bowtie2 mapping for Illumina reads (short VNTRs) and Blasr for PacBio reads (long VNTRs). Each plot shows the sensitivity of mapped/selected reads as a function of the number of repeats for different VNTRs.

total of 865 VNTR models were trained for VNTRs in coding and promoter regions of the genome, for both Illumina and PacBio. Subsequently, we tested performance for (a) read-recruitment, (b) counting of Repeat Units, and (c) detection of indels.

Running time. On PacBio WGS sequencing data at 30X coverage (aligned with Blasr), adVNTR took 10 hours to genotype all 865 VNTRs. For unaligned sequencing data of a CHB individual (20X) coverage, the time increased to 21:45 hours.

Read Recruitment. To evaluate read recruitment performance in *PacBio* sequencing, we simulated haplotypes with 30x coverage using SimLoRD(Stöcker et al., 2016) for three disease-linked VNTRs (INS, CSTB, and HIC1) in their known RU count range. To evaluate read recruitment and RU genotyping for *Illumina*, we used ART(Huang et al., 2011) to simulate haplotype WGS (shotgun 150bp) reads at 30x coverage for three disease-linked short VNTRs using known RU counts. Pairs of haplotypes were merged to get diploid samples. We compared adVNTR with BWA-MEM and Bowtie2 algorithms for assessing Illumina read recruitment (Fig. 4: A-C) and with Blasr(Chaisson and Tesler, 2012) for PacBio reads (Fig. 4: D-F). The plots show that while adVNTR works well for a range of RU counts, other mapping tools work well only when the simulated RU count matches

the reference RU count.

VNTR RU count estimation with PacBio reads. Recall that sequencing (particularly homopolymer) errors can cause lengths to change, particularly for short RU lengths and larger RU counts. To test the performance of adVNTR for RU counting on haplotype data, we compared against a naive method that estimates RU counts based on read length between the flanking regions. A total of 210 VNTRs with lengths ranging from 10bp to 90bp were selected and for each VNTR, PacBio sequence was simulated for 20 different RU counts with coverage varying from 1-40X. Fig. 5A shows RU count performance on VNTRs in *INS*, *CSTB*, and *HIC1* genes for varying RU counts and sequence coverage, while Fig. 5B shows RU counting performance on all 210 VNTRs as a function of RU lengths. adVNTR estimates are uniformly good except at low sequence coverage. We tested RU counting on diploid samples by simulating different RU counts on individuals at 3 VNTRs (Table S1). adVNTR RU counts showed 100% accuracy in each of the 52 different samples tested.

To test performance on real data where the true VNTR genotype is not known, we checked for Mendelian inheritance consistency at four disease-linked VNTRs in the AJ trio from Genome in a Bottle (GIAB)(Zook et al., 2016) and a Chinese Han trio from NCBI SRA (accession PRJEB12236). Our predictions were consistent in each case (Fig. 5C). We extended this analysis to 865 VNTRs in the coding region of human genome with RU count ≥ 10 . At a posterior probability threshold of 0.99, 99.1% of the calls in the AJ trio, and 98.5% of the calls in the Chinese trio are consistent with Mendelian inheritance. The few discrepancies can be attributed mainly to low coverage and missing data (Fig. 5E,F). Increasing sequence coverage threshold from $5\times$ to $10\times$ increased the average posterior probability from 0.91 to 0.98 resulting in improved RU count accuracy (Suppl. Fig. S1).

We also performed a long range (LR)PCR experiment on the individual NA12878 to assess the accuracy of the adVNTR genotypes. In Fig. 5D, black bands correspond to the observed PCR product lengths, and they match up with the computational product lengths based on estimated RU count (red arrows), while being different from the hg19 reference RU count. For each VNTR, there are two arrows for the predicted heterozygous RU counts, and a single arrow for the SLC6A4 VNTR that was predicted to be homozygous.

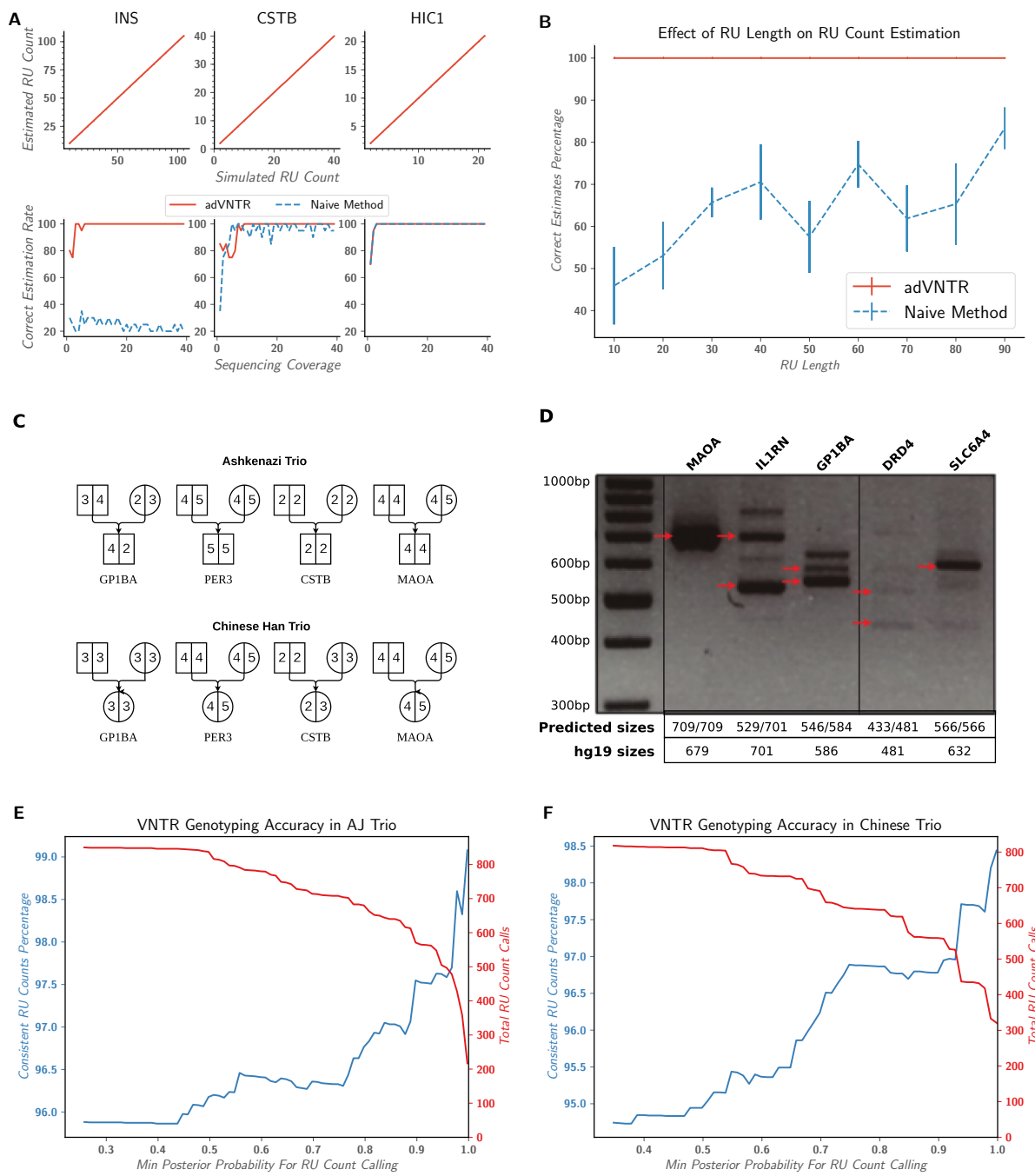


Figure 5: VNTR genotyping using PacBio data. (A) RU count estimation on simulated PacBio reads and effect of sequencing coverage on estimation. (B) Effect of repeating unit length on naive method. For shorter repeating units, insertion and deletion errors have more impact on RU count estimation. (C) Mendelian consistency of genotypes at 4 VNTR loci in the Chinese Han and Ashkenazi trios. (D) LR-PCR based validation of genotypes at disease-linked VNTRs in NA12878. (E) RU count calls and consistent calls ratio in AJ trio. (F) RU count calls and consistent calls ratio in Chinese trio.

While we could not get the VNTR discovery tool VNTRseek(Gelfand et al., 2014) to run on our machine (personal communication), we observed that the authors had predicted 125 VNTRs in the Watson sequenced genome(Wheeler et al., 2008), and 75 VNTRs in two trios as being polymorphic. In contrast, analysis of the PacBio sequencing data identified >500 examples of polymorphic VNTRs with RU counts ≥ 10 that overlap with coding regions. The results suggest that variation in RU counts of VNTRs and their role in influencing phenotypes might be greater than previously estimated.

RU counting with Illumina. We tested adVNTR RU counting performance on Illumina reads simulated using ART on three VNTRs and compared it to ExpansionHunter(Dolzhenko et al., 2017) which is designed mainly for STRs (Supp. Table S1). Of the 52 samples tested, adVNTR predicted the correct genotype in all but 6 cases, with erroneous calls in the case of high RU counts where the read length does not span the VNTR perfectly. In contrast, ExpansionHunter could not predict a majority of cases as it makes the assumption that the different RUs are mostly identical in sequence (valid for STRs but not for many VNTRs).

For short VNTRs, adVNTR can be an effective tool for larger population-scale studies of VNTR genotypes using WGS data replacing labor intensive gel electrophoresis(Byrd and Manuck, 2014; Cervera et al., 2007). Fig: 6 shows the RU count frequencies for two disease-linked VNTRs (in the coding region of GP1BA and promoter of MAOA), using 150 PCR-free WGS data obtained from 1000 genomes project(Consortium, 2015). The 2R/3R genotypes in GP1BA are associated with Aspirin Treatment failure for stroke prevention(Cervera et al., 2007). Notably, our results suggest that the 2R genotype is absent in African populations suggesting that this shorter allele arose after the out of Africa transition.

VNTR mutation/indel detection. To test indel detection, we simulated Illumina reads from 20 whole genomes after introducing a single insertion or deletion in the middle of the VNTR region in the CEL gene. As a negative control, we simulated 10 WGS experiments with a range of sequence coverage values. We ran adVNTR, Samtools mpileup(Li, 2011), and GATK Haplotype-Caller(DePristo et al., 2011) which uses GATK IndelRealigner, to identify frameshifts in each of the simulated datasets, and the 10 control datasets. On the control data, none of the tools found any variant. On the simulated indels, adVNTR made the correct prediction in each case (Suppl.

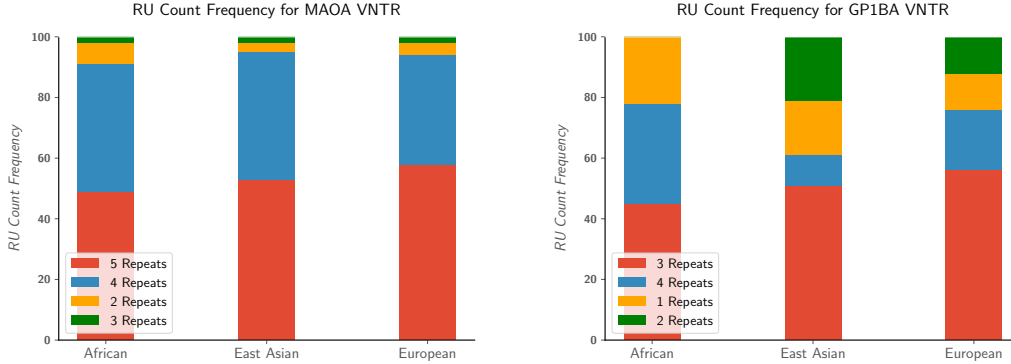


Figure 6: Population-scale genotyping of VNTRs. (A) RU count frequencies for MAOA VNTR. In more diverse populations, the dominant genotype has lower frequency than other populations which leads to having more diversity in RU count of the VNTR. (B) RU count frequencies for GP1BA VNTR. Lack of 2R genotype in African population suggests it may have occurred outside of Africa.

Table S2), while Samtools and GATK were unable to predict a single insertion or deletion. This result is not surprising as the reads have poor alignment scores, and the indel can be mapped to multiple locations (Suppl. Fig. S2)(Robinson et al., 2011).

As frameshifts in the VNTR region of the CEL gene have been linked to a monogenic form of diabetes(Ræder et al., 2006), we tested for frameshifts in CEL using whole Exome sequencing (WES) data from 2,081 cases with Type 2 Diabetes (Fuchsberger et al., 2016) and compared the numbers to 2,090 control individuals. WES data analysis is challenging as high GC-content makes it difficult to PCR-amplify this VNTR. adVNTR found that while none of the controls had any evidence of a frameshift, 8 of the 2,081 diabetes cases showed a frameshift in this VNTR region (Suppl. Fig. S3).

6 Discussion

The main contribution of our paper is the separation of VNTR discovery from VNTR genotyping, and allowing to genotype VNTRs using sequencing data instead of reference genome. The problem of genotyping VNTRs (determining diploid RU counts and mutations) is important for clinical pipelines seeking to find the genetic mechanisms of Mendelian disorders, as well as association studies where we could associate a disease with VNTR variation. In this paper, we presented adVNTR to genotype VNTRs using major next generation sequencing technologies. Our method train a model for each specific target, allowing us to tailor the approach for complex VNTRs. It solves the problem of alignment by mapping the reads to the pattern instead of a static reference

genome, and resolves the problem of mapping reads for indel detection by collapsing all RU copies inside the read. Like other STR genotyping tools, adVNTR works best when reads can span the VNTR, but (a) indel detection is possible for long VNTRs; (b) RU counting lower bounds can separate pathogenic cases from normal cases for some targets; and, (c) the increasing popularity of long read sequencing (esp. PacBio, and Nanopore) makes it possible to genotype VNTRs with higher confidence. Future research will focus on increasing the number of target VNTRs, and algorithmic strategies to speed up VNTR genotyping across the genome.

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

I would like to thank Melissa Gymrek, Sharona Shleizer-Burko, Vikas Bansal, and Vineet Bafna who are co-authors on a related manuscript on Biorxiv and helped me doing this project. Finally, I would like to thank Vineet Bafna for proof reading this work.

References

- Altschul, S. F., Gish, W., Miller, W., Myers, E. W., and Lipman, D. J., 1990. Basic local alignment search tool. *Journal of molecular biology*, **215**(3):403–410.
- Au, K. F., Underwood, J. G., Lee, L., and Wong, W. H., 2012. Improving PacBio long read accuracy by short read alignment. *PloS one*, **7**(10):e46679.
- Bakhtiari, M., Shleizer-Burko, S., Gymrek, M., Bansal, V., and Bafna, V., 2018. Targeted genotyping of variable number tandem repeats with advntr. *Genome Research*, **28**(11):1709–1719.
- Benedetti, F., Dallaspezia, S., Colombo, C., Pirovano, A., Marino, E., and Smeraldi, E., 2008. A length polymorphism in the circadian clock gene *Per3* influences age at onset of bipolar disorder. *Neuroscience letters*, **445**(2):184–187.
- Benson, G., 1999. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic acids research*, **27**(2):573.
- Berlin, K., Koren, S., Chin, C.-S., Drake, J. P., Landolin, J. M., and Phillippy, A. M., 2015. Assembling large genomes with single-molecule sequencing and locality-sensitive hashing. *Nature biotechnology*, **33**(6):623–630.
- Brookes, K., 2013. The VNTR in complex disorders: The forgotten polymorphisms? A functional way forward? *Genomics*, **101**(5):273–281.
- Byrd, A. L. and Manuck, S. B., 2014. MAOA, childhood maltreatment, and antisocial behavior: meta-analysis of a gene-environment interaction. *Biological psychiatry*, **75**(1):9–17.
- Cervera, A., Tassies, D., Obach, V., Amaro, S., Reverter, J., and Chamorro, A., 2007. The BC genotype of the VNTR polymorphism of platelet glycoprotein *Ib α* is overrepresented in patients with recurrent stroke regardless of aspirin therapy. *Cerebrovascular Diseases*, **24**(2-3):242–246.
- Chaisson, M. J. and Tesler, G., 2012. Mapping single molecule sequencing reads using basic local alignment with successive refinement (BLASR): application and theory. *BMC bioinformatics*, **13**(1):238.
- Clarke, J., Wu, H.-C., Jayasinghe, L., Patel, A., Reid, S., and Bayley, H., 2009. Continuous base identification for single-molecule nanopore DNA sequencing. *Nature nanotechnology*, **4**(4):265–270.
- Consortium, . G. P., 2015. A global reference for human genetic variation. *Nature*, **526**(7571):68–74.

- DePristo, M. A., Banks, E., Poplin, R., Garimella, K. V., Maguire, J. R., Hartl, C., Philippakis, A. A., Del Angel, G., Rivas, M. A., Hanna, M., *et al.*, 2011. A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nature genetics*, **43**(5):491–498.
- Dolzhenko, E., van Vugt, J. J., Shaw, R. J., Bekritsky, M. A., van Blitterswijk, M., Narzisi, G., Ajay, S. S., Rajan, V., Lajoie, B., Johnson, N. H., *et al.*, 2017. Detection of long repeat expansions from PCR-free whole-genome sequence data. *Genome Research*, :gr-225672.
- Durinovic-Belló, I., Wu, R., Gersuk, V., Sanda, S., Shilling, H., and Nepom, G., 2010. Insulin gene VNTR genotype associates with frequency and phenotype of the autoimmune response to proinsulin. *Genes and immunity*, **11**(2):188–193.
- Eid, J., Fehr, A., Gray, J., Luong, K., Lyle, J., Otto, G., Peluso, P., Rank, D., Baybayan, P., Bettman, B., *et al.*, 2009. Real-time DNA sequencing from single polymerase molecules. *Science*, **323**(5910):133–138.
- Eser, O., Eser, B., Cosar, M., Erdogan, M., Aslan, A., Yildiz, H., Solak, M., and Haktanir, A., 2011. Short aggrecan gene repetitive alleles associated with lumbar degenerative disc disease in Turkish patients. *Genet Mol Res*, **10**(3):1923–1930.
- Franke, B., Vasquez, A. A., Johansson, S., Hoogman, M., Romanos, J., Boreatti-Hümmer, A., Heine, M., Jacob, C. P., Lesch, K.-P., Casas, M., *et al.*, 2010. Multicenter analysis of the SLC6A3/DAT1 VNTR haplotype in persistent ADHD suggests differential involvement of the gene in childhood and persistent ADHD. *Neuropsychopharmacology*, **35**(3):656.
- Fuchsberger, C., Flannick, J., Teslovich, T. M., Mahajan, A., Agarwala, V., Gaulton, K. J., Ma, C., Fontanillas, P., Moutsianas, L., McCarthy, D. J., *et al.*, 2016. The genetic architecture of type 2 diabetes. *Nature*, **536**(7614):41–47.
- Galimberti, D., Scarpini, E., Venturelli, E., Strobel, A., Herterich, S., Fenoglio, C., Guidi, I., Scalabrini, D., Cortini, F., Bresolin, N., *et al.*, 2008. Association of a NOS1 promoter repeat with Alzheimer’s disease. *Neurobiology of aging*, **29**(9):1359–1365.
- Gelfand, Y., Hernandez, Y., Loving, J., and Benson, G., 2014. VNTRseek-a computational tool to detect tandem repeat variants in high-throughput sequencing data. *Nucleic acids research*, **42**(14):8884–8894.
- Gymrek, M., Willems, T., Guilmatre, A., Zeng, H., Markus, B., Georgiev, S., Daly, M. J., Price, A. L., Pritchard, J. K., Sharp, A. J., *et al.*, 2016. Abundant contribution of short tandem repeats to gene expression variation in humans. *Nature genetics*, **48**(1):22–29.

- Hackl, T., Hedrich, R., Schultz, J., and Förster, F., 2014. proovread: large-scale high-accuracy PacBio correction through iterative short read consensus. *Bioinformatics*, **30**(21):3004–3011.
- Haddley, K., Bubb, V., Breen, G., Parades-Esquivel, U., and Quinn, J., 2011. Behavioural genetics of the serotonin transporter. In *Behavioral Neurogenetics*, pages 503–535. Springer.
- Hijikata, M., Matsushita, I., Tanaka, G., Tsuchiya, T., Ito, H., Tokunaga, K., Ohashi, J., Homma, S., Kobashi, Y., Taguchi, Y., *et al.*, 2011. Molecular cloning of two novel mucin-like genes in the disease-susceptibility locus for diffuse panbronchiolitis. *Human genetics*, **129**(2):117–128.
- Huang, W., Li, L., Myers, J. R., and Marth, G. T., 2011. ART: a next-generation sequencing read simulator. *Bioinformatics*, **28**(4):593–594.
- Kirby, A., Gnrirke, A., Jaffe, D. B., Barešová, V., Pochet, N., Blumenstiel, B., Ye, C., Aird, D., Stevens, C., Robinson, J. T., *et al.*, 2013. Mutations causing medullary cystic kidney disease type 1 lie in a large VNTR in MUC1 missed by massively parallel sequencing. *Nature genetics*, **45**(3):299–303.
- Kirchheiner, J., Nickchen, K., Sasse, J., Bauer, M., Roots, I., and Brockmöller, J., 2007. A 40-basepair VNTR polymorphism in the dopamine transporter (DAT1) gene and the rapid response to antidepressant treatment. *The pharmacogenomics journal*, **7**(1):48.
- LaHoste, G., Swanson, J., Wigal, S., Glabe, C., Wigal, T., King, N., and Kennedy, J., 1996. Dopamine D4 receptor gene polymorphism is associated with attention deficit hyperactivity disorder. *Mol Psychiatry*, **1**(2):121–124.
- Lalioti, M. D., Scott, H. S., Buresi, C., Rossier, C., Bottani, A., Morris, M. A., Malafosse, A., and Antonarakis, S. E., 1997. Dodecamer repeat expansion in cystatin B gene in progressive myoclonus epilepsy. *Nature*, **386**(6627):847.
- Lander, E. S., Linton, L. M., Birren, B., Nusbaum, C., Zody, M. C., Baldwin, J., Devon, K., Dewar, K., Doyle, M., FitzHugh, W., *et al.*, 2001. Initial sequencing and analysis of the human genome. *Nature*, **409**(6822):860–921.
- Langmead, B. and Salzberg, S. L., 2012. Fast gapped-read alignment with Bowtie 2. *Nature methods*, **9**(4):357–359.
- Lee, H., Gurtowski, J., Yoo, S., Marcus, S., McCombie, W. R., and Schatz, M., 2014. Error correction and assembly complexity of single molecule sequencing reads. *BioRxiv*, :006395.

- Lemmers, R. J., de Kievit, P., Sandkuijl, L., Padberg, G. W., van Ommen, G.-J. B., Frants, R. R., and van der Maarel, S. M., 2002. Facioscapulohumeral muscular dystrophy is uniquely associated with one of the two variants of the 4q subtelomere. *Nature genetics*, **32**(2):235.
- Li, H., 2011. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics*, **27**(21):2987–2993.
- Li, H. and Durbin, R., 2009. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics*, **25**(14):1754–1760.
- Liu, Q., Zhang, P., Wang, D., Gu, W., and Wang, K., 2017. Interrogating the unsequenceable genomic trinucleotide repeat disorders by long-read sequencing. *Genome medicine*, **9**(1):65.
- Margulies, M., Egholm, M., Altman, W. E., Attiya, S., Bader, J. S., Bemben, L. A., Berka, J., Braverman, M. S., Chen, Y.-J., Chen, Z., *et al.*, 2005. Genome sequencing in microfabricated high-density picolitre reactors. *Nature*, **437**(7057):376–380.
- Miclotte, G., Heydari, M., Demeester, P., Rombauts, S., Van de Peer, Y., Audenaert, P., and Fostier, J., 2016. Jabba: hybrid error correction for long sequencing reads. *Algorithms for Molecular Biology*, **11**(1):10.
- Okazaki, S., Schirripa, M., Loupakis, F., Cao, S., Zhang, W., Yang, D., Ning, Y., Berger, M. D., Miyamoto, Y., Suenaga, M., *et al.*, 2017. Tandem repeat variation near the HIC1 (hypermethylated in cancer 1) promoter predicts outcome of oxaliplatin-based chemotherapy in patients with metastatic colorectal cancer. *Cancer*, .
- Pritchard, A. L., Pritchard, C. W., Bentham, P., and Lendon, C. L., 2007. Role of serotonin transporter polymorphisms in the behavioural and psychological symptoms in probable Alzheimer disease patients. *Dementia and geriatric cognitive disorders*, **24**(3):201–206.
- Pugliese, A., Zeller, M., Fernandez, A., Zalcberg, L. J., Bartlett, R. J., Ricordi, C., Pietropaolo, M., Eisenbarth, G. S., Bennett, S. T., and Patel, D. D., *et al.*, 1997. The insulin gene is transcribed in the human thymus and transcription levels correlate with allelic variation at the INS VNTR-IDD3 susceptibility locus for type 1 diabetes. *Nature genetics*, **15**(3):293–297.
- Ræder, H., Johansson, S., Holm, P. I., Haldorsen, I. S., Mas, E., Sbarra, V., Nerøen, I., Eide, S. Å., Grevle, L., Bjørkhaug, L., *et al.*, 2006. Mutations in the CEL VNTR cause a syndrome of diabetes and pancreatic exocrine dysfunction. *Nature genetics*, **38**(1):54.

- Robinson, J. T., Thorvaldsdóttir, H., Winckler, W., Guttman, M., Lander, E. S., Getz, G., and Mesirov, J. P., 2011. Integrative genomics viewer. *Nature biotechnology*, **29**(1):24–26.
- Salmela, L. and Rivals, E., 2014. LoRDEC: accurate and efficient long read error correction. *Bioinformatics*, **30**(24):3506–3514.
- Salmela, L., Walve, R., Rivals, E., and Ukkonen, E., 2016. Accurate self-correction of errors in long reads using de Bruijn graphs. *Bioinformatics*, **33**(6):799–806.
- Shriver, M. D., Jin, L., Chakraborty, R., and Boerwinkle, E., 1993. VNTR allele frequency distributions under the stepwise mutation model: a computer simulation approach. *Genetics*, **134**(3):983–993.
- Stöcker, B. K., Köster, J., and Rahmann, S., 2016. SimLoRD: Simulation of Long Read Data. *Bioinformatics*, **32**(17):2704–2706.
- Tyner, C., Barber, G. P., Casper, J., Clawson, H., Diekhans, M., Eisenhart, C., Fischer, C. M., Gibson, D., Gonzalez, J. N., Guruvadoo, L., *et al.*, 2016. The UCSC Genome Browser database: 2017 update. *Nucleic acids research*, **45**(D1):D626–D634.
- Ummat, A. and Bashir, A., 2014. Resolving complex tandem repeats with long reads. *Bioinformatics*, **30**(24):3491–3498.
- Viswanath, B., Purushottam, M., Kandavel, T., Reddy, Y. J., Jain, S., *et al.*, 2013. DRD4 gene and obsessive compulsive disorder: do symptom dimensions have specific genetic correlates? *Progress in Neuro-Psychopharmacology and Biological Psychiatry*, **41**:18–23.
- Wheeler, D. A., Srinivasan, M., Egholm, M., Shen, Y., Chen, L., McGuire, A., He, W., Chen, Y.-J., Makhijani, V., Roth, G. T., *et al.*, 2008. The complete genome of an individual by massively parallel DNA sequencing. *nature*, **452**(7189):872–876.
- Willems, T., Zielinski, D., Yuan, J., Gordon, A., Gymrek, M., and Erlich, Y., 2017. Genome-wide profiling of heritable and de novo STR variations. *Nature Methods*, .
- Worrall, B. B., Brott, T. G., Brown, R. D., Brown, W. M., Rich, S. S., Arepalli, S., Wavrant-De Vrièze, F., Duckworth, J., Singleton, A. B., Hardy, J., *et al.*, 2007. IL1RN VNTR polymorphism in ischemic stroke. *Stroke*, **38**(4):1189–1196.
- Wright, J. M., 1994. Mutation at vntrs: Are minisatellites the evolutionary progeny of microsatellites? *Genome*, **37**(2):345–347.

Zook, J. M., Catoe, D., McDaniel, J., Vang, L., Spies, N., Sidow, A., Weng, Z., Liu, Y., Mason, C. E., Alexander, N., *et al.*, 2016. Extensive sequencing of seven human genomes to characterize benchmark reference materials. *Scientific data*, **3**.

Supplementary Material

A. Appendix

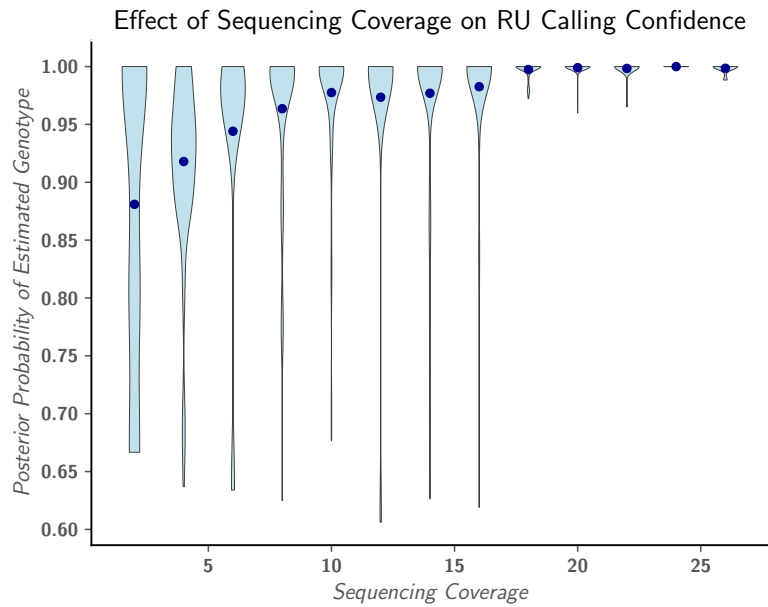


Figure S1: **Association of sequencing coverage in VNTR region and posterior probability of RU count calling.** The figure shows posterior probability of RU count estimation in AJ trio. Most of calls with low posterior probability (low confidence calls) result from low coverage in VNTR region. With at least 10 reads that span the VNTR, we will get 0.98 posterior probability for estimated genotype.

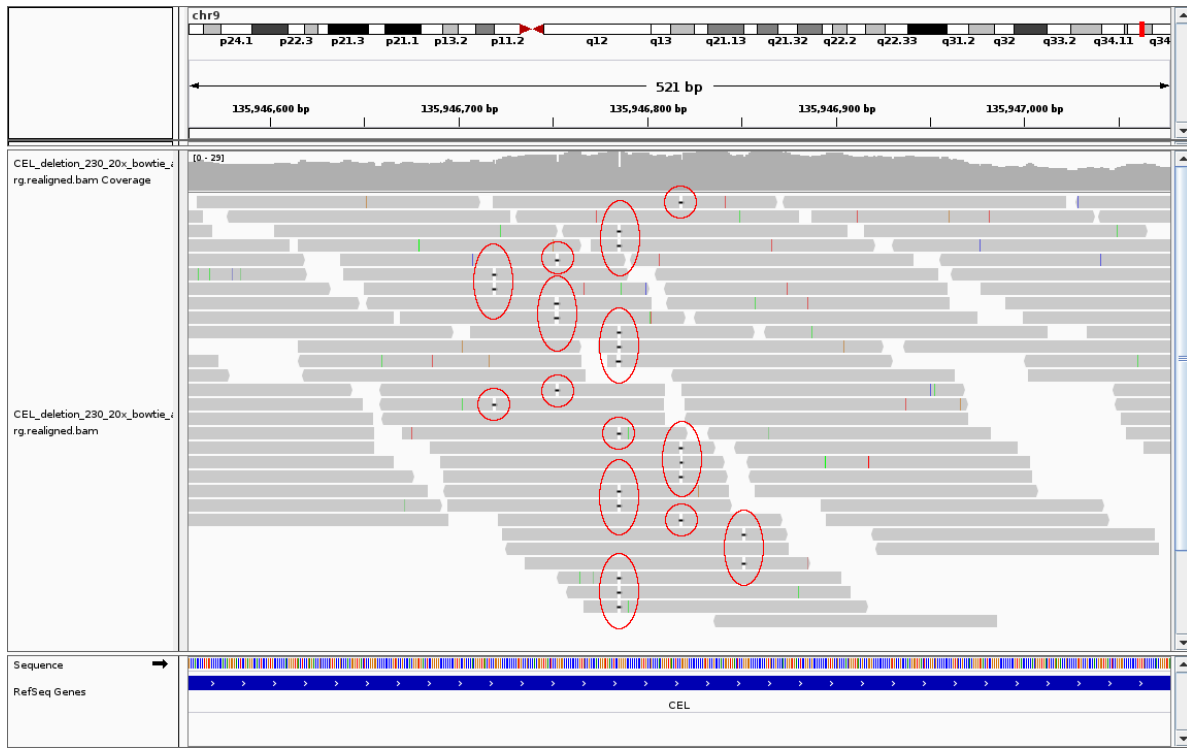


Figure S2: **Alignment stats with frameshift.** Alignment of a simulated data after running GATK IndelRealigner, when there is a deletion. With a sequencing mean of 30X, 25 reads contain the deletion but even after running realigner, deletions are mapped to five different repeating units.

Read1	GGCCACCCTGTG-CCCCACAGGGGACTCCGA
Read2	GGCCACCCTGTG-CCCCACAGGGGACTCCGA
Read3	GGCCACCCTGTG-CCCCACAGGGGACTCCGA
Read4	GGCCACCCTGTG-CCCCACAGGGGACTCCGA
Read5	GGCCACCCTGTG-CCCCACAGGGGACTCCGA
Read6	GGCCACCCTGTG-CCCCACAGGGGACTCCGA
ReferenceRepeatingUnit	GGCCACCCTGTGCCCCACAGGGGACTCCGA *****

Figure S3: **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.

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 S1: **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.

		# 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 S2: Comparison of indel finding with Samtools