

# Profiling of Short-Tandem-Repeat Disease Alleles in 12,632 Human Whole Genomes

Haibao Tang,<sup>1</sup> Ewen F. Kirkness,<sup>2</sup> Christoph Lippert,<sup>1</sup> William H. Biggs,<sup>2</sup> Martin Fabani,<sup>2</sup> Ernesto Guzman,<sup>2</sup> Smriti Ramakrishnan,<sup>1</sup> Victor Lavrenko,<sup>1</sup> Boyko Kakaradov,<sup>2</sup> Claire Hou,<sup>2</sup> Barry Hicks,<sup>1</sup> David Heckerman,<sup>1</sup> Franz J. Och,<sup>1</sup> C. Thomas Caskey,<sup>3</sup> J. Craig Venter,<sup>2,\*</sup> and Amalio Telenti<sup>2,\*</sup>

Short tandem repeats (STRs) are hyper-mutable sequences in the human genome. They are often used in forensics and population genetics and are also the underlying cause of many genetic diseases. There are challenges associated with accurately determining the length polymorphism of STR loci in the genome by next-generation sequencing (NGS). In particular, accurate detection of pathological STR expansion is limited by the sequence read length during whole-genome analysis. We developed TREDPARSE, a software package that incorporates various cues from read alignment and paired-end distance distribution, as well as a sequence stutter model, in a probabilistic framework to infer repeat sizes for genetic loci, and we used this software to infer repeat sizes for 30 known disease loci. Using simulated data, we show that TREDPARSE outperforms other available software. We sampled the full genome sequences of 12,632 individuals to an average read depth of approximately 30× to 40× with Illumina HiSeq X. We identified 138 individuals with risk alleles at 15 STR disease loci. We validated a representative subset of the samples (n = 19) by Sanger and by Oxford Nanopore sequencing. Additionally, we validated the STR calls against known allele sizes in a set of GeT-RM reference cell-line materials (n = 6). Several STR loci that are entirely guanine or cytosines (G or C) have insufficient read evidence for inference and therefore could not be assayed precisely by TREDPARSE. TREDPARSE extends the limit of STR size detection beyond the physical sequence read length. This extension is critical because many of the disease risk cutoffs are close to or beyond the short sequence read length of 100 to 150 bases.

## Introduction

Microsatellites, or short tandem repeats (STRs), are stretches of simple nucleotide repetitions in the genome; typical repeat units are 1–6 bp in length. Short tandem repeats are often polymorphic as a result of strand slippage during DNA replication and are a common source of rare genetic diseases.<sup>1</sup> The mutation rates of STRs are typically on the order of  $\sim 10^{-4}$  mutations per generation per site,<sup>2</sup> as compared to point mutation rates, which are on the order of  $\sim 10^{-8}$  mutations per generation per site for single-nucleotide variants (SNVs).<sup>3</sup> Because of the higher mutation rate, STRs offer a different level of resolution at which to study kinship and trait variations among individuals.

STRs are currently used in forensics to identify suspects from DNA traces left at a crime scene. The amplification targets the 13 CODIS (Combined DNA Index System)<sup>4</sup> STR loci, and the sizes of the amplicons are analyzed by electrophoresis. The repeat number at each loci is inferred from the size of the amplicon, and a DNA profile is generated. STRs also have a role in revealing genealogy. For example, STR loci on the Y chromosome (Y-STRs) are used for defining haplotypes that predated the use of Y-SNPs. The STR data, coupled with public genealogy databases such as Y-search, can be used for “surname inference.”<sup>5</sup>

STRs have been shown to be involved in several human genetic diseases.<sup>6</sup> Several neural-degenerative disorders, known as the “polyglutamine” (PolyQ) diseases, are caused

by variable stretches of the repeated trinucleotide CAG within protein-coding exons. Examples of PolyQ diseases are Huntington disease (HD [MIM: 143100]) and several forms of spinocerebellar ataxia (SCA). Huntington disease is caused by an expansion of the CAG repeats in the first exon of the Huntingtin gene (*HTT* [MIM: 613004]). Individuals carrying an expanded allele have motor, cognitive, and psychological symptoms that typically appear at the age of 40 years old or older, depending on the number of repeats.

STRs also occur in non-coding regions and can regulate gene expression and histone modifications, affecting the expression of nearby genes in *cis* to the STR sites.<sup>7</sup> Examples of these repeat disorders include Myotonic dystrophy (DM1 [MIM: 160900]), caused by CTG repeats; Friedreich Ataxia (FRDA [MIM: 229300]), caused by GAA repeats; and Fragile X syndrome (FXS [MIM: 300624]), caused by CGG repeats. STRs that regulate gene expression (e-STRs) are mostly enriched in genes responsible for cognitive functions and autoimmune responses.<sup>8</sup>

Whole-genome-scale analysis of human STR variation in the presumably healthy 1000 Genomes Project individuals suggests potential contributions of STRs to more complex traits.<sup>8</sup> Because most of these diseases are in the form of tri-nucleotide repeats, they are termed trinucleotide-repeat diseases (TREDs). Additionally, STR mutations are known to be associated with susceptibility to cancer.<sup>9</sup> Microsatellite instability is also a well-known hypermutability event that results from impaired DNA

<sup>1</sup>Human Longevity, Mountain View, CA 94041, USA; <sup>2</sup>Human Longevity, San Diego, CA 92121, USA; <sup>3</sup>Baylor College of Medicine, Houston, TX 77030, USA  
\*Correspondence: [jcventer@jcv.org](mailto:jcventer@jcv.org) (J.C.V.), [atelenti@humanlongevity.com](mailto:atelenti@humanlongevity.com) (A.T.)  
<https://doi.org/10.1016/j.ajhg.2017.09.013>

© 2017 The Author(s). This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

mismatch repair, as often occurs in colorectal and gastric cancer.<sup>10</sup>

Most STR loci are not usually included in routine analyses of the genetics of traits and diseases. Consequently, STRs might contribute to the “missing heritability” of complex diseases and traits.<sup>11</sup> One critical bottleneck in assaying STR loci by conventional experimental approaches is a dependency on targeting predefined sites for amplification and subsequently directly sequencing or measuring the size of the amplicons. Large-scale experimental assays are limited because it is difficult to measure multiplex loci within the same assay. Loci might also fail to amplify if they are highly expanded or if there are additional variants on the flanking regions. With whole-genome shotgun sequencing (WGS), it is now possible to type many STR loci on the basis of a single comprehensive sequencing run without the need to design separate genetic assays.

Using WGS, one can use the sequencing reads that map to STR loci to predict allele lengths. Variant calling software (for example,<sup>12,13</sup>) can identify some short indels in reads that span STRs. Other tools seek to identify STR variants by specifically examining the sequencing reads that are piled around a target STR region.<sup>14,15</sup> A popular caller, lobSTR, uses three separate steps: sensing, alignment, and allelotyping, which explicitly model two possible alleles (diploid) as well as sequencing errors typically associated with STRs (because of stutter noise).<sup>15</sup> However, lobSTR only considers reads that fully span an STR locus. Owing to the short length of Illumina reads (100–150 bases), this imposes a major limitation on the length of STR alleles that can be identified. A more recently developed STR caller, ExpansionHunter, incorporates additional evidence beyond spanning reads.<sup>16</sup>

It is also possible to estimate length variation at an STR by combining information from a prior estimate and the observed sizes of paired-end sequence fragments spanning the STR, as in STRViper.<sup>14</sup> However, STRViper assumes a single allele at each site, and this approach is inadequate for diploid human calls. Using long sequence reads such as PacBio<sup>17</sup> or Oxford Nanopore (ONP)<sup>18,19</sup> could potentially help to increase both the precision and the range of detectable variants.<sup>20,21</sup> However, the per-base cost of the long-read technologies is greater than for short-read technologies for whole-genome sequencing, limiting its widespread use for typing STRs. Indeed, few human genomes have been sequenced with PacBio or ONP because of the prohibitive costs associated with long-read sequencing platforms.

Despite recent progress,<sup>14,15,20</sup> high-throughput genotyping of STRs remains limited as a result of low effective coverage, sequencing stutters, and a lack of robust models with which to perform both haploid and diploid calls while distinguishing true variation from technical artifacts.<sup>11</sup> We built TREDPARSE to assess multiple sequence signatures suggested or implemented by previous methods.<sup>14,15,20</sup> Our testing on both simulated datasets, and more than 10,000 sequenced full human genomes,<sup>22</sup>

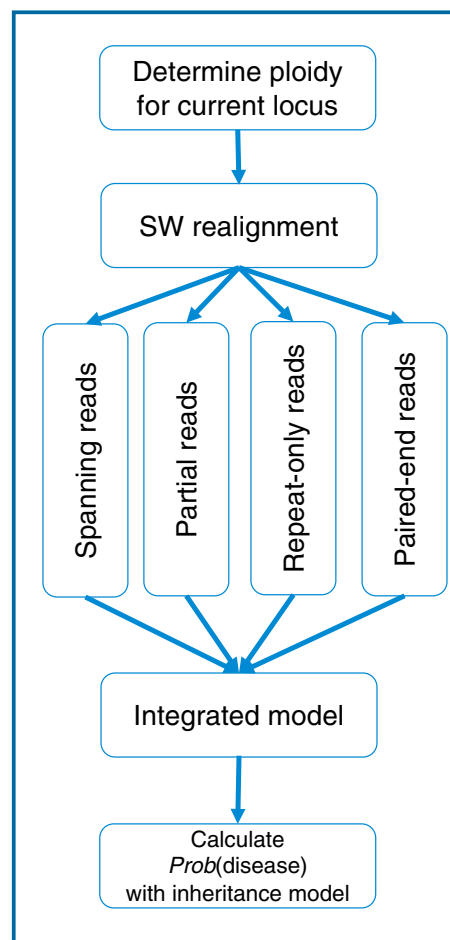
demonstrates that TREDPARSE yields highly accurate typing of many disease-related STRs.

## Material and Methods

### TREDPARSE Overview

The goal of TREDPARSE is to identify each allele length at predefined STR loci by using Illumina WGS sequence data that are sampled at sufficient depth (discussion on the sequencing depth is provided in later sections). Given a set of observed reads that are mapped around a particular STR locus, our goal is to estimate up to two haplotypes  $h_1$  and  $h_2$ , where  $1 \leq h_1 \leq h_2 \leq h_{max}$ , that represent the number of an individual's repeat units that maximize the likelihood in our model.

The TREDPARSE workflow involves a number of key steps—ploidy inference for a given locus, realignment of reads near the STR region, classification of the reads into four key types of evidence, and the deployment of a full probabilistic framework (Figure 1). The input for the workflow is typically a BAM file that contains mapped WGS reads, and the output is the maximum-likelihood size estimates, distributions over the number of repeats, and the associated probability of having each of the 30 STR-related diseases. The full probabilistic model is partitioned



**Figure 1. TREDPARSE Workflow for Calling STR-Related Genetic Disease**

The workflow includes ploidy inference, read realignment, and integration of various types of evidence in a probabilistic model.

متن کامل مقاله

دریافت فوری ←

**ISI**Articles

مرجع مقالات تخصصی ایران

- ✓ امکان دانلود نسخه تمام متن مقالات انگلیسی
- ✓ امکان دانلود نسخه ترجمه شده مقالات
- ✓ پذیرش سفارش ترجمه تخصصی
- ✓ امکان جستجو در آرشیو جامعی از صدها موضوع و هزاران مقاله
- ✓ امکان دانلود رایگان ۲ صفحه اول هر مقاله
- ✓ امکان پرداخت اینترنتی با کلیه کارت های عضو شتاب
- ✓ دانلود فوری مقاله پس از پرداخت آنلاین
- ✓ پشتیبانی کامل خرید با بهره مندی از سیستم هوشمند رهگیری سفارشات