



## STRait Razor: A length-based forensic STR allele-calling tool for use with second generation sequencing data



David H. Warshauer<sup>a</sup>, David Lin<sup>b</sup>, Kumar Hari<sup>b</sup>, Ravi Jain<sup>b</sup>, Carey Davis<sup>a</sup>, Bobby LaRue<sup>a</sup>, Jonathan L. King<sup>a</sup>, Bruce Budowle<sup>a,c,\*</sup>

<sup>a</sup> Institute of Applied Genetics, Department of Forensic and Investigative Genetics, University of North Texas Health Science Center, 3500 Camp Bowie Boulevard, Fort Worth, TX 76107, USA

<sup>b</sup> cBio, Inc., 37869 Abraham Street, Fremont, CA 94536, USA

<sup>c</sup> Center of Excellence in Genomic Medicine (CEGMR), King Abdulaziz University, Jeddah, Saudi Arabia

### ARTICLE INFO

#### Article history:

Received 8 March 2013

Received in revised form 4 April 2013

Accepted 15 April 2013

#### Keywords:

STR

Bioinformatics

Software

Second generation sequencing

MiSeq™

GAllx™

### ABSTRACT

Recent studies have demonstrated the capability of second generation sequencing (SGS) to provide coverage of short tandem repeats (STRs) found within the human genome. However, there are relatively few bioinformatic software packages capable of detecting these markers in the raw sequence data. The extant STR-calling tools are sophisticated, but are not always applicable to the analysis of the STR loci commonly used in forensic analyses. STRait Razor is a newly developed Perl-based software tool that runs on the Linux/Unix operating system and is designed to detect forensically-relevant STR alleles in FASTQ sequence data, based on allelic length. It is capable of analyzing STR loci with repeat motifs ranging from simple to complex without the need for extensive allelic sequence data. STRait Razor is designed to interpret both single-end and paired-end data and relies on intelligent parallel processing to reduce analysis time. Users are presented with a number of customization options, including variable mismatch detection parameters, as well as the ability to easily allow for the detection of alleles at new loci. In its current state, the software detects alleles for 44 autosomal and Y-chromosome STR loci. The study described herein demonstrates that STRait Razor is capable of detecting STR alleles in data generated by multiple library preparation methods and two Illumina® sequencing instruments, with 100% concordance. The data also reveal noteworthy concepts related to the effect of different preparation chemistries and sequencing parameters on the bioinformatic detection of STR alleles.

© 2013 Elsevier Ireland Ltd. All rights reserved.

## 1. Introduction

The majority of genetic analyses used in forensic casework involve the detection of short tandem repeats (STRs), which are relatively small sequences of DNA made up of repeating units of 2–6 nucleotides [1–3]. Currently, the accepted means of detecting these markers is size separation by capillary electrophoresis (CE) [4–7]. In recent years, however, second generation sequencing (SGS) technology has advanced to the point that it can be considered a viable platform for forensic DNA analyses, including STR detection [8–11]. In addition, sequencing has long been regarded as an effective means of revealing individual base variations in DNA known as single nucleotide polymorphisms (SNPs) [11–15]. While the traditional CE-based method of STR detection reveals only the length of alleles, SGS can increase resolution and detect nucleotide variation within the repeat

regions and in proximal flanking regions [16]. Furthermore, the expansive genetic coverage and read length associated with current SGS methods allow for the potential capture of genetic information related to far more forensic STR markers than are possible with conventional multiplex-based CE kits. With the ability to sequence gigabases of DNA [17,18], a properly designed assay could yield STR information in a single analysis which surpasses that of all the currently available commercial CE-based kits combined, and do so on multiple samples simultaneously. However, while the extant sequencing instruments are capable of providing such extensive data, currently available software tools for identifying forensic STR alleles within the data are only beginning to address the task.

One such tool, lobSTR [8], uses an algorithm specifically designed to identify STR alleles within SGS data. First, the software analyzes a raw FASTA/FASTQ or BAM input file, detecting reads that contain an STR sequence and identifying the repeat motif. Next, lobSTR aligns the regions that flank the STR sequence to a modified reference sequence. Finally, the algorithm determines the identity of the allele(s) based on the number of detected repeat

\* Corresponding author. Tel.: +1 8177352979.

E-mail address: [Bruce.Budowle@unthsc.edu](mailto:Bruce.Budowle@unthsc.edu) (B. Budowle).

units between the two flanking regions, applying statistical corrections to produce the most likely allelotype. While this software is certainly a refined and reasonably accurate method, it is somewhat limited in that lobSTR, by default, identifies only a single simple repeat motif. To allow the software to detect alleles at STRs that have longer, complex repeats, such as those within the D21S11 locus, for example, the user must determine the distinct simple repeats that comprise the complex motif and instruct lobSTR to identify each of these repeats individually. The resulting data must then be interpreted altogether in order to draw conclusions. This necessity makes lobSTR less applicable for the analysis of forensic STR markers, which display varying repeat motif complexity.

Recently, a method was introduced by Bornman et al. [9] that allows for the detection of STR alleles in SGS data using a different strategy. This method uses the Bowtie short read aligner [19] to align raw SGS reads to an “in silico reference,” which is a user-generated FASTA file containing the full sequence of each allele at each STR locus. To reduce erroneous allele calls, reads are filtered so that only those encompassing the entire repeat region defined in the reference file are used for allelotyping. Allele calls are made using a heuristic decision model based on Fisher’s exact test, and probability values are given for each allele call. This software also is effective for identifying STR alleles in sequence data, but requires prior knowledge of allelic sequence information. As a result, novel alleles or allelic variants, or those for which there are no documented sequence data, are a limitation for this system.

STRait Razor (the STR allele identification tool – Razor) is a Perl script designed for the Linux/Unix platform that identifies alleles at forensic STR loci based on the length of the repeat sequence, a method that is conceptually similar to the length-based allele detection offered by CE. This software is capable of handling repeat motifs ranging from simple to complex, and it does not require a reference composed of extensive allelic sequence data. As a result, the allele call results are consistent with those of current CE-based methods, and it is not confounded by unexpected sequence variation within repeats. In its first iteration of development, STRait Razor is capable of detecting alleles at 44 forensically-relevant STR loci, and others can be configured readily.

Fordyce et al. [10] independently developed a software tool that functions similarly to STRait Razor, isolating the repeat region of interest and performing length-based allelotyping. However, the algorithm was designed for use with the Roche® Genome Sequencer FLX™ and is only able to analyze FASTA files consisting of sequence data that contain Roche® Molecular Identifier (MID) tags. Currently this software only is able to identify alleles at 5 STR loci, compared with the 44 STR loci detected by STRait Razor. While STRait Razor was only tested on raw FASTQ files output by Illumina® instruments in this study, the software should, in theory, maintain compatibility with the raw read files generated by any second generation sequencing platform.

## 2. Materials and methods

To test the efficiency and accuracy of STRait Razor, an initial concordance study was performed wherein allele calls made by the software were compared with CE results. Following the University of North Texas Health Science Center Institutional Review Board approval, DNA from five Caucasian blood samples was used for this trial.

### 2.1. Sample preparation and CE typing

Blood samples were extracted using the Qiagen® QIAamp® DNA Mini kit (Qiagen Inc., Valencia, CA), following the manufacturer’s recommendations. The quantity of extracted DNA was determined

using the Applied Biosystems® Quantifiler™ Human DNA Quantification Kit (Life Technologies, Carlsbad, CA) on an Applied Biosystems® 7500 Real-Time PCR System (Life Technologies), according to the manufacturer’s protocol. Amplification was performed using the reagents from the Applied Biosystems® AmpFISTR® Identifier® PCR Amplification Kit (Life Technologies) and the Promega® PowerPlex® 16 HS, ESI 17 Pro, and Y23 Systems (Promega Corporation, Madison, WI), on an Applied Biosystems® GeneAmp® PCR System 9700 thermal cycler (Life Technologies), according to the manufacturer’s recommendations. These various kits allowed for the typing of the following STR loci: CSF1PO, D13S317, D16S539, D18S51, D19S433, D21S11, D2S1338, D3S1358, D5S818, D7S820, D8S1179, FGA, TH01, TPOX, vWA, PENTA D, PENTA E, D10S1248, D12S391, D1S1656, D22S1045, D2S441, SE33, DYS19, DYS385, DYS389I/II, DYS390, DYS391, DYS392, DYS393, DYS437, DYS438, DYS439, DYS448, DYS456, DYS458, DYS481, DYS533, DYS549, DYS570, DYS576, DYS635, DYS643, and GATA H4. Capillary electrophoresis was performed on an Applied Biosystems® 3130xl Genetic Analyzer (Life Technologies) using POP-4™ polymer (Life Technologies) and analyzed using Applied Biosystems® GeneMapper® ID v3.2 software (Life Technologies), according to the manufacturer’s recommended protocol.

### 2.2. Sample preparation, SGS, and STRait Razor typing

The quantity of extracted DNA from the blood sample was determined using a Qubit® 2.0 Fluorometer (Life Technologies). Library preparation prior to sequencing was performed using either the Illumina® TruSeq™ Custom Enrichment protocol (Illumina, Inc., San Diego, CA) or the Agilent Technologies HaloPlex™ Target Enrichment protocol (Agilent Technologies, Inc., Santa Clara, CA). Using the DesignStudio (Illumina, Inc.) and SureDesign (Agilent Technologies, Inc.) software, respectively, custom probes were designed to target the forensically-relevant STR loci. Paired-end sequencing was carried out on the GAIIX™ and MiSeq™ sequencing platforms (Illumina, Inc.). For these trials, the read lengths employed by these instruments were 2 × 146 and 2 × 251, respectively. Sample 1 was prepared using the HaloPlex™ chemistry and sequenced on both the GAIIX™ and MiSeq™ instruments. The sample also was prepared using the TruSeq™ chemistry, and subsequently sequenced on the MiSeq™. Sample 2 was prepared using the HaloPlex™ chemistry and sequenced on the GAIIX™. Samples 3, 4, and 5 were prepared using the TruSeq™ chemistry and sequenced on the MiSeq™. Following sequencing, the GAIIX™ output *bcl* files were demultiplexed and converted to a single FASTQ file using CASAVA v1.8.2 [20]. The MiSeq™ output was automatically converted to FASTQ format by the MiSeq™ Reporter software [21]. These FASTQ files served as the input for STRait Razor. The software was designed to detect the following forensic STR loci: CSF1PO, TPOX, D2S441, D3S1358, D5S818, D13S317, D18S51, D16S539, D7S820, D8S1179, TH01, vWA, D21S11, FGA, D2S1338, D19S433, PENTA D, PENTA E, D10S1248, D12S391, D1S1656, D22S1045, DYS389I/II, DYS390, DYS456, DYS19, DYS458, DYS437, DYS438, DYS448, GATA H4, DYS391, DYS392, DYS393, DYS439, DYS481, DYS533, DYS549, DYS570, DYS576, DYS643, DYS385, and DYS635. For these trials, allelicalling was performed using STRait Razor’s default flank recognition settings (1 allowable substitution and no allowable insertions or deletions). The server used for STRait Razor analysis was a Dell™ PowerEdge™ R900 blade server, with 64 GB DDR3 ECC RAM and 4 Quad-Core Intel® Xeon® E7430 CPUs (2.13 GHz each).

### 2.3. STRait Razor information

The algorithm employed by STRait Razor is relatively simple (Fig. 1). First, reads containing both a leading and trailing flanking



**Table 1**  
Autosomal STR loci detected by STRait Razor. Leading and trailing flanking sequences for both forward and reverse complement reads are listed. All flanking sequences are directly adjacent to the repeat region, except for those labeled with an asterisk (\*), which were modified to allow for increased specificity or more efficient detection. Additional alleles can be added at the discretion of the user.

Locus	Flanking region sequences	Detectable alleles
CSF1PO	GATAGATAGATT---AGGAAGTACTTA TAAGTACTTCCT---AATCTATCTATC	5, 6, 6.3, 7, 7.3, 8, 8.3, 9, 9.1, 10, 10.1, 10.2, 10.3, 11, 11.1, 11.3, 12, 12.1, 13, 14, 15, 16
D10S1248	TATTGTCTTCAT---ACTCACTCATTT AAATGAGTGAGT---ATGAGACAATA	8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19
D12S391	AAATCCCCTCTC---ACCTATGCATCC GGATGCATAGGT---GAGAGGGGATT	15, 16, 17, 17.3, 18, 18.3, 19, 19.3, 20, 21, 22, 23, 24, 25, 26
D13S317	AGATGATTGATT---ATGTATTTGTAA TTACAAATACAT---AATCAATCATCT	5, 6, 7, 7.1, 8, 8.1, 9, 10, 11, 11.1, 11.3, 12, 13, 13.3, 14, 14.3, 15, 16, 17
D16S539	GACAGACAGGTG---TCATTGAAGAC GTCTTTCAATGA---CACCTGTCTGTC	4, 5, 6, 7, 8, 8.3, 9, 9.3, 10, 11, 11.3, 12, 12.1, 12.2, 13, 13.1, 13.3, 14, 14.3, 15, 16
D18S51	TCCTCTCTCTTT---GAGACAAGGTCT AGACCTGTCTC---AAAGAGAGAGGA	7, 8, 9, 9.2, 10, 10.2, 11, 11.1, 11.2, 12, 12.2, 12.3, 13, 13.1, 13.2, 13.3, 14, 14.2, 15, 15.1, 15.2, 15.3, 16, 16.1, 16.2, 16.3, 17, 17.1, 17.2, 17.3, 18, 18.1, 18.2, 19, 19.2, 20, 20.1, 20.2, 21, 21.1, 21.2, 22, 22.1, 22.2, 23, 23.1, 23.2, 24, 24.2, 25, 26, 27, 28.1, 28.3, 39.2
D19S433	AAGATTCTGTG---AGAGAGGTAGAA TTCTACCTCTCT---CAACAGAATCTT	5.2, 6.2, 7, 8, 9, 10, 11, 11.1, 12, 12.1, 12.2, 13, 13.1, 13.2, 13.3, 14, 14.1, 14.2, 14.3, 15, 15.2, 16, 16.2, 17, 17.2, 18, 18.2, 19, 19.2, 20
D1S1656	TAAACACACACA---CATCATAACAGTT* AACTGTATGATG---TGTGTGTGTTTA*	9, 10, 11, 12, 13, 13.3, 14, 14.3, 15, 15.3, 16, 16.3, 17, 17.1, 17.3, 18, 18.3, 19, 19.3, 20, 20.3, 21
D21S11	ATAGATAGACGA---AGGCAATTCACT AGTGAATTGCTC---TCGTCTATCTAT	12, 24, 24.2, 24.3, 25, 25.1, 25.2, 25.3, 26, 26.1, 26.2, 27, 27.1, 27.2, 27.3, 28, 28.1, 28.2, 28.3, 29, 29.1, 29.2, 29.3, 30, 30.1, 30.2, 30.3, 31, 31.1, 31.2, 31.3, 32, 32.1, 32.2, 32.3, 33, 33.1, 33.2, 33.3, 34, 34.1, 34.2, 34.3, 35, 35.1, 35.2, 35.3, 36, 36.1, 36.2, 36.3, 37, 37.2, 38, 38.2, 39, 39.2, 40.2, 41.2
D22S1045	TATTTTTATAAC---GAGACTACTATC GATAGTAGTCTC---GTTATAAAAATA	8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20
D2S1338	GGATTGCAGGAG---AGGCCAAGCCAT ATGGCTTGGCCT---CTCCTGCAATCC	11, 12, 13, 14, 15, 16, 17, 18, 19, 19.3, 20, 21, 22, 23, 23.2, 23.3, 24, 25, 26, 27, 28
D2S441	TCTATGAAAAC---TATCATAACACC GGTGTTATGATA---AGTTTTATAGA	8, 9, 10, 11, 11.3, 12, 12.3, 13, 13.3, 14, 14.3, 15, 16, 17
D3S1358	AGGCTTGCATGT---ATGAGACAGGGT ACCCTGTCTCAT---ACATGCAAGCCT	8, 8.3, 9, 10, 11, 12, 13, 14, 14.3, 15, 15.1, 15.2, 15.3, 16, 16.2, 17, 17.1, 17.2, 18, 18.1, 18.2, 18.3, 19, 20
D5S818	ATTTATACCTCT---TCAAAATATTAC GTAATATTTTGA---AGAGGTATAAAT	6, 7, 8, 9, 10, 10.1, 11, 11.1, 12, 12.3, 13, 14, 15, 16, 17, 18
D7S820	GAACGAACTAAC---GACAGATTGATA TATCAATCTGTC---GTTAGTTCGTTT	5, 5.2, 6, 6.2, 6.3, 7, 7.1, 7.3, 8, 8.1, 8.2, 8.3, 9, 9.1, 9.2, 9.3, 10, 10.1, 10.3, 11, 11.1, 11.3, 12, 12.1, 12.2, 12.3, 13, 13.1, 14, 14.1, 15, 16
D8S1179	CACTGTGGGGAA---TACGAATGTACA TGTACATTCGTA---TTCCCCACAGTG	7, 8, 9, 10, 10.2, 11, 12, 12.3, 13, 14, 15, 15.1, 15.3, 16, 17, 17.1, 18, 19, 20
FGA	GAAAGGAAGAAA---CTAGCTTGTAATA TTTACAAGCTAG---TTTCTCTCTTC	12.2, 13, 13.2, 14, 14.3, 15, 15.3, 16, 16.1, 16.2, 17, 17.1, 17.2, 18, 18.1, 18.2, 19, 19.1, 19.2, 19.3, 20, 20.1, 20.2, 20.3, 21, 21.1, 21.2, 21.3, 22, 22.1, 22.2, 22.3, 23, 23.1, 23.2, 23.3, 24, 24.1, 24.2, 24.3, 25, 25.1, 25.2, 25.3, 26, 26.1, 26.2, 26.3, 27, 27.1, 27.2, 27.3, 28, 28.1, 28.2, 29, 29.1, 29.2, 30, 30.2, 31, 31.2, 32, 32.1, 32.2, 33.1, 33.2, 34.1, 34.2, 35.2, 41.1, 41.2, 42, 42.1, 42.2, 43.1, 43.2, 44, 44.2, 44.3, 45, 45.1, 45.2, 46, 46.1, 46.2, 47, 47.2, 48, 48.2, 49, 49.1, 49.2, 50.2, 50.3, 51, 51.2
Penta D	TTTTATGATTCTC---TTGAGATGGTGT* ACACCATCTCAA---GAGAATCATAAA*	1.1, 1.2, 2.2, 3.2, 4, 5, 6, 6.4, 7, 7.1, 7.4, 8, 8.1, 9, 9.1, 9.4, 10, 10.1, 10.2, 10.3, 11, 11.1, 11.2, 12, 12.1, 12.2, 12.3, 12.4, 13, 13.2, 13.3, 13.4, 14, 14.1, 14.4, 15, 15.1, 16, 17, 18
Penta E	TCCTTACAATTT---GAGACTGAGTCT AGACTCAGTCTC---AAATTGTAAGGA	5, 6, 7, 8, 9, 9.1, 9.4, 10, 10.2, 11, 11.4, 12, 12.1, 12.2, 12.3, 13, 13.2, 13.4, 14, 14.4, 15, 15.2, 15.4, 16, 16.4, 17, 17.4, 18, 18.4, 19, 19.4, 20, 20.2, 20.3, 21, 22, 23, 23.4, 24, 26
TH01	CCCTTATTTCCC---TCACCATGGAGT ACTCCATGGTGA---GGGAAATAAGGG	3, 4, 5, 5.3, 6, 6.1, 6.3, 7, 7.1, 7.3, 8, 8.3, 9, 9.1, 9.3, 10, 10.3, 11, 12, 13.3, 14
TPOX	GAACCCTCACTG---TTTGGCAAATA TATTTACCCAAA---CAGTGGGGTTC	4, 5, 6, 7, 7.3, 8, 9, 10, 10.1, 10.3, 11, 12, 13, 13.1, 14, 15, 16
vWA	GACTTGGATTGA---TCCATCCATCCT AGGATGGATGGA---TCAATCCAAGTC	10, 11, 12, 13, 14, 15, 15.2, 16, 16.1, 17, 18, 18.1, 18.2, 18.3, 19, 19.2, 20, 21, 22, 23, 24, 25

#### 2.4. Allele call comparison

The alleles detected by the CE method were compared with the allele call output files generated by STRait Razor. To be considered concordant, alleles detected via CE had to be detected by STRait Razor, based on the presence of the allelic sequence data in the input FASTQ file. It should be noted that alleles not detected by STRait Razor due to lack of respective sequence data, whether the result of kit chemistry limitations or inadequate sequencing read length, were not considered discordant.

### 3. Results

Based on the various combinations of library preparation methods and sequencing platforms utilized, a total of 7 overall comparisons with CE data were made for the 5 samples. The assortment of multiplex-based CE kits used for these samples allowed for the comparison of STR alleles at all loci that are

currently detectable with STRait Razor: 22 autosomal STR loci and 22 Y-STR loci, with DYS385 designated as a single locus. Across all trials, a total of 427 alleles were compared. The allele calls made by STRait Razor were completely concordant with the genotype results generated by the CE method. Of the 427 alleles compared, 403 alleles were detected, with 100% concordance (Tables 3 and 4). For the 24 alleles not detected by the software, a subsequent manual analysis of the FASTQ input files revealed that there were no sequence reads for these alleles in which the full repeat regions (including the surrounding flanking sequences) were present. These undetected alleles were not represented in the sequencing data because of the library preparation method (e.g., random shearing of genomic DNA) and/or the read length used, and thus could not be recognized by STRait Razor. These instances do not reflect inadequacy of the software and should not be considered evidence of discordance. Intra-repeat variation was occasionally observed in the data output by STRait Razor. For example, an examination of the



**Table 2**

Y-chromosome STR loci detected by STRait Razor. Leading and trailing flanking sequences for both forward and reverse complement reads are listed. All flanking sequences are directly adjacent to the repeat region, except for those labeled with an asterisk (\*), where the proximities were modified to allow for increased specificity or more efficient detection. Additional alleles can be added at the discretion of the user.

Locus	Flanking region sequences	Detectable alleles
DYS19	TATATAGTGTTT---TATAGTGACACT AGTGTCACTATA---AAACACTATATA	9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19
DYS385	GAGAAAGAAAGG---GGAGGACTATGT* ACATAGTCCTCC---CCTTCTTTCTC*	7, 8, 9, 9.2, 10, 10.2, 11, 11.2, 11.3, 12, 12.1, 12.2, 12.3, 13, 13.1, 13.2, 13.3, 14, 14.2, 14.3, 15, 15.1, 15.2, 15.3, 16, 16.2, 16.3, 17, 17.1, 17.2, 17.3, 18, 18.1, 18.2, 19, 19.2, 19.3, 20, 21, 22, 23, 24, 25, 28
DYS389I	ATTATCTATGTA---TCCCTCCCTCTA TAGAGGGAGGGA---TACATAGATAAT	9, 10, 11, 12, 13, 14, 15, 16, 17
DYS389II	TCTATGTGTGTG---TCCCTCCCTCTA TAGAGGGAGGGA---CACACACATAGA	24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36
DYS390	ATATTCTATCTA---TCATCTATCTAT ATAGATAGATGA---TAGATAGAATAT	17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29
DYS391	TCTGTCTGTCTG---TCTGCCTATCTG CAGATAGGCAGA---CAGACAGACAGA	5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16
DYS392	TCACCATTTAAT---TTACTAAGGAAT ATTCTTAGTAA---AATAAATGGTGA	4, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 20
DYS393	TTGTGCAATAC---GAGACATACCTC* GAGGTATGTCTC---GTATTGACACAA*	7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18
DYS437	ATGCCCATCCGG---TCATCTATCATC GATGATAGATGA---CCGGATGGGCAT	10, 11, 12, 13, 14, 15, 16, 17, 18, 19
DYS438	GTAACACAGTATA---TATTTGAAATGG CCATTTCAAATA---TATACTGTTAC	6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 18
DYS439	TGATAAATAGAA---GAAAGTATAAGT ACTTATACTTTC---TTCTATTATCA	5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19
DYS448	AGACATGGATAA---AGAGAGGTAAG CTTTACCTCTCT---TTATCCATGTCT	14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24
DYS456	TGTGATAATGTA---ATTCATTAGT AACTAATGGAAT---TACATTATCACA	9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 22, 23
DYS458	AAGAAAAGGAAG---GGAGGTTGGGCG CGCCACCCTCC---CTTCTTTTCTT	11, 12, 13, 14, 15, 16, 16.2, 17, 17.2, 18, 18.2, 19, 19.2, 20, 21, 22, 23, 24
DYS481	TTCAGCATGCTG---GAGTCTTGCAAC* GTTGCAAGACTC---CAGCATGCTGAA*	16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32
DYS533	TAGCTAGCTATC---ATCATCTATCAT ATGATAGATGAT---GATAGCTAGCTA	7, 8, 9, 10, 11, 12, 13, 14, 15
DYS549	GATTAGAAAGAT---GAAAAAATCTAC GTAGATTTTTC---ATCTTTTAATC	9, 10, 11, 12, 13, 14, 15
DYS570	CTCCAAGTTCTC---TTTTGTAGATA TATCTACAAAA---AGGAACITGGAG	12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24
DYS576	CATCTCTGAATA---AAAAGCCAAGA TCTTGGCTTTTT---TATTCAGAGATG	11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22
DYS635	ATAGAATCTCTA---TCACATTTTCTT* AAGAAAATGTGA---TAGAGATTCTAT*	13, 16, 17, 18, 19, 20, 21, 21.3, 22, 23, 24, 25, 26, 27, 28, 29, 30
DYS643	AAACTACTGTGC---CTTCTTTTAA TTAAAAAGAAAG---GCACAGTAGTTT	7, 8, 9, 10, 11, 12, 13, 14, 15
GATA H4	TAGGTAGGTAGG---ATGGATAGATTA TAATCTATCCAT---CCTACTACCTA	8, 9, 10, 11, 12, 13, 14, 15

sequence data output file for Sample 1 revealed that the “15” allele at locus D3S1358 consisted of both the “TCTA(TCTG)<sub>2</sub>(TCTA)<sub>12</sub>” variant and the “TCTA(TCTG)<sub>3</sub>(TCTA)<sub>11</sub>” variant, at a ratio of approximately 1:1. Conversely, the sequence data output file for Sample 2, which also displays a “15” allele at this locus, revealed that only the “TCTA(TCTG)<sub>2</sub>(TCTA)<sub>12</sub>” variant was present. This difference would not be detected by traditional CE methods. The time required for analysis, using a paired-end analysis method that detected both autosomal and Y chromosome STR alleles, ranged from 16 min (for dual 395 MB input files) to 285 min (for dual 7.9 GB input files) on the 16-core server used for this study.

#### 4. Discussion

The results of this study show the efficiency and accuracy of STR allele detection with STRait Razor. They also reveal the close relationship between the software, the library preparation chemistries, and the sequencing platforms used to produce the sequence information.

Read length is arguably the most important factor (followed by coverage) that impacts the allelic detection capability of

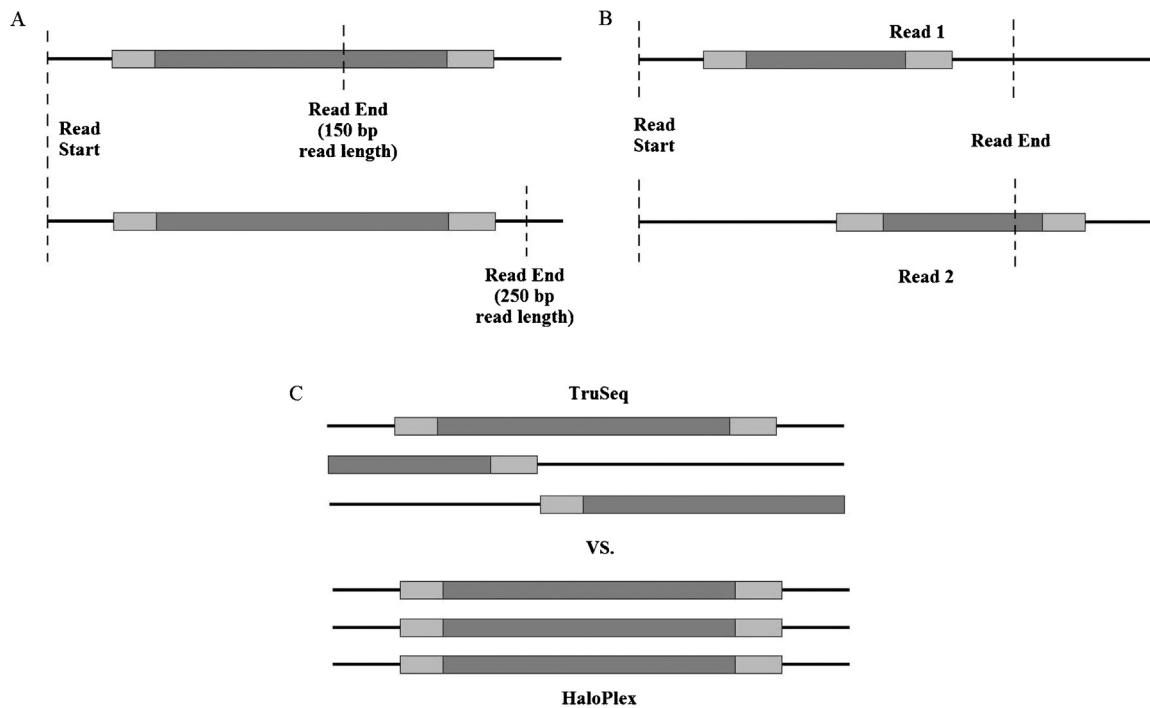
STRait Razor. The HaloPlex<sup>TM</sup> chemistry relies on enzymatic cleavage [24], which creates fragments with consistent start and end points. Depending on the length of the allele in question and the position of the repeat region within the resulting fragment(s), it is possible for sequencing reads to be produced that only partially span the repeat region and its associated flanking sequences (Fig. 2a). Since STRait Razor requires reads that contain all of this information, alleles may go undetected without complete repeat region traversal. An example of this phenomenon is locus D2S1338 in Sample 1 (HaloPlex<sup>TM</sup> preparation, GAIIX<sup>TM</sup> sequencing). For this locus, the “18” allele was called, but the “25” allele was too long to be covered completely by the sequencing read, and thus went undetected. When this same sample subsequently was sequenced on the MiSeq<sup>TM</sup> platform using a longer read length [17,18], the “25” allele was detected. In addition, if a repeat region is situated toward the beginning of a HaloPlex<sup>TM</sup> fragment, the allele is likely to be detected in one direction of a paired-end analysis. However, when the reads are sequenced from the opposite direction, the repeat region is oriented toward the end of the read and may not be completely encompassed (Fig. 2b). This situation can be seen in loci such as D7S820 and vWA in Sample

**Table 3**  
Comparison of CE allele calls and STRait Razor results – Autosomal STRs. Alleles detected by both CE and STRait Razor analysis of SGS data are shown in bold in the columns for each sample. The numbers of reads in which an allele was detected by STRait Razor are listed in parentheses next to the respective allele. The first number in parentheses represents the abundance of the allele in Read 1 of the paired-end sequencing run, while the second number represents the abundance of the allele in Read 2. Alleles not detected by STRait Razor due to lack of relevant sequence data are denoted by “[–]”.

		Detection method											
		Sample 1				Sample 2		Sample 3		Sample 4		Sample 5	
		CE	STRait Razor (TruSeq™ prep, GAIIx™)	STRait Razor (HaloPlex™ prep, GAIIx™)	STRait Razor (HaloPlex™ prep, MiSeq™)	CE	STRait Razor (HaloPlex™ prep, GAIIx™)	CE	STRait Razor (TruSeq™ prep, MiSeq™)	CE	STRait Razor (TruSeq™ prep, MiSeq™)	CE	STRait Razor (TruSeq™ prep, MiSeq™)
S T R loci	CSF1PO	<b>12</b>	<b>12</b> (332, 315)	<b>12</b> (1135, 976)	<b>12</b> (324, 322)	<b>12</b>	<b>12</b> (2445, 4155)	<b>10, 11</b>	<b>10</b> (89, 84), <b>11</b> (54, 75)	<b>12</b>	<b>12</b> (304, 310)	<b>10, 12</b>	<b>10</b> (152, 189), <b>12</b> (162, 145)
	D13S317	<b>12</b>	<b>12</b> (104, 96)	<b>12</b> (639, 0)	<b>12</b> (217, 214)	<b>8, 11</b>	<b>8</b> (2513, 4888), <b>11</b> (1995, 4377)	<b>11, 12</b>	<b>11</b> (49, 61), <b>12</b> (50, 48)	<b>12</b>	<b>12</b> (220, 209)	<b>8, 11</b>	<b>8</b> (101, 107), <b>11</b> (86, 76)
	D16S539	<b>10, 12</b>	<b>10</b> (86, 70), <b>12</b> (64, 43)	<b>10</b> (886, 2896), <b>12</b> (686, 1546)	<b>10</b> (181, 135), <b>12</b> (174, 138)	<b>9, 12</b>	<b>9</b> (1733, 7394), <b>12</b> (1165, 3209)	<b>8, 13</b>	<b>8</b> (62, 69), <b>13</b> (51, 56)	<b>9, 12</b>	<b>9</b> (121, 110), <b>12</b> (82, 77)	<b>11, 13</b>	<b>11</b> (117, 98), <b>13</b> (91, 73)
	D18S51	<b>15, 21</b>	<b>15</b> (25, 29), <b>21</b> (16, 17)	<b>15</b> (598, 3258), <b>21</b> (341, 541)	<b>15</b> (123, 123), <b>21</b> (106, 105)	<b>15, 16</b>	<b>15</b> (957, 5711), <b>16</b> (819, 4734)	<b>14, 16</b>	<b>14</b> (20, 16), <b>16</b> (13, 17)	<b>15, 18</b>	<b>15</b> (34, 22), <b>18</b> (17, 17)	<b>12</b>	<b>12</b> (94, 79)
	D19S433	<b>12, 14</b>	<b>12</b> (39, 19), <b>14</b> (23, 13)	<b>12</b> (589, 2620), <b>14</b> (542, 1543)	<b>12</b> (141, 138), <b>14</b> (126, 133)	<b>13, 14</b>	<b>13</b> (1634, 6001), <b>14</b> (1423, 4149)	<b>13, 14</b>	<b>13</b> (16, 5), <b>14</b> (10, 11)	<b>13</b>	<b>13</b> (44, 33)	<b>12, 15</b>	<b>12</b> (12, 19), <b>15</b> (28, 17)
	D21S11	<b>29, 30</b>	[–], [–]	[–], [–]	<b>29</b> (70, 17), <b>30</b> (57, 14)	<b>29, 31</b>	[–], [–]	<b>29</b>	<b>29</b> (17, 9)	<b>29</b>	<b>29</b> (38, 25)	<b>28, 29</b>	<b>28</b> (16, 11), <b>29</b> (16, 11)
	D2S1338	<b>18, 25</b>	<b>18</b> (19, 18), <b>25</b> (9, 0)	<b>18</b> (80, 1), [–]	<b>18</b> (54, 52), <b>25</b> (29, 25)	<b>20, 25</b>	<b>20</b> (137, 1), [–]	<b>18</b>	<b>18</b> (53, 53)	<b>17, 24</b>	<b>17</b> (74, 67), <b>24</b> (35, 35)	<b>20, 23</b>	<b>20</b> (70, 70), <b>23</b> (55, 37)
	D3S1358	<b>15</b>	<b>15</b> (109, 110)	<b>15</b> (449, 4079)	<b>15</b> (439, 395)	<b>15, 16</b>	<b>15</b> (1046, 4579), <b>16</b> (852, 3995)	<b>16, 18</b>	<b>16</b> (32, 29), <b>18</b> (36, 29)	<b>16, 18</b>	<b>16</b> (71, 74), <b>18</b> (78, 66)	<b>14, 17</b>	<b>14</b> (99, 93), <b>17</b> (72, 84)
	D5S818	<b>11</b>	<b>11</b> (76, 66)	[–]	<b>11</b> (66, 67)	<b>11</b>	[–]	<b>11, 12</b>	<b>11</b> (20, 23), <b>12</b> (30, 26)	<b>11</b>	<b>11</b> (101, 86)	<b>12, 13</b>	<b>12</b> (56, 62), <b>13</b> (45, 34)
	D7S820	<b>9, 12</b>	<b>9</b> (1, 2), <b>12</b> (5, 4)	<b>9</b> (0, 1754), <b>12</b> (0, 1650)	<b>9</b> (70, 70), <b>12</b> (57, 57)	<b>10</b>	<b>10</b> (0, 7733)	<b>10, 11</b>	<b>10</b> (4, 3), <b>11</b> (3, 4)	<b>9, 13</b>	<b>9</b> (3, 7), <b>13</b> (8, 8)	<b>11</b>	<b>11</b> (17, 15)
	D8S1179	<b>13</b>	<b>13</b> (98, 87)	<b>13</b> (1722, 5068)	<b>13</b> (335, 220)	<b>13</b>	<b>13</b> (3941, 10527)	<b>10, 12</b>	<b>10</b> (36, 38), <b>12</b> (24, 25)	<b>12, 14</b>	<b>12</b> (52, 47), <b>14</b> (48, 49)	<b>13, 16</b>	<b>13</b> (68, 78), <b>16</b> (58, 41)
	FGA	<b>20, 21</b>	<b>20</b> (36, 25), <b>21</b> (26, 20)	<b>20</b> (770, 3819), <b>21</b> (581, 3419)	<b>20</b> (168, 129), <b>21</b> (146, 134)	<b>19, 21</b>	<b>19</b> (1019, 5066), <b>21</b> (890, 4849)	<b>23, 25</b>	<b>23</b> (12, 16), <b>25</b> (12, 13)	<b>22, 24</b>	<b>22</b> (40, 30), <b>24</b> (18, 23)	<b>20</b>	<b>20</b> (84, 71)
	TH01	<b>9, 9.3</b>	<b>9</b> (160, 146), <b>9.3</b> (132, 178)	<b>9</b> (3172, 5571), <b>9.3</b> (2893, 5493)	<b>9</b> (260, 255), <b>9.3</b> (297, 297)	<b>8, 9.3</b>	<b>8</b> (5701, 8471), <b>9.3</b> (4963, 8350)	<b>7</b>	<b>7</b> (150, 162)	<b>9.3</b>	<b>9.3</b> (287, 292)	<b>9.3</b>	<b>9.3</b> (393, 390)
	TPOX	<b>8, 9</b>	<b>8</b> (90, 96), <b>9</b> (99, 80)	<b>8</b> (4832, 5208), <b>9</b> (4488, 4710)	<b>8</b> (527, 479), <b>9</b> (475, 428)	<b>11</b>	<b>11</b> (11043, 15943)	<b>8, 11</b>	<b>8</b> (35, 34), <b>11</b> (32, 29)	<b>8, 12</b>	<b>8</b> (113, 105), <b>12</b> (84, 86)	<b>8</b>	<b>8</b> (176, 185)
	vWA	<b>16, 17</b>	<b>16</b> (53, 37), <b>17</b> (39, 24)	<b>16</b> (299, 0), <b>17</b> (213, 0)	<b>16</b> (55, 55), <b>17</b> (36, 36)	<b>15, 20</b>	<b>15</b> (669, 0), <b>20</b> (0, 3)	<b>17</b>	<b>17</b> (49, 56)	<b>14, 17</b>	<b>14</b> (60, 60), <b>17</b> (57, 52)	<b>15, 17</b>	<b>15</b> (63, 65), <b>17</b> (56, 61)
	Penta D	<b>10</b>	<b>10</b> (24, 35)	<b>10</b> (388, 0)	<b>10</b> (180, 0)	<b>14, 15</b>	<b>14</b> (214, 0), [–]	<b>9, 11</b>	<b>9</b> (11, 14), <b>11</b> (12, 10)	<b>12, 14</b>	<b>12</b> (26, 23), <b>14</b> (25, 20)	<b>9, 12</b>	<b>9</b> (23, 24), <b>12</b> (29, 30)
	Penta E	<b>11, 12</b>	<b>11</b> (6, 7), <b>12</b> (9, 5)	<b>11</b> (98, 121), <b>12</b> (123, 104)	<b>11</b> (105, 87), <b>12</b> (107, 88)	<b>10, 11</b>	<b>10</b> (290, 993), <b>11</b> (243, 191)	<b>5, 7</b>	<b>5</b> (10, 7), <b>7</b> (6, 5)	<b>11, 12</b>	<b>11</b> (9, 4), <b>12</b> (11, 6)	<b>7, 14</b>	<b>7</b> (20, 19), <b>14</b> (8, 7)
	D10S1248	<b>13, 14</b>	<b>13</b> (106, 94), <b>14</b> (69, 67)	<b>13</b> (672, 4537), <b>14</b> (329, 3613)	<b>13</b> (192, 199), <b>14</b> (159, 171)	<b>13</b>	<b>13</b> (3545, 17065)	<b>12, 13</b>	<b>12</b> (81, 73), <b>13</b> (71, 69)	<b>13, 14</b>	<b>13</b> (134, 127), <b>14</b> (152, 134)	<b>16</b>	<b>16</b> (211, 211)
	D12S391	<b>15, 17</b>	<b>15</b> (50, 50), <b>17</b> (35, 31)	<b>15</b> (1464, 548), <b>17</b> (1101, 447)	<b>15</b> (167, 117), <b>17</b> (162, 106)	<b>21</b>	<b>21</b> (2123, 0)	<b>15, 24</b>	<b>15</b> (37, 31), <b>24</b> (23, 19)	<b>17, 21</b>	<b>17</b> (90, 82), <b>21</b> (84, 79)	<b>18, 20</b>	<b>18</b> (68, 73), <b>20</b> (58, 63)
	D1S1656	<b>16, 18.3</b>	<b>16</b> (63, 47), <b>18.3</b> (37, 52)	<b>16</b> (499, 0), <b>18.3</b> (468, 0)	<b>16</b> (59, 75), <b>18.3</b> (78, 41)	<b>11, 12</b>	<b>11</b> (2262, 2848), <b>12</b> (1872, 2284)	<b>16.3, 18.3</b>	<b>16.3</b> (33, 24), <b>18.3</b> (25, 29)	<b>15, 17.3</b>	<b>15</b> (79, 75), <b>17.3</b> (67, 52)	<b>12, 16</b>	<b>12</b> (98, 83), <b>16</b> (76, 67)
	D22S1045	<b>16, 17</b>	<b>16</b> (10, 16), <b>17</b> (7, 16)	<b>16</b> (901, 2637), <b>17</b> (718, 1639)	<b>16</b> (107, 105), <b>17</b> (108, 104)	<b>15, 16</b>	<b>15</b> (2042, 5133), <b>16</b> (1415, 3535)	<b>15, 16</b>	<b>15</b> (15, 14), <b>16</b> (13, 15)	<b>16</b>	<b>16</b> (35, 38)	<b>11, 16</b>	<b>11</b> (94, 94), <b>16</b> (40, 40)
	D2S441	<b>11, 15</b>	<b>11</b> (72, 49), <b>15</b> (32, 37)	<b>11</b> (291, 3863), <b>15</b> (0, 3685)	<b>11</b> (124, 145), <b>15</b> (113, 167)	<b>11.3, 14</b>	<b>11.3</b> (615, 7807), <b>14</b> (0, 7024)	<b>11, 14</b>	<b>11</b> (74, 71), <b>14</b> (68, 64)	<b>10, 11.3</b>	<b>10</b> (106, 100), <b>11.3</b> (149, 137)	<b>10, 11</b>	<b>10</b> (131, 125), <b>11</b> (130, 138)
	Total alleles	38	36	34	38	37	32	40	40	37	37	38	38

**Table 4**  
Comparison of CE allele calls and STRait Razor results – Y-Chromosome STRs. Alleles detected by both CE and STRait Razor analysis of SGS data are shown in bold in the columns for each sample. The numbers of reads in which an allele was detected by STRait Razor are listed in parentheses next to the respective allele. The first number in parentheses represents the abundance of the allele in Read 1 of the paired-end sequencing run, while the second number represents the abundance of the allele in Read 2. Alleles not detected by STRait Razor due to lack of relevant sequence data are denoted by “[–]”.

S T R loci		Detection method											
		Sample 1				Sample 2		Sample 3		Sample 4		Sample 5	
		CE	STRait Razor (TruSeq™ prep, GAIIx™)	STRait Razor (HaloPlex™ prep, GAIIx™)	STRait Razor (HaloPlex™ prep, MiSeq™)	CE	STRait Razor (HaloPlex™ prep, GAIIx™)	CE	STRait Razor (TruSeq™ prep, MiSeq™)	CE	STRait Razor (TruSeq™ prep, MiSeq™)	CE	STRait Razor (TruSeq™ prep, MiSeq™)
DYS19	<b>14</b>	<b>14</b> (7, 11)	<b>14</b> (231, 1221)	<b>14</b> (61, 62)	<b>15</b>	<b>15</b> (6, 2341)	<b>14</b>	<b>14</b> (8, 18)	<b>16</b>	<b>16</b> (13, 13)	<b>14</b>	<b>14</b> (14, 15)	
DYS385	<b>11, 13</b>	<b>11</b> ( <b>12, 12</b> ), <b>13</b> ( <b>12, 11</b> )	[–], [–]	[–], [–]	<b>11, 14</b>	[–], [–]	<b>11, 14</b>	<b>11</b> (5, 3), <b>14</b> (4, 5)	<b>11, 14</b>	<b>11</b> (22, 29), <b>14</b> (8, 19)	<b>11, 14</b>	<b>11</b> (12, 13), <b>14</b> (24, 16)	
DYS389I	<b>13</b>	<b>13</b> (102, 99)	<b>13</b> (423, 1)	<b>13</b> (26, 0)	<b>13</b>	<b>13</b> (1605, 0)	<b>13</b>	<b>13</b> (54, 47)	<b>13</b>	<b>13</b> (137, 137)	<b>13</b>	<b>13</b> (120, 103)	
DYS389II	<b>29</b>	[–]	[–]	<b>29</b> (23, 0)	<b>30</b>	[–]	<b>29</b>	<b>29</b> (9, 3)	<b>28</b>	<b>28</b> (25, 20)	<b>29</b>	<b>29</b> (13, 17)	
DYS390	<b>24</b>	<b>24</b> (14, 16)	<b>24</b> (115, 3495)	<b>24</b> (115, 126)	<b>25</b>	<b>25</b> (259, 3386)	<b>24</b>	<b>24</b> (16, 19)	<b>23</b>	<b>23</b> (37, 43)	<b>23</b>	<b>23</b> (40, 37)	
DYS391	<b>10</b>	<b>10</b> (175, 167)	<b>10</b> (952, 39)	<b>10</b> (99, 24)	<b>10</b>	<b>10</b> (3179, 1431)	<b>10</b>	<b>10</b> (73, 80)	<b>10</b>	<b>10</b> (180, 182)	<b>12</b>	<b>12</b> (166, 165)	
DYS392	<b>13</b>	<b>13</b> (3, 8)	<b>13</b> (885, 1965)	<b>13</b> (82, 78)	<b>11</b>	<b>11</b> (1466, 2850)	<b>13</b>	<b>13</b> (7, 8)	<b>13</b>	<b>13</b> (7, 8)	<b>13</b>	<b>13</b> (11, 10)	
DYS393	<b>13</b>	<b>13</b> (9, 2)	<b>13</b> (0, 360)	<b>13</b> (14, 13)	<b>14</b>	<b>14</b> (0, 1023)	<b>13</b>	<b>13</b> (2, 7)	<b>13</b>	<b>13</b> (2, 3)	<b>13</b>	<b>13</b> (10, 10)	
DYS437	<b>15</b>	<b>15</b> (85, 70)	<b>15</b> (0, 4020)	<b>15</b> (247, 238)	<b>14</b>	<b>14</b> (0, 11,064)	<b>15</b>	<b>15</b> (77, 77)	<b>15</b>	<b>15</b> (148, 133)	<b>14</b>	<b>14</b> (141, 146)	
DYS438	<b>12</b>	<b>12</b> (42, 36)	<b>12</b> (324, 285)	<b>12</b> (62, 32)	<b>11</b>	<b>11</b> (884, 871)	<b>12</b>	<b>12</b> (48, 49)	<b>12</b>	<b>12</b> (79, 68)	<b>12</b>	<b>12</b> (96, 93)	
DYS439	<b>12</b>	[–]	<b>12</b> (428, 2296)	<b>12</b> (134, 78)	<b>10</b>	<b>10</b> (1789, 6072)	<b>12</b>	<b>12</b> (2, 0)	<b>11</b>	<b>11</b> (2, 2)	<b>13</b>	<b>13</b> (3, 1)	
DYS448	<b>19</b>	[–]	[–]	<b>19</b> (17, 5)	<b>19</b>	[–]	<b>19</b>	<b>19</b> (11, 3)	<b>19</b>	<b>19</b> (21, 16)	<b>18</b>	<b>18</b> (12, 11)	
DYS456	<b>15</b>	<b>15</b> (10, 13)	<b>15</b> (523, 1402)	<b>15</b> (80, 54)	<b>14</b>	<b>14</b> (2723, 6296)	<b>15</b>	<b>15</b> (13, 10)	<b>16</b>	<b>16</b> (11, 9)	<b>15</b>	<b>15</b> (22, 21)	
DYS458	<b>17</b>	<b>17</b> (10, 6)	<b>17</b> (56, 258)	<b>17</b> (31, 21)	<b>15</b>	<b>15</b> (152, 1300)	<b>19</b>	<b>19</b> (11, 9)	<b>17</b>	<b>17</b> (13, 12)	<b>17</b>	<b>17</b> (17, 24)	
DYS481	<b>22</b>	<b>22</b> (19, 21)	<b>22</b> (227, 1943)	<b>22</b> (150, 101)	<b>23</b>	<b>23</b> (663, 3830)	<b>25</b>	<b>25</b> (20, 18)	<b>23</b>	<b>23</b> (30, 43)	<b>22</b>	<b>22</b> (50, 45)	
DYS533	<b>12</b>	<b>12</b> (26, 36)	<b>12</b> (184, 0)	<b>12</b> (30, 10)	<b>12</b>	<b>12</b> (511, 0)	<b>12</b>	<b>12</b> (37, 37)	<b>12</b>	<b>12</b> (34, 37)	<b>14</b>	<b>14</b> (41, 27)	
DYS549	<b>13</b>	<b>13</b> (44, 49)	<b>13</b> (743, 0)	<b>13</b> (151, 101)	<b>12</b>	<b>12</b> (2649, 0)	<b>13</b>	<b>13</b> (39, 38)	<b>13</b>	<b>13</b> (46, 58)	<b>13</b>	<b>13</b> (60, 62)	
DYS570	<b>17</b>	<b>17</b> (44, 51)	<b>17</b> (646, 0)	<b>17</b> (76, 28)	<b>20</b>	<b>20</b> (777, 0)	<b>18</b>	<b>18</b> (64, 69)	<b>17</b>	<b>17</b> (73, 66)	<b>17</b>	<b>17</b> (145, 134)	
DYS576	<b>19</b>	<b>19</b> (45, 43)	<b>19</b> (0, 341)	<b>19</b> (3, 11)	<b>17</b>	<b>17</b> (0, 512)	<b>19</b>	<b>19</b> (34, 22)	<b>18</b>	<b>18</b> (74, 74)	<b>18</b>	<b>18</b> (65, 57)	
DYS635	<b>24</b>	<b>24</b> (10, 5)	<b>24</b> (220, 71)	<b>24</b> (18, 20)	<b>25</b>	<b>25</b> (774, 700)	<b>23</b>	<b>23</b> (18, 15)	<b>24</b>	<b>24</b> (33, 29)	<b>23</b>	<b>23</b> (30, 27)	
DYS643	<b>10</b>	<b>10</b> (43, 25)	<b>10</b> (2392, 989)	<b>10</b> (221, 111)	<b>10</b>	<b>10</b> (3314, 1407)	<b>11</b>	<b>11</b> (10, 15)	<b>11</b>	<b>11</b> (34, 30)	<b>10</b>	<b>10</b> (46, 43)	
GATA H4	<b>12</b>	<b>12</b> (23, 21)	<b>12</b> (290, 2196)	<b>12</b> (97, 93)	<b>12</b>	<b>12</b> (511, 3795)	<b>11</b>	<b>11</b> (21, 27)	<b>11</b>	<b>11</b> (34, 47)	<b>11</b>	<b>11</b> (33, 40)	
Total alleles	23	20	19	21	23	19	23	23	23	23	23	23	



**Fig. 2.** Read length-related issues. In this figure, the dark gray bars represent the repeat region, while the light gray bars represent the flanking regions. The bold black lines represent surrounding sequence data. (A) Two identical fragments, such as those produced by the HaloPlex™ chemistry, are sequenced with 2 different read lengths. (B) A repeat region situated within a fragment is sequenced toward the beginning of the read in Read 1 of a paired-end sequencing run, but is sequenced toward the end of the read in Read 2 of the paired-end run. (C) A comparison of repeat region location consistency between the TruSeq™ and HaloPlex™ library preparation chemistries. This illustration is an example of a small subset of sequence reads, and does not reflect the actual proportions of reads produced by each library preparation method.

1, where the alleles are detected only in one set of paired-end reads and not the other. Some library preparation redesign may overcome the truncated repeat region reads.

The TruSeq™ chemistry [25] is less prone to these issues because the random fragmentation of DNA allows for a much more diverse positioning of repeat regions within the resulting fragments (Fig. 2c). Therefore, there is a higher likelihood of at least some reads encompassing the entire repeat region. This design explains why the majority of the alleles that were not detected by STRait Razor for Sample 1 following HaloPlex™ preparation and GAIIX™ sequencing were detected normally by the program for this sample following TruSeq™ preparation on the same instrument. Despite this advantage, the non-enzymatic random fragmentation employed by the TruSeq™ chemistry may result in lower read counts for some alleles in comparison with HaloPlex™, due to the fewer resulting fragments containing the complete repeat region. In some cases, the random fragmentation method simply may not create any fragments that contain the repeat region of interest. This limitation may explain the undetected alleles in Sample 1 at loci DYS439 and DYS448 following TruSeq™ preparation and GAIIX™ sequencing. Coverage depth differences in the results of this study, however, also may be explained by the fact that the regions targeted by the TruSeq™ kit for these trials were, by design, approximately 100 times larger than those targeted by the HaloPlex™ kit. Finally, it should be noted that in cases where the larger allele of a heterozygous pair at a given locus is not detected due to read length issues, a detected stutter allele may give the impression of a true heterozygous allele. Such issues may eventually be overcome by a combination of secondary statistical allelotyping software that takes into account the allele coverage ratios, and customization of the allelic information modules utilized by STRait Razor for the specific chemistry and instrumentation used for sequencing.

It should be noted that sequence variations (e.g., insertions or deletions) that reside outside of the flanking regions used by STRait Razor but within the primer-binding sites used by commercial STR/CE kits may result in discordant results between SGS and CE analyses. While this discordance was not observed during the course of this study, most likely due to limited sampling, it is not unique to SGS. Such discrepancies can occur due to different primer-binding site locations between various constructs of CE-based kits, as well.

The small portions of repeat sequence that are exactly matched to reads during the filtering portion of the algorithm were deliberately kept short (~3 simple repeat units or less). This exact matching step is important, as it removes irrelevant sequence data that may have been captured due to chance flanking sequence homology. However, the use of exact matching can potentially remove repeat regions with intra-repeat variation. Thus, sequences of repeats that are shorter than the smallest known allele at each locus are used, so that they may align at various points throughout the captured reads and still allow reads with intra-repeat variation to be retained. The Penta D locus, however, requires special attention, because its smallest known allele is “1.1” [26], albeit an extremely rare allele. The repeat sequence used by STRait Razor for exact match filtering at this locus is only one repeat unit long, since using a longer sequence would not allow for intra-repeat variability detection in the smaller alleles. The use of such a small portion of repeat sequence results in less thorough filtering of the reads, but this does not appear to affect the concordance of the resulting allelic data.

While the results of this study are encouraging, it should be noted that all software has its strengths and limitations. STRait Razor is simply an alternative method of approaching an increasingly important facet of SGS data analysis. The software is being offered freely so that those working with SGS and STR detection have an additional tool to facilitate their analyses. STRait



Razor is a fairly straightforward program, and the authors encourage enhancements to this software. With this phase of STRait Razor, probabilities of allele calls are not provided, nor does the program perform any stutter filtering or attempt to differentiate between heterozygosity and homozygosity. The software only reports the allele calls and related coverage at each locus; the analyst makes the final decision in an informed manner. A separate program is currently in development that reads the colon-delimited text file(s) output by STRait Razor and detects the two allele calls with the highest count values (if any), comparing them based on a user-defined “abundance ratio” to make homozygous and heterozygous allelotype calls with confidence. Presently, STRait Razor does not filter reads based on FASTQ quality information. However, the flanking region detection and small-repeat match verification performed by STRait Razor can be considered inherent filtering steps. Quality-based filtering may be incorporated into future versions of the program. In its current format, STRait Razor is designed to detect all the known alleles at each of the tested loci, according to the allelic information listed in STRBase and the Y-Chromosome Haplotype Reference Database (YHRD) [27,28]. These alleles are defined in the modules used by the software and may be modified by the user to include other rare or undocumented variants. In the future, the software may be altered to allow for the intuitive calling of alleles based on repeat length alone, without the need for allelic definitions.

## 5. Conclusion

In its current state of development, STRait Razor offers forensic DNA analysts a simple and effective means of detecting STR alleles in SGS data. Its ability to do so is limited primarily by the capabilities of the library preparation chemistries and sequencing platforms used to generate the raw data. The software has been shown to properly and efficiently analyze data resulting from a variety of commercially available library preparation chemistries and sequencers. The results of this study also indicate that the software is capable of identifying alleles in SGS sequence data with 100% concordance. STRait Razor provides multiple options for customization, and new loci can be added for detection at the discretion of the analyst. The program is freely available (see Supplementary Materials), and the scientific community is encouraged to make improvements to STRait Razor to increase its utility as a forensic tool.

Supplementary data associated with this article can be found, in the online version, at <http://dx.doi.org/10.1016/j.fsi-gen.2013.04.005>.

## Acknowledgments

This project was supported in part by Award No. 2012-DN-BX-K033, awarded by the National Institute of Justice, Office of Justice Programs, U.S. Department of Justice. The opinions, findings, and conclusions or recommendations expressed in this publication are those of the author(s) and do not necessarily reflect those of the U.S. Department of Justice.

The authors would like to thank Illumina, Inc. and Agilent Technologies, Inc. for providing their technical expertise and contributing a portion of the sample preparation and sequencing chemistries used in this study.

## References

- [1] A. Edwards, A. Civitello, H.A. Hammond, C.T. Caskey, DNA typing and genetic mapping with trimeric and tetrameric tandem repeats, *Am. J. Hum. Genet.* 49 (1991) 746–756.
- [2] A. Edwards, H.A. Hammond, L. Jin, C.T. Caskey, R. Chakraborty, Genetic variation at five trimeric and tetrameric repeat loci in four human population groups, *Genomics* 12 (1992) 241–253.
- [3] H. Ellegren, Microsatellites: simple sequences with complex evolution, *Nat. Rev. Genet.* 5 (2004) 435–445.
- [4] K. Lazaruk, P.S. Walsh, F. Oaks, D. Gilbert, B.B. Rosenblum, S. Menchen, D. Scheibler, H.M. Wenz, C. Holt, J. Wallin, Genotyping of forensic short tandem repeat (STR) systems based on sizing precision in a capillary electrophoresis instrument, *Electrophoresis* 19 (1998) 86–93.
- [5] P. Gill, D.J. Werrett, B. Budowle, R. Guerrieri, An assessment of whether SNPs will replace STRs in national DNA databases – joint considerations of the DNA working group of the European Network of Forensic Science Institutes (ENFSI) and the Scientific Working Group on DNA Analysis Methods (SWGDM), *Sci. Justice.* 44 (2004) 51–53.
- [6] P.J. Collins, L.K. Hennessy, C.S. Leibel, R.K. Roby, D.J. Reeder, P.A. Foxall, Developmental validation of a single-tube amplification of the 13 CODIS STR loci, D2S1338, D19S433, and amelogenin: the AmpFISTR® Identifier® PCR Amplification Kit, *J. Forensic Sci.* 49 (2004) 1265–1277.
- [7] K. Oostdik, J. French, D. Yet, B. Smalling, C. Nolde, P.M. Vallone, E.L. Butts, C.R. Hill, M.C. Kline, T. Rinta, A.M. Gerow, S.R. Allen, C.K. Huber, J. Teske, B. Krenke, M. Ensenberger, P. Fulmer, C. Sprecher, Developmental validation of the PowerPlex® 18D system, a rapid STR multiplex for analysis of reference samples, *Forensic Sci. Int. Genet.* 7 (2013) 129–135.
- [8] M. Gymrek, D. Golan, S. Rosset, Y. Erlich, lobSTR: a short tandem repeat profiler for personal genomes, *Genome Res.* 22 (2012) 1154–1162.
- [9] D.M. Bornman, M.E. Hester, J.M. Schuetter, M.D. Kasoji, A. Minard-Smith, C.A. Barden, S.C. Nelson, G.D. Godbold, C.H. Baker, B. Yang, J.E. Walther, I.E. Tornes, P.S. Yan, B. Rodriguez, R. Bundschuh, M.L. Dickens, B.A. Young, S.A. Faith, Short-read, high-throughput sequencing technology for STR genotyping, *Biotechniques – Rapid Dispatches* (2012) 1–6.
- [10] S.L. Fordyce, M.C. Avila-Arcos, E. Rockenbauer, C. Børsting, R. Frank-Hansen, F.T. Petersen, E. Willerslev, A.J. Hansen, N. Morling, M.T. Gilbert, High-throughput sequencing of core STR loci for forensic genetic investigations using the Roche Genome Sequencer FLX platform, *Biotechniques* 51 (2011) 127–133.
- [11] M.M. Holland, M.R. McQuillan, K.A. O’Hanlon, Second generation sequencing allows for mtDNA mixture deconvolution and high resolution detection of heteroplasmy, *Croat. Med. J.* 52 (2011) 299–313.
- [12] R. Nielsen, J.S. Paul, A. Albrechtsen, Y.S. Song, Genotype and SNP calling from next-generation sequencing data, *Nat. Rev. Genet.* 12 (2011) 443–451.
- [13] D.W. Craig, J.V. Pearson, S. Szlinger, A. Sekar, M. Redman, J.J. Corneveaux, T.L. Pawlowski, T. Laub, G. Nunn, D.A. Stephan, N. Homer, M.J. Huentelman, Identification of genetic variants using bar-coded multiplexed sequencing, *Nat. Methods* 5 (2008) 887–893.
- [14] D.C. Koboldt, K. Chen, T. Wylie, D.E. Larson, M.D. McLellan, E.R. Mardis, G.M. Weinstein, R.K. Wilson, L. Ding, VarScan: variant detection in massively parallel sequencing of individual and pooled samples, *Bioinformatics* 25 (2009) 2283–2285.
- [15] A. McKenna, M. Hanna, E. Banks, A. Sivachenko, K. Cibulskis, A. Kernysky, K. Garimella, D. Altshuler, S. Gabriel, M. Daly, M.A. DePristo, The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data, *Genome Res.* 20 (2010) 1297–1303.
- [16] E.C. Berglund, A. Kiialainen, A.C. Syvänen, Next-generation sequencing technologies and applications for human genetic history and forensics, *Investig. Genet.* 2 (2011) 1–15.
- [17] Illumina® GAlx™ Specifications: [http://www.illumina.com/documents/products/datasheets/datasheet\\_genome\\_analyzerllx.pdf](http://www.illumina.com/documents/products/datasheets/datasheet_genome_analyzerllx.pdf).
- [18] Illumina® MiSeq™ Specifications: [http://www.illumina.com/documents/products/datasheets/datasheet\\_miseq.pdf](http://www.illumina.com/documents/products/datasheets/datasheet_miseq.pdf).
- [19] B. Langmead, C. Trapnell, M. Pop, S.L. Salzberg, Ultrafast and memory-efficient alignment of short DNA sequences to the human genome, *Genome Biol.* 10 (2009) R25.
- [20] CASAVA v.1.8.2: [http://support.illumina.com/downloads/casava\\_182.ilmn](http://support.illumina.com/downloads/casava_182.ilmn).
- [21] MiSeq Reporter: [http://support.illumina.com/sequencing/sequencing\\_software/miseq\\_reporter/downloads.ilmn](http://support.illumina.com/sequencing/sequencing_software/miseq_reporter/downloads.ilmn).
- [22] AGREP: <http://laurikari.net/tre/download/>.
- [23] PPSS: <http://code.google.com/p/ppss/>.
- [24] Agilent Technologies® HaloPlex™ Specifications: <http://www.genomics.agilent.com/GenericB.aspx?pagetype=Custom&subpagetype=Custom&pageid=3081>.
- [25] Illumina TruSeq Specifications: [http://www.illumina.com/Documents/%5Cproducts%5Cdatasheets%5Cdatasheet\\_truseq\\_custom\\_enrichment\\_kit.pdf](http://www.illumina.com/Documents/%5Cproducts%5Cdatasheets%5Cdatasheet_truseq_custom_enrichment_kit.pdf).
- [26] Penta D Facts Sheet (STRBase): [http://www.cstl.nist.gov/strbase/str\\_Penta\\_D.htm](http://www.cstl.nist.gov/strbase/str_Penta_D.htm).
- [27] STRBase: [http://www.cstl.nist.gov/strbase/str\\_fact.htm](http://www.cstl.nist.gov/strbase/str_fact.htm).
- [28] Y-Chromosome Haplotype Reference Database: <http://www.yhrd.org/Research/Loci>.