

STR Sequence Diversity in Population Samples and Nomenclature Guidance for the “Next Generation”

Katherine B. Gettings¹, Seth A. Faith², Brian Young², Esley Heizer Jr.², Kevin M. Kiesler¹, Elizabeth Montano², Christine Baker², Angela Minard-Smith², Richard Guerrieri² and Peter M. Vallone¹

¹National Institute of Standards and Technology, 100 Bureau Drive, Gaithersburg, MD 20899

²Battelle, 505 King Avenue, Columbus, OH 43201

ABSTRACT

As STR loci were being identified in the 1990s, various nomenclature systems were developed for different loci, with the primary variation being whether or not to “count” non-repeat bases interspersed in the repeat motif. In 1997, the ISFG issued guidelines on STR nomenclature, in an attempt to provide a common currency for information exchange. Historical precedent already existed for some loci, and this was maintained to avoid confusion, resulting in several commonly used forensic loci having complicated and contradictory nomenclature systems. This has not been an issue within the forensic community, as the capillary electrophoresis (CE)-length analyses are kit-based, with corresponding computer programs that automatically count repeats in a standardized manner. Now, as the costs associated with next-generation sequencing (NGS) methods decline, forensic research laboratories are beginning to explore the increase in information sequencing STR loci may provide. As a new generation of scientists begins interrogating these loci on a deeper level, an understanding of historical nomenclature is needed to achieve bioinformatic concordance with existing CE data. In the work presented here, NGS results from population samples exemplify the sequence variation that exists in forensic STR loci (SNPs and InDels within and outside of STR allele regions and repeat motif changes) as well as the complexity and inconsistency of the current nomenclature. This experimental sequence data gives an indication of the level of diversity expected in the larger population and provides examples of how sub-alleles can improve discrimination and mixture deconvolution in forensic casework. The different purposes of nomenclature—manual comparisons, forensic reports, database searching, court explanations—are discussed and examples of possible NGS-compatible nomenclature systems that may meet the needs of the forensic community are shown.

MATERIALS and METHODS



Figure 1. Overview of experimental design.

NIST population samples (N=183 consisting of 70 Caucasian, 68 African American, and 45 Hispanic individuals) were amplified twice in 96-well plates, with 0.5 ng input DNA per sample in 25 µL reaction volumes. Duplicate amplicons were combined during the clean-up step, prior to library generation. Sequencing template libraries were prepared in 96-well format with the TruSeq DNA PCR-Free Sample Preparation Kit HS (Illumina, San Diego CA, USA). Sequencing was performed in two runs (96 samples/run) on the MiSeq system (Illumina) using the 600 cycle MiSeq Reagent Kit v3 (Illumina).

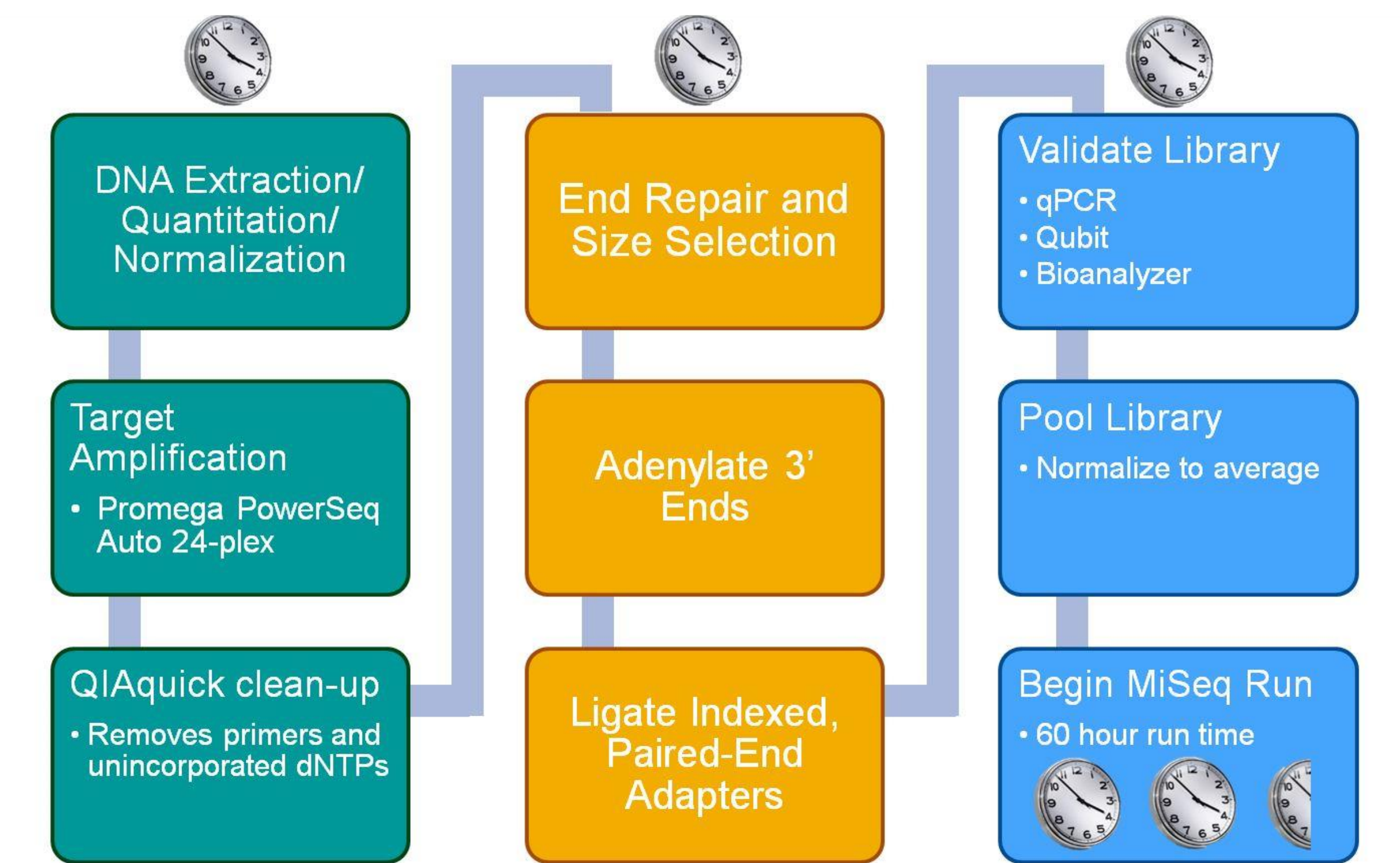


Figure 2. Overview of sample preparation workflow. Each clock represents a minimum of one day. All steps other than the actual MiSeq run were performed manually in a 96-well format.

Analysis of .fastq files to produce STR allele calls was performed with two different bioinformatic pipelines: ExactID (Battelle Memorial Institute, Columbus OH, USA, see ISHI 2014 poster #69 for more information), and STRait Razor [1]. Allelic balance based on coverage was evaluated to determine zygosity. Only majority sequences (two for heterozygotes or one for homozygotes) were considered as evidence supporting allele calls, and only the repeat regions of the majority sequences were analyzed further (e.g. sequences that were consistent with stutter, and sequences that did not match the majority sequence within the repeat region were excluded from further analysis). Genotypes from both ExactID and STRaitRazor were independently analyzed for concordance to CE based genotypes (generated previously with PowerPlex Fusion (Promega)). Discordances were evaluated further to determine the true genotype/sequence.

RESULTS

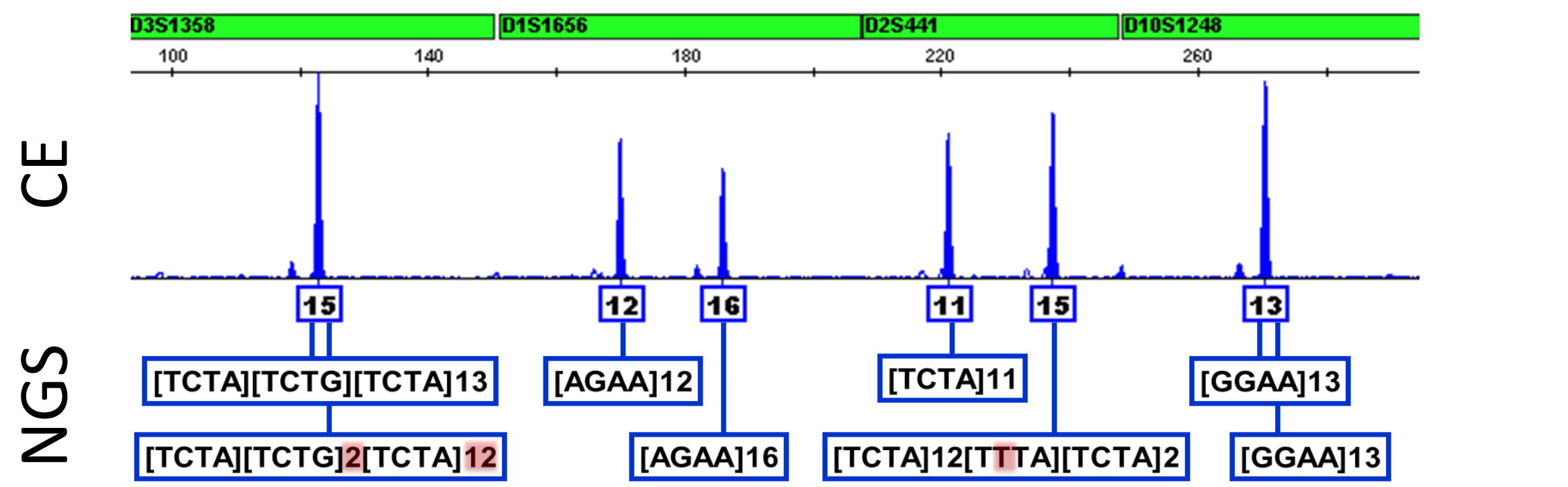


Figure 3. CE (PowerPlex Fusion, four loci shown) electropherogram for one sample of the 183 tested. Below the CE-based allele designations are the sequences obtained. D3S1358 is homozygous by length but heterozygous by sequence, D151656 is a simple repeat heterozygote, D2S441 has a simple repeat 11 allele and a different motif caused by a C→T SNP at the 15 allele, and D151248 is a simple repeat homozygote.

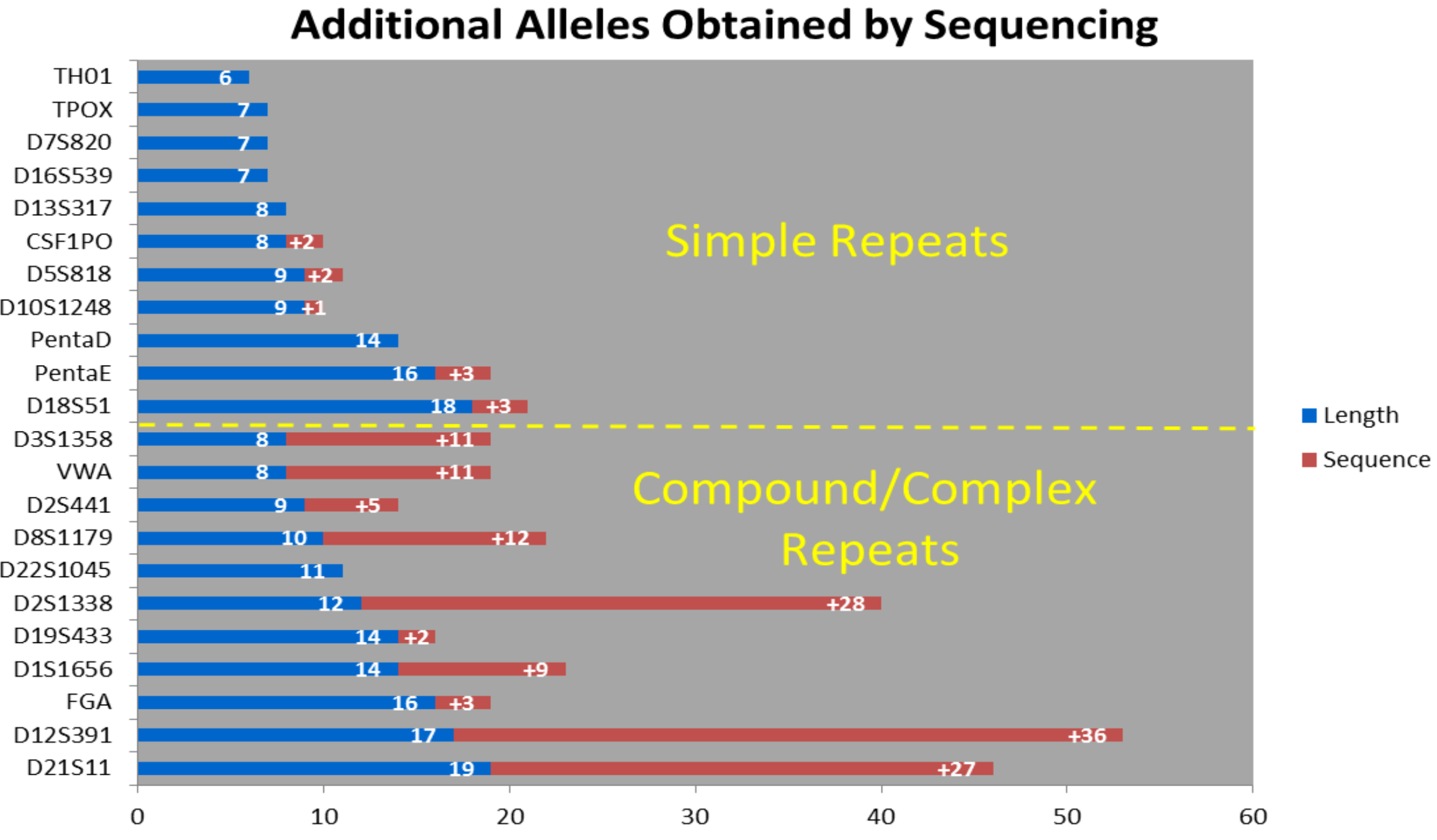


Figure 4. In blue are the number of different length-based alleles observed in this dataset (N=183), and in red are the number of additional sequenced-based alleles observed. Loci are grouped by repeat motif type (simple vs compound/complex) and sorted within each group by number of length-based alleles, smallest to largest.

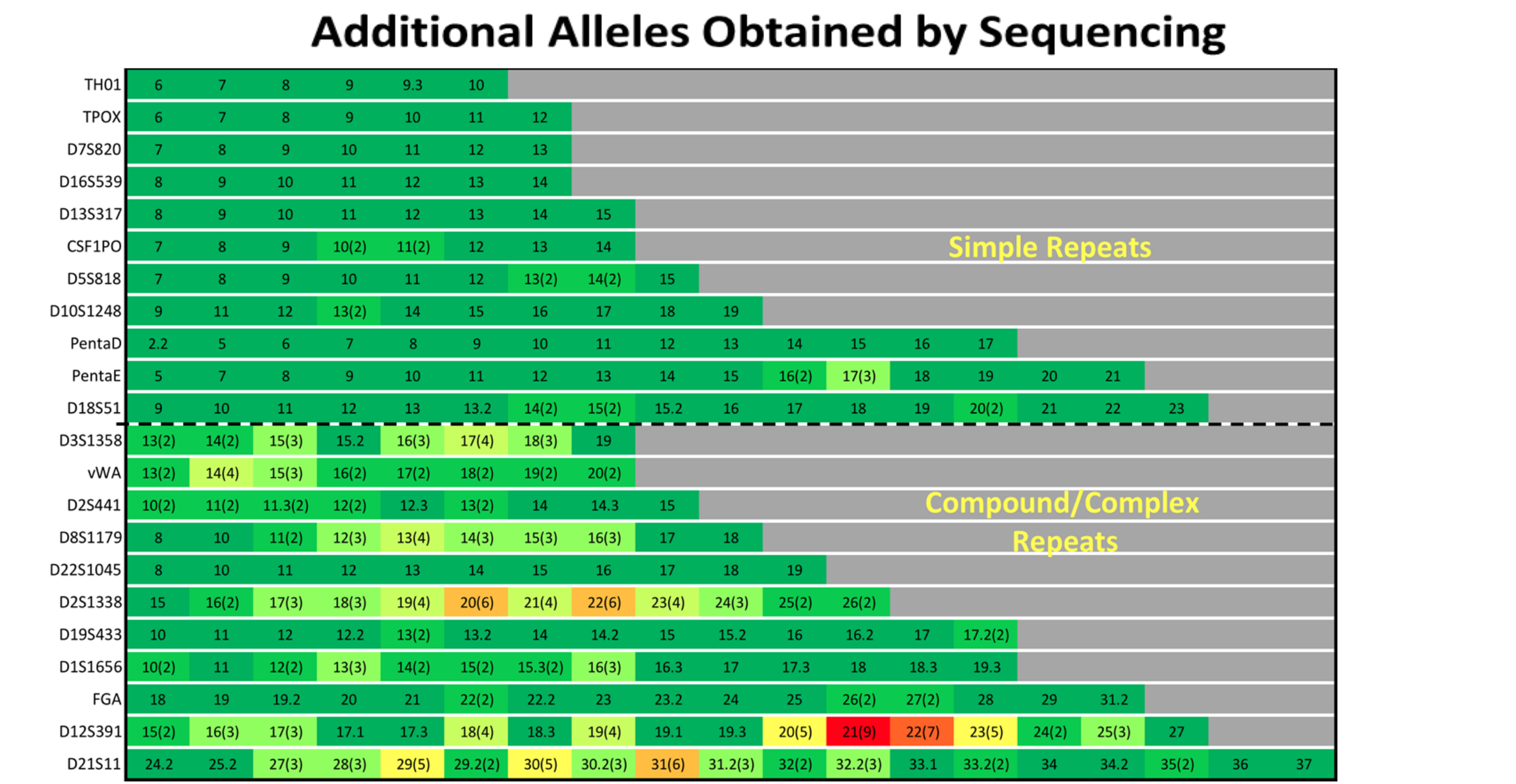


Figure 5. Heatmap showing the variant count for each allele (counts in parentheses, alleles with no parenthetical notation show no sequence variants in this dataset). Alleles are color coded with the darkest green shading representing no sequence variation; shading changes to yellow-orange-red with increasing sequence variation.

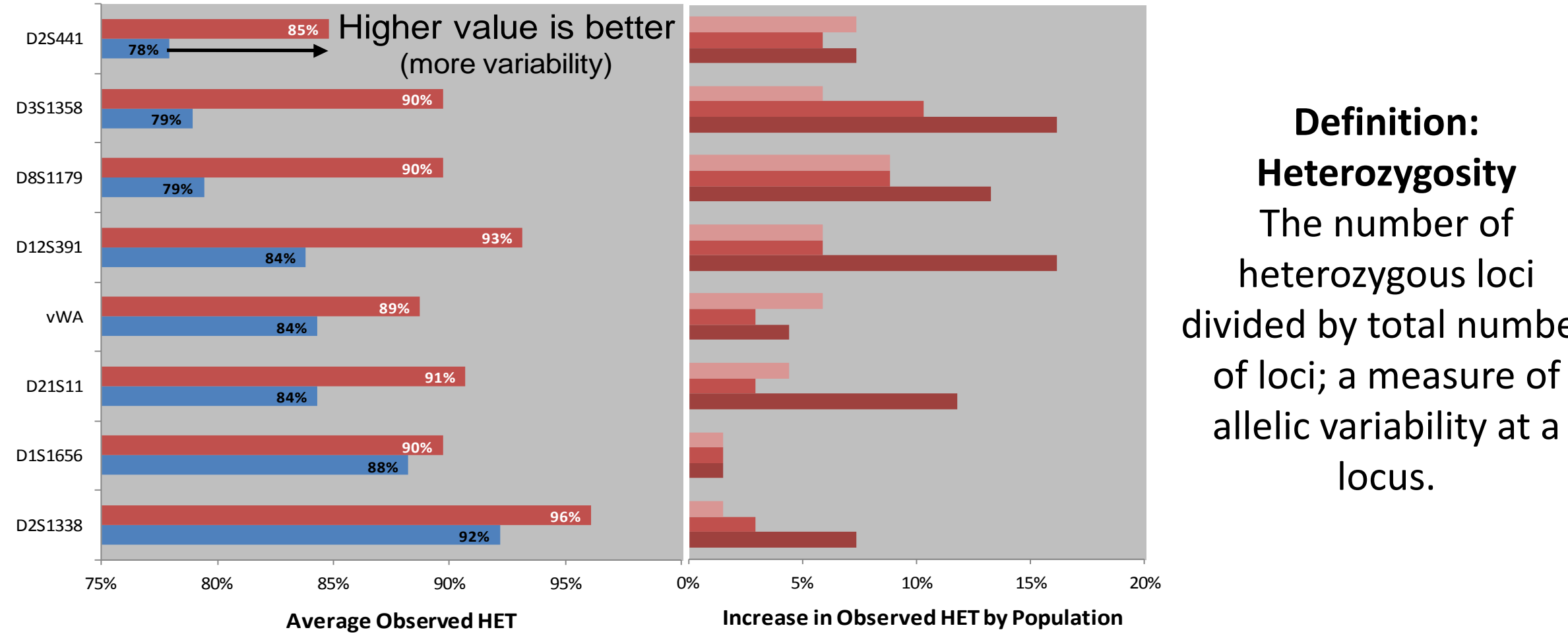


Figure 6. (left) Percent heterozygosity by length (blue) and by sequence (red), averaged from three populations (N=183) rank ordered by HET length, top eight loci shown. (right) The increase in heterozygosity by sequence, broken out by population.

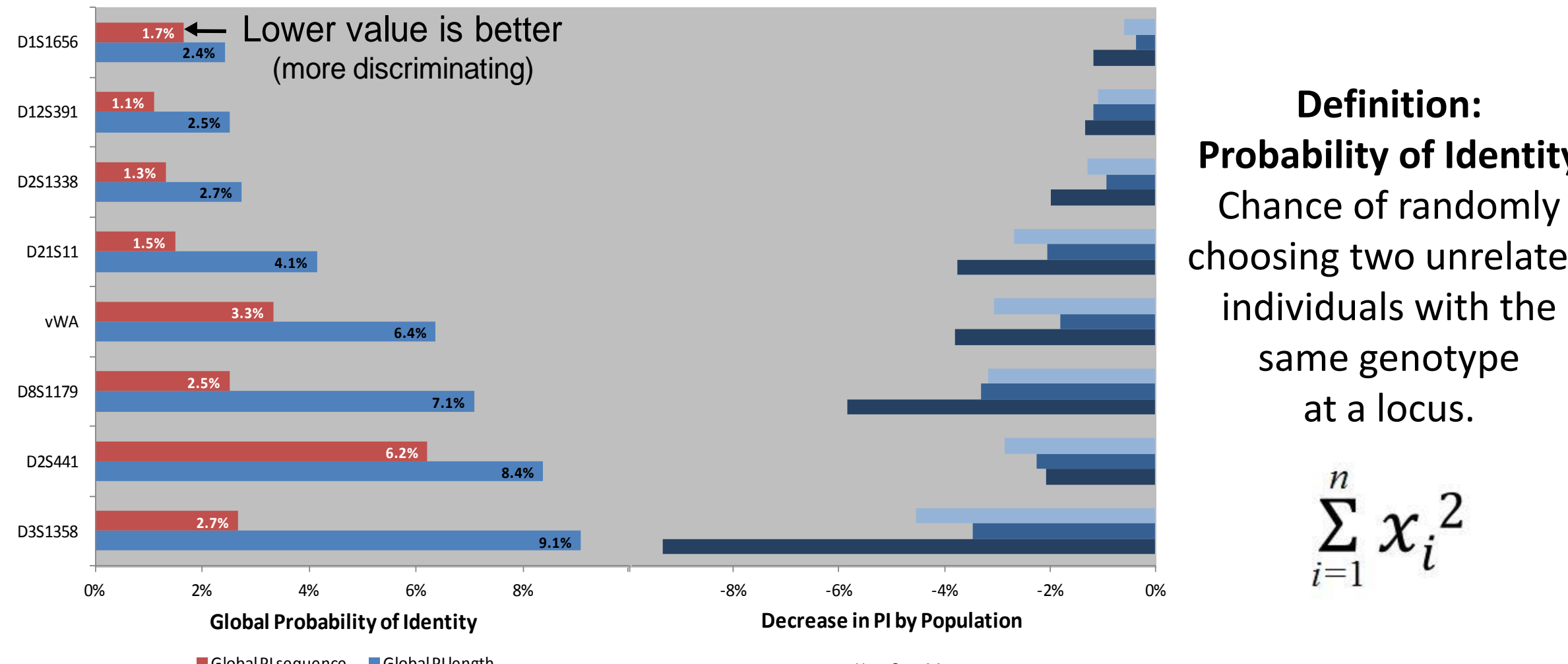


Figure 7. (left) Percent probability of identity by length (blue) and by sequence (red), based on average allele frequency from these three populations (N=183) rank ordered by PI length, top eight loci shown. (right) The decrease in probability of identity by sequence, broken out by population.

EXAMPLE LOCUS

The D21S11 locus provides an example of the sequence complexity that exists and the need for a system of nomenclature that is meaningful and expandable as more unique alleles are sequenced. For the set of population samples sequenced in this project (N=183), at the D21S11 locus, 16 sequences were found which had not previously been reported (equaling 4.4% of chromosome 21 sequences), see Table 1 rows in red.

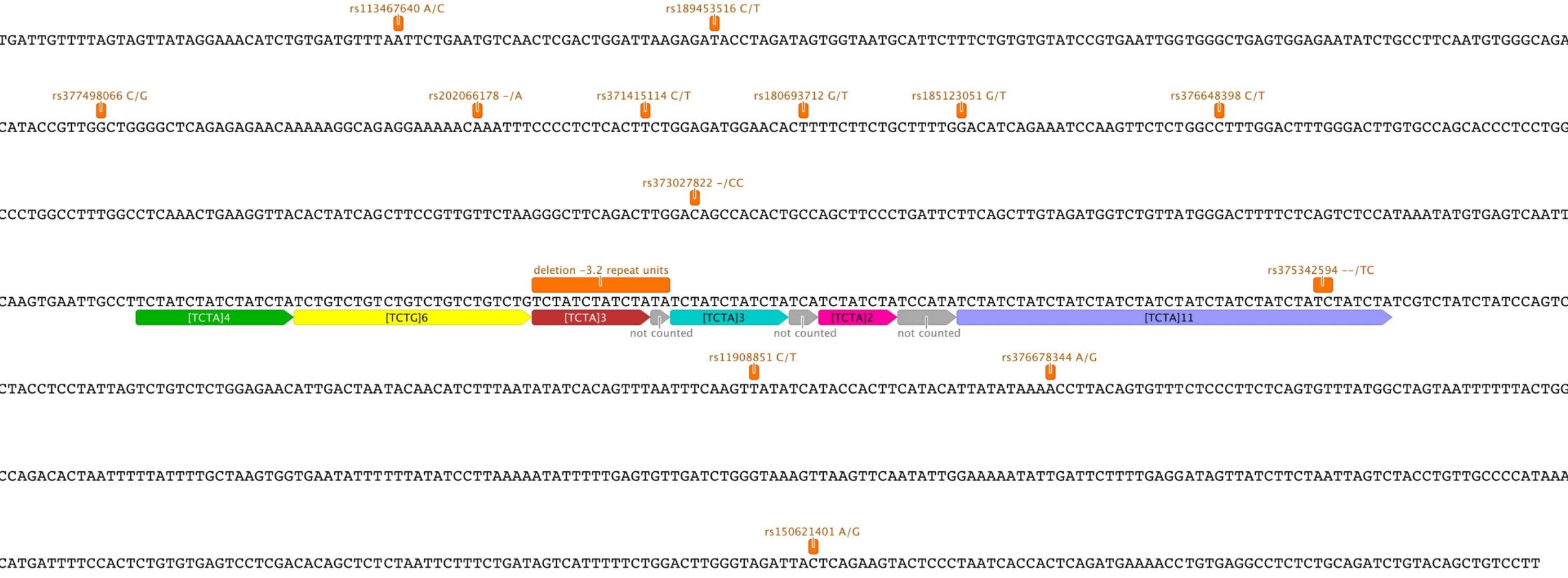


Figure 8. GRCh38 sequence (GenBank) at the D21S11 locus (repeat region and 500 bases upstream and downstream), annotated with information relevant to forensic NGS. The subunits of the repeat motif are shown in green, yellow, red, aqua, pink and purple (these regions are “counted”, resulting in a 29 allele), while the interspersed non-repeat sequences are shown in gray (these regions are “not counted”). Two different regions where deletions result in 2 alleles, as well as the locations of numerous SNPs (from dbSNP) that have been observed in the flanking regions are shown in orange. Annotation created with Geneious v7.1.7.

D21S11	[TCTA]4	[TCTG]4	[TCTA]3	TA	[TCTA]1	TCA	[TCTA]2	TCCATA	[TCTA]3n		
allele	[TCTA]4-13	[TCTG]3-13	[TCTA]3	TA	[TCTA]2-3	TCA	[TCTA]2	TCCATA	[TCTA]6-18		
24	[TCTA]4	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]6	Sanger	Griffiths et al. (1998)
25	[TCTA]3	[TCTA]3	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]10	Sanger	Schwartz et al. (1996)
26	[TCTA]4	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]8	Sanger	Moller et al. (1994)
26	[TCTA]6	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]9	Sanger	Wang et al. (2004)
27	[TCTA]4	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]9	Sanger	Moller et al. (1994)
27	[TCTA]5	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]9	Sanger	Griffiths et al. (1998)
27	[TCTA]6	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]8	Schwartz et al. (1996)	
28	[TCTA]4	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]10	isa	Gelardi et al. (2014)
28	[TCTA]4	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]10	Sanger	Moller et al. (1994)
28	[TCTA]5	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]10	Sanger	Wang et al. (2004)
28	[TCTA]5	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]10	Msseq	NIST 183
28	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]9	Sanger	Zhou et al. (1997)
28	[TCTA]6	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]9	isa	Gelardi et al. (2014)
29	[TCTA]4	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]11	Sanger	Griffiths et al. (1998)
29	[TCTA]4	[TCTG]7	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]10	Sanger	Plante et al. (2012)
29	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]11	Msseq	NIST 183
29	[TCTA]6	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]10	isa	Gelardi et al. (2014)
29	[TCTA]6	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]10	Sanger	Zhou et al. (1997)
30	[TCTA]4	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]12	Sanger	Schwartz et al. (1996)
30	[TCTA]4	[TCTG]7	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]11	isa	Gelardi et al. (2014)
30	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]11	Sanger	Zhou et al. (1997)
30	[TCTA]6	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]11	Sanger	Griffiths (1998)
30	[TCTA]6	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]10	Sanger	Brinkmann et al. (1996a)
30	[TCTA]7	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]10	isa	Gelardi et al. (2014)
31	[TCTA]4	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]13	isa	Gelardi et al. (2014)
31	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]12	Sanger	Griffiths et al. (1998)
31	[TCTA]6	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]12	Sanger	Moller et al. (1994)
31	[TCTA]6	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]11	Sanger	Zhou et al. (1997)
31	[TCTA]7	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]11	Sanger	Schwartz et al. (1996)
31	[TCTA]7	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]10	Msseq	NIST 183
31	[TCTA]4	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]13	Sanger	Wang et al. (2004)
32	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]13	Sanger	Zhou et al. (1997)
32	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]13	isa	Gelardi et al. (2014)
32	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]13	Sanger	Griffiths et al. (1998)
32	[TCTA]6	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]12	Msseq	NIST 183
32	[TCTA]6	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]11	isa	Gelardi et al. (2014)
33	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]14	Sanger	Zhou et al. (1997)
33	[TCTA]6	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]14	isa	Gelardi et al. (2014)
33	[TCTA]6	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]13	Sanger	Wang et al. (2004)
34	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]15	Sanger	Zhou et al. (1997)
34	[TCTA]6	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]13	Sanger	Wang et al. (2004)
34	[TCTA]10	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]11	Sanger	Brinkmann et al. (1996a)
35	[TCTA]6	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]11	Sanger	Wang et al. (2004)
35	[TCTA]10	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]12	Sanger	Griffiths et al. (1998)
35	[TCTA]10	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]10	Msseq	NIST 183
35	[TCTA]11	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]11	Sanger	Brinkmann et al. (1996a)
36	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]14	Sanger	Wang et al. (2004)
36	[TCTA]10	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]13	Sanger	Brinkmann et al. (1996a)
36	[TCTA]10	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]12	Sanger	Brinkmann et al. (1996a)
36	[TCTA]10	[TCTG]7	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]11	Msseq	NIST 183
36	[TCTA]11	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]12	Sanger	Griffiths et al. (1998)
37	[TCTA]9	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]15	Sanger	Wang et al. (2004)
37	[TCTA]10	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]13	Sanger	Griffiths et al. (1998)
37	[TCTA]11	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]11	Msseq	NIST 183
38	[TCTA]13	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]12	Sanger	Griffiths et al. (1998)

allele	[TCTA]4-6	[TCTG]5-6	[TCTA]2-3	TA	[TCTA]2-3	TCA	[TCTA]2	TCCATA	[TCTA]8-16	TA	TCTA	Platform	Reference
28.2	[TCTA]4	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]8	TA	TCTA	Sanger	Zhou et al. (1997)
29.2	[TCTA]5	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]10	TA	TCTA	Sanger	Zhou et al. (1997)
29.2	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]9	TA	TCTA	MsSeq	NIST 183
30.2	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]11	TA	TCTA	Sanger	Shewhart et al. (1996)
30.2	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]10	TA	TCTA	MsSeq	NIST 183
30.2	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]11	TA	TCTA	Sanger	Griffiths et al. (1998)
31.2	[TCTA]5	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]12	TA	TCTA	MsSeq	NIST 183
31.2	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]12	TA	TCTA	MsSeq	NIST 183
31.2	[TCTA]6	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]11	TA	TCTA	Sanger	Griffiths et al. (1998)
32.2	[TCTA]4	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]13	TA	TCTA	Sanger	Brinkmann et al. (1996b)
32.2	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]12	TA	TCTA	Sanger	Griffiths et al. (1998)
32.2	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]13	TA	TCTA	Sanger	Brinkmann et al. (1996a)
32.2	[TCTA]6	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]12	TA	TCTA	MsSeq	NIST 183
32.2	[TCTA]6	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]11	TA	TCTA	isa	Gelardi et al. (2014)
33.2	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]13	TA	TCTA	Sanger	Griffiths et al. (1998)
33.2	[TCTA]6	[TCTG]5	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]13	TA	TCTA	Sanger	Brinkmann et al. (1996a)
33.2	[TCTA]6	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]12	TA	TCTA	Sanger	Brinkmann et al. (1996a)
34.2	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]14	TA	TCTA	Sanger	Griffiths et al. (1998)
35.2	[TCTA]5	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]15	TA	TCTA	Sanger	Zhou et al. (1997)
36.2	[TCTA]6	[TCTG]6	[TCTA]3	TA	[TCTA]3	TCA	[TCTA]2	TCCATA	[TCTA]16	TA	TCTA	Sanger	Zhou et al. (1997)