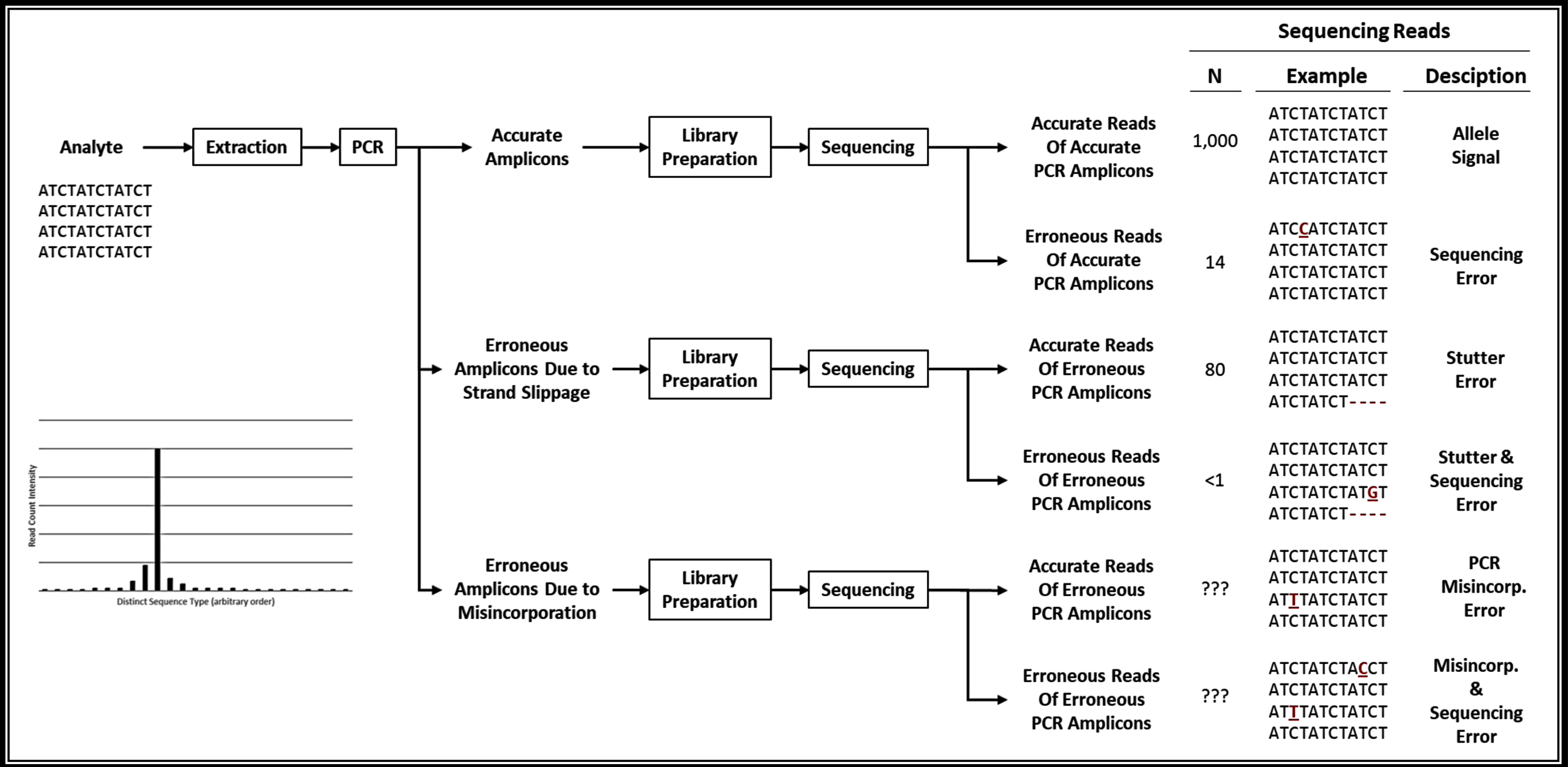


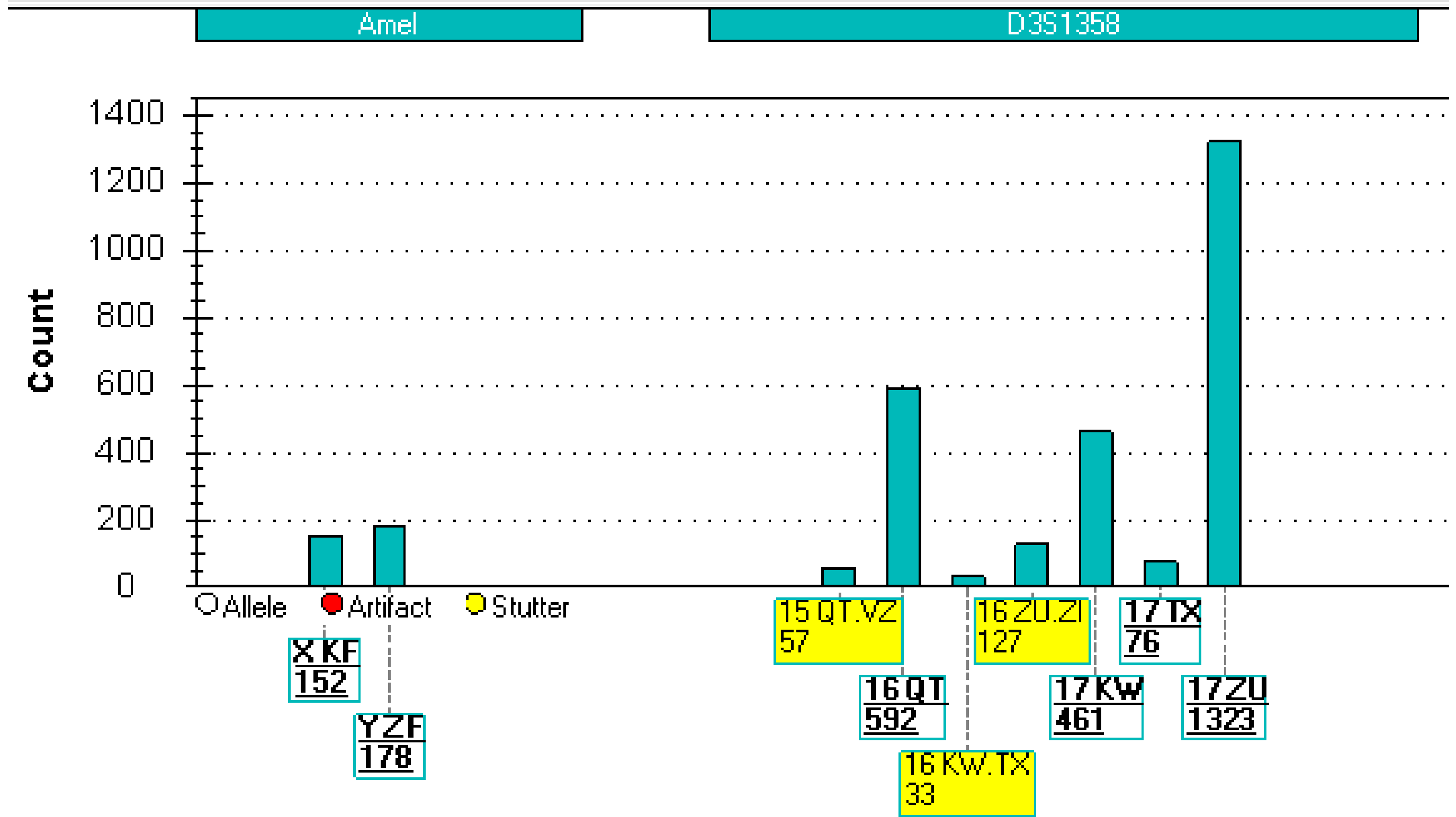
Filtering Artifacts from Forensic PCR-MPS Data for Mixture Analysis

Brian Young, Tom Faris, Jeff Smith, Luigi Armogida
NicheVision Forensics brian@nichevision.com

Model of PCR-MPS Allele and Artifact Generation



Stutter Artifacts



Non-Stutter Artifacts

Table 1. Example of an allele, non-stutter artifact, and stutter artifact in reads targeting D2S1338. Artifacts can be attributed to parent alleles using SID labels in a 'parent-tic-child' format.

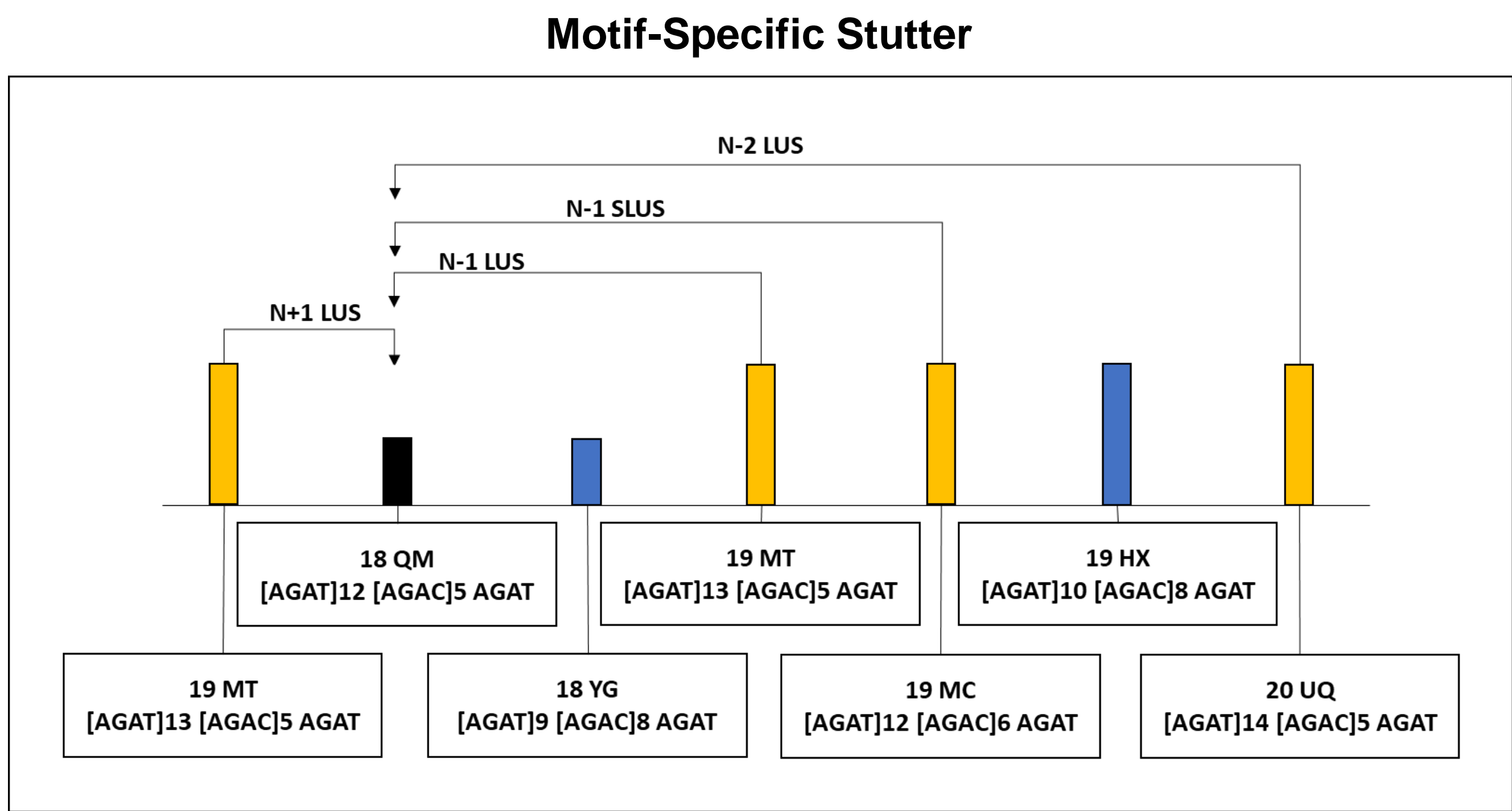
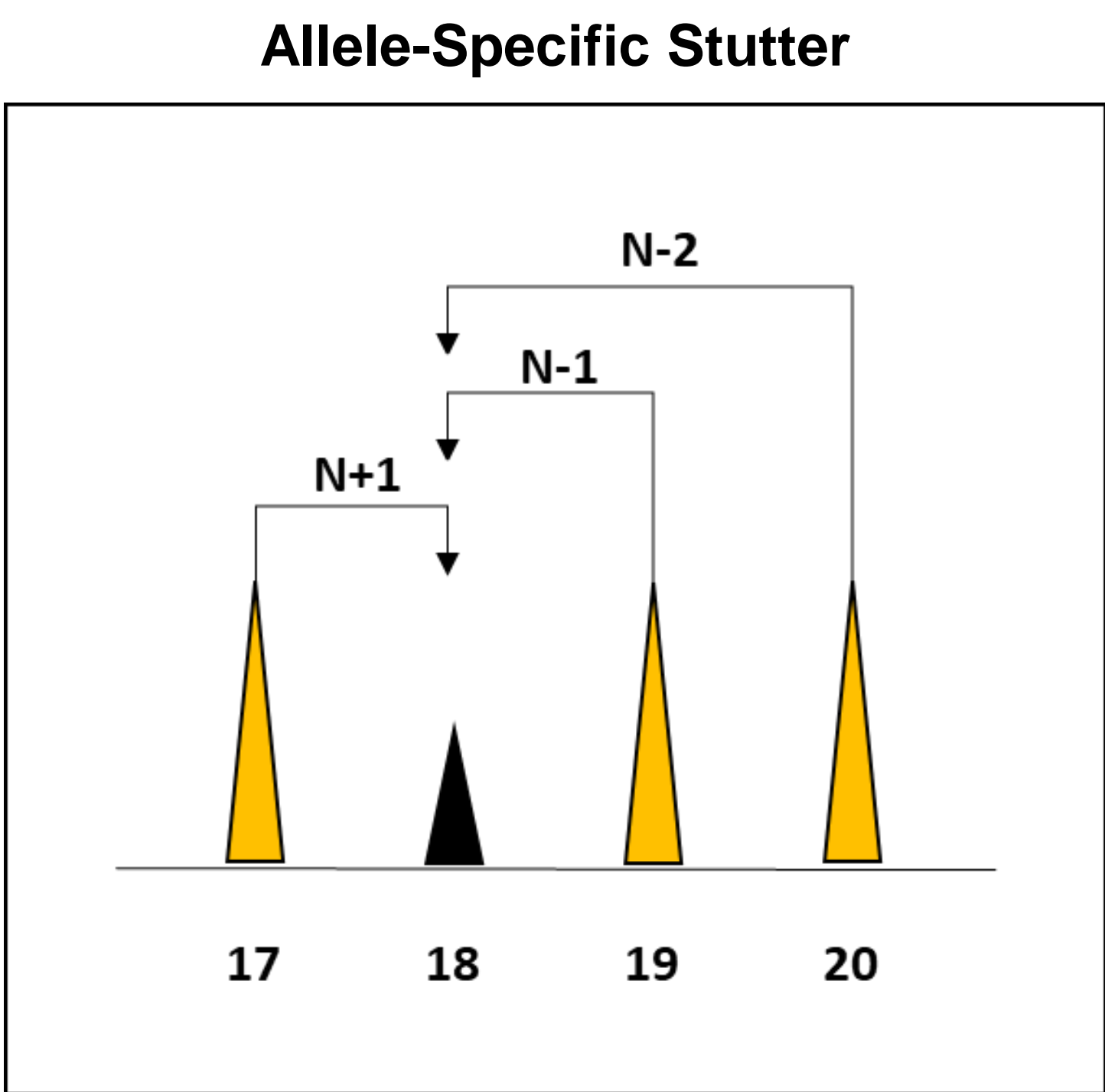
Category	Label	Bracketed Sequence of Haplotype Allele
Allele	25 BJ	GAG [GGAA]2 GGAC [GGAA]15 [GGCA]7 AGGCCAAGCCATT
Non-Stutter	25 BJ AS	GAG [GGAA]2 GGAC [GGAA]15 [GGCA]7 AGGCCAAGCCATGT
Stutter	24 BJ.SS	GAG [GGAA]2 GGAC [GGAA]14 [GGCA]7 AGGCCAAGCCATT

PCR-MPS Read Sequences Are 'Haplotype Alleles'

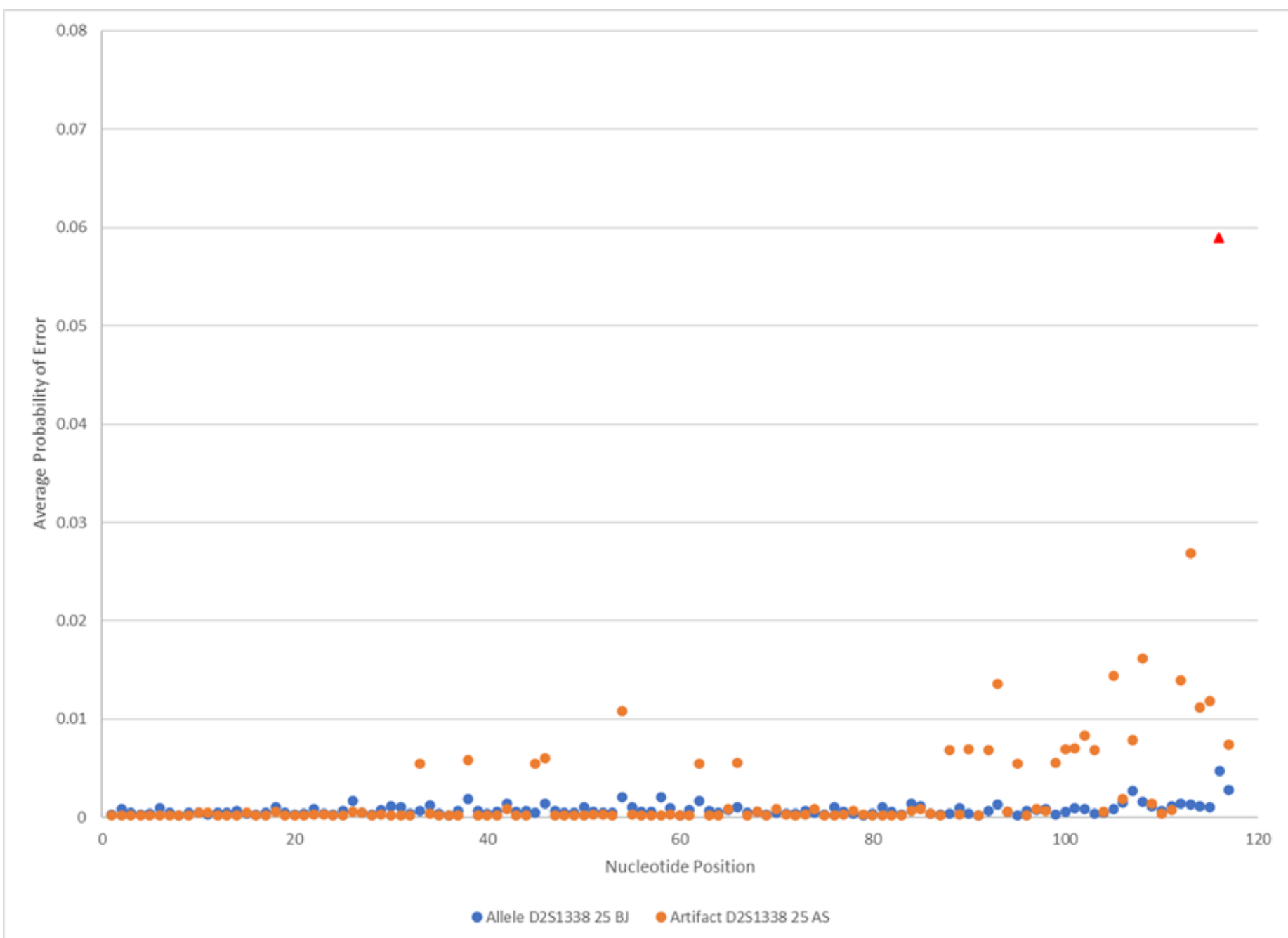
Table 2. Trimmed read sequences of PCR-MPS amplicons are 'haplotype alleles' possibly containing multiple STR, SNP, and DIP markers. The entire sequence string is required to describe the allele, or to discriminate alleles from artifacts. Three isometric alleles targeting D16S539 are shown. Raw string and hash digests uniquely define haplotype alleles. SID labels using two letters of the hash digest generate pronounceable labels for labeling distinct sequence types within a locus profile.

Raw Sequence	SID Hash Digest	SID Label	Bracket Notation
TCCTCTTCCCTAGATCAATACAGACAGACAGACAGGTGGATAGATAGATAGATAGATAGATAGATAGATAGATATCATTGAAAGACAAAC	QBCNTZTPIQYLRLCQCCIBFFRWAQZSBBVXQQGGXXTTFEWIJZKXYGMXPLC	10 QB	[GATA]10 rs11642858
CAGAGATGGATGATAGATAC	XIVKFZSKIMSNTGFYEINGIJOGBCEIOLTLUEEGHOWJDBNKLZYJFMYQHE	10 XI	[GATA]10
TCCTCTTCCCTAGATCAATACAGACAGACAGGTGGATAGATAGATAGATAGATAGATAGATAGATAGATATCATTGAAAGACAAAC	ERZHTOJAVBFTVGHWKTWALRKAMDZPEKJSVYQQOVVYKGBYPWVNOZLZJD	10 ER	[GATA]5 GACA [GATA]4

Stutter Artifact Discrimination

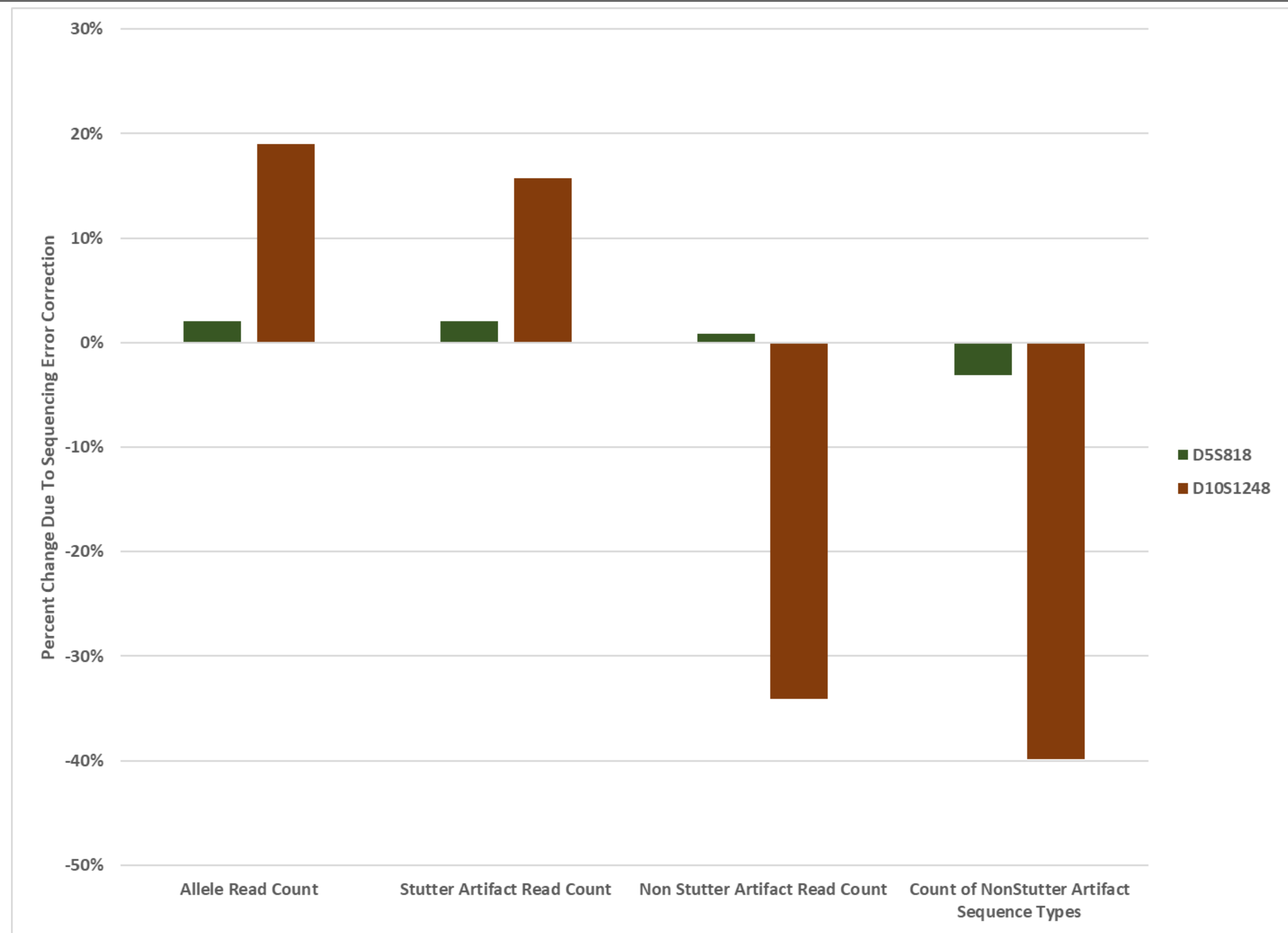


Non-Stutter Artifact Detection



Read Pair	Base Calls and Q Scores
Original Forward Read	G A T C A C A G G T 30 32 28 30 30 34 40 38 30 30
Original Reverse Read	G A T C T C A G G T 28 30 30 28 20 32 40 36 32 34
Forward Read	G A T C A C A G G T 30 32 28 30 30 34 40 38 30 30
Corrected Reverse Read	G A T C A C A G G T 28 30 30 28 30 32 40 36 32 34

Figure 6. Illustration of overlapping paired-end reads that differ in base call at a position (highlighted rectangle). Both forward and reverse reads are shown in forward orientation. When overlapping paired-end reads differ by sequence, they cannot both be correct (top panel). Quality scores can be used to infer which is the correct sequence. Optionally, the inferred incorrect base call can be corrected (bottom panel).



(1) Sample H2_3P-C_10-30-60_1_R702-A508_S17. <https://data.nist.gov/od/id/mds2-2157>
(2) Young B, Faris T, Armogida L. A nomenclature for sequence-based forensic DNA analysis. Forensic Sci. Intl. Genet. 42(2019)14-20.