



STRSeq: A resource for sequence-based STR analysis

Abstract - The STR Sequencing Project (STRSeq)¹ was initiated to facilitate the description of sequence-based alleles at Short Tandem Repeat (STR) loci targeted in human identification assays. STRSeq data are maintained as GenBank records at the U.S. National Center for Biotechnology Information (NCBI). Each GenBank record contains: 1) observed sequence of an STR region, 2) annotation of the repeat region (“bracketing” consistent with the guidance of the International Society for Forensic Genetics) and flanking region polymorphisms, 3) information regarding the sequencing assay and data quality, and 4) backward compatible length-based allelic designation. STRSeq GenBank records are organized within a BioProject at NCBI www.ncbi.nlm.nih.gov/bioproject/380127, which is sub-divided by:

- Commonly used autosomal STR Loci
- Alternate autosomal STR Loci
- Y-chromosomal STR loci
- X-chromosomal STR loci

The BioProject will initially contain aggregate alleles across **4,612 samples** submitted by four laboratories: National Institute of Standards and Technology (NIST, the project organizer), Kings College London (KCL), University of North Texas Health Sciences Center (UNT), and University of Santiago de Compostela (USC)²⁻⁶.

Currently, **1145** records have been uploaded into STRSeq. These submitted records represent the evaluation of approximately 25% of the current data set for 28 autosomal STR loci. **The plot to the right** is a representation of 35 autosomal STR loci from the four data sets representing **1739** unique sequences, the additional **594** autosomal sequences will potentially be uploaded to the STRSeq Project after additional evaluation and confirmation. STRSeq is meant to continually grow as new data sets are developed and new sequences are reported and confirmed.

Current participating laboratories and samples provided for evaluation



1786 Samples



1043 Samples

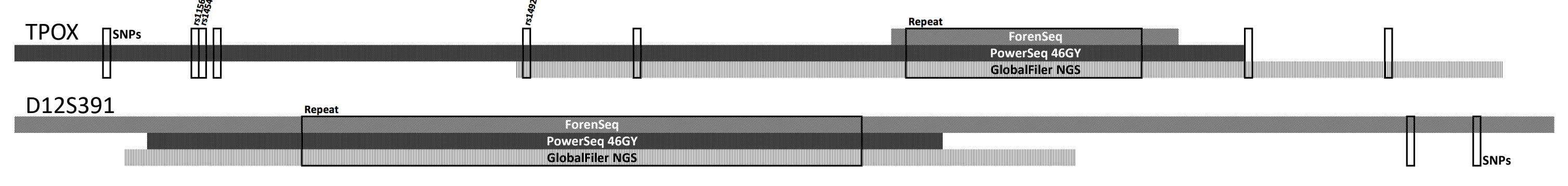


839 Samples



944 Samples

Method – Data from three Forensic STR sequencing kits are currently being evaluated for inclusion in STRSeq: 1) **ForenSeq DNA Signature Prep Kit** (Verogen), 2) **PowerSeq 46GY** (beta version, Promega), and 3) **GlobalFiler NGS** (Thermo Fisher). All of the samples currently represented in STRSeq have been sequenced with ForenSeq. NIST has sequenced subsets of >600 samples with PowerSeq (included in the current STRSeq records) and >200 samples with GlobalFiler NGS, which are currently being reviewed. **Below** are examples of sequence overlap at the TPOX and D12S391 loci across the sequencing kits.



TPOX allele 8

Homo sapiens microsatellite TPOX 8 [AATG] sequence
167 bp linear DNA
Accession: MF042428.1 GI: 119790699
BioProject: PubMed Taxonomy
GenBank FASTA Graphics

Homo sapiens microsatellite TPOX 8 [AATG] sequence
167 bp linear DNA
Accession: MF042428.1 GI: 119790699
BioProject: PubMed Taxonomy
GenBank FASTA Graphics

Homo sapiens microsatellite TPOX 8 [AATG] sequence
167 bp linear DNA
Accession: MF042428.1 GI: 119790699
BioProject: PubMed Taxonomy
GenBank FASTA Graphics

D12S391 allele 27

Homo sapiens microsatellite D12S391 27 [AGATT] sequence
241 bp linear DNA
Accession: MH167199.1 GI: 138124258
BioProject: PubMed Taxonomy
GenBank FASTA Graphics

Homo sapiens microsatellite D12S391 27 [AGATT] sequence
241 bp linear DNA
Accession: MH167200.1 GI: 138124259
BioProject: PubMed Taxonomy
GenBank FASTA Graphics

Homo sapiens microsatellite D12S391 27 [AGGT] sequence
241 bp linear DNA
Accession: MH167201.1 GI: 138124260
BioProject: PubMed Taxonomy
GenBank FASTA Graphics

NCBI houses the GenBank records. **To the left** are examples of records in the TPOX and D12S391 BioProjects. Orange boxes highlight differences in the records. All four TPOX allele “8” sequences have the same repeat but vary in the flanking regions, whereas D12S391 allele “27” varies in the repeat region but not in the flanking region.

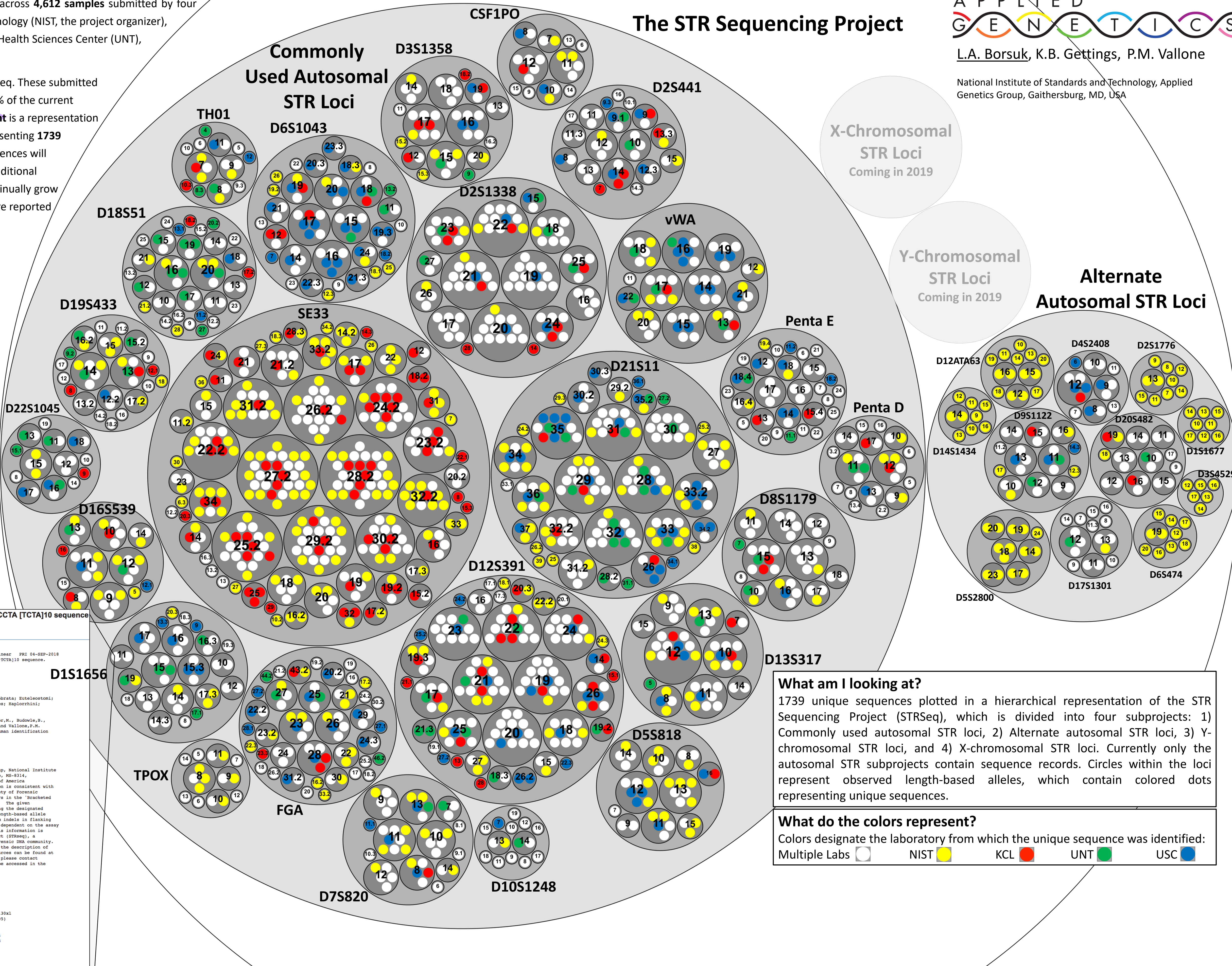
To the right is an example of a D1S1656 allele “11” record. The goal is to provide the user with maximal information about the reported sequence. This includes confirmed length-based allele size, sequence coordinates corresponding to the GRCh38 reference sequence, and flanking region variants compared to the GRCh38 reference sequence.

Homo sapiens microsatellite D1S1656 11 CCTA [TCTA]10 sequence

```

GenBank: MH174236.1
FASTA
GSI
LOCUS MH14836 126 bp DNA linear PK2 04-SEP-2018
DEFINITION Homo sapiens microsatellite D1S1656 11 CCTA [TCTA]10 sequence.
ACCESSION MH14836
VERSION MH14836.1
DBLINK BioProject: FBX038553
REFERENCE Strategy STR D1S1656.
SOURCE Homo sapiens (human)
ORGANISM Homo sapiens
Bovines; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates; Haplorhina; Catarrhini; Hominoidea; Homo.
REFERENCE 1 (bases 1 to 126)
AUTHORS Gettings,K.B., Borsuk,L.A., Ballard,D., Bodner,M., Budowle,B., Devesse,L., King,J., Parson,W., Phillips,C., and P.M. Vallone. STRSeq: A catalog of sequence diversity at human identification Short Tandem Repeat loci. Forensic Science International: Genetics 31, 111-117 (2017).
JOURNAL FORENSIC SCIENCE INTERNATIONAL: GENETICS 31
PUBMED 28185123
REFERENCES 2 (bases 1 to 126)
AUTHORS Bodner,M., Budowle,B., Devesse,L., King,J., Parson,W., Phillips,C., and P.M. Vallone. STRSeq: A catalog of sequence diversity at human identification Short Tandem Repeat loci. Forensic Science International: Genetics 31, 111-117 (2017).
JOURNAL FORENSIC SCIENCE INTERNATIONAL: GENETICS 31
PUBMED 28185123
COMMENT
##humanSTR-STRID#
STR locus name || D1S1656
Length-based allele || 11
Bracketed repeat || CCTA [TCTA]10
Sequencing technology || MiSeq FGx; MiSeq
Coverage || X30X
Length-based tech. || PowerPlex Fusion; 3130x1
Assay || Qubit3 (GC_00001405)
Chromosome || 1
RefSeq Accession || NC_000001.11
Chrom. Location || 23079553..23079574
Repeat Location || 23079616..23079683
Cytogenetic Location || 1q44.2
##humanSTR-ID#
Source/Qualifiers
1..126
/organism="Homo sapiens"
/mol_type="genomic DNA"
/db_xref="taxon:9606"
1..126
/contig="Promega PowerSeq DNA Signature Prep Kit"
6..126
/contig="Promega PowerSeq 46G System"
62..115
/ftp_upload=
/submitter="NIST"
/submitter_contact="NIST"
/submitter_email="nists@nist.gov"
ORIGINS
1: ttagagaa tagatcaatc agggacaa atatataatc ataacattaa cacacacaa
6: acatatactc atctatactc atctatactc atctatactc atctatactc cacagtgtac
121: cctaga
  
```

The STR Sequencing Project



What am I looking at?
1739 unique sequences plotted in a hierarchical representation of the STR Sequencing Project (STRSeq), which is divided into four subprojects: 1) Commonly used autosomal STR loci, 2) Alternate autosomal STR loci, 3) Y-chromosomal STR loci, and 4) X-chromosomal STR loci. Currently only the autosomal STR subprojects contain sequence records. Circles within the loci represent observed length-based alleles, which contain colored dots representing unique sequences.

What do the colors represent?
Colors designate the laboratory from which the unique sequence was identified:
Multiple Labs NIST KCL UNT USC

Discussion –The next sets of autosomal, X, and Y STRSeq records representing the remainder of the current samples will be uploaded to the STRSeq BioProject over the next year. Future plans for this NIJ-funded effort include a pathway for researchers to submit additional alleles and customized interface tools, which would allow users to easily search the STRSeq data set. In addition to providing a framework for communication among laboratories, the ability to search the BioProject can be leveraged as QC for rare sequences encountered in forensic casework.

Acknowledgements – The authors express gratitude to the NCBI staff who have facilitated development of the BioProject: Drs. Lori Black, Melissa Landrum, Ilene Mizrahi, Kim Pruitt, George Riley, and Steven Sherry. The authors also acknowledge the input of the European Commission project DNASEQEX (HOME/2014/ISFP/AG/LAWX/400007135) and the support of the ENFSI DNA Working Group and thank the many practitioners and researchers who provided valuable feedback.

References

- Gettings, K.B., Borsuk, L.A., Ballard, D., Bodner, M., Budowle, B., Devesse, L., King, J., Parson, W., Phillips, C., and P.M. Vallone. **STRSeq: A catalog of sequence diversity at human identification Short Tandem Repeat loci**. Forensic Science International: Genetics 31, 111-117 (2017).
- Gettings, K.B., Borsuk, L.A., Steffen, C.R., Klesler, K.M., and P.M. Vallone. **Sequence-based U.S. population data for 27 autosomal STR loci**. Forensic Science International: Genetics 37, 106-115 (2018).
- Borsuk, L.A., Gettings, K.B., Steffen, C.R., Klesler, K.M., and P.M. Vallone. **Sequence-based U.S. population data for the SE33 locus**. Electrophoresis 21, (2018).
- Devesse, L., Ballard, D., Davenport, L., Riethorst, L., Mason-Buck, G., and D. Syndercombe Court. **Congordance of the ForenSeq system and characterization of sequence-specific autosomal STR alleles across two major population groups**. Forensic Science International: Genetics 34, 57-61 (2018).
- Novroski, N.M.M., King, J.L., Churchill, J.D., Seah, L.H., and B. Budowle. **Characterization of genetic sequence variation of 58 STR loci in four major population groups**. Forensic Science International: Genetics 25, 214-226 (2016).
- Phillips, C., Devesse, L., Ballard, D., van Weert, L., de la Puente, M., Melis, S. Alvarez Iglesias, V., Freire-Aradas, A., Oldroyd, N., Holt, C., Syndercombe Court, D., Carracedo, A., and M.V. Lareu. **Global patterns of STR sequence variation: Sequencing the CEPH human genome diversity panel for 58 forensic STRs using the Illumina ForenSeq DNA Signature Prep Kit**. Electrophoresis 21, (2018).

RawGraphic/plot ref
(Circle Packing Plot) Mauri, M., Elli, T., Caviglia, G., Uboldi, G., and M. Azzi. **RAWGraphs: A Visualisation Platform to Create Open Outputs**. In *Proceedings of the 12th Biannual Conference on Italian SIGCHI Chapter* (p. 28:1–28:5). New York, NY, USA: ACM. (2017).

(ref) Parson, W., Ballard, D., Budowle, B., Butler, J.M., Gettings, K.B., Gill, P., Gumsjö, L., Hares, D., Irwin, J., King, J., de Knijff, P., Morling, N., Prinz, M., Schneider, P.M., Van Neste, C., Willuweit, S., and C. Phillips. **Massively Parallel Sequencing of Forensic STRs: Considerations of the DNA Commission of the International Society of Forensic Genetics (ISFG) on minimal nomenclature requirements**. Forensic Science International: Genetics 37, 54–63 (2016).

NIST Funding Sources – This work was funded in part by the National Institute of Justice (NIJ) interagency agreement 1609-602-18NIJ: “Forensic DNA Applications of Next Generation Sequencing”.

Disclaimer – Points of view in this document are those of the authors and do not necessarily represent the official position or policies of the U.S. Department of Commerce. **Product:** Certain commercial equipment, instruments, and materials are identified in order to specify experimental procedures as completely as possible. In no case does such identification imply a recommendation or endorsement by NIST, nor does it imply that any of the materials, instruments, or equipment identified are necessarily the best available for the purpose. **Data:** The National Institute of Standards and Technology (NIST) uses its best efforts to deliver high quality data that have been selected on the basis of sound scientific judgement. However, NIST makes no warranties to that effect, and NIST shall not be liable for any damage that may result from errors or omissions in the data set.

Poster # 61 at:
29th ISHI, Phoenix, AZ
September 24-27, 2018
Email: strseq@nist.gov