

STR Nomenclature Meeting

April 11-12, 2019 London



5' to 3':

Walther Parson, Lisa Borsuk, Peter Schneider, Brian Young, Rebecca Just, Jodi Irwin, David Ballard, Sascha Willuweit, Cydne Holt, Chris Phillips, Jonathan King, Tunde Huszar, Peter Gill, Christian Sell, Kris Van der Gaag, Laurence Devesse, Claus Borsting, Doug Hares, Katherine Gettings, Rob Lagace, Jerry Hoogenboom, Martin Bodner, Peter deKnijff, Sebastian Ganschow, Pedro Barrio, Teresa Gross

Agenda: April 11, 2019

8:30am - 8:45am	Arrival
8:45am - 9:30am	Welcome and Opening Remarks - STRAND WG
9:30am - 10:00am	Tunde Huszar
10:00am - 10:15am	Coffee Break
10:15am - 10:45am	Pedro Barrio
10:45am - 11:15am	Claus Borsting
11:15am - 11:45am	Brian Young
11:45am - 12:00pm	Discussion
12:00pm - 1:30pm	Lunch on your own
1:30pm - 2:00pm	Peter deKnijff
2:00pm - 2:30pm	Kris Van der Gaag / Jerry Hoogenboom
2:30pm - 3:00pm	Sascha Willuweit
3:00pm - 3:15pm	Coffee Break
3:15pm - 3:45pm	Rebecca Just
3:45pm - 4:15pm	Peter Gill
4:15pm - 4:45pm	Sebastian Ganschow
4:45pm - 5:00pm	Discussion and Day 1 Closing Remarks

Agenda: April 12, 2019

8:30am - 8:45am	Arrival
8:45am - 9:00am	Welcome and Opening Remarks
9:00am - 11:00am	STRAND WG Facilitated Discussion <ul style="list-style-type: none">● Reference genomes, existing databases● Quality control● Bioinformatics● Implementation
11:00am - 11:15am	Coffee Break
11:15am - 12:00pm	Summary, Path Forward, and Closing Remarks

STRAND *working group*
align | name | define

STR sequence nomenclature

NGS of STRs: Nomenclature Panel at ISFG Conference Krakow 2015



NGS of STRs: Considerations of the ISFG (2016)

Forensic Science International: Genetics 22 (2016) 54–63



Contents lists available at ScienceDirect

Forensic Science International: Genetics

journal homepage: www.elsevier.com/locate/fsig



Massively parallel sequencing of forensic STRs: Considerations of the DNA commission of the International Society for Forensic Genetics (ISFG) on minimal nomenclature requirements



Walther Parson^{a,b,*}, David Ballard^c, Bruce Budowle^{d,e}, John M. Butler^f, Katherine B. Gettings^f, Peter Gill^{g,h}, Leonor Gusmão^{i,j,k}, Douglas R. Hares^l, Jodi A. Irwin^l, Jonathan L. King^d, Peter de Knijff^m, Niels Morlingⁿ, Mechthild Prinz^o, Peter M. Schneider^p, Christophe Van Neste^q, Sascha Willuweit^r, Christopher Phillips^s

the **full sequence** (sequence string),
the **alignment of sequences** relative to a reference sequence
the **annotation** of alleles

3.1. Most Cited Articles, 2018 (Published IF Window 2016-2017)

Citations	Citations (lifetime)	Article Title	Authors	Publication Year	Document Type
32	66	Massively parallel sequencing of forensic STRs: Considerations of the DNA commission of the International Society for Forensic Genetics (ISFG) on minimal nomenclature requirements	Parson W.,Ballard D.,Budowle B.,Butler J.M.,Gettings K.B.,Gill P.,Gusmao L.,Hares D.R.,Irwin J.A.,King J.L.,Knijff P.D.,Morling N.,Prinz M.,Schneider P.M.,Neste C.V.,Willuweit S.,Phillips C.	2016	Article
29	81	Evaluation of the Illumina® Beta Version ForenSeq™ DNA Signature Prep Kit for use in genetic profiling	Churchill J.D.,Schmedes S.E.,King J.L.,Budowle B.	2016	Article
27	55	Sequence variation of 22 autosomal STR loci detected by next generation sequencing	Gettings K.B.,Kiesler K.M.,Faith S.A.,Montano E.,Baker C.H.,Young B.A.,Guerrieri R.A.,Vallone P.M.	2016	Article
26	39	Developmental validation of the MiSeq FGx Forensic Genomics System for Targeted Next Generation Sequencing in Forensic DNA Casework and Database Laboratories	Jager A.C.,Alvarez M.L.,Davis C.P.,Guzman E.,Han Y.,Way L.,Walichewicz P.,Silva D.,Pham N.,Caves G.,Bruand J.,Schlesinger F.,Pond S.J.K.,Varlaro J.,Stephens K.M.,Holt C.L.	2017	Article
23	35	Massively parallel sequencing of short tandem repeats - Population data and mixture analysis results for the PowerSeq™ system	Van Der Gaag K.J.,De Leeuw R.H.,Hoogenboom J.,Patel J.,Storts D.R.,Laros J.F.J.,De Knijff P.	2016	Article

Forensic Science International: Genetics 24 (2016) 97–102

Contents lists available at ScienceDirect

Forensic Science International: Genetics

journal homepage: www.elsevier.com/locate/fsig



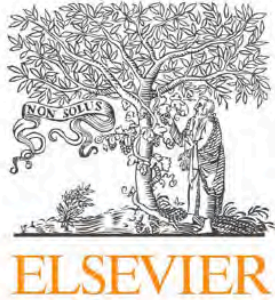
ELSEVIER



Recommendations of the DNA Commission of the International Society for Forensic Genetics (ISFG) on quality control of autosomal Short Tandem Repeat allele frequency databasing (STRidER)



Martin Bodner^a, Ingo Bastisch^b, John M. Butler^c, Rolf Fimmers^d, Peter Gill^{e,f}, Leonor Gusmão^{g,h,i}, Niels Morling^j, Christopher Phillips^k, Mechthild Prinz^l, Peter M. Schneider^m, Walther Parson^{a,n,*}



Contents lists available at [ScienceDirect](#)

Forensic Science International: Genetics

journal homepage: www.elsevier.com/locate/fsigen



Research paper

STRSeq: A catalog of sequence diversity at human identification Short Tandem Repeat loci



Katherine Butler Gettings^{a,*}, Lisa A. Borsuk^a, David Ballard^b, Martin Bodner^c, Bruce Budowle^{d,e}, Laurence Devesse^b, Jonathan King^d, Walther Parson^{c,f}, Christopher Phillips^g, Peter M. Vallone^a

NCBI BioProject—STRseq and STRidER
Collaboration in QC and exchange of data

“The devil’s in the detail”: Release of an expanded, enhanced and dynamically revised forensic STR Sequence Guide

(2018)



C. Phillips^{a,*}, K. Butler Gettings^b, J.L. King^c, D. Ballard^d, M. Bodner^e, L. Borsuk^b, W. Parson^{e,f}

^a Forensic Genetics Unit, Institute of Forensic Sciences, University of Santiago de Compostela, Spain

^b National Institute of Standards and Technology, Biomolecular Measurement Division, Gaithersburg, MD, USA

^c Institute of Applied Genetics, Department of Molecular and Medical Genetics, University of North Texas Health Science Center, Fort Worth, TX, USA

^d King’s Forensics, King’s College London, Franklin-Wilkins Building, London, UK

^e Institute of Legal Medicine, Medical University of Innsbruck, Innsbruck, Austria

^f Forensic Science Program, The Pennsylvania State University, University Park, PA, USA, USA



+ **revised** STR Sequence Guide as **dynamic** document at STRidER

HOME QUERY BATCH QUERY ABOUT FREQUENCIES FORMULAE QUALITY CONTROL STR SEQ NOMENCLATURE

STR Sequence Nomenclature

The ‘Forensic STR Sequence Structure’ file is an updated set of forensic STR sequences that was originally published as *Supplementary File S1* in the article:

The most recent version of this permanently curated and updated Forensic STR sequence structure file containing updated information is available for download [here](#). The updates since the last version are reported in a change log contained in the file. To receive information on new releases of the Forensic STR sequence structure file and to stay updated about STRidER, [register here](#) for the STRidER newsletter.

Goals

Continue collection of STR sequence information to understand variation

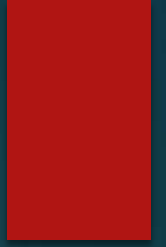
Update STRSeq @NCBI and STR sequence guide @STRidER

Harmonize efforts to develop a common STR nomenclature system

Propose a common STR nomenclature system to the community through the ISFG

This is a nomenclature panel

This is **NOT** a panel to focus on technical and analytical problems



Katherine Gettings - NIST Applied Genetics Group

- ▶ NIST population sample sequencing
- ▶ Reference materials for STR sequencing

STRAND *working group*

align | name | define

NIST population sample sequencing

Forensic Science International: Genetics 37 (2018) 106–114

Contents lists available at ScienceDirect

Forensic Science International: Genetics

journal homepage: www.elsevier.com/locate/FSIG

Research paper

Sequence-based U.S. population data for 27 autosomal STR loci

Katherine Butler Gettings^a, Lisa A. Borsuk, Carolyn R. Steffen, Kevin M. Kistler^a

^a U.S. National Institute of Standards and Technology, Biomolecular Measurement Division, 100 Bureau Drive, Gaithersburg, MD 20899, USA

ARTICLE INFO

ABSTRACT

Keywords: STR, Allele frequency, Sequence

This manuscript reports Short Tandem Repeat (STR) sequence data for 27 autosomal STR loci: D1S1656, TPOX, D2S441, D16S1043, D7S820, D8S1179, D9S1122, D10S1248, D17S1301, D18S51, D19S433, D20S482, D21S11, Penta E, and Penta D. All samples have been evaluated using the same bioinformatic pipelines and all samples have been evaluated at all loci. Each reported sequence includes by the most recent guidance of the International Commission on Forensic Genetics (ICFG) for each population and accession numbers are reported for each sequence, and a link to the STR allele frequency project (STRAF) is provided. The increase in heterozygosity across populations (approximately 10% to 5% per population) and the increase in heterozygosity from 10 to 5 percent per population point increase in average heterozygosity. Direct traditional CE results, such as informing stutter model parameters for population sampling strategies in light of sequence-based STR loci. This NIST 1036 data set is expected to be a valuable resource for forensic casework by providing high-confidence sequence data already the basis for population statistics in many forensic laboratories.

1. Introduction

In forensic casework, Short Tandem Repeat (STR) allele frequency data is used to calculate statistical weight when a person of interest cannot be excluded as a possible contributor of genetic material recovered from an item of evidence. This statistical weight should be derived from the same level of information which was used for comparison; therefore, implementation of STR sequencing into forensic casework necessitates the development of appropriate allele frequency data sets.

Several recent publications have reported sequence-based allele frequency data for autosomal STR loci [1–6]. This manuscript reports high-confidence autosomal STR sequence-based allele frequencies for N = 1036 across 27 autosomal STR loci: D1S1656, TPOX, D2S441, D2S1338, D3S1358, D4S2408, FGA, D5S818, CSF1PO, D6S1043, D7S820, D8S1179, D9S1122, D10S1248, TH01, vWA, D12S391, D13S317, Penta E, D16S539, D17S1301, D18S51, D19S433, D20S482, D21S11, Penta E, and Penta D.

The preceding allele frequencies for each locus by population in this dataset; analysis with quality flanking sequence at every loci confirmation of all null same length but different repeat units.

The preceding allele frequencies for each locus by population in this dataset; analysis with quality flanking sequence at every loci confirmation of all null same length but different repeat units.

^{*} Corresponding author at: National Institute of Standards and Technology, Biomolecular Measurement Division, E-mail addresses: katherine.gettings@nist.gov (K.B. Gettings), lisa.borsuk@nist.gov (L.A. Borsuk), becky.kevin.kistler@nist.gov (K.M. Kistler), ceter.vallone@nist.gov (P.M. Vallone).

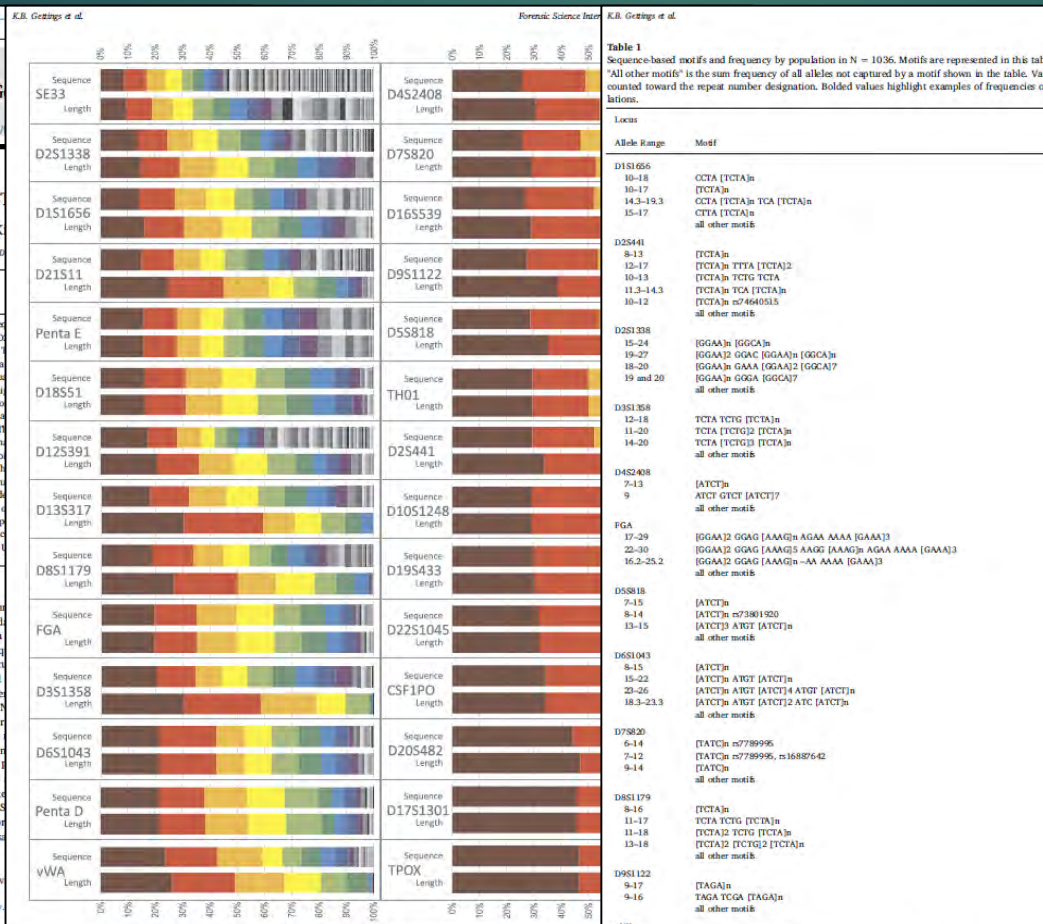


Fig. 2. Across-population allele frequency distribution per locus, by sequence and by length in N = 1036. Loci are sorted in descending order of the most common allele at each locus (first column top to bottom followed by second column top to bottom). The color coding facilitates comparisons within and across loci, with any remaining alleles shown in grayscale. Sequence data for this reported in [9].

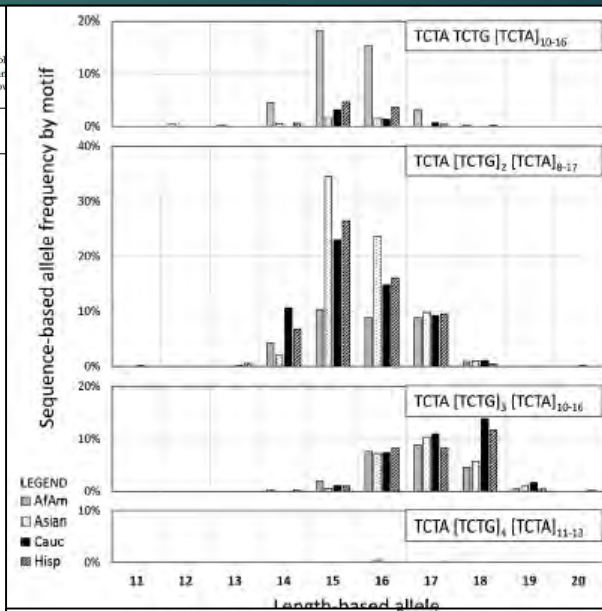


Fig. 3. D3S1358 frequency distribution among the primary motifs by length-based allele and population in N = 1036. The motif is defined as: the first subunit is fixed TCTA, the second subunit is definitive of the motif with TCTG varying from one to four, and the third subunit contains a widely varying number of TCTA repeats. For simplicity, seven additional rare motif alleles present in the data set have been excluded from this figure.

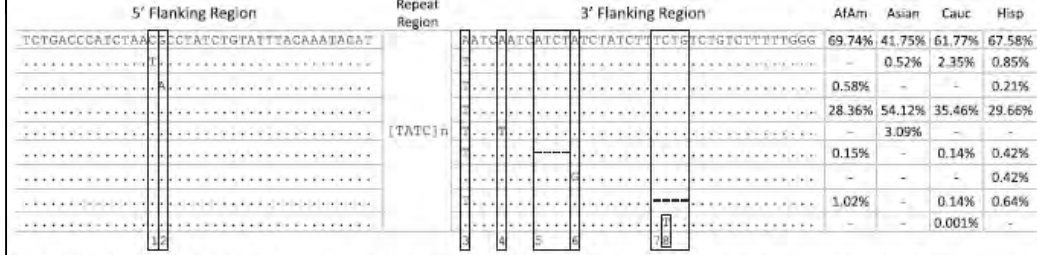


Fig. 4. D13S317 frequency distribution by population of the nine flanking region motifs identified in N = 1036. The first row of 5' and 3' flanking sequence is consistent with GRCh38, and is the most common sequence found in this data set. Dots in subsequent rows represent bases matching the first row. Flanking polymorphisms are identified by numbers one through eight in the bottom row: 1) rs73250432 C > T, 2) rs146621667 G > A, 3) rs9546005 A > T, 4) rs202043589 A > T, 5) rs1442523705 delATCT, 6) rs2137543825 A > G, 7) rs561167308 delTCTG, and 8) rs768323113 C > T. Variation in repeat region length, combined with these flanking region polymorphisms, results in 32 sequence-based alleles at this locus. Three additional D13S317 alleles in this data set result from repeat region sequence variants, not observed once, and have been excluded from this figure.

STRAND *working group*

align | name | define

NIST Public Data Repository

1.1.0-beta

<https://doi.org/10.18434/T4/1500024>

Public Data Resource

Sequence-based U.S. population data for 27 autosomal STR loci

Contact: [Katherine Gettings](#)

Identifier: [doi:10.18434/T4/1500024](https://doi.org/10.18434/T4/1500024)

Version: 1.0.1 Last modified: 2018-06-14

Description

This information and data are supplemental files associated with: K.B. Gettings, L.A. Borsuk, C.R. Steffen, K.M. Genetics 37 (2018) 106-115. The primary data consists of sequence-based allele frequencies for N=1036 and D2S1338, D3S1358, D4S2408, FGA, D5S818, CSF1PO, D6S1043, D7S820, D8S1179, D9S1122, D10S1248, TH01. This information is expected to support the implementation of sequence-based STR analysis in forensic applications. The 42 sequencing runs performed to generate the allele frequency data (S1 - Run Metrics); coverage per locus (S2 - Coverage); allele frequency data (S3 - Frequencies); GRCh38 reference coordinates for genomic regions reported (S5 - Flank Polymorph); number of alleles, expected and observed heterozygosity, and p-values associated with disequilibrium (S7 - LD p-values); and pairwise Fst values by population for the 27 auSTR loci (Supp Table 8 - P quality control of the data.

Subject Keywords: STR, forensic, sequence, population, allele frequency

Data Access

These data are public.

Files Click on the file/row in the table below to view more details.

Name

[NIST1036_auSTR_Seq_SuppFile1.pdf](#)

[NIST1036_auSTR_Seq_SuppFile1.pdf.sha256](#)

[NIST1036_auSTR_Seq_SuppTables.xlsx](#)

[NIST1036_auSTR_Seq_SuppTables.xlsx.sha256](#)

References

This data is referenced in:

<https://doi.org/10.1016/j.fsigen.2018.07.013>

11 pages of methods, including:

Population, Sample Type	Locus	CE	UAS	SR	Source
African American, Blood	Penta D	2.2,13.4	2.2,14	2.2,13.4	1 bp deletion rs536566765
African American, Buccal	D5S818	7,12*	8,12	8,12*	Assumed 4 bp deletion outside ForenSeq amplicon
Hispanic, Blood	D7S820	10.3,11	11,11	10.3,11	1 bp deletion, rs897512434
Caucasian, Blood	D9S1122	(12),14	12,14	11.2,14	2 bp deletion rs754976988, overlaps with CE primer binding site

Table D. Discordance between CE and sequence data observed in N=1036. Bolded genotypes used in allele frequency calculations. *8,12 used in 1036 sequence-based frequencies, 7,12 used in 1036 CE-based frequencies [2].

STRAND *working group*

align | name | define

NIST Public Data Repository

1.1.0-beta

<https://doi.org/10.18434/T4/1500024>

Public Data Resource
Sequence-based U.S. population data for 27 autosomal STR loci

Contact: [Katherine Gettings](#)

Identifier: [doi:10.18434/T4/1500024](https://doi.org/10.18434/T4/1500024)

[Visit Home Page](#)

Table with columns: Locus, Allele, Frequency, Count, and Sequence. The table lists 27 autosomal STR loci and their corresponding alleles and frequencies. A large watermark text '~800 unique sequences' is overlaid on the table.

NIST1036_auSTR_Seq_SuppTables.xlsx

application/vnd.openxmlformats-officedocument.spreadsheetml.sheet

NIST1036_auSTR_Seq_SuppTables.xlsx sha256

text/plain

References

This data is referenced in:

<https://doi.org/10.1016/j.fsigen.2018.07.013>

STRAND *working group*

align | name | define

Controls for STR Sequencing

- ▶ NIST SRM 2391d
- ▶ Expected release Summer 2019
- ▶ STR sequence data:
 - ▶ ForenSeq
 - ▶ Precision ID GlobalFiler NGS STR Panel v2
 - ▶ PowerSeq 46GY (prototype)



STRAND *working group*

align | name | define

Controls for STR Sequencing

- ▶ NIST GiaB
 - ▶ 7 Coriell cell lines
 - ▶ One individual and two trios
 - ▶ PCR-free prep, HiSeq 150 and 250, PacBio
- ▶ STR project
 - ▶ Any “novel” marker can be characterized
 - ▶ ISFG Poster targeting ~500 “novel” STRs



STRAND *working group*

align | name | define

Lisa Borsuk - NIST Applied Genetics Group



- ▶ BioProject Structure
- ▶ Record Structure
- ▶ Current Sample Sets
- ▶ BioProject Status

STRAND *working group*

align | name | define



The STR Sequencing Project (human)

Umbrella project

National Institute of Standards and Technology

Accession: PRJNA380127



U.S. National Library of Medicine
National Center for Biotechnology Information

<https://www.ncbi.nlm.nih.gov/bioproject/380127>



Homo sapiens

STRSeq Commonly Used Autosomal STR Loci

Umbrella project

National Institute of Standards and Technology

Accession: PRJNA380345 ID: 380345



Homo sapiens

STRSeq Alternate Autosomal STR Loci

Umbrella project

National Institute of Standards and Technology

Accession: PRJNA380346 ID: 380346



Homo sapiens

STRSeq X-Chromosomal STR Loci

Umbrella project

National Institute of Standards and Technology

Accession: PRJNA380348 ID: 380348



Homo sapiens

STRSeq Y-Chromosomal STR Loci

Umbrella project

National Institute of Standards and Technology

Accession: PRJNA380347 ID: 380347



STRAND *working group*

align | name | define

Homo sapiens Accession: PRJNA380345 ID: 380345

STRSeq Commonly Used Autosomal STR Loci

This sub-project of the STR Sequencing Project encompasses the data pertaining to 24 autosomal STR loci which are commonly targeted in human identification assays.

Accession	PRJNA380345
Type	Umbrella project
Publications (total 5)	1. Borsuk LA <i>et al.</i> , "Sequence-based US population data for the SE33 locus.", <i>Electrophoresis</i> , 2018 Jun 1;39(21):2694-2701 More...
Submission	Registration date: 24-Mar-2017 National Institute of Standards and Technology
Related Resources	<ul style="list-style-type: none"> • STRSeq • STRidER
Relevance	Human Identification

NAVIGATE UP

This project is a component of the The STR Sequencing Project (human)

NAVIGATE ACROSS

3 additional projects are components of the The STR Sequencing Project (human).

Project Data:

Resource Name	Number of Links
SEQUENCE DATA	
Nucleotide (Genomic DNA)	1251
PUBLICATIONS	
PubMed	5

Homo sapiens encompasses the following 24 sub-projects:

Project Type	Number of Projects	
targeted loci	24	
BioProject accession	Organism	Title
PRJNA380553	Homo sapiens	STRSeq D1S1656 Sequence-Based Alleles (National Institute of Standards...)
PRJNA380554	Homo sapiens	STRSeq TPOX Sequence-Based Alleles (National Institute of Standards...)
PRJNA380555	Homo sapiens	STRSeq D2S441 Sequence-Based Alleles (National Institute of Standards...)
PRJNA380556	Homo sapiens	STRSeq D2S1338 Sequence-Based Alleles (National Institute of Standards...)
PRJNA380558	Homo sapiens	STRSeq D3S1358 Sequence-Based Alleles (National Institute of Standards...)
List all 24 'targeted loci' projects...		

Links to BioProjects

The STR Sequencing Project (human)

Accession: PRJNA380127 ID: 380127

The purpose of STRSeq is to facilitate the description of sequence-based alleles at the Short Tandem Repeat (STR) loci targeted in human identification assays. This collaborative effort of the international forensic DNA community, which has been endorsed by the executive board of the ISFG (International Society of Forensic Genetics), provides a framework for communication among laboratories. Each record contains: (a) observed sequence of an STR region, (b) annotation of the repeat region ("bracketing") and flanking region polymorphisms, (c) information regarding the sequencing assay and data quality, and (d) backward compatible length-based allelic designation. Data within the umbrella project is organized into locus sub-projects, and can be accessed by browsing, BLAST searching, or ftp download at NCBI. For comments or questions, please contact strseq@nist.gov.

Accession	PRJNA380127
Type	Umbrella project
Publications (total 5)	1. Borsuk LA <i>et al.</i> , "Sequence-based US population data for the SE33 locus.", <i>Electrophoresis</i> , 2018 Jun 1;39(21):2694-2701 More...
Submission	Registration date: 22-Mar-2017 National Institute of Standards and Technology
Related Resources	<ul style="list-style-type: none"> • STRSeq • STRidER
Relevance	Human Identification

Project Data:

Resource Name	Number of Links
SEQUENCE DATA	
Nucleotide (Genomic DNA)	1320
PUBLICATIONS	
PubMed	5

The STR Sequencing Project (human) encompasses the following 4 sub-projects:

Project Type	Number of Projects	
Umbrella project	4	
BioProject accession	Name	Title
PRJNA380345	Homo sapiens	STRSeq Commonly Used Autosomal STR Loci (National Institute of Standards...)
PRJNA380346	Homo sapiens	STRSeq Alternate Autosomal STR Loci (National Institute of Standards...)
PRJNA380347	Homo sapiens	STRSeq Y-Chromosomal STR Loci (National Institute of Standards...)
PRJNA380348	Homo sapiens	STRSeq X-Chromosomal STR Loci (National Institute of Standards...)

STRAND *working group*

align | name | define



GenBank → Send to: -

Homo sapiens microsatellite rs1019813099 sequence

GenBank: MH174843.1 Graphics →

FEATURES

source	Location/Qualifiers
	1..138
	/organism="Homo sapiens"
	/mol_type="genomic DNA"
	/db_xref="taxon:9606"
misc feature	1..117
	/note="Illumina ForenSeq DNA Signature Prep Kit"
misc feature	6..138
	/note="Promega PowerSeq 46GY System"
variation	38
	/note="C/T SNP"
	/db_xref="dbSNP:rs1019813099"
repeat region	62..117
	/rpt_type=tandem
	/satellite="microsatellite:D1S1656"

ORIGIN

```

1 ttcagagaaa tagaatcact agggaaccaa atatatacac atacaattaa acacacacac
61 acctatctat ctatctatct atctatctat ctatctatct atctatctat ctatctacat
121 cacacagttg acccttga
//
    
```

62..117
/rpt_type=tandem
/satellite="microsatellite:D1S1656"

repeat_region ▾ [Feature](#) 1 of 1 MH174843: 1 segment [Details](#) [Display: FASTA](#) [GenBank](#) [Help](#) ✕

of Standards and Technolo
Gaithersburg, Maryland 20

collaborative effort of the international forensic DNA community. The purpose of this project is to facilitate the description of sequence-based STR alleles. Additional resources can be found at strseq.nist.gov. For questions or feedback, please contact strseq@nist.gov. Allele frequency data can be accessed in the strider.online database.

GenBank → Send to: -

Homo sapiens microsatellite D1S1656 14 CCTA [TCTA]13 rs1019813099 sequence

GenBank: MH174843.1 [FASTA](#) [Graphics](#)

Go to:

LOCUS MH174843 138 bp DNA linear PRI 04-SEP-2018

DEFINITION Homo sapiens microsatellite D1S1656 14 CCTA [TCTA]13 rs1019813099 sequence.

ACCESSION MH174843

VERSION MH174843.1

DBLINK BioProject: [PRJNA380553](#)

KEYWORDS STRseq; STR; D1S1656.

SOURCE Homo sapiens (human)

ORGANISM [Homo sapiens](#)
Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini; Catarrhini; Hominidae; Homo.

REFERENCE 1 (bases 1 to 138)

AUTHORS Gettings,K.B., Borsuk,L.A., Ballard,D., Bodner,M., Budowle,B., Devesse,L., King,J., Parson,W., Phillips,C. and Vallone,P.M.

TITLE STRseq: A catalog of sequence diversity at human identification Short Tandem Repeat loci

JOURNAL Forensic Sci Int Genet 31, 111-117 (2017)

PUBMED [2888135](#)

REFERENCE 2 (bases 1 to 138)

AUTHORS NIST,A.G.G.

TITLE Direct Submission

JOURNAL Submitted (06-APR-2018) Applied Genetics Group, National Institute of Standards and Technology, 100 Bureau Drive, MS-8314, Gaithersburg, Maryland 20899, United States of America

COMMENT Annotation ('bracketing') of the repeat region is consistent with the guidance of the ISFG (International Society of Forensic Genetics), PMID: 26844919. Lower case letters in the 'Bracketed repeat' region below denote uncounted bases. The given length-based allele value was determined using the designated length-based technology. Variation in the length-based allele between individuals or assays can result from indels in flanking regions. The length of reported sequence is dependent on the assay and the quality of the flanking sequence. This information is provided as part of the STR Sequencing Project (STRseq), a collaborative effort of the international forensic DNA community. The purpose of this project is to facilitate the description of sequence-based STR alleles. Additional resources can be found at strseq.nist.gov. For questions or feedback, please contact strseq@nist.gov. Allele frequency data can be accessed in the strider.online database.

```

##HumanSTR-START##
STR locus name      :: D1S1656
Length-based allele :: 14
Bracketed repeat   :: CCTA [TCTA]13
Sequencing technology :: MiSeq FGx; MiSeq
Coverage           :: >30X
Length-based tech.  :: PowerPlex Fusion, 3130xl
Assembly           :: GRCh38 (GCF_000001405)
Chromosome         :: 1
RefSeq Accession   :: NC_000001.11
Chrom. Location    :: 230769555..230769704
Repeat Location    :: 230769616..230769683
Cytogenetic Location :: 1q42.2
##HumanSTR-END##
    
```

FEATURES

source	Location/Qualifiers
	1..138
	/organism="Homo sapiens"
	/mol_type="genomic DNA"
	/db_xref="taxon:9606"
misc feature	1..117
	/note="Illumina ForenSeq DNA Signature Prep Kit"
misc feature	6..138
	/note="Promega PowerSeq 46GY System"
variation	38
	/note="C/T SNP"
	/db_xref="dbSNP:rs1019813099"
repeat region	62..117
	/rpt_type=tandem
	/satellite="microsatellite:D1S1656"

ORIGIN

```

1 ttcagagaaa tagaatcact agggaaccaa atatatacac atacaattaa acacacacac
61 acctatctat ctatctatct atctatctat ctatctatct atctatctat ctatctacat
121 cacacagttg acccttga
//
    
```

STRAND *working group*

align | name | define




1786 Samples

1043 Samples

839 Samples

944 Samples

ForenSeq Verogen,
PowerSeq 46GY (prototype)
Promega, and GlobalFiler
NGS Thermo Fisher (and Sanger
()
CE Information

ForenSeq Verogen,
CE Information

ForenSeq Verogen,
CE Information

ForenSeq Verogen,
CE Information

Working towards including more sets of data
Collaborating with STRidER

STRAND *working group*

align | name | define



BioProject Status

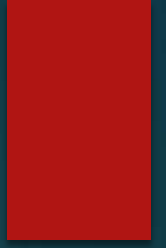
Project Data:

Resource Name	Number of Links
SEQUENCE DATA	
Nucleotide (Genomic DNA)	1320
PUBLICATIONS	
PubMed	5

- ▶ Currently sequences for autosomal loci have been submitted
 - ▶ Additional autosomal sequences continue to be loaded
- ▶ Y and X data is going to be submitted in the next few months
 - ▶ 2019 Summer - Fall

STRAND *working group*

align | name | define



Jonathan King

STRAND *working group*

align | name | define

▶ Current stable release

▶ v3

▶ fastq processing (C++ script)



▶ Data visualization and processing (Microsoft Excel)



▶ Development version

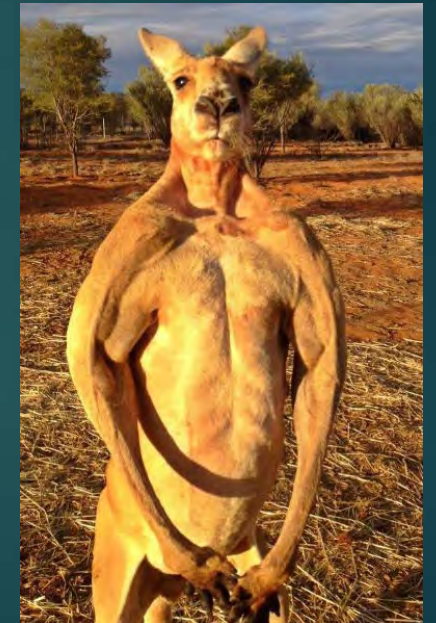
▶ R package

▶ Combination of C++ script and functions

▶ Full UI

▶ Data visualization

▶ Allele-calling



thekangaroosanctuary

STRAND *working group*

align | name | define

- ▶ **Default settings** (ForenSeq-STRs only)
 - ▶ ~2s per sample
 - ▶ (~15s from fastq)
 - ▶ ~97.9% of alleles called automatically
 - ▶ (>99% excluding D22S1045...)
 - ▶ Stutter Assessment
 - ▶ ~6s per sample
- ▶ **Population-level processing**

STRAND *working group*

align | name | define

► In Development UI

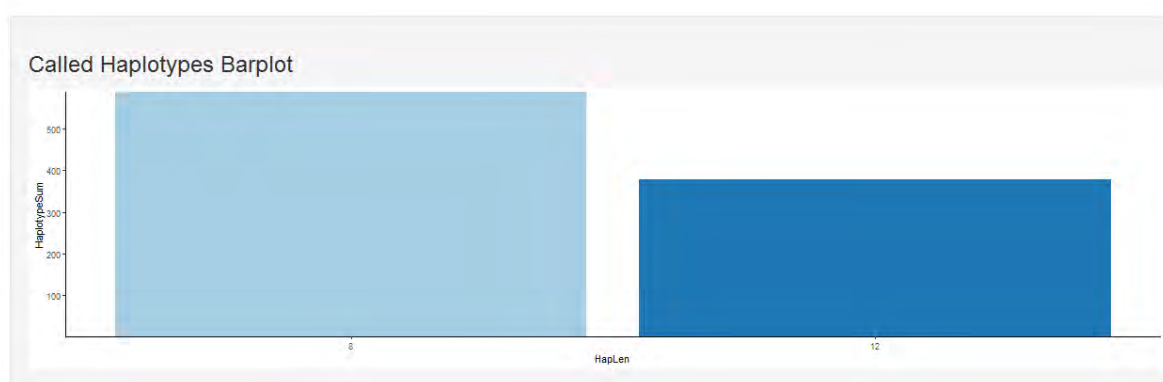
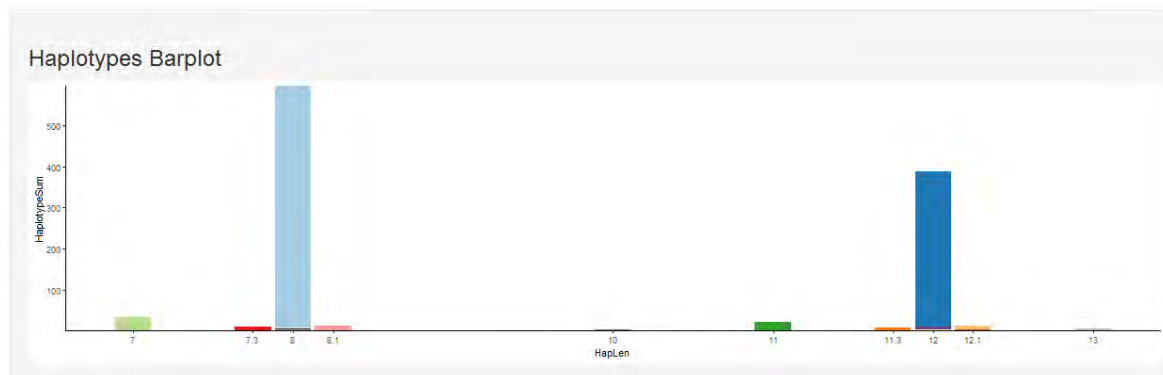
Edit and save Locus Table

Shiny app for analysis of STRait Razor data.

	Status	LocusRank	Locus	Allele	RAP	AR	SB	HaplotypeSur
1	<input checked="" type="checkbox"/>	1.00	D7S820	8.00	0.51	0.64	0.00	588.01
2	<input checked="" type="checkbox"/>	2.00	D7S820	12.00	0.33	0.64	0.00	378.01
3	<input type="checkbox"/>	3.00	D7S820	7.00	0.03		0.00	32.01
4	<input type="checkbox"/>	4.00	D7S820	11.00	0.02		0.00	22.01
5	<input type="checkbox"/>	5.00	D7S820	8.10	0.01		0.00	13.01
6	<input type="checkbox"/>	6.00	D7S820	7.30	0.01		0.00	10.01
7	<input type="checkbox"/>	7.00	D7S820	12.10	0.01		0.00	10.01
8	<input type="checkbox"/>	8.00	D7S820	11.30	0.01		0.00	7.01
9	<input type="checkbox"/>	9.00	D7S820	13.00	0.00		0.00	5.01
10	<input type="checkbox"/>	10.00	D7S820	12.00	0.00		0.00	5.01
11	<input type="checkbox"/>	11.00	D7S820	12.00	0.00		0.00	4.01
12	<input type="checkbox"/>	12.00	D7S820	12.10	0.00		0.00	3.01
13	<input type="checkbox"/>	13.00	D7S820	8.00	0.00		0.00	2.01
14	<input type="checkbox"/>	14.00	D7S820	8.00	0.00		0.00	2.01
15	<input type="checkbox"/>	15.00	D7S820	8.00	0.00		0.00	2.01
16	<input type="checkbox"/>	16.00	D7S820	10.00	0.00		0.00	2.01
17	<input type="checkbox"/>	17.00	D7S820	8.00	0.00		0.00	2.01

Save table

Save Output Table



STRAND *working group*

align | name | define



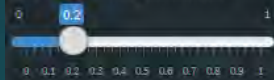
STRait Razor Analysis



Relative Allele Percentage Threshold



Strand Balance Threshold



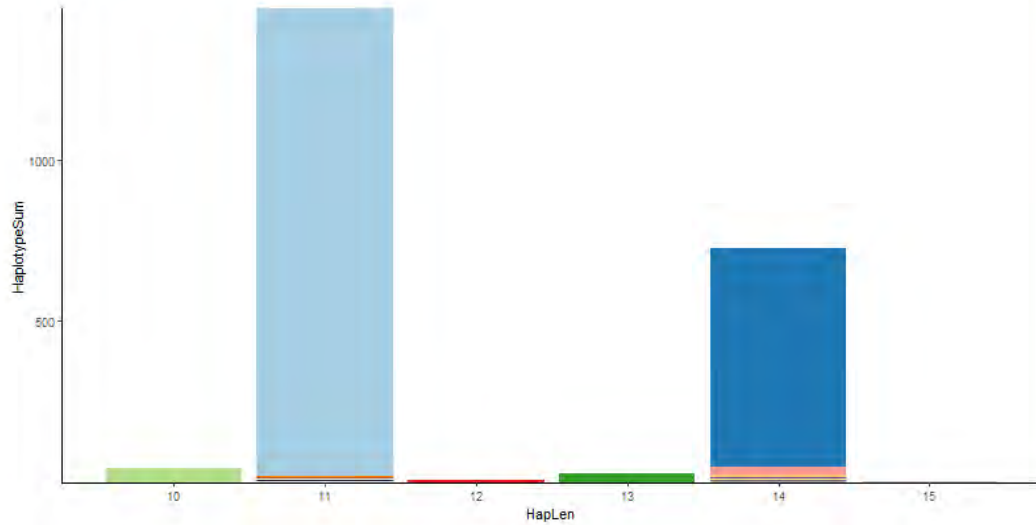
Haplotype Depth Threshold



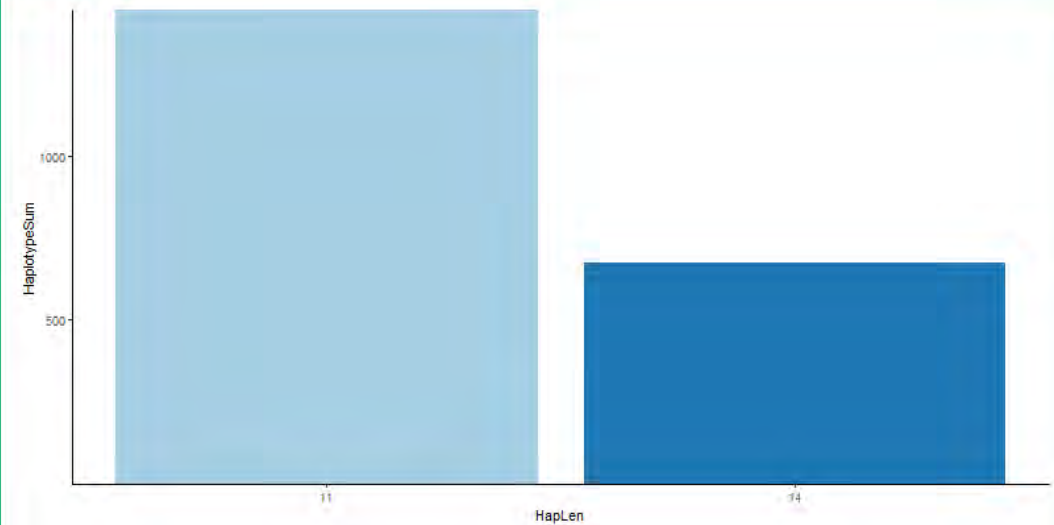
Save Output Table

Status	LocusRank	Locus

D2S441



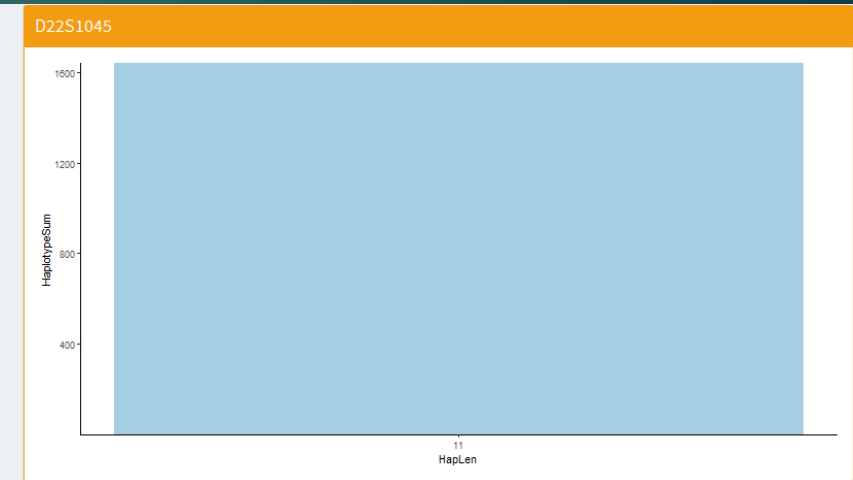
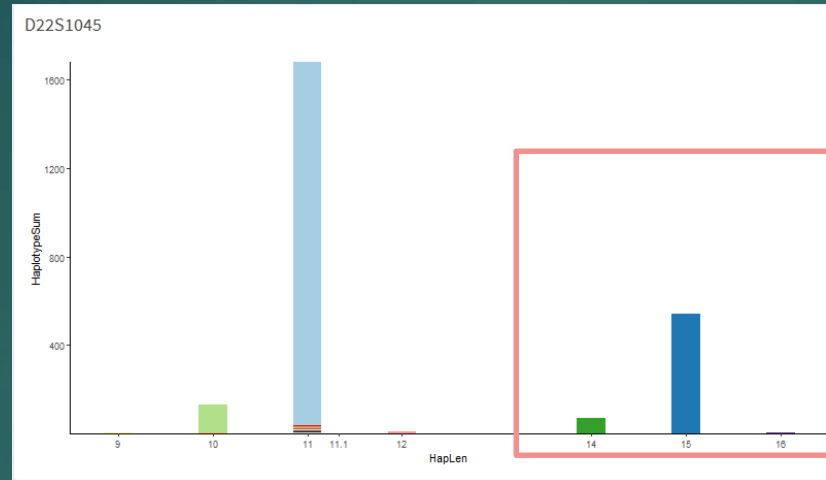
D2S441



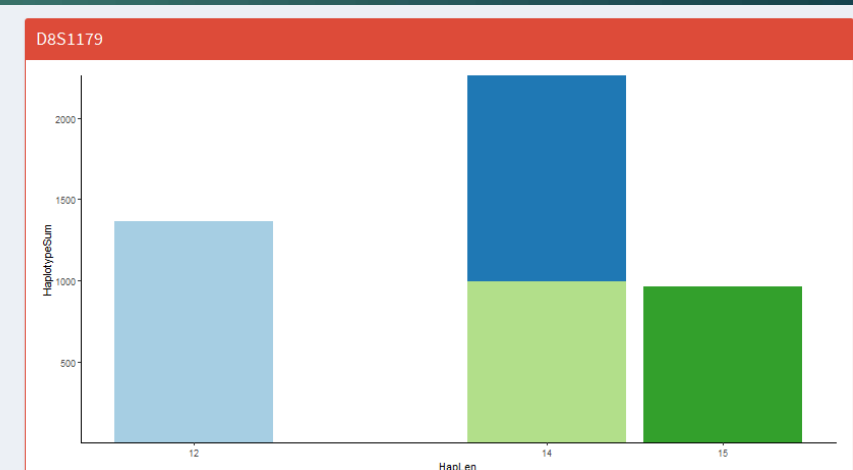
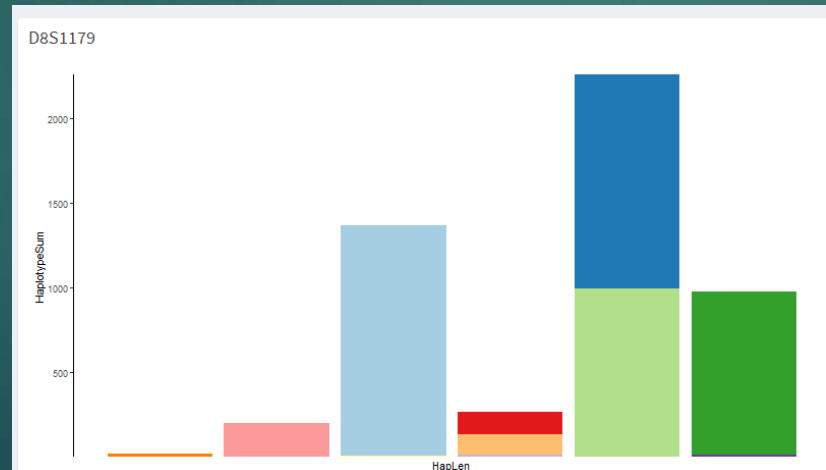
STRAND *working group*

align | name | define

Caution



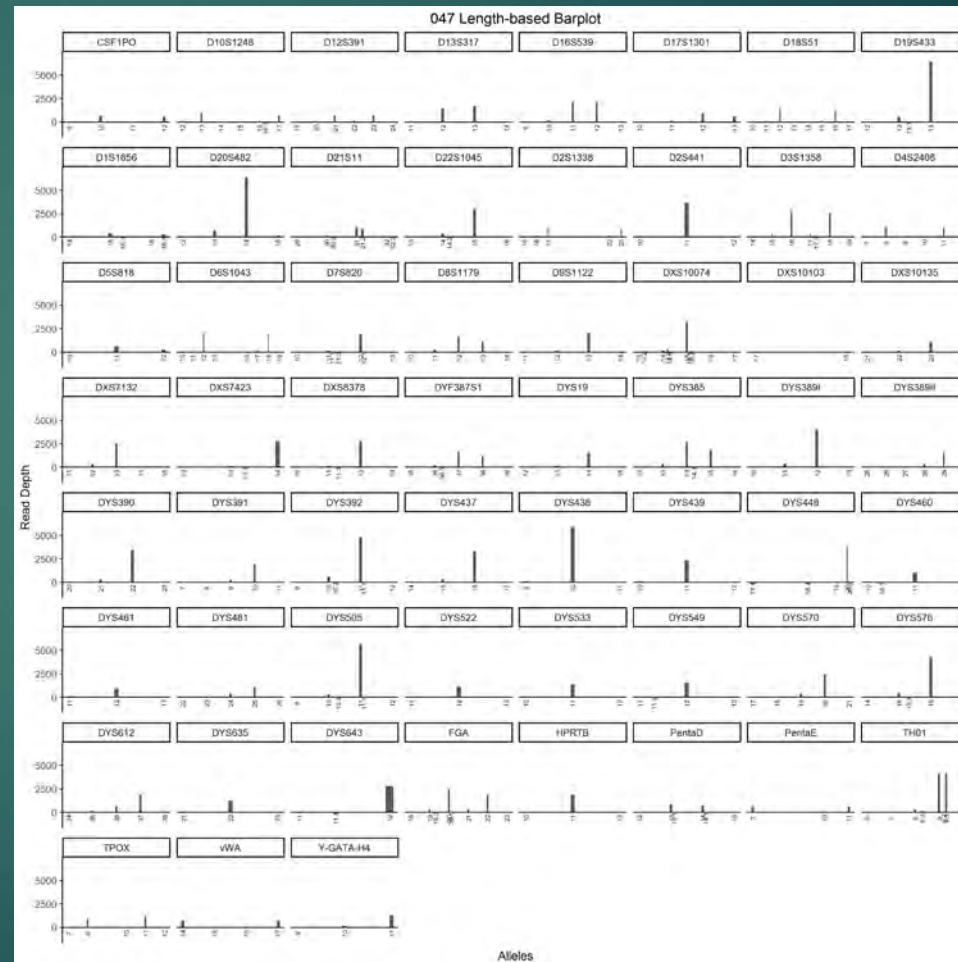
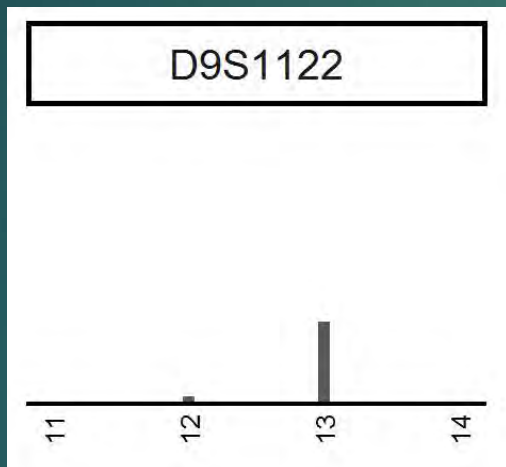
Warning



STRAND *working group*

align | name | define

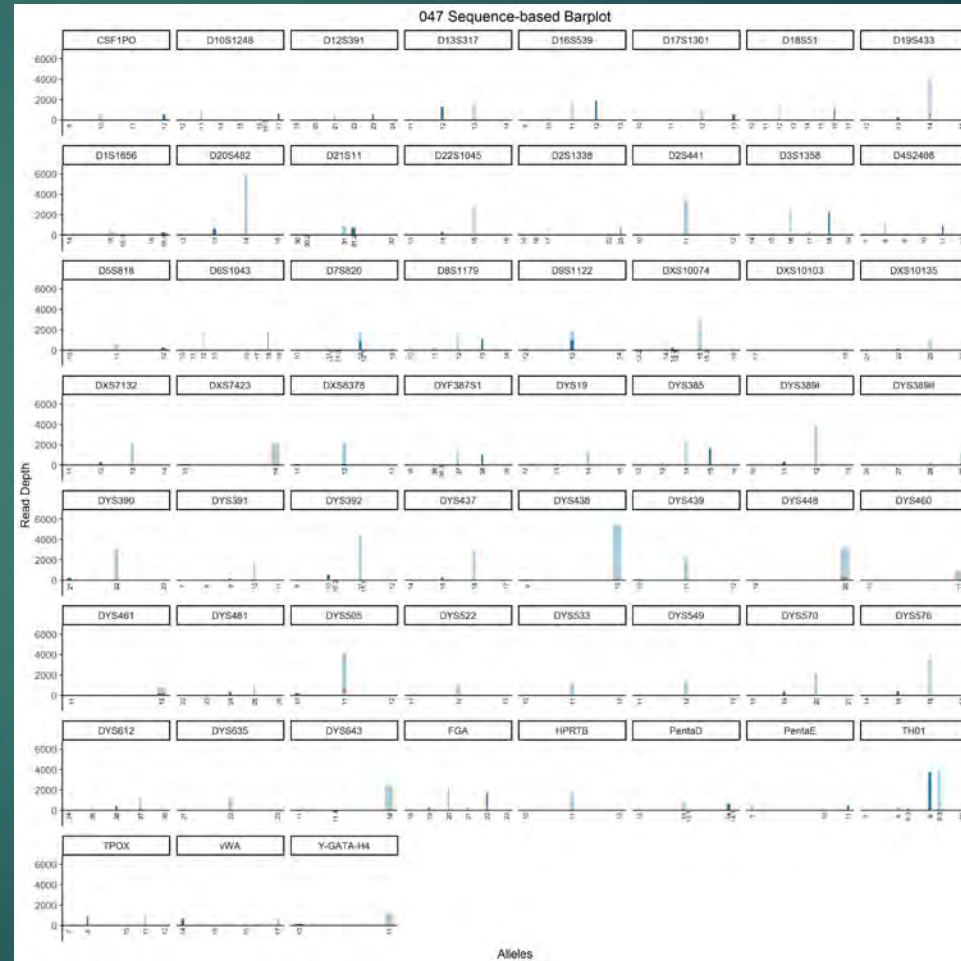
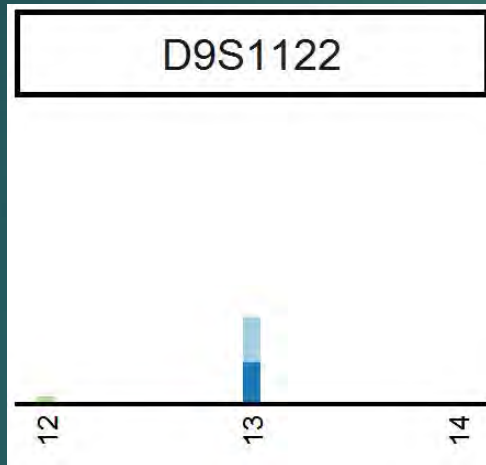
► Length-based Allele Calling



STRAND *working group*

align | name | define

► Sequence-based



STRAND *working group*

align | name | define

- ▶ Stutter Assessment
- ▶ Refine Levels
 - ▶ Locus
 - ▶ Most widely implemented
 - ▶ Length-based allele
 - ▶ Most widely studied
 - ▶ Haplotype
 - ▶ Most useful ultimately...(probably)

STRAND *working group*

align | name | define

▶ D21S11

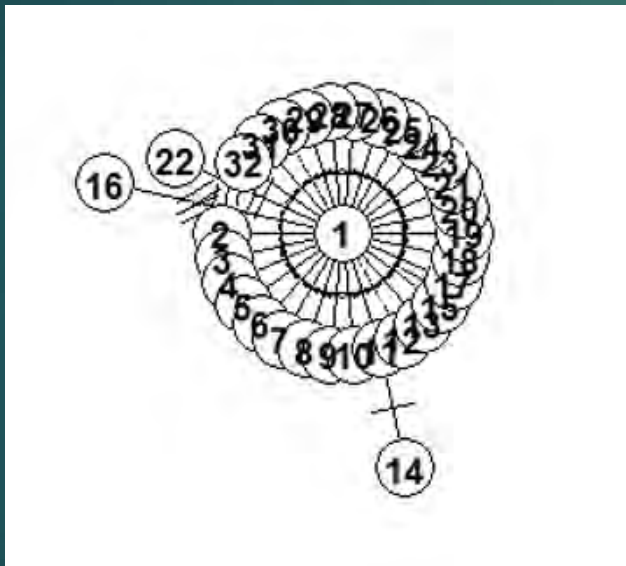
- ▶ PA: [TCTA]4 [TCTG]6 [TCTA]3 TA [TCTA]3 TCA [TCTA]2 TCCA TA [TCTA]9
 - ▶ [TCTA]4 [TCTG]6 [TCTA]3 TA [TCTA]3 TCA [TCTA]2 TCCA TA [TCTA]8
 - ▶ >6%
 - ▶ [TCTA]4 [TCTG]5 [TCTA]3 TA [TCTA]3 TCA [TCTA]2 TCCA TA [TCTA]9
 - ▶ >1%
 - ▶ [TCTA]3 [TCTG]6 [TCTA]3 TA [TCTA]3 TCA [TCTA]2 TCCA TA [TCTA]9
 - ▶ <1%
 - ▶ [TCTA]4 [TCTG]6 [TCTA]3 TA [TCTA]2 TCA [TCTA]2 TCCA TA [TCTA]9
 - ▶ <1%

STRAND *working group*

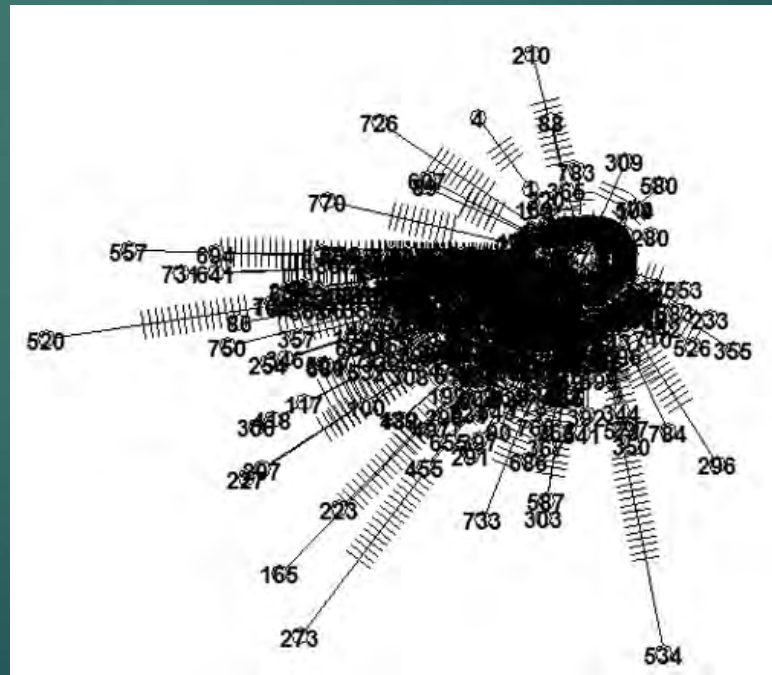
align | name | define

► Understanding Error via Haplotype Networks

CSF1PO



D19S433



STRAND *working group*

align | name | define

- ▶ Future Directions
 - ▶ Long-term vision of STRait Razor
 - ▶ Alignment-based
 - ▶ Dynamic flanking region processing/anchor assignment
 - ▶ Web interface

STRAND *working group*

align | name | define

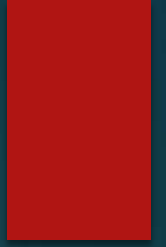
UNT Center for Human ID Research and Development Unit



STRAND *working group*

align | name | define

Chris Phillips

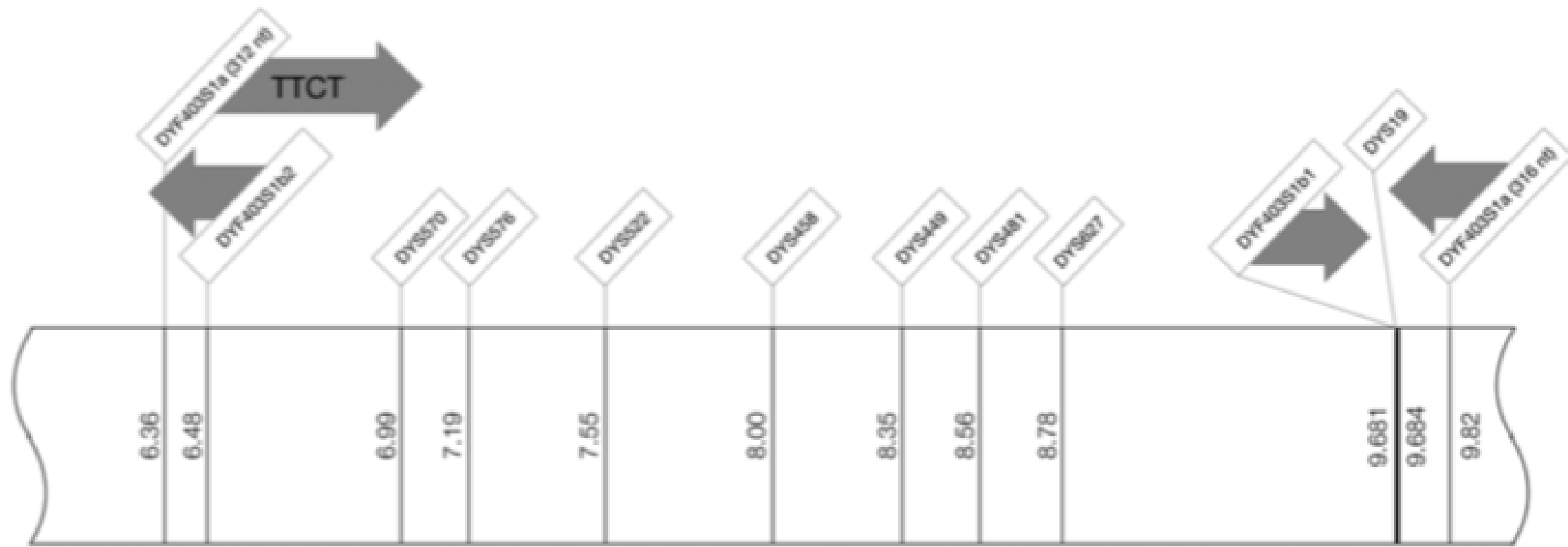


Reference genomes, existing databases

GRCh38 has regular updates but no 'audits' reveal the sequence build or alignments for forensic STRs have changed

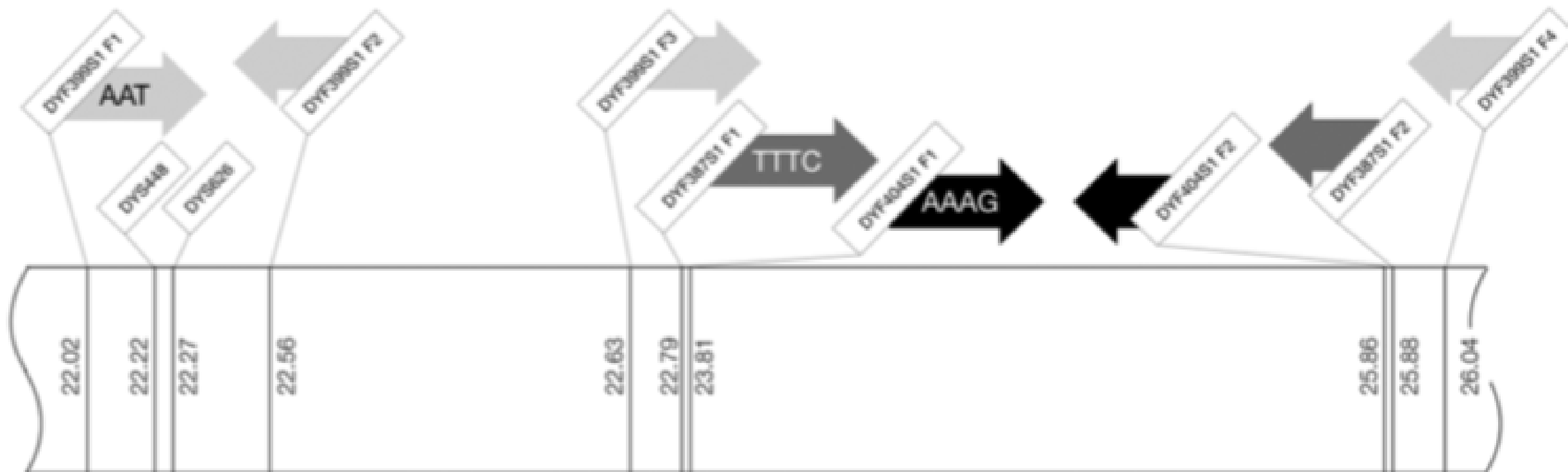
GRCh37 and GRCh 38 do have different sequence builds in one or two loci - most complex comparison of sequences is DXS10146

No distinction can be made between the fragments of duplication-inversion Y-STRs when we move away from the single strand RefSeq framework



6 Mb

10 Mb



21.8 Mb

26.2 Mb

Reference genomes, existing databases

GRCh38 has regular updates but no 'audits' reveal the sequence build or alignments for forensic STRs have changed

GRCh37 and GRCh 38 do have different sequence builds in one or two loci - most complex comparison of sequences is DXS10146

No distinction can be made between the fragments of duplication-inversion Y-STRs when we move away from the single strand RefSeq framework

Handling of insertions is quite low key when these are above a certain length

Short Indels can be insertions or deletions of the sequence element
- this is rarely fixed in 1000 Genomes / dbSNP annotations

[A/-] is rarely given as [-/A]

[G/GA] rarely [GA/G]

longer sequence elements are given as insertions and this includes STR alleles when these are compiled in 1000 Genomes (not always)

		GRCh 38															
minima reported		1	2	3	4	5	6	7	8	9	10	11	12	13	Repeat structure	Rpt	
		T C T A	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A			
8	1	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]8	8	
9	2	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]9	9	
10	3	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]10	10	
11	4	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]11	11	
11	5	T C T A	T C T A	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]2 TCTG [TCTA]8	11	
11	6	T C T A	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	TCTA TCTG [TCTA]9	11	
12	7	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]12	12	
12	8	T C T A	T C T A	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]2 TCTG [TCTA]9	12	
12	9	T C T A	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	TCTA TCTG [TCTA]10	12	
13	10	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]13	13	
13	11	T C T A	T C T A	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]2 TCTG [TCTA]10	13	
13	12	T C T A	T C T A	T C T G	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]2 [TCTG]2 [TCTA]9	13	
13	13	T C T A	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	TCTA TCTG [TCTA]11	13	
14	14	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]14	14	
14	15	T C T A	T C T A	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]2 TCTG [TCTA]11	14	
14	16	T C T A	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	TCTA TCTG [TCTA]12	14	
15	17	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]15	15	
15	18	T C T A	T C T A	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]2 TCTG [TCTA]12	15	
15	19	T C T A	T C T A	T C T G	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]2 [TCTG]2 [TCTA]11	15	
15	20	T C T A	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	TCTA TCTG [TCTA]13	15	
16	21	T C T A	T C T A	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]2 TCTG [TCTA]13	16	
16	22	T C T A	T C T A	T C T G	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]2 [TCTG]2 [TCTA]12	16	
16	23	T C T A	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	TCTA TCTG [TCTA]14	16	
16	24	T C T A	T C T G	T C T G	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	TCTA [TCTG]3 [TCTA]12	16	
17	25	T C T A	T C T A	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]2 TCTG [TCTA]14	17	
17	26	T C T A	T C T A	T C T G	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]2 [TCTG]2 [TCTA]13	17	
18	27	T C T A	T C T A	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]2 TCTG [TCTA]15	18	
18	28	T C T A	T C T A	T C T G	T C T G	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	T C T A	[TCTA]2 [TCTG]2 [TCTA]14	18	
al framework																	
8	1	A	
9	2	A	
10	3	A	
11	4	A	
11	5	A	.	G	
11	6	
12	7	A	
12	8	A	.	G	
12	9	
13	10	A	
13	11	A	.	G	
13	12	A	.	G	.	G	
13	13	
14	14	A	
14	15	A	.	G	
14	16	
15	17	A	
15	18	A	.	G	
15	19	A	.	G	.	G	
15	20	
16	21	A	.	G	
16	22	A	.	G	.	G	
16	23	
16	24	G	.	G	
17	25	A	.	G	
17	26	A	.	G	.	G	
18	27	A	.	G	
18	28	A	.	G	.	G	

CEPH
D8S1179

Reference genomes, existing databases

The two sets of co-ordinates and sequence linked to them won't change substantially

		D9S2157									
		1	2	3	4	5	6	7	8	9	10
Reference sequence	C T G T C T C A A T	A T A A T A A T A A T A A T A A T A A T A A T A A T A A T A A T A									
Flanking SNP IUPAC codes		A T A A T A A T A A T A A T A A T A A T A A T A A T A A T A A T A									
GRCh38 coordinates	133160272 133160273 133160274 133160275 133160276 133160277 133160278 133160279 133160280 133160281	133160282 133160283 133160284 133160285 133160286 133160287 133160288 133160289 133160290 133160291 133160292 133160293 133160294 133160295 133160296 133160297 133160298 133160299 133160300 133160301 133160302 133160303 133160304 133160305 133160306 133160307 133160308 133160309 133160310 133160311									
GRCh37 coordinates	136035659 136035660 136035661 136035662 136035663 136035664 136035665 136035666 136035667 136035668	136035669 136035670 136035671 136035672 136035673 136035674 136035675 136035676 136035677 136035678 136035679 136035680 136035681 136035682 136035683 136035684 136035685 136035686 136035687 136035688 136035689 136035690 136035691 136035692 136035693 136035694 136035695 136035696 136035697 136035698									
Distance from repeat region	10 9 8 7 6 5 4 3 2 1	1 2 3 4 5 6 7 8 9 10									

We can accommodate whatever framework for defining the repeat region we agree to

Reference genomes, existing databases

Databases of genomic details pertinent to a variant

Sequence databases generally compile good quality data for SNPs/Indels

STR structure data is poor - description and detail varies between loci

There may be a need to systematise some forensic STR locus names (e.g. FGA) - less regularised than SNPs/Indels, so locating novel loci very difficult

Databases of the population variation observed in a variant

In last two years >10K genome projects have found many rare variants

No STR sequence variation database exists - we will have to build our own

Forensic STR databases of CE-alleles have functioned well up to now

STR.Base popSTR ALFRED

Databases of genomic details pertinent to a variant

Databases of the population variation observed in a variant

Maps of clustered variants:

STRs and their accompanying flanking SNPs

Multiple STRs in close proximity

Indels that might influence CE size estimates

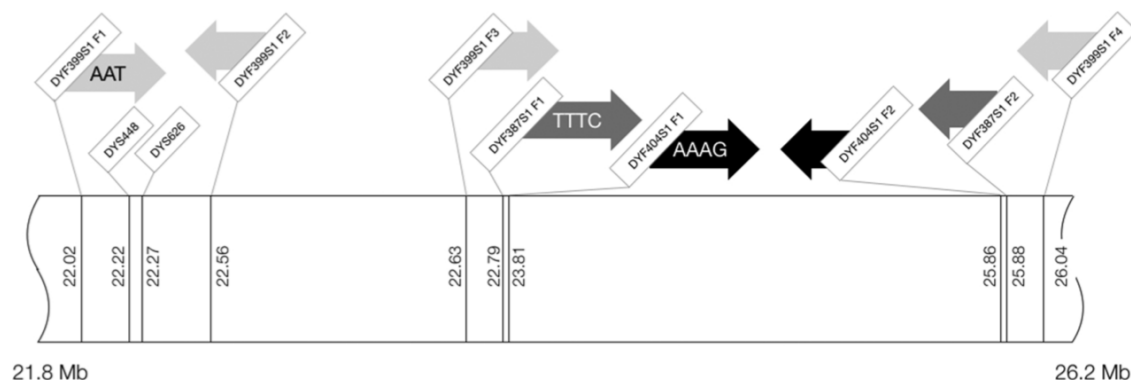
Arrangement of multiple fragment STRs

Ensembl

NCBI Genbank

STR-specific

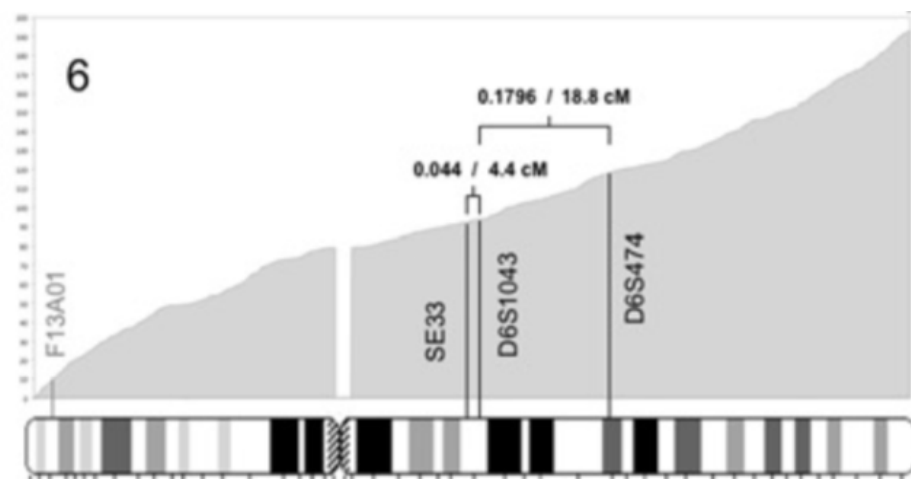
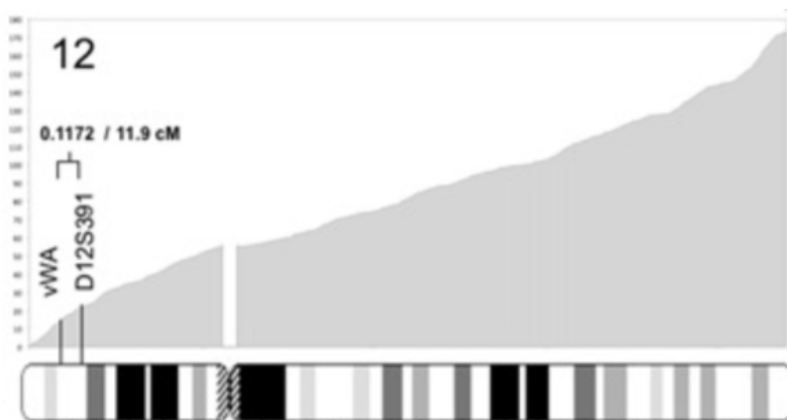
BioProjects



Gauging linkage between close STRs:

HapMap

cM estimates of syntenic loci



Identifying less well established STRs as unique

NCBI Probe

D5S2500 (in a CE kit) vs D5S2800 (in an MPS kit)

Chrom.	Core STRs	dbSNP rs-number identifier for STR	Chrom.	Supplementary STRs	Kit	dbSNP rs-number identifier for STR
C1	D1S1656	rs113633160	C1	F13B	Promega CS7	rs10643350
C2	TPOX	rs113475620	C2	D2S1360	Qiagen HD-plex	rs113680434
C2	D2S1338	rs112111672				
C2	D2S441	rs10203882				
C3	D3S1358	rs111694514	C3	D3S1744	Qiagen HD-plex	rs113865588
C4	FGA	rs67296980	C4	D4S2366	Qiagen HD-plex	rs113820309
C5	D5S818	rs112497490	C5	D5S2500	Qiagen HD-plex	rs111362704
C5	CSF1PO	rs113729910				
C6	SE33	rs71021371	C6	D6S474	Qiagen HD-plex	rs113991233
C7	D7S820	rs112714641	C7	D7S1517	Qiagen HD-plex	rs112397288
C8	D8S1179	rs67563232	C8	D8S1132	Qiagen HD-plex	rs71307053
			C9	Penta C	Promega CS7	rs72398274
C10	D10S1248	rs113518246	C10	D10S2325	Qiagen HD-plex	no SNPs found
C11	TH01	rs71029110				
C12	D12S391	rs113002069				
C12	vWA	rs10579907				
C13	D13S317	rs111980288				
C15	Penta E	rs8036258	C15	FES-FPS	Promega CS7	rs6229
C16	D16S539	rs112689398				
C18	D18S51	rs10560567				
C19	D19S433	rs113951851				
C21	D21S11	rs113145752	C21	D21S2055	Qiagen HD-plex	rs113225349
C21	Penta D	rs7279663				
C22	D22S1045	rs112790319				

Databases of genomic details pertinent to a variant

Databases of the population variation observed in a variant

Maps of clustered variants:

STRs and their accompanying flanking SNPs

Multiple STRs in close proximity

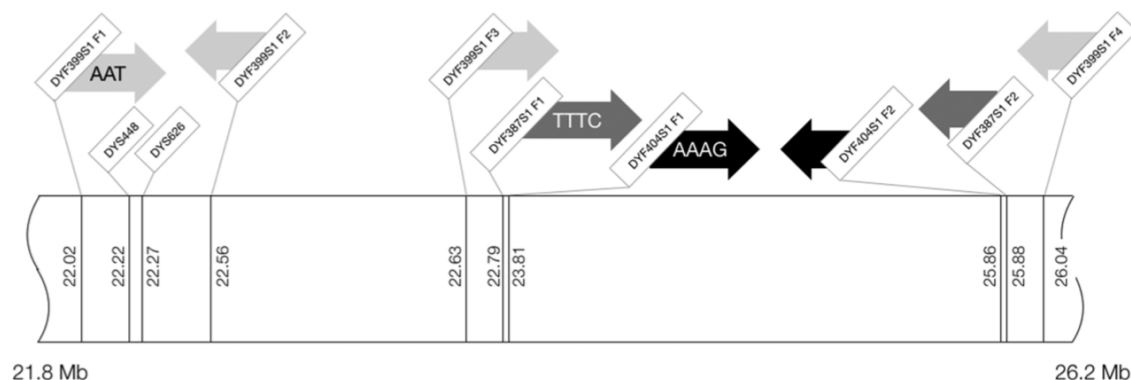
Indels that might influence CE size estimates

Arrangement of multiple fragment STRs

Ensembl

NCBI Genbank

*STR-specific
BioProjects*



Gauging linkage between close STRs:

HapMap

cM estimates of syntenic loci

Identifying less well established STRs as unique

NCBI Probe

D5S2500 (in a CE kit) vs D5S2800 (in an MPS kit)

Agreement to use the RefSeq reference strand:

5' to 3' single sequence with fixed coordinates for both current builds

Consensus amongst MPS users for loci and level of detail

Consensus amongst STRAND members for each STR's START-END nucleotides and repeat region structures

*STR
Sequence
Guide*

*Guide
Wiki
Database*

*'Red-point' SNPs
Mobility-shift SNPs
Rare 4-nt Indels*

✓ Flanking variants even at low frequency are easy to compile

STRSeq compiling full sequences including flanks

dbSNP is accelerating the uptake of Indels and assignment of rs-numbers

1000 genomes is now supplemented by gnomAD / TOPMed - tens of thousands of samples

Databases of genomic details pertinent to a variant

Databases of the population variation observed in a variant

dbSNP is undergoing a transition in presentation and scope to accommodate many more rare SNPs identified by >10,000 sample projects

NIH U.S. National Library of Medicine
National Center for Biotechnology Information

dbSNP Short Genetic Variations

Search for rs Search
Example: rs268

Reference SNP (rs) Report Download Facebook Twitter YouTube

[Switch to classic site](#)

rs12913832 **Current Build 152**
Released October 2, 2018

Organism	<i>Homo sapiens</i>	Clinical Significance	Reported in ClinVar
Position	chr15:28120472 (GRCh38.p12) ?	Gene : Consequence	HERC2 : Intron Variant
Alleles	A>G	Publications	92 citations
Variation Type	SNV Single Nucleotide Variation	Genomic View	See rs on genome
Allele Frequency	G=0.45329 (56919/125568, TOPMED) A=0.4419 (13667/30926, GnomAD) G=0.177 (888/5008, 1000G) (+ 3 more)		

gnomAD
genome aggregation database

No results found

gnomAD Genome Aggregation Database

NHLBI Trans-Omics for Precision Medicine

Centers Projects/Studies Working Groups Data Publications EAP ELSI Workshops

About TOPMed

- Contents
- Overview
 - Study Characteristics
 - Study Designs
 - Participant Diversity
 - Whole Genome Sequencing
 - Resources for the Scientific Community

TOPMed Trans-Omics for Precision Medicine

rs1051822965 SNP

Most severe consequence

[Intron variant](#) | [See all predicted consequences](#)

Alleles

C/G | Ancestral: C | Highest population MAF: < 0.01

Change tolerance

CADD: G:3.952 | GERP: -3.35

Location

Chromosome 11:2171074 (forward strand) | VCF: 11 2171074 rs1051822965 C G

Evidence status



HGVS names

This variant has 10 HGVS names - [Show](#)

Synonyms

ClinGen Allele Registry [CA216229061](#) (G)

Original source

Variants (including SNPs and indels) imported from dbSNP (release 151) | [View in dbSNP](#)

About this variant

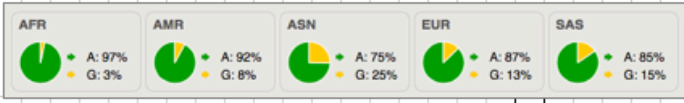
This variant overlaps [6 transcripts](#).

gnomADg:ALL	C: 0.999899112 (29733)	G: 0.000100888 (3)
gnomADg:af	C: 0.999635569 (8229)	G: 0.000364431 (3)
gnomADg:amr	C: 1.000 (832)	G: 0.000
gnomADg:asj	C: 1.000 (264)	G: 0.000
gnomADg:eas	C: 1.000 (1412)	G: 0.000
gnomADg:fin	C: 1.000 (3432)	G: 0.000
gnomADg:nfe	C: 1.000 (14546)	G: 0.000
gnomADg:oth	C: 1.000 (1018)	G: 0.000

[AATG]n ATG [AATG]n	TH01
above does not describe the reference sequence	Reference sequence
	Flanking SNP IUPAC codes
Chr: 11	GRCh38 coordinates
	GRCh37 coordinates
	Distance from repeat region

1	2	3	4	5	6	7
A A T G	A A T G	A A T G	A A T G	A A T G	A A T G	A A T G
A A T G	A A T G	A A T G	A A T G	A A T G	A A T G	A A T G
2171088	2171089	2171090	2171091	2171092	2171093	2171094
2192318	2192319	2192320	2192321	2192322	2192323	2192324
1	2	3	4	5	6	7
1	2	3	4	5	6	7

TOPMed and 'red-point' (very rare) SNPs



[TAGA]n [CAGA]n TAGA	VWA
(reverse sequence listed in STRbase)	Reference sequence
	Flanking SNP IUPAC codes
Chr: 12	GRCh38 coordinates
	GRCh37 coordinates
	Distance from repeat region

1	2	3	4	5	6	7	8	9
T A G A	T A G A	T A G A	T A G A	T A G A	T A G A	T A G A	T A G A	T A G A
T A G A	T A G A	T A G A	T A G A	T A G A	T A G A	T A G A	T A G A	T A G A
5983936	5983937	5983938	5983939	5983940	5983941	5983942	5983943	5983944
6093102	6093103	6093104	6093105	6093106	6093107	6093108	6093109	6093110
1	2	3	4	5	6	7	8	9
1	2	3	4	5	6	7	8	9

rs771794429 SNP

gnomAD Genome Aggregation Database

Most severe consequence

[Intron variant](#) | [See all predicted consequences](#)

Alleles

G/A | Ancestral: G | MAF: 0.01 (A) | Highest population MAF: < 0.01

Change tolerance

CADD: A:1.395 | GERP: -2.86

Location

Chromosome 12:5983962 (forward strand) | VCF: 12 5983962 rs771794429 G A

Evidence status



HGVS names

This variant has 4 HGVS names - [Show](#)

Synonyms

ClinGen Allele Registry [CA232291691](#) (A)

Original source

Variants (including SNPs and indels) imported from dbSNP (release 151) | [View in dbSNP](#)

About this variant

gnomADg:af	G: 0.999 (8041)	A: 0.001 (5)
------------	-----------------	--------------

rs199970098 SNP

TOPMed Trans-Omics for Precision Medicine Program

Most severe consequence

[Intron variant](#) | [See all predicted consequences](#)

Alleles

G/A | Ancestral: A | Highest population MAF: 0.01

Change tolerance

CADD: A:0.041 | GERP: -3.74

Location

Chromosome 12:5983974 (forward strand) | VCF: 12 5983974 rs199970098 G A

Evidence status



HGVS names

This variant has 4 HGVS names - [Show](#)

Synonyms

ClinGen Allele Registry [CA232291705](#) (A)

Original source

Variants (including SNPs and indels) imported from dbSNP (release 151) | [View in dbSNP](#)

About this variant

TOPMed	G: 0.985	A: 0.015
--------	----------	----------

STR Flanking Region Variation

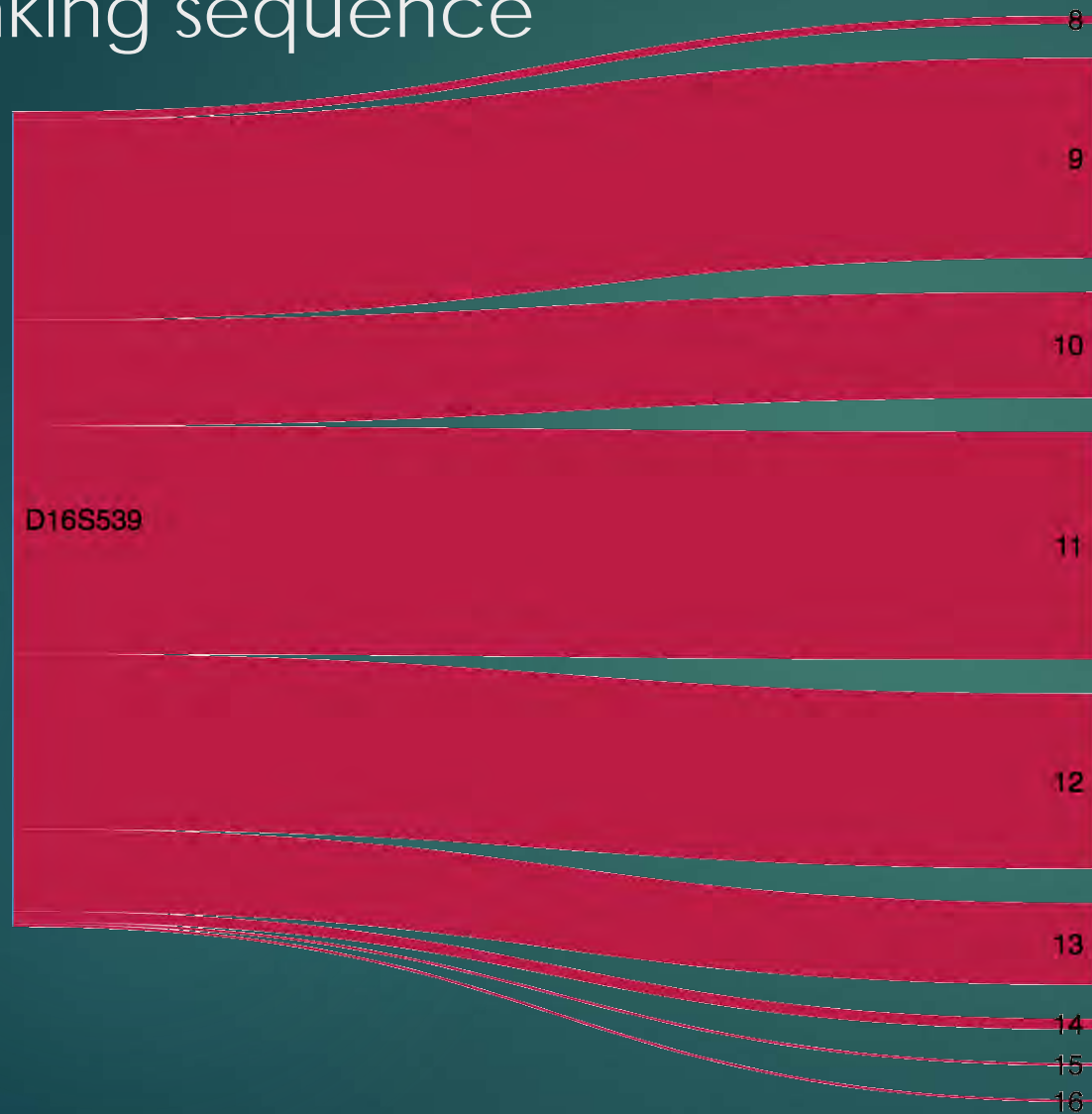
David Ballard – King's College London



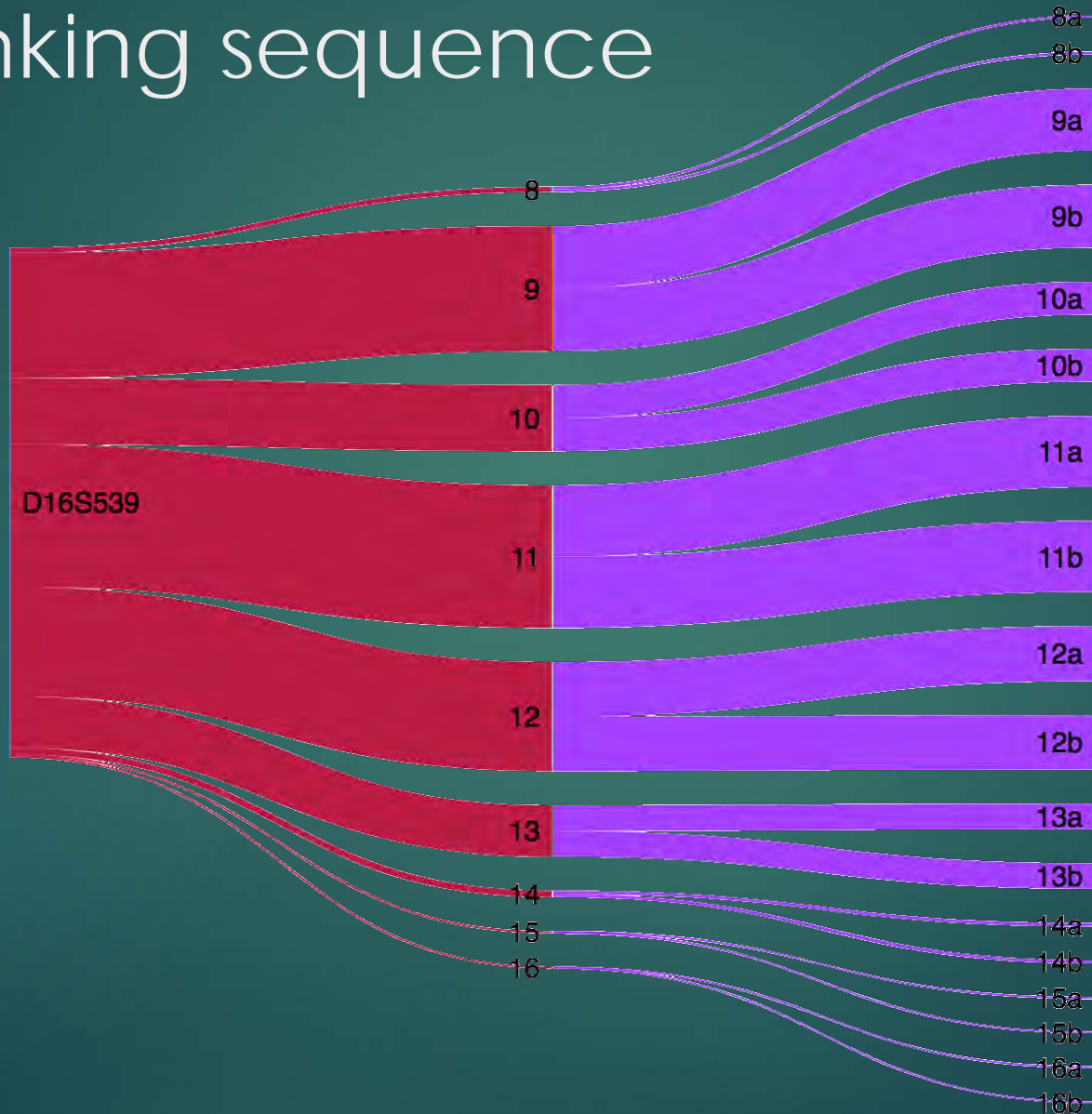
D16S539 – Known Flanking SNP Variation

																AFR					AMR					ASN					EUR					SAS														
																9					10					11																								
																A T A G A T A					G A T A					G A T A																								
																A T A G A T A					G A T A					G A T A																								
86386337	86386338	86386339	86386340	86386341	86386342	86386343	86386344	86386345	86386346	86386347	86386348	86386349	86386350	86386351	86386352	86386353	86386354	86386355	86386356	86386357	86386358	86386359	86386360	86386361	86386362	86386363	86386364	86386365	86386366	86386367	86386368	86386369	86386370	86386371	86386372	86386373	86386374	86386375	86386376	86386377	86386378	86386379	86386380	86386381	86386382	86386383				
																T	C	A	T	T	G	A	A	A	G	A	C	A	A	A	A	C	A	G	A	G	A	T	G	G	A	T	G	A	T	A	G			
																M																																		
																rs11642858																																		

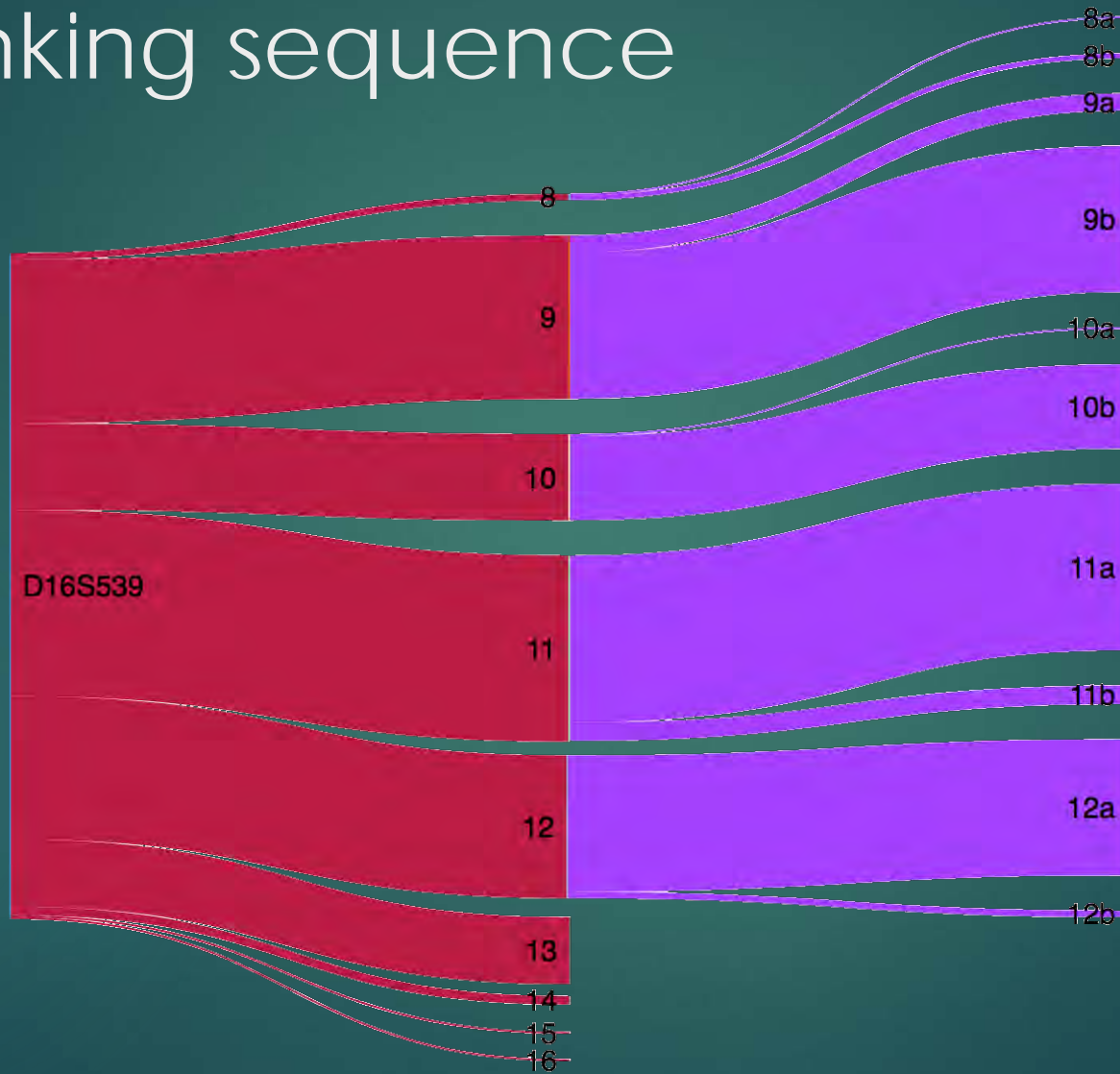
D16S539 – Expected frequency of alleles using flanking sequence



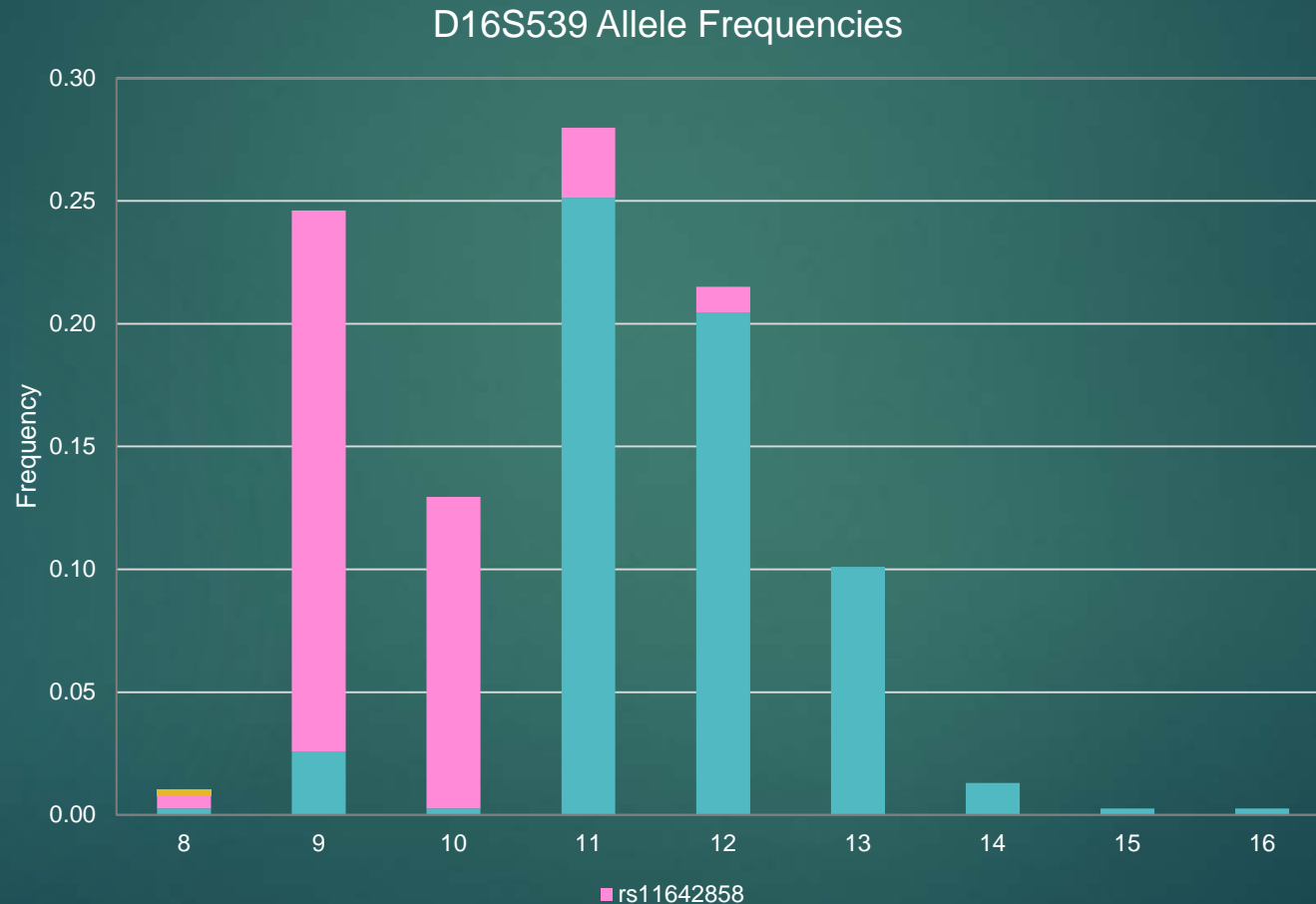
D16S539 – Expected frequency of alleles using flanking sequence



D16S539 – Actual frequency of alleles using flanking sequence



D16S539 – Why isn't the flanking region helping?



STR nomenclature



To consider?

- For discrimination purposes, extended flanking region variation is of limited use due to linkage with specific STR repeat sequence alleles
- A nomenclature system just describing the repeat region (or a defined region around the repeat) would be simple and capture almost all of the useful variation
- A collection of all common STR sequence allele variation is already available



UNIVERSITY OF
LEICESTER

STR Sequence Nomenclature – The view from Leicester

Tunde Huszar

University of Leicester

Alec Jeffreys Forensic Genomics Unit

Department of Genetics and Genome Biology

Leicester, UK

th201@leicester.ac.uk



STR Sequence Nomenclature

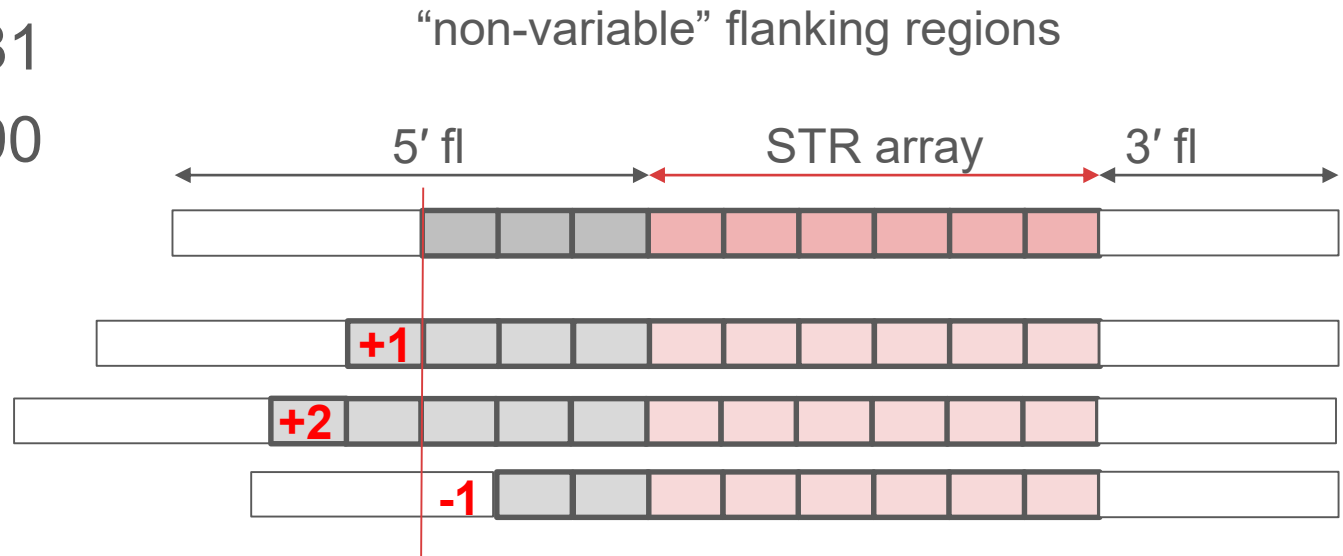
- ✓ Lessons from a phylogenetic framework for Y-STRs
- ✓ Existing databases – using STRSeq
- ✓ Local database – LeiceSTRSeq

Phylogenetic framework



✓ variable array limits:

- DYS385a,b
- DYS481
- DYS390



Huszar et al. FSI: Genetics (2018), 35, 97-106

A phylogenetic framework facilitates Y-STR variant discovery and classification via massively parallel sequencing

Phylogenetic framework - DYS385a,b



conventionally : (STRBase, strbase.nist.gov)

DYS385a,b: $[GAAA]_n$

2016 ISFG recommendation: (Parson et al. 2016 FSI Gen)

DYS385a : $[TTTC]_n$

DYS385b: $[GAAA]_n$

no distinction between a/b in kits – suggestion:

DYS385a,b: $[GAAA]_n$

Huszar et al. FSI: Genetics (2018), 35, 97-106

A phylogenetic framework facilitates Y-STR variant discovery and classification via massively parallel sequencing

Phylogenetic framework - DYS385a,b



✓ Variable flanking regions:

recognise repeat structure, rather than calling several SNPs

Allele	Observed #	General structure of alleles including variable flanking sequences	CE allele name designation
canonical	193	AAGG[6] GAAA[n]	n
variant	4	AAGG[5] GAAA[n]	n - 1
variant	2	AAGG[7] GAAA[n]	n + 1
variant	2	AAGG[8] GAAA[n]	n + 2
variant	*	AAGG[9] GAAA[n]	n + 3

* observed in Novroski et al. 2016, FSI Gen

Huszar et al. FSI: Genetics (2018), 35, 97-106

A phylogenetic framework facilitates Y-STR variant discovery and classification via massively parallel sequencing

Phylogenetic framework



- ✓ high sequence variability – flexible software

FDSTools (Hoogenboom et al. 2017, FSI Gen)

non-standard populations, new variants, non-human STRs

- ✓ Multiple software/analysis – against bioinformatic nulls

UAS / STRaitRazor / FDSTools /

/ commercial software / in-house scripts

Huszar et al. FSI: Genetics (2018), 35, 97-106

A phylogenetic framework facilitates Y-STR variant discovery and classification via massively parallel sequencing

Existing databases

✓ STRSeq: (Gettings et al. 2017, FSI Gen)

GenBank records at NCBI - (unique Acc#)

BioProjects (**Auto**, Auto+, Y-, X-STRs)

“Project Data:

No public data is linked to this project. Any recently released data that cites this project will be linked to it within a few days.”

✓ STRidER: (Bodner et al. 2016, FSI Gen)

STR sequence guide v4 (Phillips et al. 2018, FSI Gen)

user interface, pathway for submission, QC

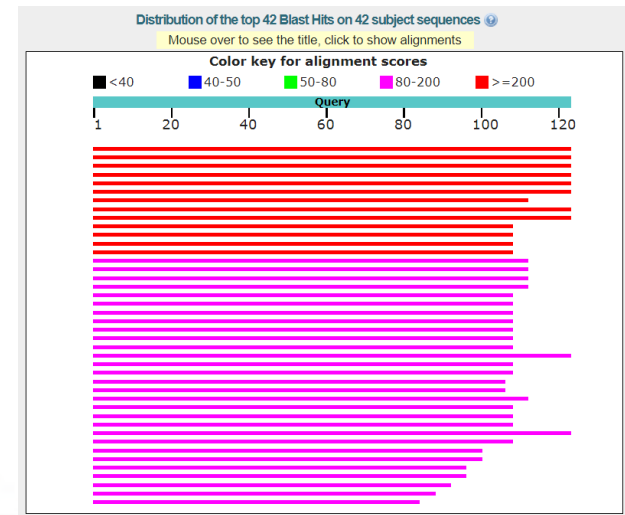
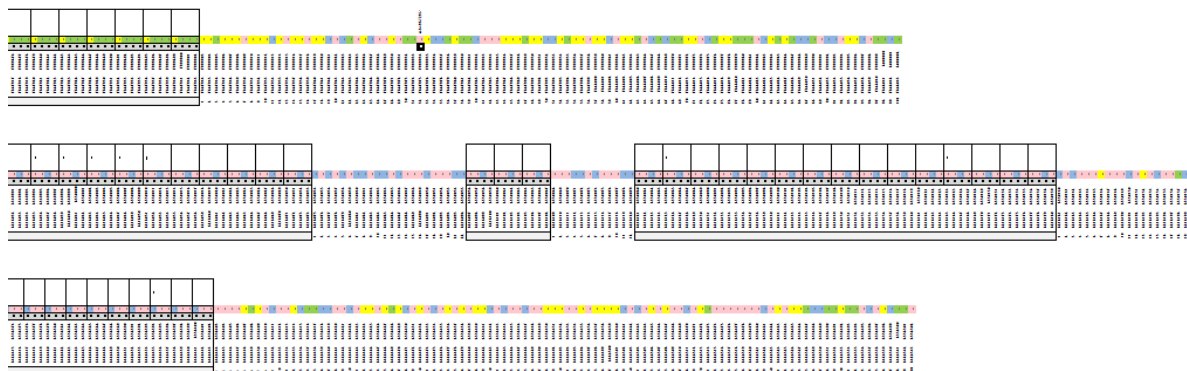
Existing databases

✓ STRSeq: (**Auto**, Auto+, Y-, X-STRs)

search by BLASTn
detailed, comparable,
sequence identifiers

✓ STRidER:

compendium for variation



Sequences producing significant alignments:

Select: All None Selected: 0

Alignments Download Compare Clustal DistanceMatrix All results

Description	Max Score	Total Score	Query Cover	E value	Per. Ident	Accession
<input type="checkbox"/> Homo sapiens microsatellite D1S1856.13 (TCTA)13 sequence	228	228	100%	4e-56	100.00%	MH174842.1
<input type="checkbox"/> Homo sapiens microsatellite D1S1856.13 TCTA GCTA (TCTA)11 sequence	222	222	100%	2e-54	99.19%	MH174841.1
<input type="checkbox"/> Homo sapiens microsatellite D1S1856.13 CCTA (TCTA)12 sequence	222	222	100%	2e-54	99.19%	MH174840.1
<input type="checkbox"/> Homo sapiens microsatellite D1S1856.14 (TCTA)14 sequence	209	209	100%	1e-50	96.85%	MH174845.1
<input type="checkbox"/> Homo sapiens microsatellite D1S1856.14 CCTA (TCTA)13 sequence	209	209	100%	1e-50	96.85%	MH174844.1
<input type="checkbox"/> Homo sapiens microsatellite D1S1856.14 CCTA (TCTA)13 rs1019813099 sequence	204	204	100%	6e-49	96.06%	MH174843.1

Existing databases - Issues

✓ STRSeq: not yet applicable for Auto+, Y-, X-STRs

instead: literature search
 redundant task
 no unique sequence ID
 LeiceSTRSeq – local DB

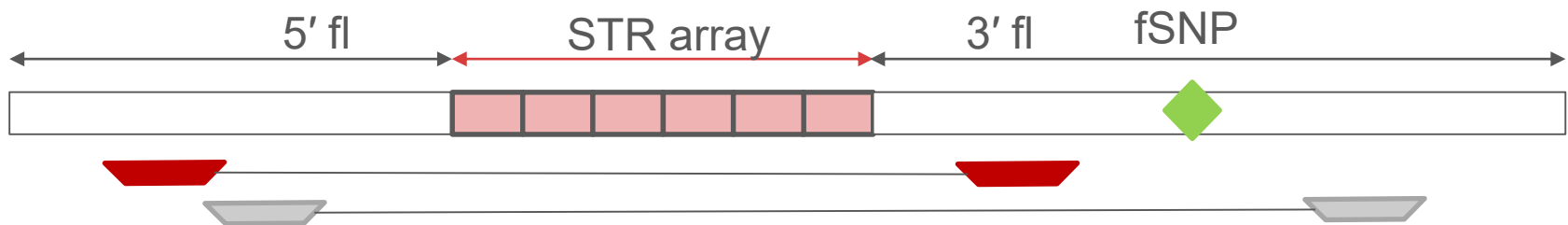
Y-STR	sequence	mutant count	variant type (SNP, indel, STR)	artibeans.in.gov	Zhao et al. 2015 reports all alleles	Wron et al. 2016 reports all alleles	Wendt et al. 2016 reports all alleles	Jure et al. 2017 reports all alleles	Novroski et al. 2017 reports all alleles	Waghaizer et al. 2015 reports novels	Churhill et al. 2016 reports novels only	Wendt et al. 2017 reports novels only	Forster et al. 1998 specific	Reid et al. 2002 specific	D'Amano et al. 2010 specific	Lee et al. 2016 specific	novel?	
	CE12-17_TCTA[9-14]CCTA[0-1]TCTA[3]																	
DYS19	CE12_TCTA[13]	1	ISNP	-	-	-	-	-	-				NA	NA	NA	NA	novel	
	CE12_TCTA[9]CCTA[1]TCTA[3]	2		+	-	-	-	-	-				NA	NA	NA	NA	-	
	CE13_TCTA[10]CCTA[1]TCTA[3]	11		+	+	-	-	-	-				NA	NA	NA	NA	-	
	CE14_TCTA[11]CCTA[1]TCTA[3]	32		+	+	+	-	-	-				NA	NA	NA	NA	-	
	CE15_TCTA[12]CCTA[1]TCTA[3]	36		+	+	+	+	-	-				NA	NA	NA	NA	-	
	CE16_TCTA[13]CCTA[1]TCTA[3]	15		+	+	+	+	+	-				NA	NA	NA	NA	-	
	CE17_TCTA[14]CCTA[1]TCTA[3]	4		+	+	+	+	+	+				NA	NA	NA	NA	-	
	CE9-22_AAGG[5-8]GAAA[9-22]																	
DYS385a,b	CE9_AAGG[5]GAAA[10]	1		-	NA	-	-	-	-				NA	NA	NA	NA	novel	
	CE9_AAGG[6]GAAA[9]	3		+	NA	-	-	-	-				NA	NA	NA	NA	-	
	CE10_AAGG[6]GAAA[10]	7		+	NA	-	-	-	-				NA	NA	NA	NA	-	
	CE11_AAGG[6]GAAA[11]	18		+	NA	+	-	-	-				NA	NA	NA	NA	-	
	CE12_AAGG[6]GAAA[12]	17		+	NA	+	+	-	-				NA	NA	NA	NA	-	
	CE13_AAGG[5]GAAA[14]	1		-	NA	-	-	-	-				NA	NA	NA	NA	novel	
	CE13_AAGG[6]GAAA[13]	28		+	NA	+	+	-	-				NA	NA	NA	NA	-	
	CE14_AAGG[6]GAAA[14]	38		+	NA	+	+	+	-				NA	NA	NA	NA	-	
	CE15_AAGG[5]GAAA[16]	1		-	NA	-	-	-	-				NA	NA	NA	NA	novel	
	CE15_AAGG[6]GAAA[15]	20		+	NA	+	+	+	-				NA	NA	NA	NA	-	
	CE15_AAGG[8]GAAA[13]	1		-	NA	-	-	-	-				NA	NA	NA	NA	novel	
	CE16_AAGG[6]GAAA[16]	17		+	NA	+	+	+	-			+	NA	NA	NA	NA	-	
	CE16_AAGG[6]GAAA[2]TAAA[1]GAAA[13]	2	ISNP	-	NA	-	+	+	-			+	NA	NA	NA	NA	novel	

✓ STRidER:

great QC for CE-based submission,
 but no MPS submission pathway yet

Existing databases – Potential issues

✓ STRSeq: GenBank Acc# - unique sequence ID



Kit1 STR only – unique ID#1

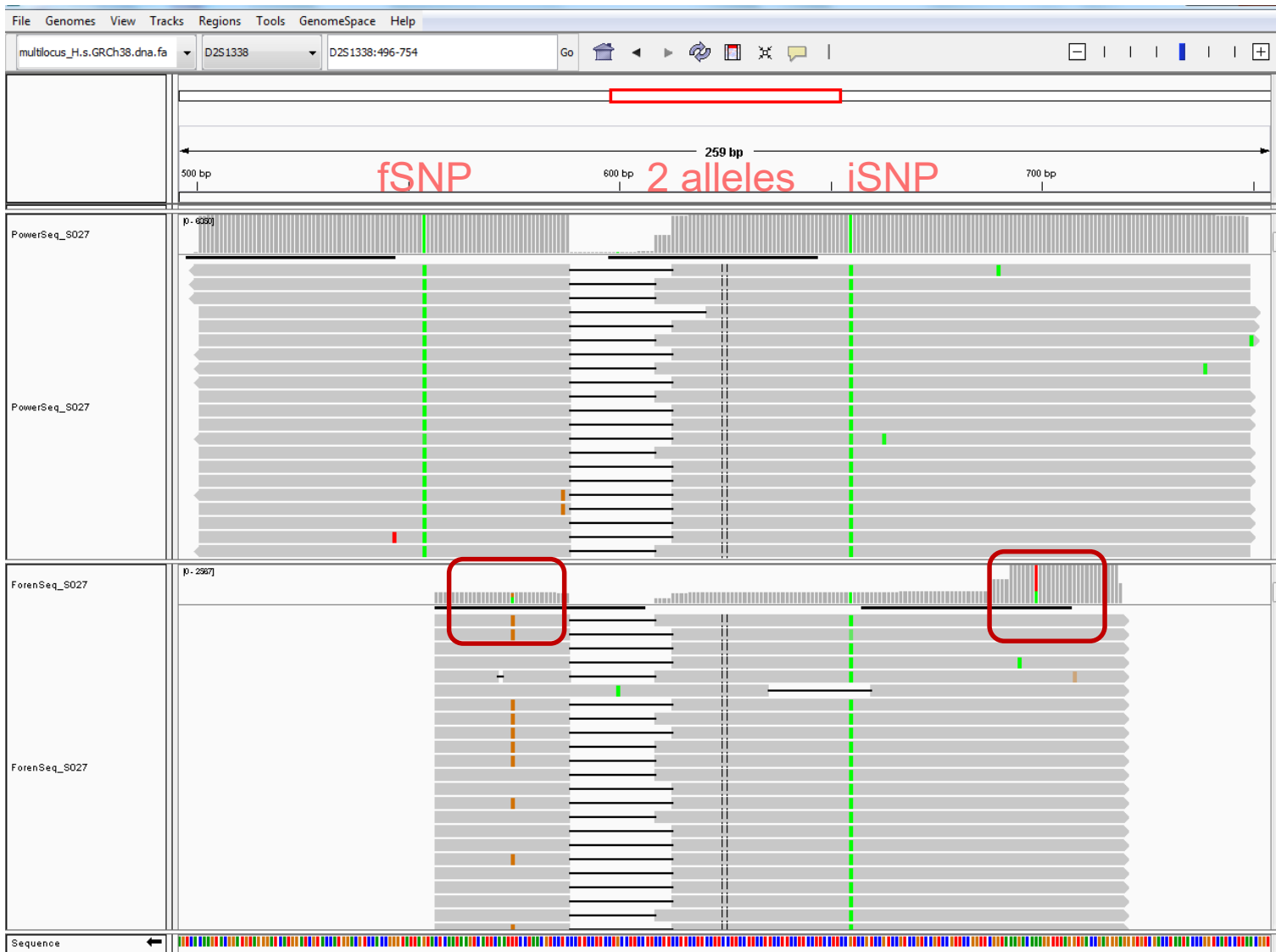
Kit2 STR + **flanking SNP** – unique ID#2

D16S539 – same sample on GRCh38, chr16

PowerSeq:CE8_GATA[8] – in STRSeq as MH167241.1

ForenSeq:CE8_GATA[8]_rs11642858 – in STRSeq as MK570017.1

Existing databases – Potential issues

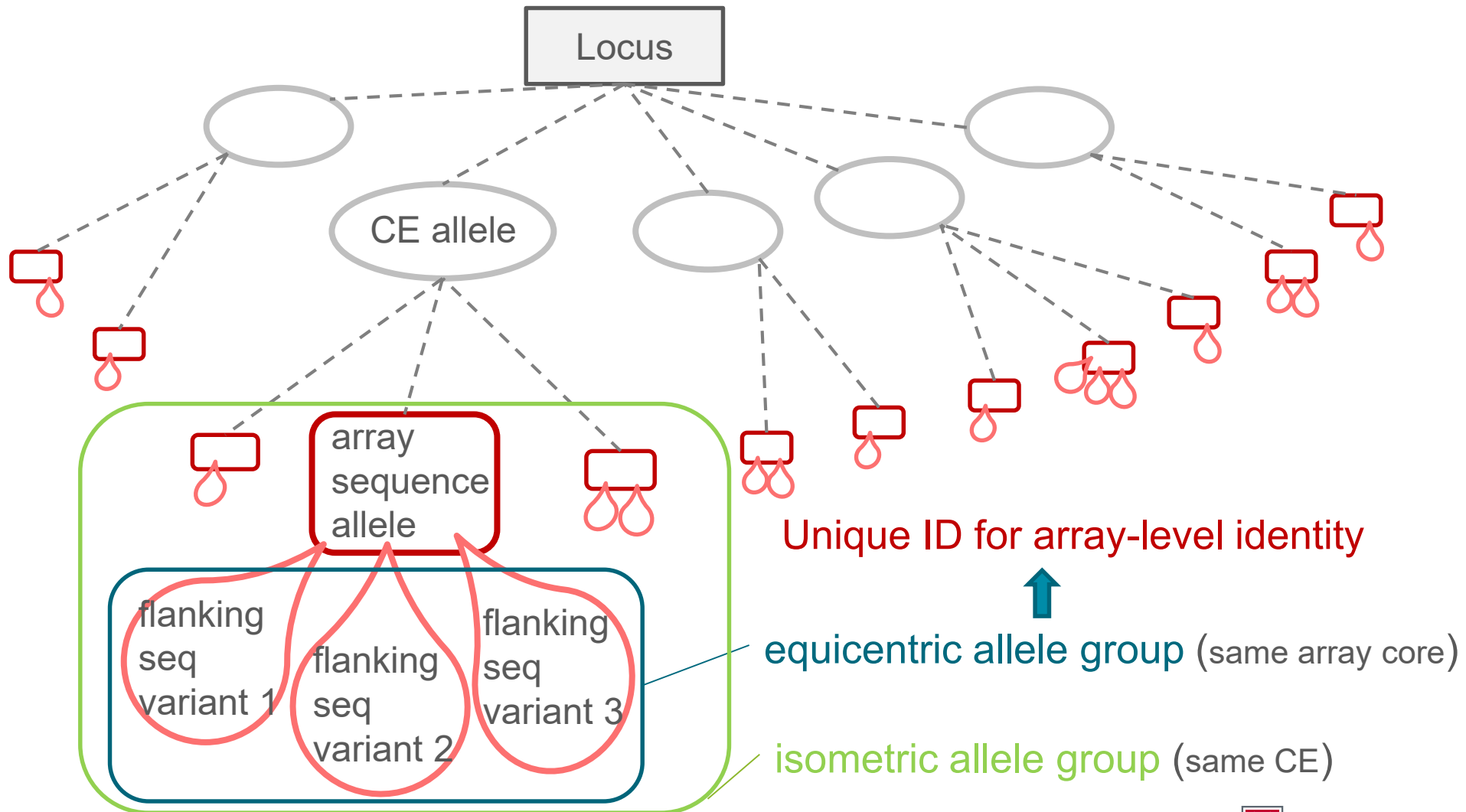


Primer interference – trimming!
(Huszar et al. 2019, FSI Gen.)

Existing databases – Potential issues

- ✓ STRSeq: GenBank Acc# - unique sequence ID
 - one DNA type – several unique sequence IDs
 - sequence ID groups – based on STR array seq identity
 - flanking region difference – exclusion / no match
 - real OR kit / software / reported region difference
 - unique ID currently at flanking region variants level
 - not ideal for automated matches

Unique ID for comparable reporting



Local database - LeiceSTRSeq

- ✓ To help MPS-based projects with redundant tasks:
 - literature search for allele variants
 - allele IDs from STRSeq – by array-level identity groups
 - extra: kit, software, amplified/reported coordinates
 - STRSeq ID/publications for reference, annotated string
- ✓ Difficulties:
 - lack of personnel: constant screening, input, cleaning data and development
 - Excel-based user's copy, Access-based background DB

Summary

- ✓ more flexible/sensible array definition – repeat vs SNPs
- ✓ unique ID for STR array-level identity – comparable match
- ✓ flanking variants - with clear reporting coordinates
(reference genome, kit type, software)
- ✓ Current state: doable, but not user-friendly
 - central QC
 - curation of submitted new alleles
 - build and maintain database
 - cross-platform searchable interface

Acknowledgements

Mark Jobling

Jon Wetton



High Performance Computing cluster
NUCLEUS Genomics Lab

Other Leicester MPS projects:





CE-MPS Discordances

in a study of 31 autosomal STR loci
from 498 Spanish individuals

Pedro A. Barrio¹, Pablo Martin¹, Antonio Alonso¹, The DNASEQEX Consortium

¹ Servicio de Biología del Instituto Nacional de Toxicología y Ciencias Forenses (INTCF),
Departamento de Madrid
pedro.barrio@justicia.es



Comparison of two MPS platforms:
Ion S5 (*Thermo Fisher Scientific*) MiSeq FGx™ (*Illumina*)

STRs Standardization typing by MPS

International Exchange of MPS data

Population Studies



Madrid



Innsbruck



Berlin



498 samples



Spanish ancestries representing all the 17 Autonomous Communities of Spain (i.e. "regions")



Precision ID GlobalFiler® NGS STR Panel v2 on Ion S5 System



Concordance Study CE/MPS



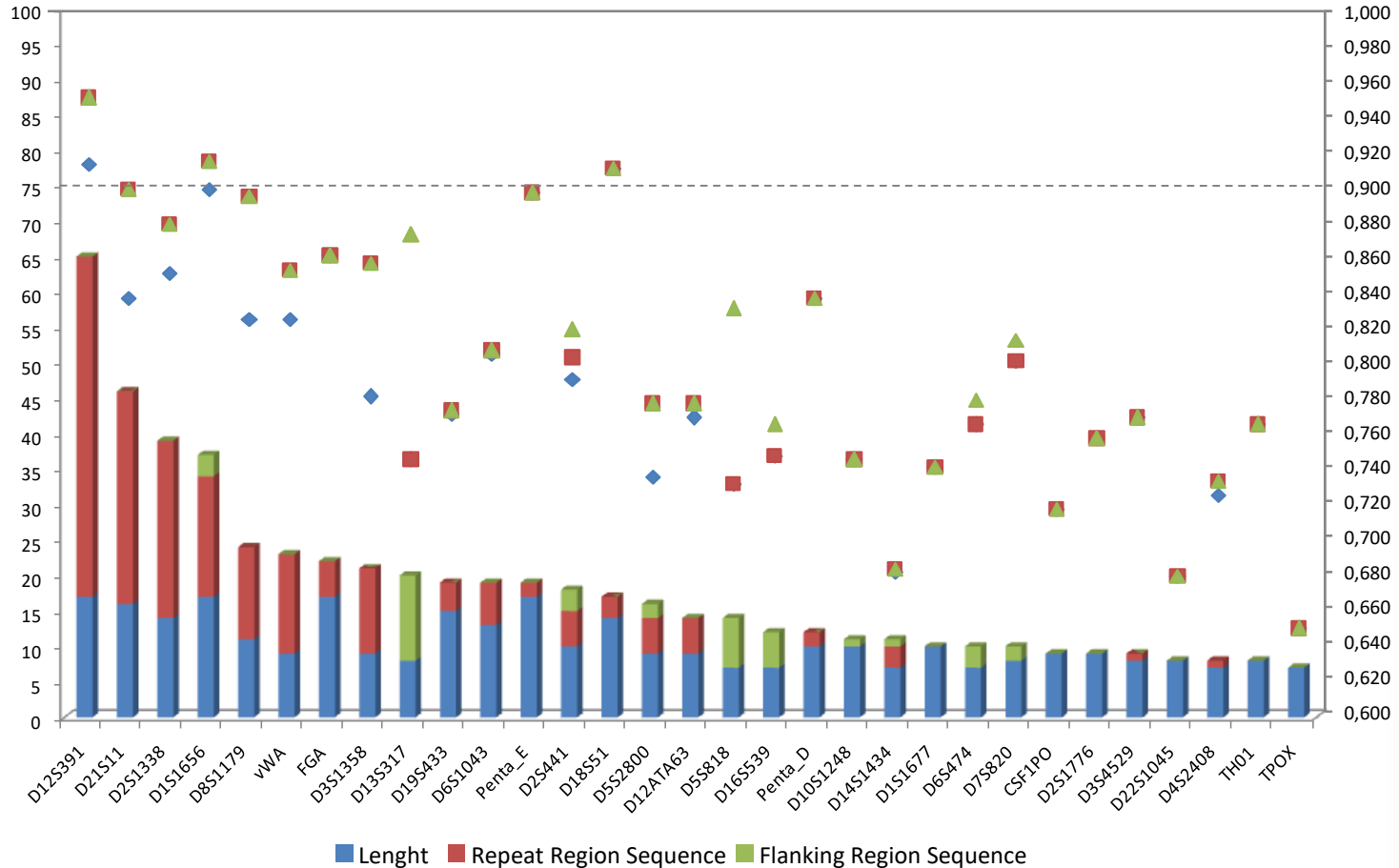
221 samples

PowerPlex Fusion 6C System (Promega, Madison, WI, USA)



STR allelic gains by sequence:

Number of alleles compared to heterozygosity observed for the 31 auSTR loci



Putative mismatches:

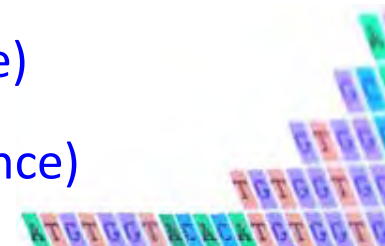
3 loci: Penta D, D2S441 and D19S433

Penta D		D2S441		D19S433	
Reference sequence	Flanking SNP IUPAC codes	Reference sequence	Flanking SNP IUPAC codes	Reference sequence	Flanking SNP IUPAC codes
GRCh38 coordinates	GRCh37 coordinates	GRCh38 coordinates	GRCh37 coordinates	GRCh38 coordinates	GRCh37 coordinates
Distance from repeat region		Distance from repeat region		Distance from repeat region	

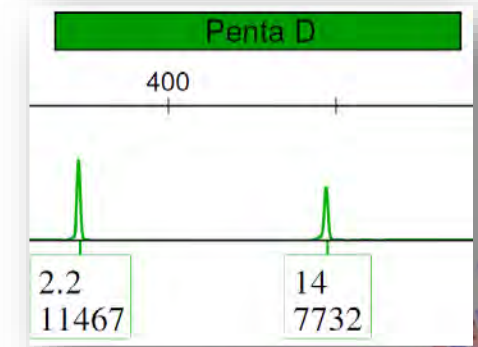
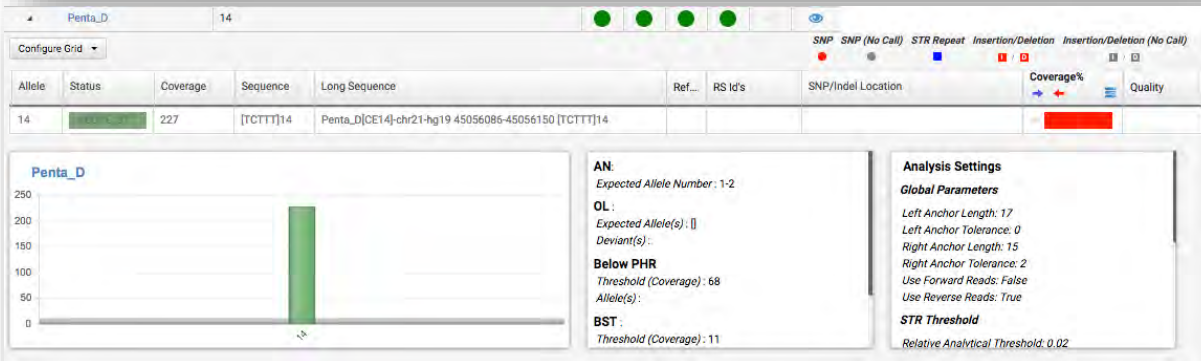
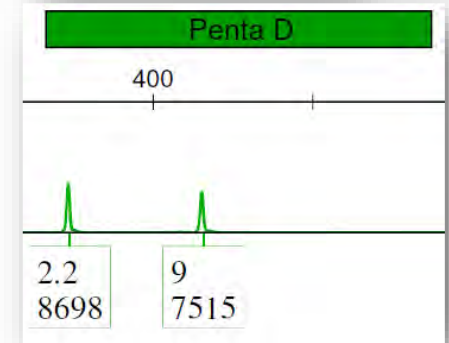
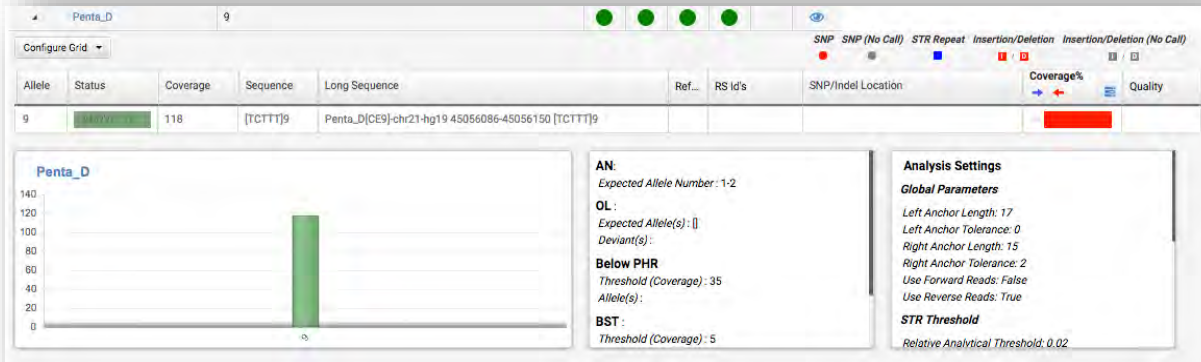
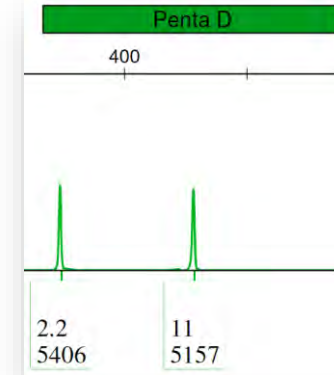
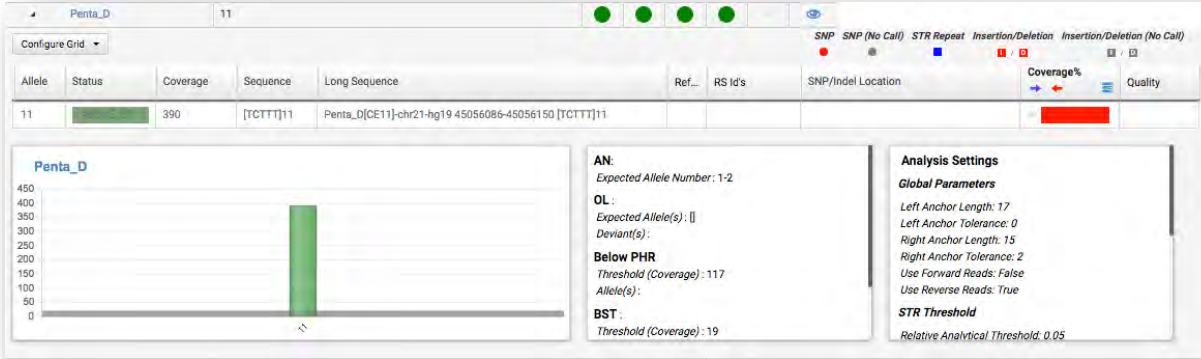
5 samples: 5 samples out of 221 (97.73 % sample concordance)

1 locus/sample: 5 loci out of 5083 (99.90 % locus concordance)

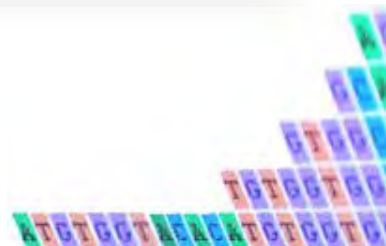
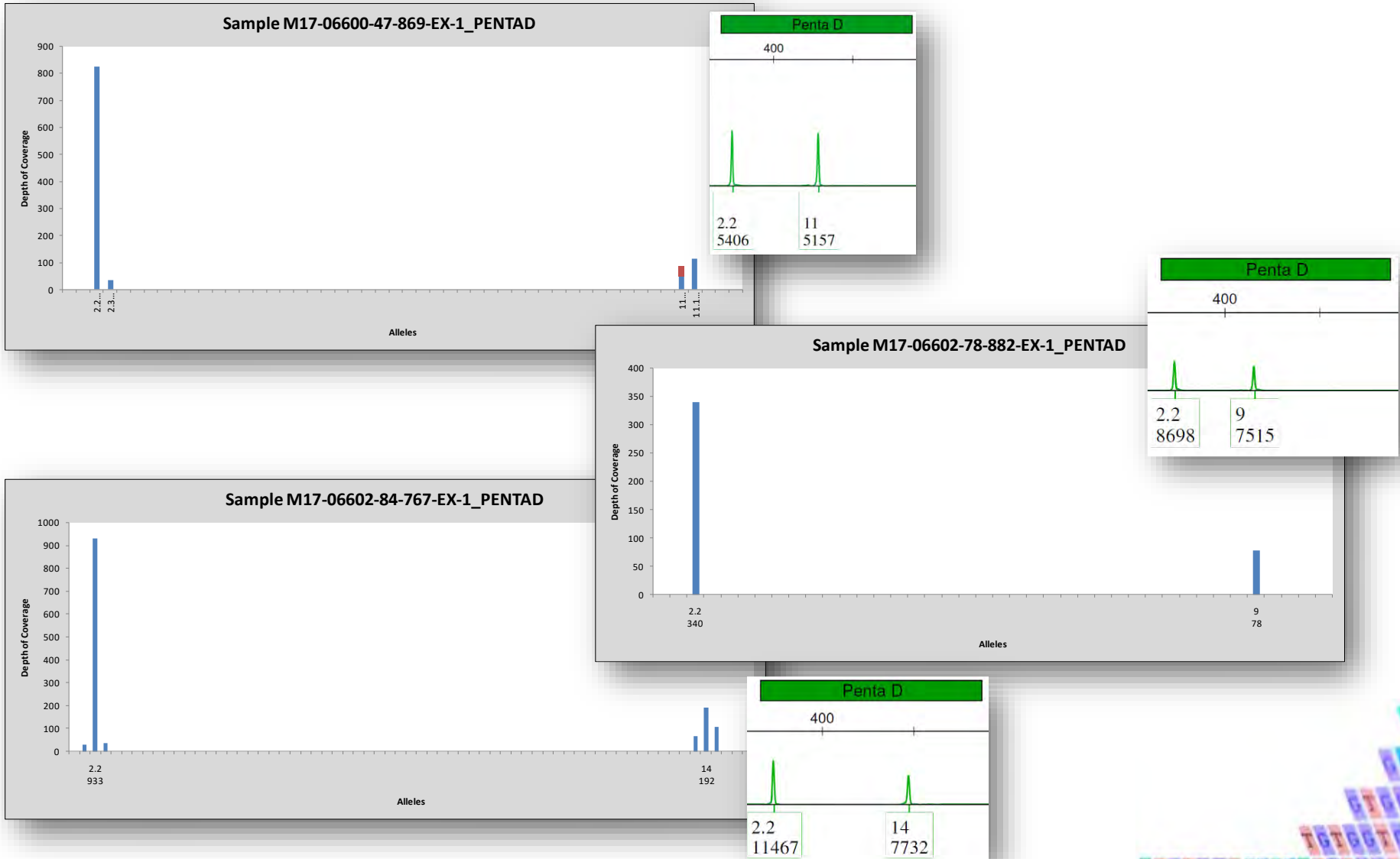
1 allele/locus: 5 alleles out of 10166 (99.95 % allele concordance)



3 samples: M17-06600-47-869-EX-1, M17-06602-78-882-EX-1 and M17-06602-84-767-EX-1



Additional analysis: *STRait Razor v3* (Woerner et al., 2017)



1 sample: M17-06603-60-608-EX-1

Allel...	Status	Coverag...	Sequence	Long Sequence	Re...	RS Id's	SNP/Indel Location	Coverage%	Quality
12	STUTTER	36	[TCTA]9 TTTA[TCTA]2	D2S441[CE12]-chr2-hg19 68239079-68239126 [TCTA]9 TTTA[TCTA]2					
13	90195.811	594	[TCTA]10 TTTA[TCTA]2	D2S441[CE13]-chr2-hg19 68239079-68239126 [TCTA]10 TTTA[TCTA]2					
14	90195.811	492	[TCTA]11 TTTA[TCTA]2	D2S441[CE14]-chr2-hg19 68239079-68239126 [TCTA]11 TTTA[TCTA]2					



AN:
Expected Allele Number: 1-2

OL:
Expected Allele(s): [8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 9.1, 11.3, 12.3, 13.3, 14.3]
Deviant(s):

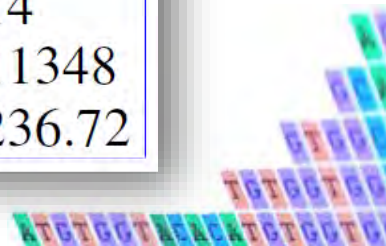
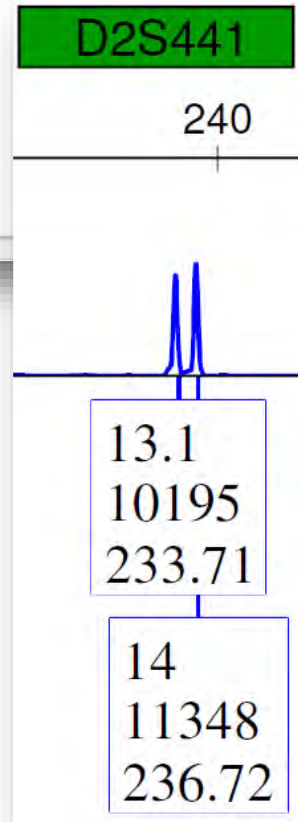
Below PHR
Threshold (Coverage): 178
Allele(s):

BST:

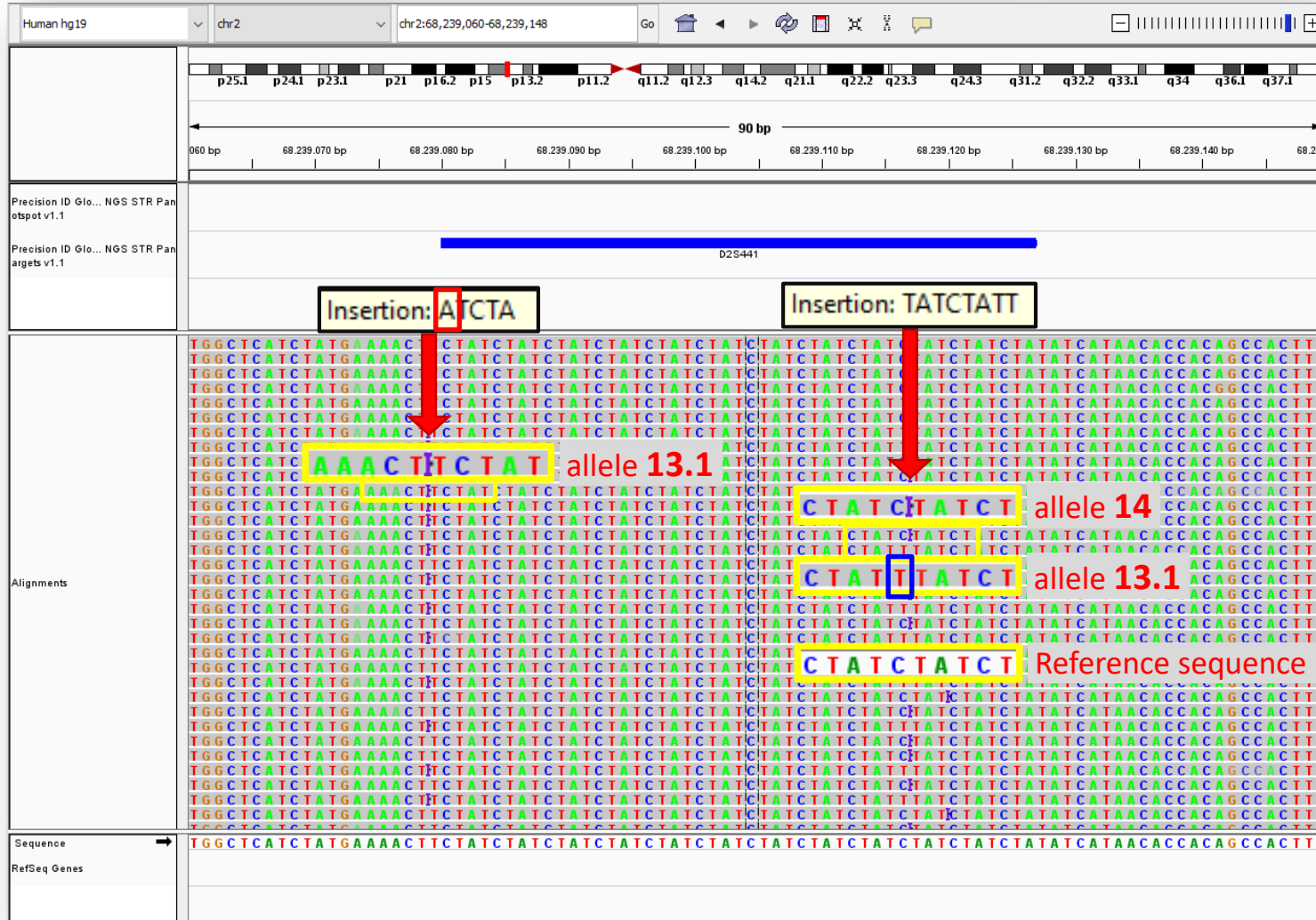
Analysis Settings

Global Parameters
Left Anchor Length: 15
Left Anchor Tolerance: 2
Right Anchor Length: 15
Right Anchor Tolerance: 2
Use Forward Reads: True
Use Reverse Reads: True

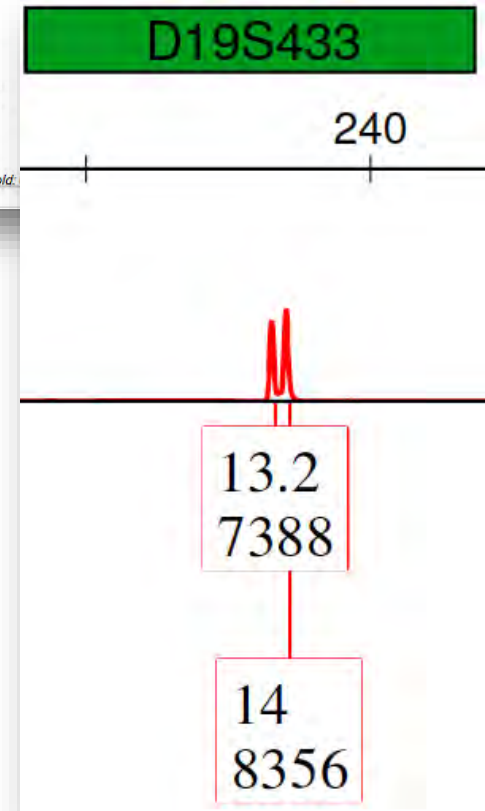
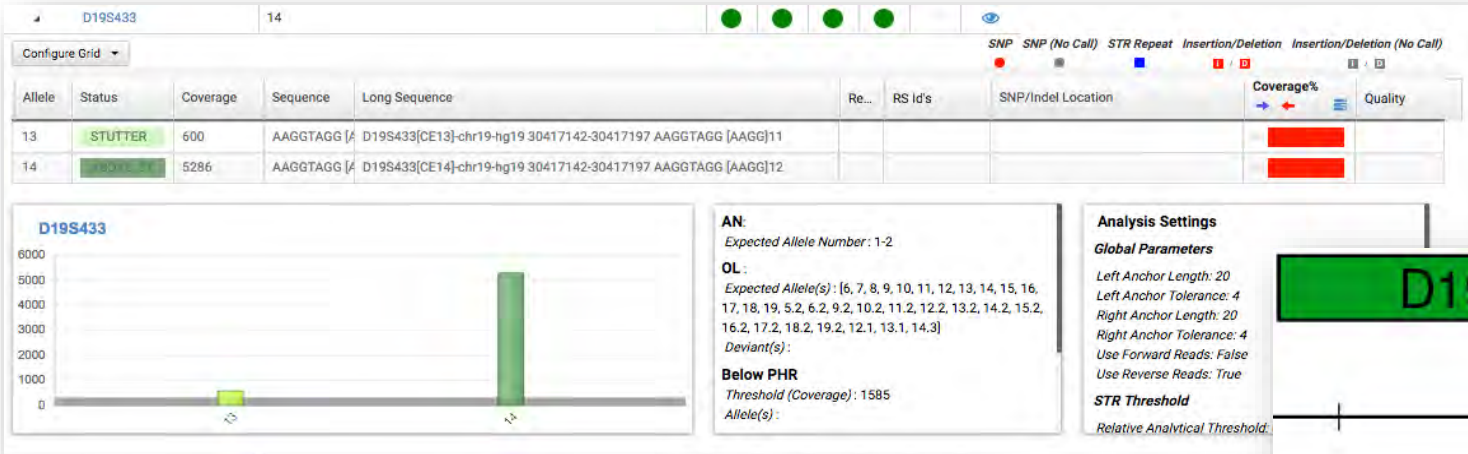
STR Threshold
Relative Analytical Threshold: 0.02



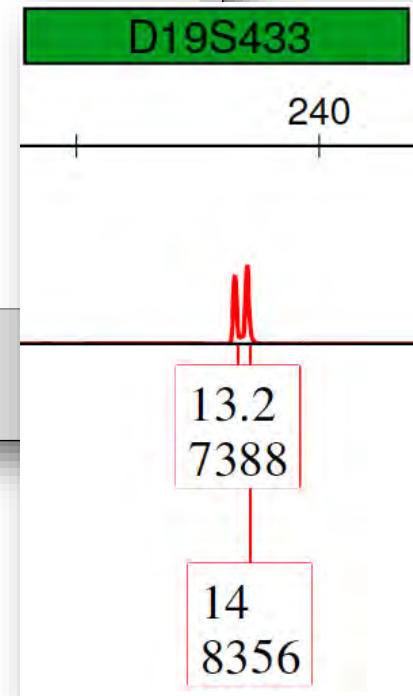
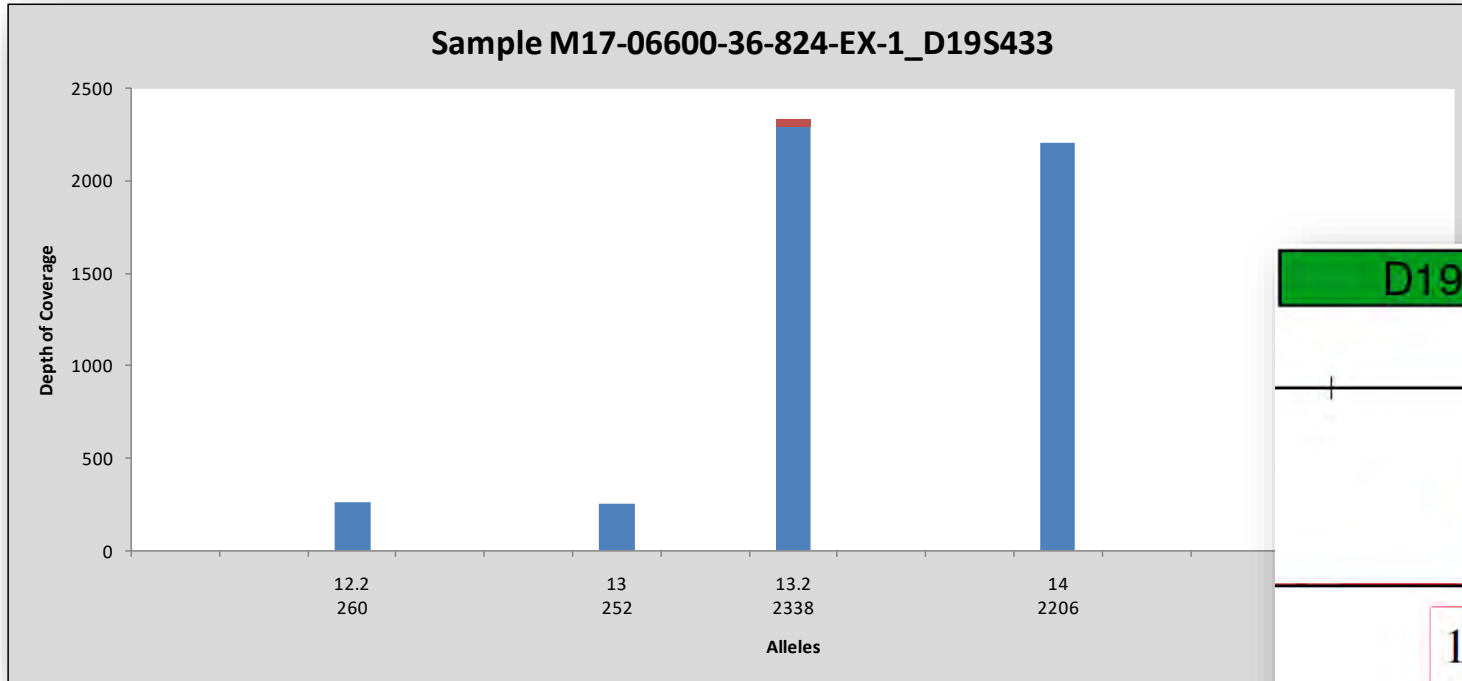
Additional analysis: IGV v2.4.16 (Robinson et al., 2011; Thorvaldsdottir et al., 2013)



1 sample: M17-06600-36-824-EX-1



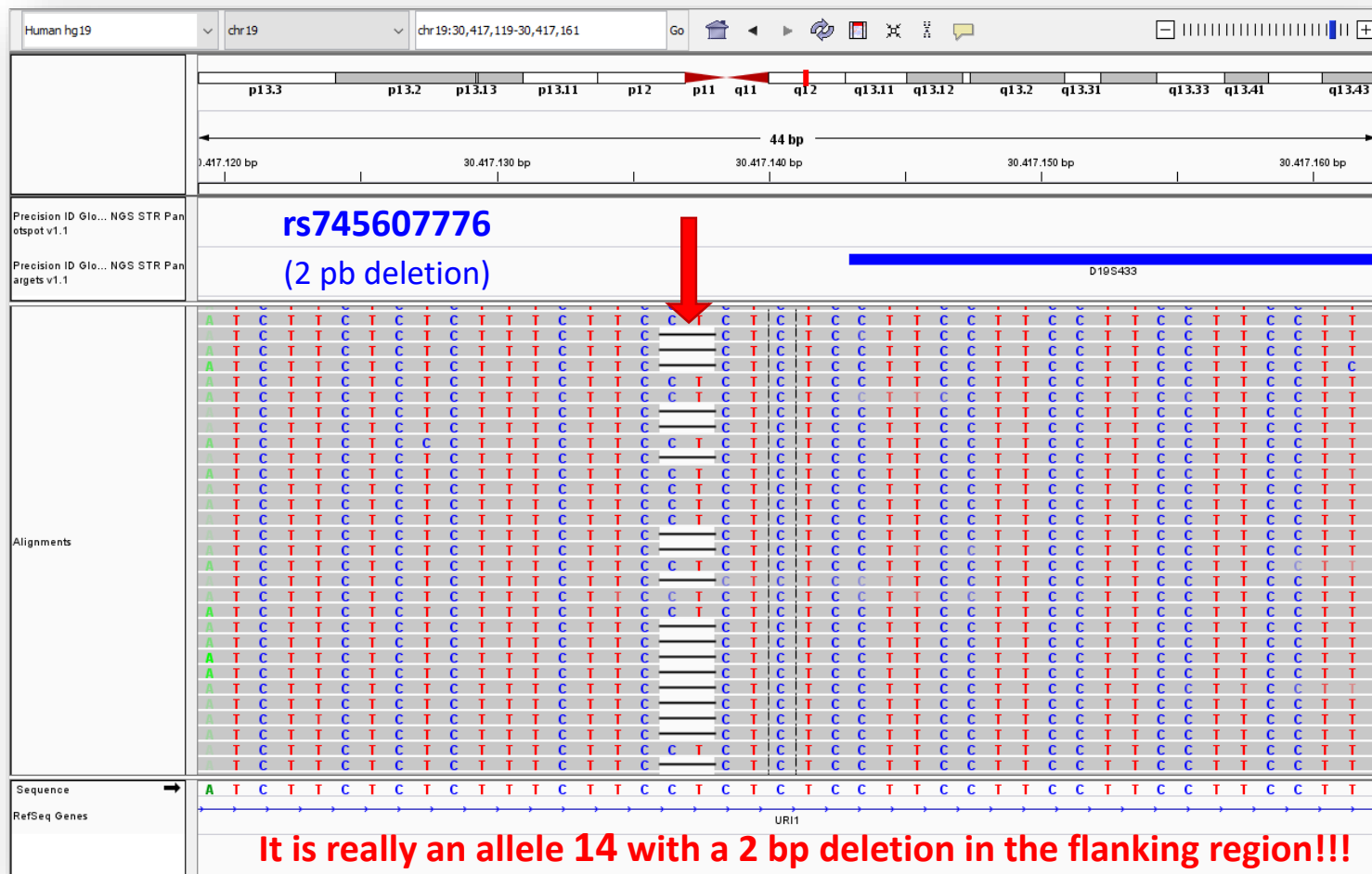
Additional analysis: *STRait Razor v3* (Woerner et al., 2017)



Additional analysis: *STRSeq catalog* (Gettings et al., 2017)

BioProject: [PRJNA380574](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA380574) D19S433[CE13.2]  [CCTT]12 ccta CCTT cttt CCTT rs745607776

Integrative Genomics Viewer - *IGV v2.4.16* (Robinson et al., 2011; Thorvaldsdottir et al., 2013)



* Findings based on the **CE** (*Fusion 6C*) - **MPS** (*Global NGS*) concordance study:

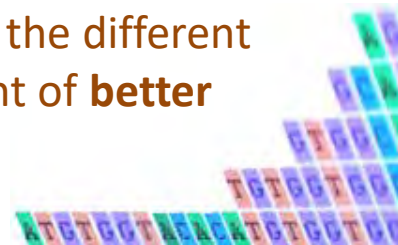
- Deletions in the flanking region (13 bp deletion in Penta D associated to CE allele 2.2) may cause “**bioinformatic null alleles**” if the selected region matches the anchor region of sequence recognition by the software.
- Deletions and insertions in the flanking regions also generated discrepancies between the **CE micro-variant alleles** and **MPS** data for Penta D, D2S441, and D19S433.

Penta_D[CE2.2]-chr21-hg19 45056086-45056150 **[AAAGA]5** 45056081-45056093-**delAAAAGAAAGAAAA**

D2S441[CE13.1]-chr2-hg19 68239079-68239126 **A [TCTA]10 TTTA [TCTA]2**

D19S433[CE13.2]-chr19-hg19 30417142-30417197 **[CCTT]12** ccta **CCTT** cttt **CCTT** 30417137-3041713-**delGA**

- The knowledge of the MPS sequence demonstrated that **these microvariants were erroneously called by CE**, as they were alleles with complete repeat units plus insertions or deletions in the flanking region.
- The identification and classification of these small discrepancies in the different STR markers will allow software improvement and the development of **better comparison tools between CE and MPS data**.





Institute of Legal Medicine, Medical University of Innsbruck. Austria



Institute of Legal Medicine and Forensic Science, Charité - Universitätsmedizin Berlin. Germany



Center for Human Identification at the University of North Texas Health Science Center. USA



National Institute of Toxicology and Forensic Sciences. Madrid Department. Spain



Thanks
for your attention!!



Acknowledgements



This project has been funded with support from the European Commission (grant **HOME/2014/ISFP/AG/LAWX/4000007135** under the Internal Security Funding Police programme of the European Commission---Directorate General Justice and Home Affairs). This publication reflects the views only of the authors, and the European Commission cannot be held responsible for any use which may be made of the information contained therein.



The authors would like to thank members of LIMS Administrators Team of the General Subdirectorate of New Technologies of Justice (SGNTJ) of the Ministry of Justice (Spain) for their helpful technical support.



Thermo Fisher Scientific Inc. Provider of *Precision ID GlobalFiler® NGS STR Panel v2* and *Converge v2.1* software. The authors would like to thank Matt Phipps for technical support.



Faculty of Health and Medical Sciences

STR sequencing in Copenhagen

Claus Børsting MSc, PhD
Senior researcher

Section of Forensic Genetics
Department of Forensic Medicine
Faculty of Health and Medical Sciences
University of Copenhagen
Denmark



STR sequencing in Copenhagen



A bit of history

Case reporting

NGS nomenclature issues



STR sequencing in Copenhagen - history



2011: 454 GS Junior (Roche)

Fordyce et al., (2011) *Biotechniques* 51, 127-133.
Rockenbauer et al., (2014) *FSI genet.* 8, 68-72.
Dalsgaard et al., (2014) *FSI genet.* 8, 195-199.
Gelardi et al., (2014) *FSI genet.* 12, 38-41.

2014: Ion PGM (Thermo Fisher Scientific)

Fordyce et al., (2015) *FSI genet.* 14, 132-140.
Friis et al., (2016) *FSI genet.* 21, 68-75.
Vilsen et al., (2017) *FSI genet.* 28, 82-89.

2015: MiSeq FGx (Verogen)

Hussing et al., (2018) *Foren. Sci. Res.* 3, 111-123.
Vilsen et al., (2018) *FSI genet.* 35, 107-112.
Hussing et al., (2019) *Int. J. Legal Med.* 133, 325-334.
Simayijiang et al., (2019) in preparation.
DNASEQEX 29 Y-STR panel test (2019).

2016: Ion S5 System (Thermo Fisher Scientific)

Precision ID Globalfiler mixture ID test (2017).
Project Iceberg: Precision ID Globalfiler NGS STR panel test (2018-19).



STR sequencing in Copenhagen - history



Nomenclature*:

- Locus name (used in forensic genetics)
- Length of repeat region/length of subunit
- Sequence(s) of subrepeat(s) followed by the number of repeats
- Variation in the flanking regions

Examples:

TH01[9] AATG[9]

D5S818[12] AGAT[9]ACAT[1]AGAT[2]rs25768[T]rs73801920[G]

DXS10135[23] AAGA[3]GAAAGGA[1]AAGA[19]AAAG[1]del:9306454-6

Python script: STRinNGS[‡]

*Gelardi *et al.*, FSI genet. (2014) 12, 38-41

‡ Friis *et al.*, (2016) FSI genet. 21, 68-75.



Case reporting in Copenhagen - history



Forensic Genetics in Denmark

- Accredited according to ISO 17025
- University based service
- Independent of the State administration, police, and judiciary system

Report everything we detect

- Unless the quality of the results prevent it
- Example: LPL in the FFFL panel
- Example: rs7520386 and rs576261 in the Precision ID Identity Panel



STR sequencing nomenclature



Examples:

TH01[9] AATG[9]

D5S818[12] AGAT[9]ACAT[1]AGAT[2]rs25768[T]rs73801920[G]

DXS10135[23] AAGA[3]GAAAGGA[1]AAGA[19]AAAG[1]del:9306454-6

STRidER nomenclature*:

TH01 [AATG]9

D5S818 [ATCT]9 ATGT [ATCT]2 rs25768[T] rs73801920[G]

DXS10135 [AAGA]3 gaaagga [AAGA]19 AAAG del:9306454-6

*Bodner et al., (2016) FSI genet. 24, 97-102.

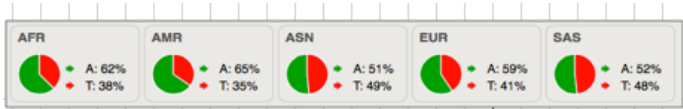
Phillips et al., (2018) FSI genet. 34, 162-169.





STR sequencing nomenclature

D13S317	
Reference sequence	C C C A T C T A A C G A C C C T A T T G T A T T T A C A A A T A C A T
Flanking SNP IUPAC codes	Y R
GRCh38 coordinates	82147991 82147992 82147993 82147994 82147995 82147996 82147997 82147998 82147999 82148000 82148001 82148002 82148003 82148004 82148005 82148006 82148007 82148008 82148009 82148010 82148011 82148012 82148013 82148014 82148015 82148016 82148017 82148018 82148019 82148020 82148021 82148022 82148023 82148024 82148025 82148026 82148027 82148028 82148029 82148030 82148031 82148032
GRCh37 coordinates	82722126 82722127 82722128 82722129 82722130 82722131 82722132 82722133 82722134 82722135 82722136 82722137 82722138 82722139 82722140 82722141 82722142 82722143 82722144 82722145 82722146 82722147 82722148 82722149 82722150 82722151 82722152 82722153 82722154 82722155 82722156 82722157 82722158 82722159 82722160 82722161 82722162 82722163 82722164 82722165 82722166 82722167
Distance from repeat region	34 33 32 31 30 29 28 27 26 25 24 23 22 21 20 19 18 17 16 15 14 13 12 11 10 9 8 7 6 5 4 3 2 1



8	9	10	11
A A A A A A A A	A A A A A A A A	A A A A A A A A	A A A A A A A A
T T T T T T T T	T T T T T T T T	T T T T T T T T	T T T T T T T T
C C C C C C C C	C C C C C C C C	C C C C C C C C	C C C C C C C C
G G G G G G G G	G G G G G G G G	G G G G G G G G	G G G G G G G G
A A A A A A A A	A A A A A A A A	A A A A A A A A	A A A A A A A A
T T T T T T T T	T T T T T T T T	T T T T T T T T	T T T T T T T T
C C C C C C C C	C C C C C C C C	C C C C C C C C	C C C C C C C C
T T T T T T T T	T T T T T T T T	T T T T T T T T	T T T T T T T T
A A A A A A A A	A A A A A A A A	A A A A A A A A	A A A A A A A A
T T T T T T T T	T T T T T T T T	T T T T T T T T	T T T T T T T T
C C C C C C C C	C C C C C C C C	C C C C C C C C	C C C C C C C C
T T T T T T T T	T T T T T T T T	T T T T T T T T	T T T T T T T T
A A A A A A A A	A A A A A A A A	A A A A A A A A	A A A A A A A A
T T T T T T T T	T T T T T T T T	T T T T T T T T	T T T T T T T T
C C C C C C C C	C C C C C C C C	C C C C C C C C	C C C C C C C C
T T T T T T T T	T T T T T T T T	T T T T T T T T	T T T T T T T T
A A A A A A A A	A A A A A A A A	A A A A A A A A	A A A A A A A A
T T T T T T T T	T T T T T T T T	T T T T T T T T	T T T T T T T T
C C C C C C C C	C C C C C C C C	C C C C C C C C	C C C C C C C C
T T T T T T T T	T T T T T T T T	T T T T T T T T	T T T T T T T T
A A A A A A A A	A A A A A A A A	A A A A A A A A	A A A A A A A A
T T T T T T T T	T T T T T T T T	T T T T T T T T	T T T T T T T T
C C C C C C C C	C C C C C C C C	C C C C C C C C	C C C C C C C C
T T T T T T T T	T T T T T T T T	T T T T T T T T	T T T T T T T T
A A A A A A A A	A A A A A A A A	A A A A A A A A	A A A A A A A A
T T T T T T T T	T T T T T T T T	T T T T T T T T	T T T T T T T T
C C C C C C C C	C C C C C C C C	C C C C C C C C	C C C C C C C C
T T T T T T T T	T T T T T T T T	T T T T T T T T	T T T T T T T T
A A A A A A A A	A A A A A A A A	A A A A A A A A	A A A A A A A A
T T T T T T T T	T T T T T T T T	T T T T T T T T	T T T T T T T T
C C C C C C C C	C C C C C C C C	C C C C C C C C	C C C C C C C C
T T T T T T T T	T T T T T T T T	T T T T T T T T	T T T T T T T T
A A A A A A A A	A A A A A A A A	A A A A A A A A	A A A A A A A A
T T T T T T T T	T T T T T T T T	T T T T T T T T	T T T T T T T T
C C C C C C C C	C C C C C C C C	C C C C C C C C	C C C C C C C C
T T T T T T T T	T T T T T T T T	T T T T T T T T	T T T T T T T T





STR sequencing nomenclature

converge

Upload Approve Reject Apply Changes Reset Current Data Set to Default

D5S818 11, 12

Configure Grid

Allele	Status	Coverage	Sequence	Long Sequence	Ref/...	RS Id's	SNP/indel Location	Coverage%	Quality
10	STUTTER	310	[AGAT]10	D5S818[CE10]-chr5-hg19 123111250-123111293 [AGAT]10					
11		4235	[AGAT]11	D5S818[CE11]-chr5-hg19 123111250-123111293 [AGAT]11					
11		360	[AGAT]11	D5S818[CE11]-chr5-hg19 123111250-123111293 [AGAT]11			-----■-----		
11		3867	[AGAT]11	D5S818[CE11]-chr5-hg19 123111250-123111293 [AGAT]11 123111306-G	A/G	rs25768	-----■-----●-----		17.69
12		2733	[AGAT]12	D5S818[CE12]-chr5-hg19 123111250-123111293 [AGAT]12					
12		2699	[AGAT]12	D5S818[CE12]-chr5-hg19 123111250-123111293 [AGAT]12			-----■-----		

D5S818

4500
4000
3500
3000
2500
2000
1500
1000
500
0

10 11 12

AN:
Expected Allele Number: 1-2

OL:
Expected Allele(s): [6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 10.1, 11.1, 12.1, 12.3]
Deviant(s):

Below PHR
Threshold (Coverage): 1270
Allele(s):

BST:

Analysis Settings
Global Parameters
Left Anchor Length: 15
Left Anchor Tolerance: 2
Right Anchor Length: 15
Right Anchor Tolerance: 2
Use Forward Reads: True
Use Reverse Reads: True

STR Threshold
Relative Analytical Threshold: 0.02

360 stutter reads

What is the SNP-STR haplotype



Summary



Nomenclature

STR region (clear definitions in STRidER)

- Start and end (genome build)
- Forward strand
- Subrepeat format

Flanking regions

- InDels should be incorporated for back-compatibility with CE
 - PCR primers should be made available (old and new kits)
- SNPs?
 - Global nomenclature for SNP-STR - which SNPs?
 - Restrict multiplex design
 - Challenge sequencing length

Database in place (STRSeq)





Thank you for the invitation



THE SEQUENCE IDENTIFIER (SID) OPERATIONAL NOMENCLATURE FOR MIXED-DNA CASEWORK USING MPS

STRAND Working Group Nomenclature Meeting

11 April 2019

Brian Young, Ph.D.

brian@nichevision.com

NicheVision Forensics, Inc.

Outline

1. Motivation
2. The Challenge of MPS Artifacts in Mixed DNA Interpretation
3. Sequence Identifier (SID) Nomenclature
4. Why We Need SID: Labeling of Artifacts by ‘Phylogeny’
5. Mixture Interpretation of Sequence-Based MPS Data Using MixtureAce™ Plugin to ArmedXpert™ Software

Motivation

- Casework by MPS Has Been Hindered by a Lack of Available Methods for Interpreting Mixed DNA Samples
- We Developed a Mixed DNA Interpretation Method But it Requires an ‘Operational’ Nomenclature

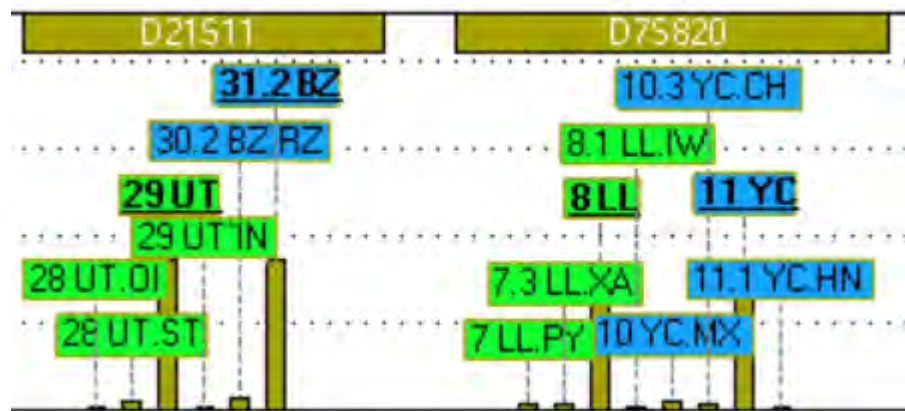
Outline

1. Motivation
2. The Challenge of MPS Artifacts in Mixed DNA Interpretation
3. Sequence Identifier (SID) Nomenclature
4. Why We Need SID: Labeling of Artifacts by ‘Phylogeny’
5. Mixture Interpretation of Sequence-Based MPS Data Using MixtureAce™ Plugin to ArmedXpert™ Software

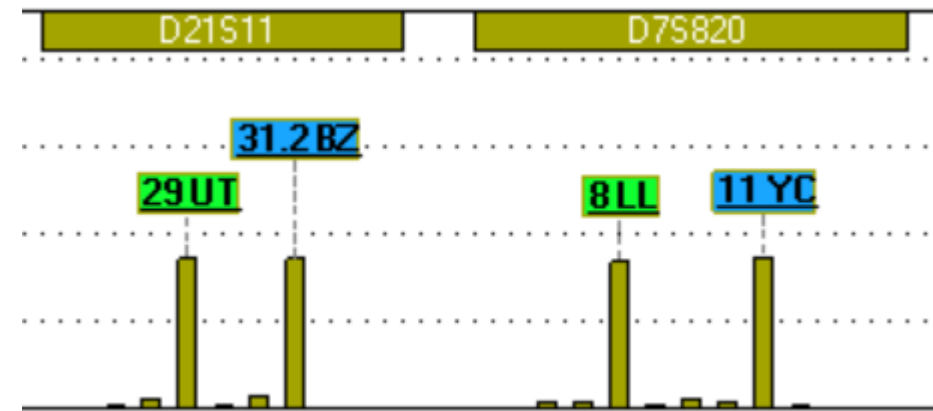
Challenges

1. Space in Computer Displays is Limited
2. Both Artifacts and Alleles Must be Labeled In Mixture Analysis
3. Artifacts are Errors, Not Polymorphisms (rs# Not Available)
4. Universal: Every Laboratory Gets Equivalent Labels for Equivalent Sequences
5. Sequence Type Labeling That is Easy to Vocalize When Discussing Profiles

Example Screen Captures from ArmedXpert/MixtureAce



Show Stutter and Non-Stutter Artifacts



Filter Stutter and Non-Stutter Artifacts

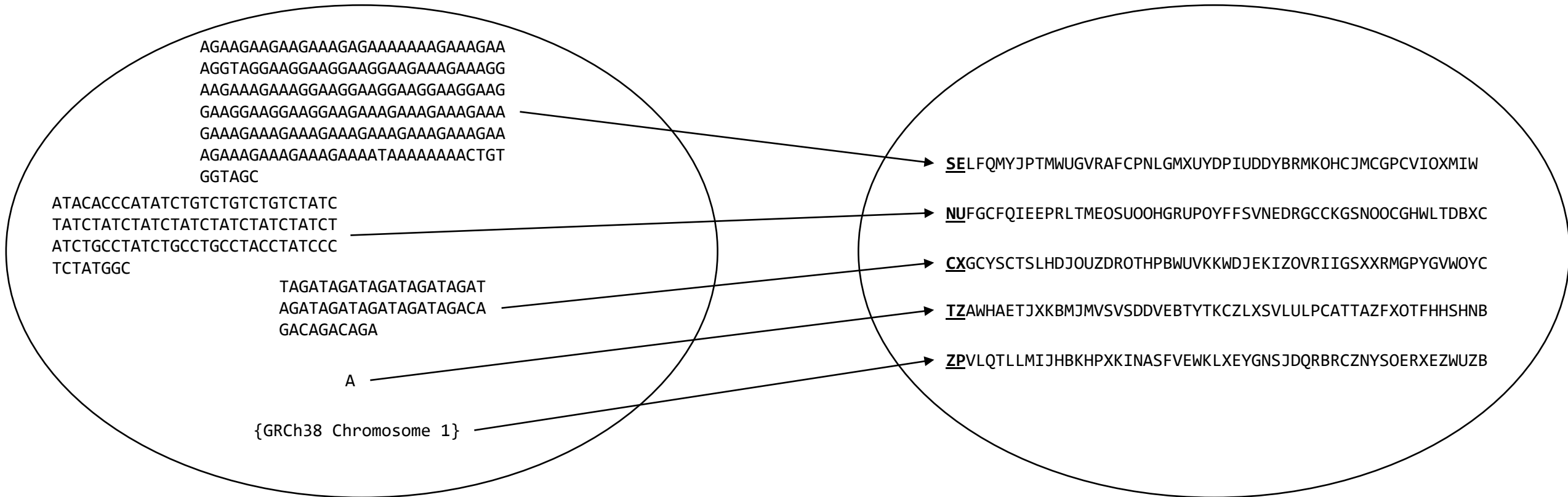
Outline

1. Motivation
2. The Challenge of MPS Artifacts in Mixed DNA Interpretation
3. Sequence Identifier (SID) Nomenclature
4. Why We Need SID: Labeling of Artifacts by ‘Phylogeny’
5. Mixture Interpretation of Sequence-Based MPS Data Using MixtureAce™ Plugin to ArmedXpert™ Software

SID: Shorthand Labels for Sequences

Domain of "All" Forensic Sequences
($< 10^8$)

Range of SID Function
($2^{256} \sim 26^{55} \sim 10^{77}$)



Any Sequence String → Mapping Function → 55-Character SID Label (Little-Endian)
 Digits Arranged Least- to Most-Significant

SIDs are Shorthand Labels for Sequences

Domain of "All" Forensic Sequences ($< 10^8$)

1,000 forensic loci
100 alleles each = 10^8
1,000 trim sites each

Number of possible sequences of a 200-nucleotide amplicon = $4^{200} \sim 10^{120}$

Number of human chromosomes on earth $\sim (8 \times 10^9) \times (2)$
 $\sim 1.6 \times 10^{10}$

Number of different sequences observed in allele surveys $< \sim 100$

Range of SID Function ($2^{256} \sim 26^{55} \sim 10^{77}$)

SELFQMYJPTMWUGVRAFCPNLGMXUYDPIUDDYBRMKOHCJMCGPCVIOXMIW
NUFGCFQIEEPRLTMEOSU00HGRUPOYFFSVNEDRGCKGKSNOOCGHWTDBXC
CXGCYSCTSLHDJOUZDROTHPBWUVKKWDJEKIZOVRIIGSXXRMGPYGVWOYC
TZAWHAETJXKBMJMVSVDVBEVTKCZLXSVLULPCATTAZFXOTFHSHNB
ZPVLQTLMLIJHBKHPXKINASFVEWKLXEYGNSSJQBRBCZNYSOERXEZUWZB

Any Sequence String \longrightarrow Mapping Function \longrightarrow 55-Character SID Label (Little-Endian)
Digits Arranged Least- to Most-Significant

SIDs are Shorthand Labels for Sequences

Domain of "All" Forensic Sequences
($< 10^8$)

1,000 forensic loci
100 alleles each
1,000 trim sites each
1,000 artifacts each

$= 10^{11}$

Range of SID Function
($2^{256} \sim 26^{55} \sim 10^{77}$)

SELFQMYJPTMWUGVRAFPCPNLGMXUYDPIUDDYBRMKOHCJMCGPCVIOXMIW
NUFGCFQIEEPRLTMEOSU00HGRUPOYFFSVNEDRGCCCKGSNOOCGHWLTDBXC
CXGCYSCTSLHDJOUZDROTHPBWUVKKWDJEKIZOVRIIGSXXRMGPYGVWOYC
TZAWHAETJXKBMJMVSVSDDVEBTYTKCZLXSVLULPCATTAZFXOTFHSHNB
ZPVLQTLLEMIJHBKHPXKINASFVEWKLXEYGNSSJQBRBCZNYSOERXEZUWZB

Any Sequence String \longrightarrow Mapping Function \longrightarrow 55-Character SID Label (Little-Endian)
 Digits Arranged Least- to Most-Significant

Algorithm is Highly Sensitive to Small Changes In Input Strings

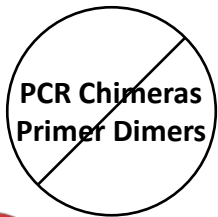
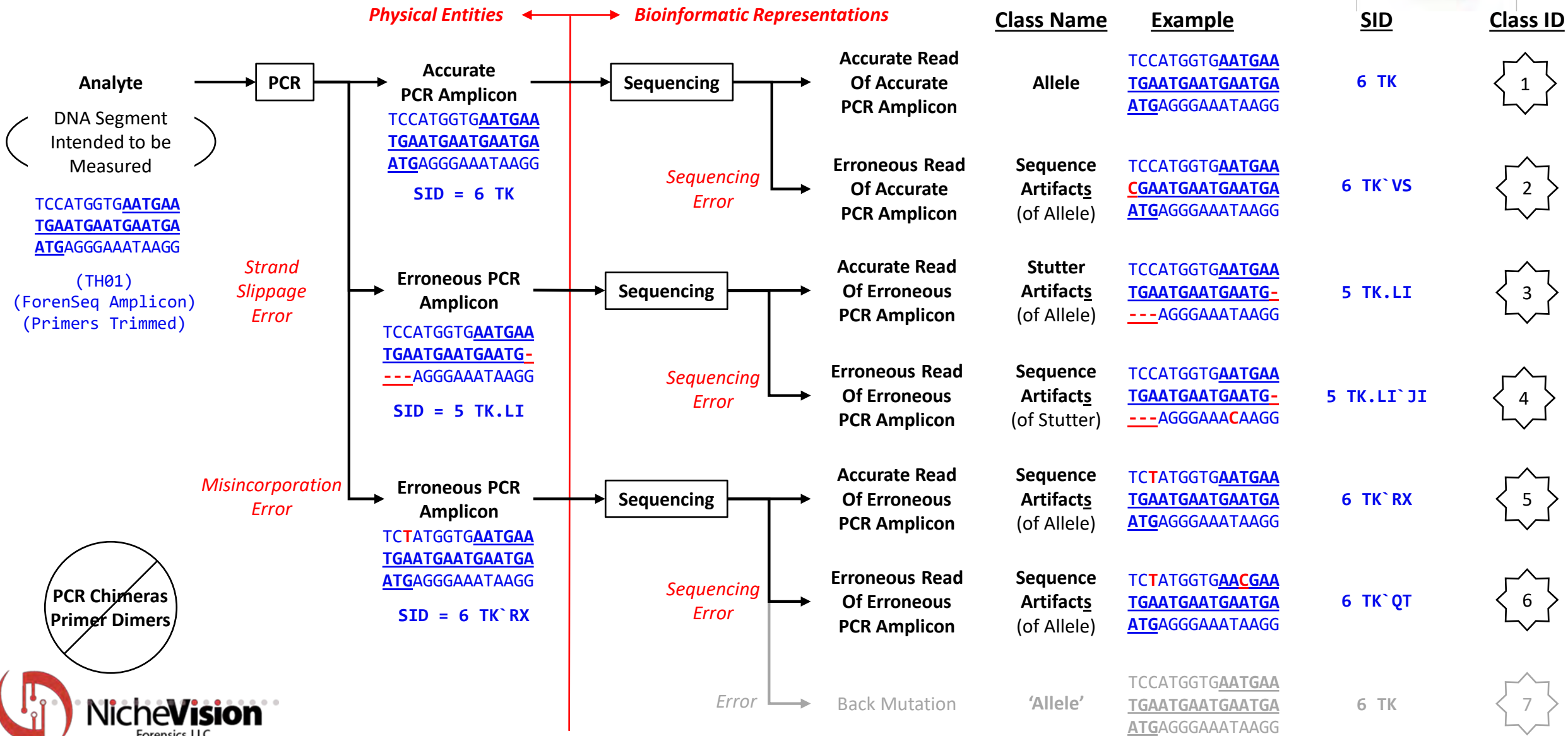
In cryptography, the “avalanche effect” means that nearly all digits change in response to even a single change in the input string.

Letter Changes	Sequence String	SID Labels
Base Sequence	TAGATAGA	EJZUKYGKYTNLXLXDTVLVGOTSMYKZJVBDKFJZLWIXEVULUGXZZDFGFJD
Substitution	TAGGTAGA	TXSUCZXTZDJWUVOCTXTTLMPNJRMUHJRUDLRFFDGWFVQQCIKYHMYFVB
Deletion	TAGATAG_	INOGHEHOUNOGOBUYWZBCLPJOHQJMLUJUZYFWTAMSIHWKZCALUGNUAC
Insertion	TAGATAGAT	NPUCXCAJHYLRPJWBFHTLJUCXXMYTKBAHZIUKJIBUFPJYFOEIOWLQCDB

Outline

1. Motivation
2. The Challenge of MPS Artifacts in Mixed DNA Interpretation
3. Sequence Identifier (SID) Nomenclature
4. **Why We Need SID: Labeling of Artifacts by ‘Phylogeny’**
5. Mixture Interpretation of Sequence-Based MPS Data Using MixtureAce™ Plugin to ArmedXpert™ Software

'Phylogenetic' Model of the PCR-MPS Method



Artifact Assignment Rules

- Shared stutter artifacts can be apportioned to multiple parent alleles
- Ambiguous non-stutter artifacts remain as candidate alleles (potential minor)
- Categories of ambiguous non-stutter artifacts:
 1. Questioned sequences that are minimally close to > 1 parent
 2. Questioned sequences with minimal distances to any parent $> \delta$

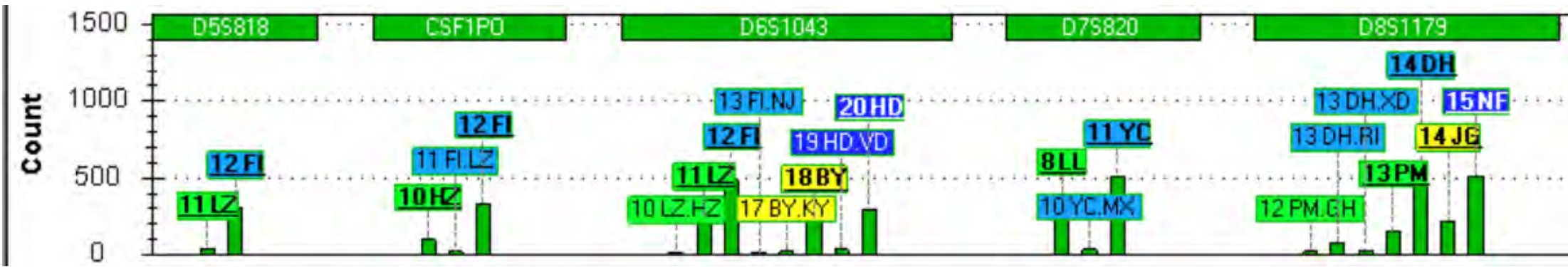
SID & Connector Conventions for MixtureAce Parent-Child Associations

PARENT . STUTTER	BT . NU	(dot)
PARENT ` SEQARTIFACT	BT ` HL	(tick)

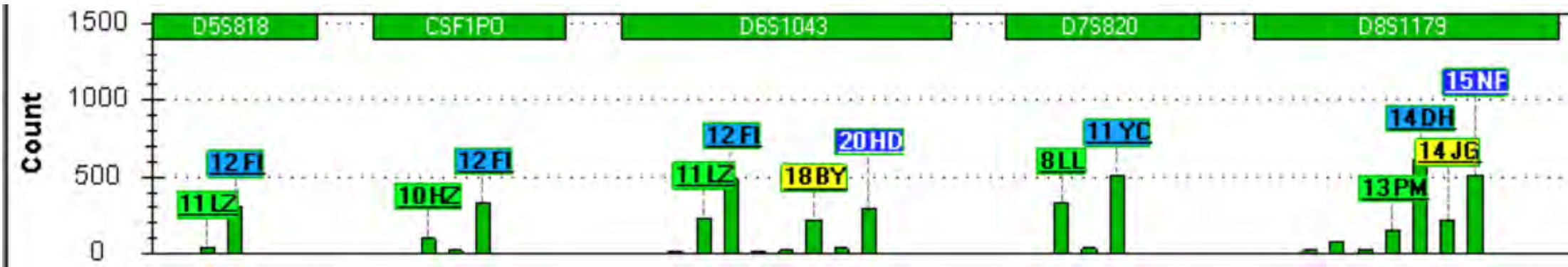
Outline

1. Motivation
2. The Challenge of MPS Artifacts in Mixed DNA Interpretation
3. Sequence Identifier (SID) Nomenclature
4. Why We Need SID: Labeling of Artifacts by ‘Phylogeny’
5. Mixture Interpretation of Sequence-Based MPS Data Using MixtureAce™ Plugin to ArmedXpert™ Software

3:1 Mixture, AT = 10

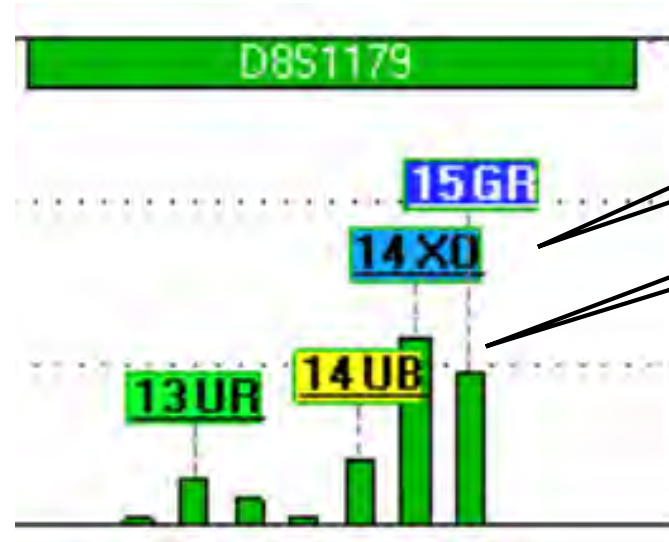


Show Artifacts



Filter Artifacts

Isoalleles Are Displayed Separately



Contributor	Allele Length	Allele Sequence
2800M (Major)	14	[TCTA]1[TCTG]1[TCTA]12
2391c-A (Minor)	14	[TCTA]2[TCTG]1[TCTA]11

ArmedXpert Match & Comparison Tool With Sequence-Based Alleles



MixtureAce v1.23

Views Data QC Checks Match & Comparison Interpretation Reporting

Set Reference Find Where Reference is Included Find Included in the Reference Find Foreign to Reference Find Where Reference Includes Foreign Match & Compare Profiles Clear Match

A: Table Reference A B: Table Reference B Compare Tables Organize

Sources Windows

C:\Users\brian\OneDrive\Desktop\OCME15Feb91\2800M_2391cA_3_1.fastq

	FGA	D5S818	CSF1PO	D6S1043	D7S820	D8S1179	D9S1122	D10S1248	TH01
2800M_2391cA_3_1.fastq	20 CV, 21 AC, 23 VR	12 DR	10 FY, 12 LY	11 EL, 12 WY, 18 TM, 20 HA	8 PM, 11 JU	13 UR, 14 XO, 14 UB, 15 GR	11 OK, 12 ZG, 12 NG	13 LR, 15 KY, 16 YK	6 TK, 8 RZ
2800M.fastq [1 Reference]	20 CV, 23 VR	12 DR	12 LY	12 WY, 20 HA	8 PM, 11 JU	14 XO, 15 GR	12 ZG, 12 NG	13 LR, 15 KY	6 TK, 9.
Exact match									
Included									
Find Where Reference Included									

Clear Selection

ArmedXpert Mixture Deconvolution With Sequence-Based Alleles



MixtureAce v1.23 - [Mixture Interpretation - DDA Interpretation 2800M_2391cA_3_1.fastq]

Views Data QC Checks Match & Comparison Interpretation Reporting

Begin Mixture Interpretation Highlight Single Source Likelihood Ratio Probability of Inclusion RMP User Defined Organize

Mixture Interpretation Frequency Calculations Windows

Setup
Pick via mouse 2800M_2391cA_3_1.fastq Operations Contributor # 2

Locus D21S11 (4) Profile 01 29 KN, 31.2 TK P. Avg(0.74)
Profile 02 28 TC, 32.2 SL P. Avg(0.26)

Alleles 28 TC, 29 KN, 31.2 TK, 32.2 SL
RFUs 110, 266, 360, 84
BPs 183, 187, 197, 201

Peaks
Apply Globally 100% Profile 01
Apply Stutter 0.09

	28 TC	29 KN	31.2 TK	32.2 SL
Profile 01	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
Profile 02	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>

Mixture Information
All combinations have: PHr >= 0, MPh >= 0, mP >= 0.1

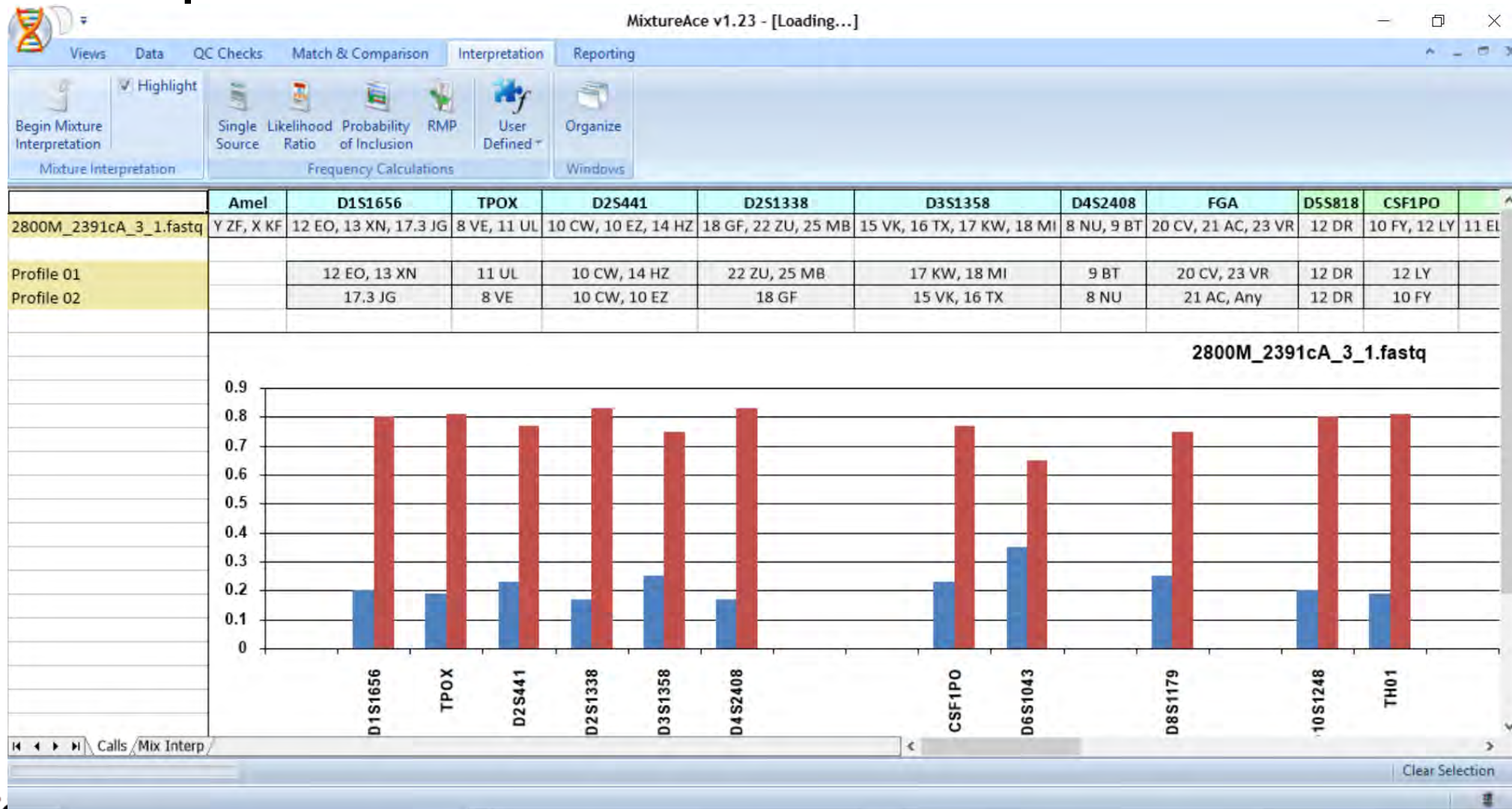
For a 2-contributor 4-allele mixture of types AB & CD: 3/3-combination(s):
 31.2 TK, 32.2 SL(phr = 0.23; p = 0.54) • 28 TC, 29 KN(phr = 0.41, p = 0.46)
 28 TC, 31.2 TK(phr = 0.31; p = 0.57) • 29 KN, 32.2 SL(phr = 0.32; p = 0.43)
 29 KN, 31.2 TK(phr = 0.74; p = 0.76) • 28 TC, 32.2 SL(phr = 0.76; p = 0.24)

Mixture Interpretation of 2800M_2391cA_3_1.fastq

	28 TC	29 KN	31.2 TK	32.2 SL
	110	266	360	84

PHr 0.00 Multi PHr
mPH 0 HT 75
mP 0.10
Popout calls View call report Add Comment

ArmedXpert Mixture Deconvolution Call Report



Acknowledgements

- SID Nomenclature
 - Nate Caldwell
 - Tom Faris
- Data
 - Dr. Elisa Wurmbach OCME
 - Dr. Michael Marciano, Syracuse U.
- Software Development
 - Luigi Armogida
 - Tom Faris
 - Abdul Alali



- **Lab-staff:** Kristiaan van der Gaag
Thirsa Kraaijenbrink
Rick van Leeuwen
Jerry Hoogenboom

Practical Issues:

- We now have to describe sequence variation.
- We have a unique opportunity to repair forward – reverse mistakes.
- There will be “repeat” confusion (N-CE \neq N-MPS).
- Have to deal with SNP variation.
- There is already a HGVS nomenclature in use.
- Describing MPS – STR variants in reports (articles) must be compatible with CE – STR results.
- Allele input in stat-progs and databases.

Two choices:

- A one size fits all solution (a code).
- A set of different solutions.

Minimal criteria:

- No online only repository / coding solution.
- Should be able to accommodate kit design variation.
- Accept FASTA – format.

MPS - STR Nomenclature

nomenclature challenges

D13S317-11	TCTAACGCCT	ATCTGTATTT	ACAAATACAT	-TATC-TATC	-TATC-TATC
D13S317 [CE12]A...	-.....-.....	-.....-.....
D13S317-11	-TATC-TATC	-TATC-TATC	-TATC-TATC	-TATC-----	-----AATC
D13S317 [CE12]	-.....-.....	-.....-.....	-.....-.....	-.....-TATC	-TATC-.....
D13S317-11	-AATC-ATCT	-ATCT-ATCT	-TT		
D13S317 [CE12]	-----.....	-.....-.....	-..		

D13S317[CE11]-Chr13-GRCh37-g.82.722.160:82.722.223-TATC[11]AATC[2]ATCT[3]

D13S317[CE12]-Chr13-GRCh37-g.82.722.160:82.722.223-TATC[13]AATC[1]ATCT[3]-g.x.136G>A

D13S317[CE12] - TATC[13]AATC[1]ATCT[3] – 24 136G>A

D13S317[CE12] – 101



Different naming approaches for each situation

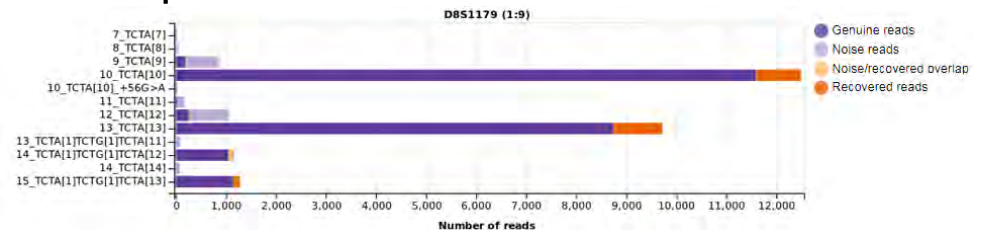
- The entire sequence
 - Exchange of data
 - Software processing and storage in database
- **Human-readable format; preserving sequence information**
 - **Publications, presentations, reports, etc.**
 - **Manual profile interpretation**
 - **Within-case manual comparisons**
- Very short code
 - Compact tables
 - Software that can't handle complete sequences

STRNaming



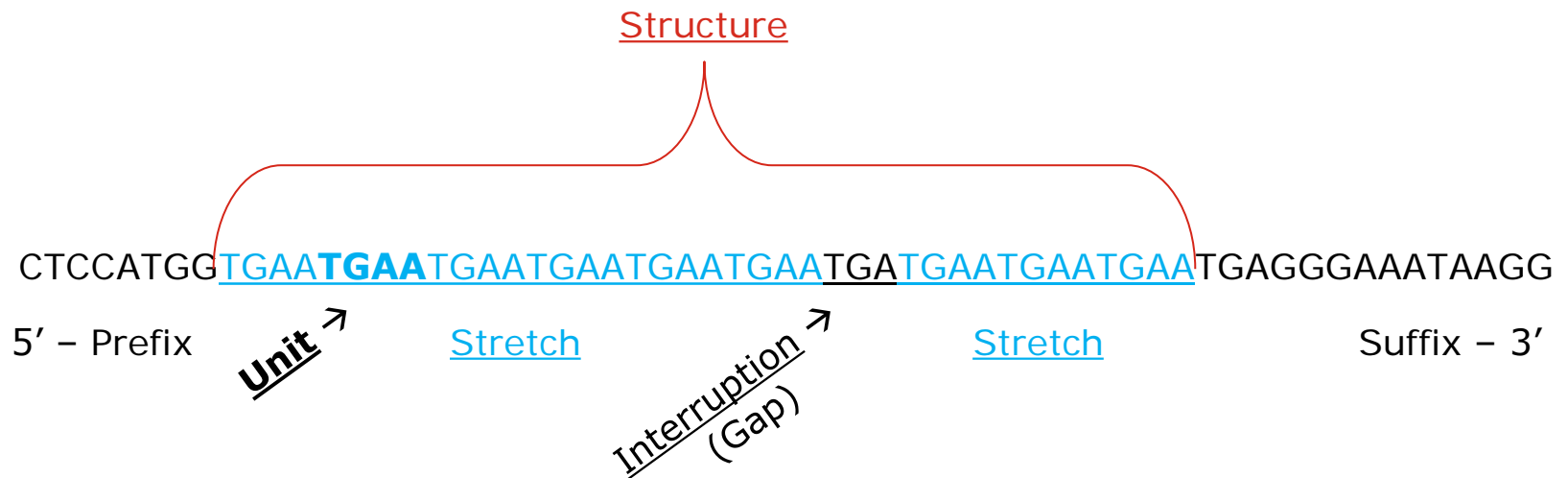
Our focus: human-interpretable names

- One unique name describing one sequence
 - One-to-one relationship
- Sequence-descriptive
 - Small change in sequence → small change in name
 - The exact sequence change is obvious from the name
 - **Essential** in interpreting stutters/hybrids/artefacts by humans!
- Standardised and automated
 - Should work the same on any locus in current or future use
 - Don't want to spend weeks working out the details of a locus
 - Publically available online & offline; open source algorithm





Definitions





Which name would you prefer?

CTCCATGGTGAATGAATGAATGAATGAATGAATGATGAATGAATGAATGAGGGAAATAAGG

CTCCATGGTGAATGAATGAATGAATGAATGAATGATGAATGAATGAATGAGGGAAATAAGG

CTCCATGGTGAATGAATGAATGAATGAATGAATTGATGAATGAATGAATGAGGGAAATAAGG

CTCCATGGTGAATGAATGAATGAATGAATGAATGATGAATGAATGAATGAGGGAAATAAGG

CTCCATGGTGAATGAATGAATGAATGAATGAATTGATGAATGAATGAATGAGGGAAATAAGG

1. **ATGA[6]**TGA[1]**ATGA[3]** Nice! Two stretches of the same repeat unit!
2. **TGAA[6]**TGA[1]**TGAA[3]** Nice! Moved 5' as far as possible!
3. **ATGA[6]****TGAA[3]** Nice! No interruptions!
4. **TGAA[6]**TGA[2]**ATGA[3]** Nice! Only repeats, no interruptions, short prefix/suffix!

→ None of these is bad!

Also, which locus is this?



Which name would you prefer?

CTCCATGGTGAATGAATGAATGAATGAATGAATGATGAATGAATGAATGAGGGAAATAAGG 9.3

CTCCATGGTGAATGAATGAATGAATGAATGAATGATGAATGAATGAATGAGGGGAAATAAGG 9.3

CTCCATGGTGAATGAATGAATGAATGAATGAATTGATGAATGAATGAATTGAGGGAAATAAGG 9.3

CTCCATGGTGAATGAATGAATGAATGAATGAATGATGAATGAATGAATGAGGGAAATAAGG 9.3

CTCCATGGTGAATGAATGAATGAATGAATGAATTGATGAATGAATGAATGAGGGAAATAAGG 9.3

CE9.3 **ATGA**[6]TGA[1]**ATGA**[3]

Nice! Two stretches of the same repeat unit!

CE9.3 **TGAA**[6]TGA[1]**TGAA**[3]

Nice! Moved 5' as far as possible!

CE9.3 **ATGA**[6]**TGAA**[3]

Nice? No interruptions!

CE9.3 **TGAA**[6]TGA[2]**ATGA**[3]

Nice? Only repeats, no interruptions, short prefix/suffix!

→ Do we care about the CE number?

→ What would allele 9 look like?



Which name would you prefer?

CTCCATGGTGAATGAATGAATGAATGAATGAATGAATGAATGAATGAGGGAAATAAGG 9

CTCCATGGTGAATGAATGAATGAATGAATGAATGAATGAATGAATGAGAGGGAAATAAGG 9

CTCCATGGTGAATGAATGAATGAATGAATGAATGAATGAATGAATGAATGAGGGAAATAAGG 9

CTCCATGGTGAATGAATGAATGAATGAATGAATGAATGAATGAATGAGGGAAATAAGG 9

CTCCATGGTGAATGAATGAATGAATGAATGAATGAATGAATGAATGAAGGGAAATAAGG 9

CE9.3 **ATGA**[6]TGA[1]**ATGA**[3]

CE9.3 **TGAA**[6]TGA[1]**TGAA**[3]

CE9.3 **ATGA**[6]**TGAA**[3]

CE9.3 **TGAA**[6]TGA[2]**ATGA**[3]

CE9 **ATGA**[9]

CE9 **TGAA**[9]

CE9 **ATGA**[8]A[1]

CE9 **TGAA**[9]TGA[1]

Nice!

Nice!

Need to account for the extra A...

Name has changed a lot?

→ We have 100s of markers more...



Our solution!

STRNaming

- Figures out what is repeated
- Figures out the nicest way of shortening that
 - The rules aren't simple, but STRNaming does it for you
- Apply some pretty notation...
- Done!

Availability

- Free online & offline STRNaming tools
- Open-source the algorithm



But how?



The reference sequence

1. Find longest repeat stretch → *extract repeat unit only*
 - Repeat until no more non-overlapping stretches are found
2. Shift the extracted units (optional step)
 - E.g., CTAT → TATC, ATCT, TCTA
3. Find all stretches of all extracted units (ignoring overlap)
4. Select the nicest combination of repeat stretches (structure)

1. CTTCTATCTATCTATCTATCTATCTATCTATCTATCTATCTAATCTATCTATCTTCA
CTAT ATCT

2. CTAT TATC ATCT TCTA

3. CTTCTATCTATCTATCTATCTATCTATCTATCTATCTATCTAATCTATCTATCTTCA
CTTCCTATCTATCTATCTATCTATCTATCTATCTATCTATCTAATCTATCTATCTTCA
CTTCCTATCTATCTATCTATCTATCTATCTATCTATCTATCTAATCTATCTATCTTCA
CTTCCTATCTATCTATCTATCTATCTATCTATCTATCTATCTAATCTATCTATCTTCA

4. CTTCTATCTATCTATCTATCTATCTATCTATCTATCTAATCTATCTATCTTCA



The reference sequence

- From the reference, we store:
 - Which repeat **units** were used in its name → TCTA, ATCT
 - The prefix and suffix sequence → CTTCCTA, TCA
 - The sequence of the longest interruption (if very long) → n/a
 - Its length and CE allele number → 58bp = CE9

Sequence:

CTTCCTATCTATCTATCTATCTATCTATCTATCTATCTATCTAATCTATCTATCTTCA

Structure:

CTTCCTA(1)TCTA(9)ATCT(3)TCA(1)

Name:

CE9_TCTA[9]ATCT[3]



Naming other sequences

- From the reference, we had stored:
 - Which repeat **units** were used in its name → TCTA, ATCT
 - The prefix and suffix sequence → CTCCTA, TCA
 - The sequence of the longest interruption (if very long) → n/a
 - Its length and CE allele number → 58bp = CE9

- Find all stretches of TCTA and ATCT (ignoring overlap)

CTTCCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTTCA
CTTCCTATCTATCTATCTATCTATCTATCTATCTATCTAATCTATCTATCTTCA

- Select the nicest combination of repeat stretches (structure)

CTTCCTATCTATCTATCTATCTATCTATCTATCTATCTAATCTATCTATCTTCA

- Find and shorten repeat stretches in interruptions

- Can use repeat units other than TCTA and ATCT

- Convert the structure into an allele name

- Simply a matter of notation, e.g.: CE9_TCTA[9]ATCT[3]



But what is nice???

Calculate a score for each STR structure

- Bonus points for:
 - Every base covered by a repeat
 - Every repeat of a unit
 - Every interruption that is exactly x repeat units long
- Penalty points for:
 - Every distinct repeat unit used
 - Every interruption introduced between stretches
 - Every base in an interruption
 - Every base inserted or deleted in the prefix or suffix

**Same for
all loci!**

How many points? → Mostly common sense, expert opinion needed for the complex alleles



Some words on notation

- Variation in suffix

GGTAAACAGTATAATTTTC[10]TATTTG**A**AATGGA
 CE10_

GGTAAACAGTATAATTTTC[10]**+**7**A>**C
 CE10_

GGTAAACAGTATAATTTTC[10]**+**9**A>**-
 CE9.3_

GGTAAACAGTATAATTTTC[10]**+**9.1**->**A
 CE10.1_



Some words on notation

- Variation in prefix

GGTA**A**CAGTATATTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTATTGAAATGGA
 CE10_TTTTC[10]

10 9 8 7 6 5 4 3 2 1
 GGTA**C**CAGTATATTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTATTGAAATGGA
 CE10_TTTTC[10]_-9A>C

10 9 7 6 5 4 3 2 1
 GGTAACAGTATATTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTATTGAAATGGA
 CE9.3_TTTTC[10]_-8A>-

10 9 8 7 6 5 4 3 2 1
 GGTA**AA**A**A**CAGTATATTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTTTCTATTGAAATGGA
 CE10.1_TTTTC[10]_-8.1->A




Some words on notation

- Large interruptions (example: DYS389)

AGATTGATAGAGGGAGGGGATAGATAGATAGATAGATAGATAGATAGATAGATAGATAGATAGACAGACAGACAGATACATAGATAATACAGA
DYS389I CE12_GATA[9]GACA[3]

AGATTGATAGAGGGAGGGGATAGATAGATAGATAGATAGATAGATAGATAGATAGATAGATAGACAGACAGACAGATACATAGATAATACAGATGAGAGTTG
GATACAGAAGTAGGTATAATGATAGATAGATAGATAGATAGATAGATAGATAGATAGATAGATAGATAGATAGACAGACAGACAGACAGACAGACAC
ACACATAGATAATACAGA

DYS389II CE29_GATA[9]GACA[3][]GATA[12]GACA[6] 

AGATTGATAGAGGGAGGGGATAGATAGATAGATAGATAGATAGATAGATAGATAGATAGATAGACAGACAGACAGATAGATAGATAATACAGATGAGAGTTG
GATACAGAAGTAGGTATAATGATAGATAGATAGATAGATAGATAGATAGATAGATAGATAGATAGATAGATAGACAGACAGACAGACAGACAGACAC
ACACATAGATAATACAGA

DYS389II CE29_GATA[9]GACA[3][3C>G]GATA[12]GACA[6]



Eager to try it out?
Just drop me a line!

j.hoogenboom@nfi.nl



Final remarks

Ideas on the to-do list:

- Auto-convert to forward strand: ISFG reconsideration
- Analyse reference sequence with more upstream and downstream sequence content
 - Stabilises names across kits
 - Auto-determine analysed range per kit
 - Provide table of analysed range

→ We are happy to share what we've got!

→ We are open for suggestions and improvements!



Acknowledgements

LUMC (FLDO)

- Rick de Leeuw
- Peter de Knijff

King's College / Verogen

- Laurence Devesse

NFI

- Titia Sijen

STRAND Working Group

*Thank
you*



Some rough thoughts

Peter Gill, Oyvind Bleka

April 11, 2019

STR Nomenclature Meeting, London, UK

Mixtures

- Probabilistic genotyping open-source software EuroForMix
- Modules that can analyse SNPs and STRs using LUS
- We are designing software that will convert sequences into LUS primary, secondary, tertiary etc.
- When we analyse mixtures, we will encounter genotypes for which a reference sample is not available.
- Also we cannot tell if a sequence that is in an apparent stutter position of a major contributor, is indeed a stutter.
 - We cannot make the assumption of stutter – but we can filter low level products.

Mixtures

- When we go to court we report the strength of evidence
 - E.g.. the probability of the evidence if it has come from Mr X and an unknown individual(s) vs. the probability of the evidence if it has come from x unknown individuals.
- With numerous loci to contend with and large numbers of case-stains, it becomes very time-consuming to suggest that the operator should examine each locus separately
- Consequently, we are looking towards complete automation of the interpretation process where the output that the court is interested in is the likelihood ratio.
- Strictly we don't need nomenclature to do this, we can use raw sequence.

Why do we need a nomenclature at all?

- It is feasible to use raw strings of sequenced bases for computers
- But for database searches there must be a way to compare strings with traditional nomenclature – i.e. back-compatibility..
- Sequence strings need to be available and tied to a nomenclature, whatever it may be.
- It is easy for a computer to turn a string into a designation using a look-up table.

Stutters

- In probabilistic genotyping it is a requirement to identify stutters.
- The LUS method will do this, but we also have to consider secondary and tertiary LUS because these will also stutter, albeit at lower levels.
- Given an allele, the expectation of a given stutter is probabilistic, and the probability of stutter, and its sequence, will be dependent upon the length of the parent allele.
- Many of the 'exotic' stutters will be very low level, and the evidential value would be low – they could be filtered. But we can anticipate that they may be important with future iterations of probabilistic genotyping. Consequently we should build this into the existing strategy in order to prevent having to revisit at a later time.

Examples using D21S11

- There are 5x type 37.2 variants

	RU	LUS1	LUS2	LUS3
[TCTA] ₇ [TCTG] ₁₄ [TCTA] ₃ TCA [TCTA] ₂ TCCA TA [TCTA] ₁₂	39.2.3	14	12	7
[TCTA] ₉ [TCTG] ₁₂ [TCTA] ₃ TCA [TCTA] ₂ TCCA TA [TCTA] ₁₂	39.2.3	12	12	9
[TCTA] ₉ [TCTG] ₁₃ [TCTA] ₃ TCA [TCTA] ₂ TCCA TA [TCTA] ₁₁	39.2.3	13	11	9
[TCTA] ₁₀ [TCTG] ₁₁ [TCTA] ₃ TCA [TCTA] ₂ TCCA TA [TCTA] ₁₂	39.2.3	12	11	10
[TCTA] ₁₁ [TCTG] ₁₁ [TCTA] ₃ TCA [TCTA] ₂ TCCA TA [TCTA] ₁₁	39.2.3	11	11	11

All the same

35.2 variant

[TCTA]₅ [TCTG]₆ [TCTA]₃ TA [TCTA]₃ TCA [TCTA]₂ TCCA TA [TCTA]₁₅ TA TCTA

36.2.2.2.3;15;6;5;3;3;2

Stutter

37.2 variant								LUS notation
[TCTA] ₇ [TCTG] ₁₄ [TCTA] ₃ TCA [TCTA] ₂ TCCA TA [TCTA] ₁₂								39.2.3;14;12;7

RU	LUS1	LUS2	LUS3	
39.2.3	14	12	7	Parent allele
38.2.2	13	12	7	-1 stutter
38.2.2	14	11	7	-1 stutter
38.2.2	14	12	6	-1 stutter
37.2.2	13	12	6	-1 stutter at 2 positions

Repeat exercise for +1 stutters, -2 stutters and multiple stutters

Type 37.2 variant where, primary, secondary, tertiary LUS are the same length

$[TCTA]_{11} [TCTG]_{11} [TCTA]_3 TCA [TCTA]_2 TCCA TA [TCTA]_{11}$ **39.2.3;11;11;11**

The notation does not tell us the order in the sequence, but this could be rectified by ordering the LUS variants as they appear e.g. 39;11;11;3;.3;2;1;.1;11
But this is a bit unwieldy

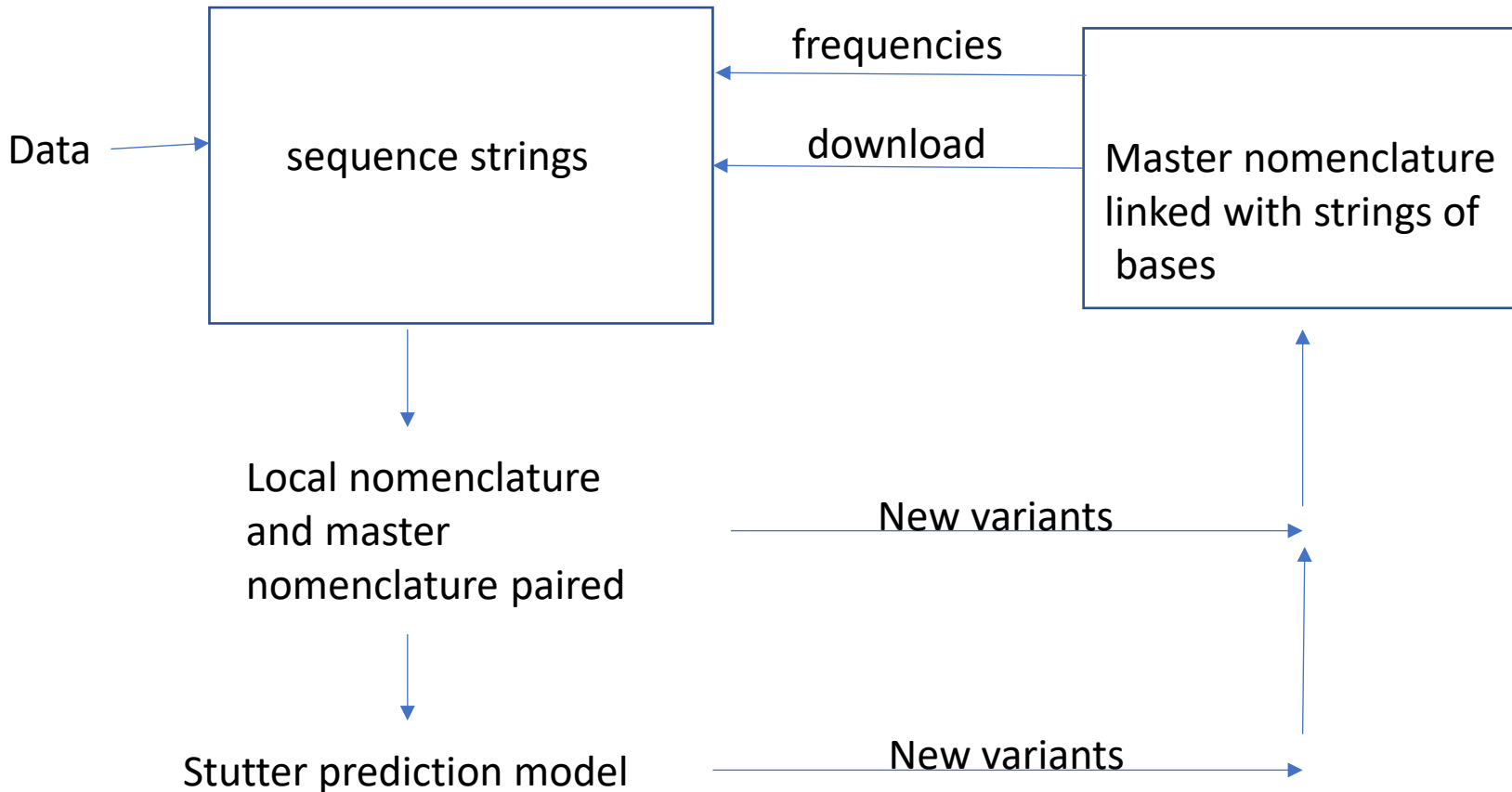
Conclusion

- Computers can take raw sequence and turn it into a local nomenclature that can be used for a particular software.
- LUS notation is probably sufficient if we use it on a per case basis – more difficult for a universal nomenclature.
- Within case a) we need to identify an allele b) we need to apply a frequency (use sequence comparison for this).
- A local nomenclature based on LUS is useful to assess stutters.
- But this can be local – per case.

Nomenclature

EuroForMix

STRidER



TOASTR

A user-friendly web app
for STR sequencing

LABCON-OWL

CORE BUSINESS

- clinical molecular diagnostics
- forensic DNA analysis

RESEARCH & DEVELOPMENT

- forensic multiplex kits for MPS
- bioinformatics solution for MPS-STR genotyping



Landscape of open-access genotyping tools



**mainstream
aligners**

- gapped alignment problem:
trade-off between runtime and accuracy
- bias due to incomplete reads
- do not report genotypes
- no stutter/noise modelling



**genome-wide
STR profilers**

- „selective“ alignment
- use whole-genome data
- legally and technically
out of scope?

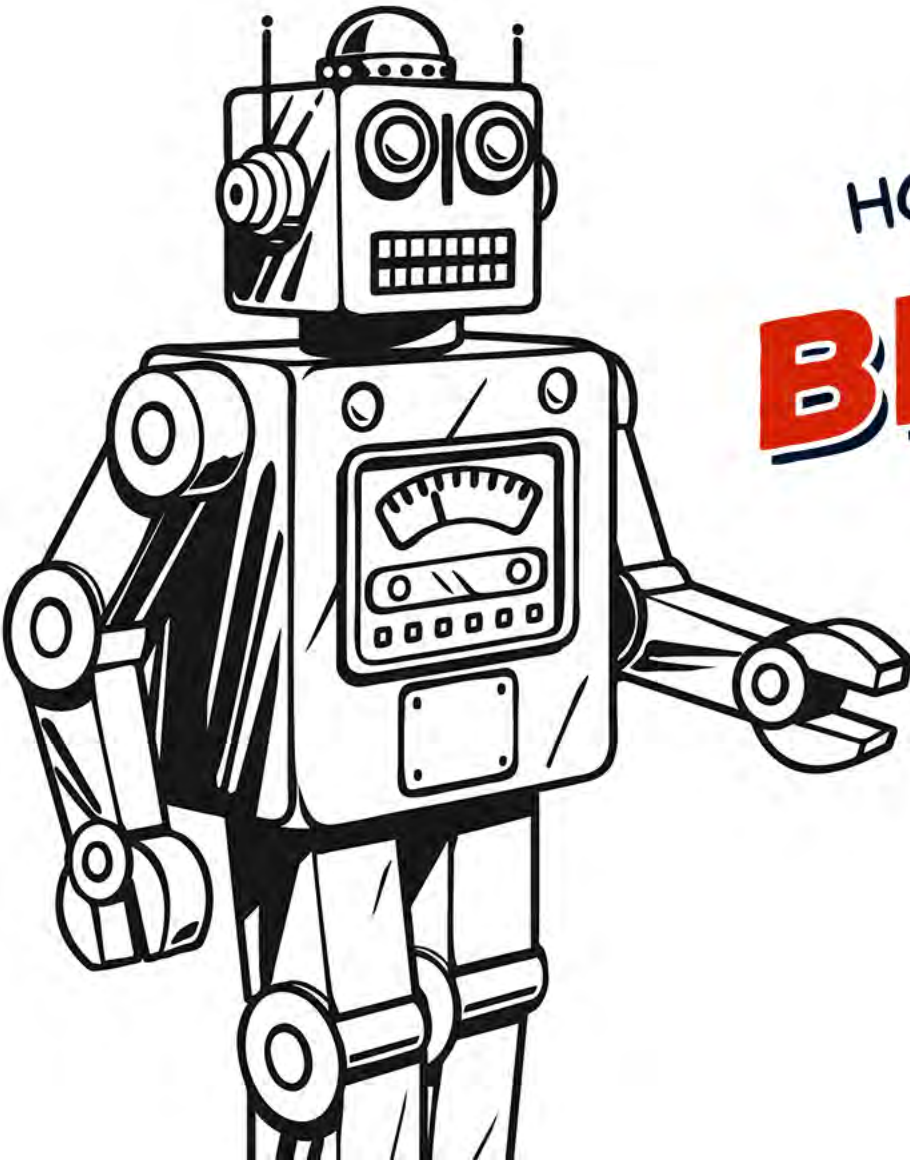


**locus-centric
STR tools**

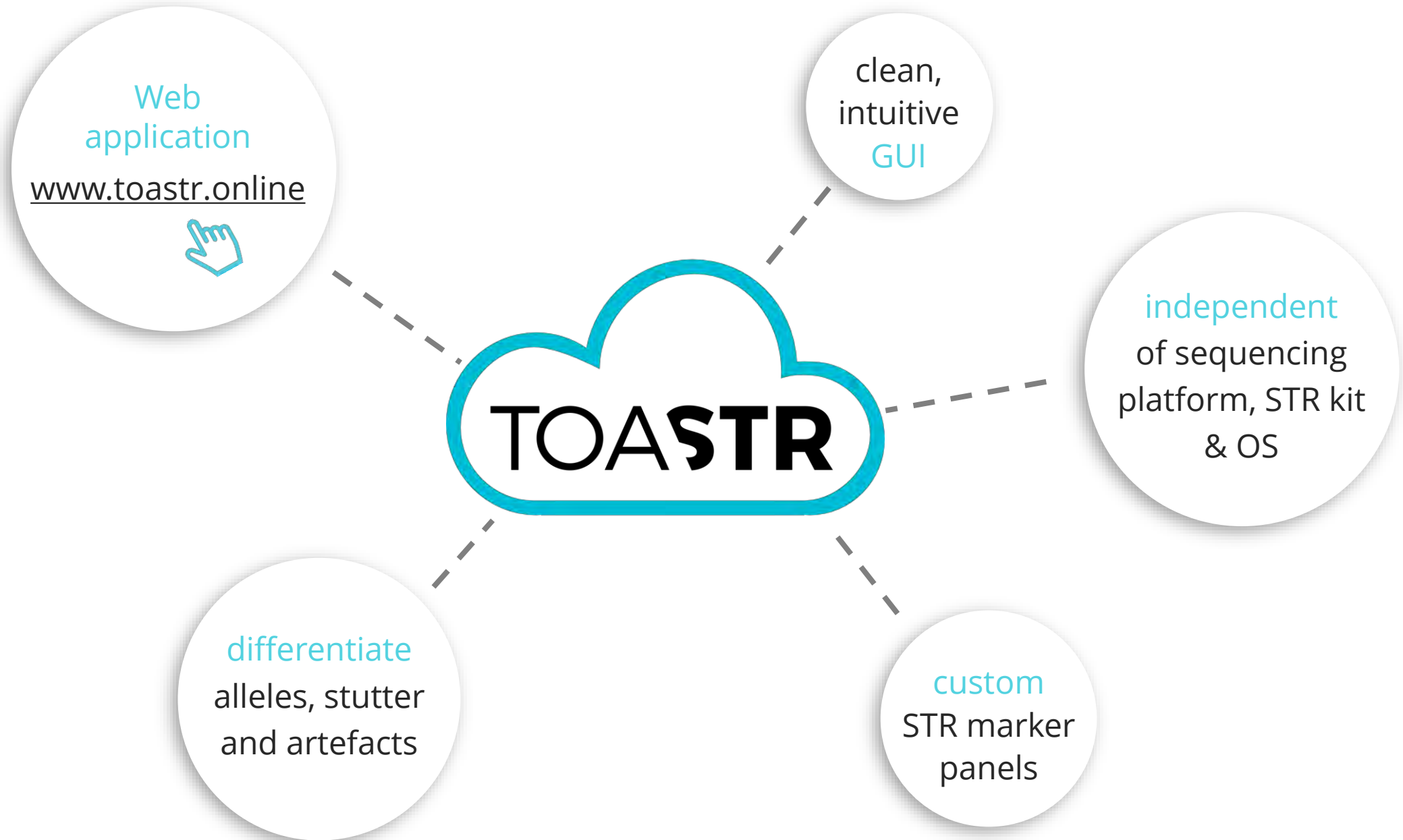
- anchor flanking sequences
- extract repeat region
- analyze length & sequence

Landscape of open-access genotyping tools

	Software	Programming language(s)	Platform	Input format	GUI	Parallel processing	Stutter modelling	Result visualization	ISFG nomenclature
genome-wide	lobSTR (Gymrek et al., 2012)	C/C++/R	Unix	FASTA, FASTQ, BAM	no	yes	yes	no	no
	RepeatSeq (Highnam et al., 2013)	C++/Python	Unix	BAM	no	yes	yes	no	no
	STR Viper (Cao et al., 2013)	Java	any	BAM, SAM	no	no	no	no	no
locus-centric	STRait Razor v3 (Woerner et al., 2017)	C++	Windows, Unix	FASTQ	no	yes	no	yes (Excel workbook)	yes (known alleles)
	MyFLq (Van Neste et al., 2015)	Python	Web, Docker, Illumina BaseSpace	FASTA, FASTQ	yes	yes	no	yes	no
	FDSTools (Hoogenboom et al., 2017)	Python	Unix	FASTA	no	no	yes	yes	no
	STRinNGS (Friis et al., 2016)	Python/R	Unix	FASTQ, BAM	no	no	no	yes	no
	SEQ Mapper (Lee et al., 2016)	.NET	Windows, Web	FASTA, FASTQ	yes	no	no	no	no
	Altius (Bailey et al., 2017)	Python	Web	FASTQ	yes	yes	no	yes	yes
	toaSTR (Ganschow et al., 2018)	Perl	Web	FASTA, FASTQ	yes	yes	yes	yes	yes



HOW CAN WE MAKE
BIOINFORMATICS
easy?



toaSTR algorithm

1

CALLING

2

CLUSTERING

3

FORMATTING

4

MODELLING

>Seq_1

```
GTGGTGTGTATTCCCTGTGCCTTTGGGGTTTTCTGTCGTTACACGCATCTCCTTTTTGGCTGTGCCTTTGGGGGCATCTCTTATACTCATGAAATC
AACAGAGGCTTGCATGTATCTGTCTGTCTGTCTGTCTGTCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTAGAGACTTTGTCTTTCCTCTG
TCTCCCCTCTTTTACTTCCGTTTTCTGTCGTTACACTTCTCTTCTCCTTTTTCTGTCGTTACACTTTTTCTGTCGTTACATTATTCCCTGTGCCT
```

>Seq_2

```
GTGGTGTGTATTCCCTGTGCCTTTGGGGTTTTCTGTCGTTACACGCATCTCCTTTTTGGCTGTGCCTTTGGGGGCATCTCTTATACTCATGAAATC
AACAGAGGCTTGCATGTATCTGTCTGTCTGTCTGTCTGTCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTAGAGACTTTGTCTTTCCTCTG
CTTTTACTTCCGTTTTCTGTCGTTACACTTCTCTTCTCCTTTTTCTGTCGTTACACTTTTTCTGTCGTTACATTATTCCCTGTGCCT
```

recognition elements

toaSTR algorithm

1

CALLING

2

CLUSTERING

3

FORMATTING

4

MODELLING

GAGGCTTGCATGTA TCTGTCTGTCTGTCTGTCTGTCTGTCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATGAGACTTTGTCTTTC
GAGGCTTGCATGTA TCTGTCTGTCTGTCTGTCTGTCTGTCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATGAGACTTTGTCTTTC
GAGGCTTGCATGTA TCTGTCTGTCTGTCTGTCTGTCTGTCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATGAGACTTTGTCTTTC

5000x **TCTGTCTGTCTGTCTGTCTGTCTGTCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTA**

GAGGCTTGCATGTA TCTGTCTGTCTGTCTGTCTGTCTGTCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATGAGACTTTGTCTTTC
GAGGCTTGCATGTA TCTGTCTGTCTGTCTGTCTGTCTGTCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATGAGACTTTGTCTTTC
GAGGCTTGCATGTA TCTGTCTGTCTGTCTGTCTGTCTGTCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATCTATGAGACTTTGTCTTTC

3000x **TCTGTCTGTCTGTCTGTCTGTCTGTCTATCTATCTATCTATCTATCTATCTATCTATCTATCTA**

→ Observations

toaSTR algorithm

1

CALLING

2

CLUSTERING

3

FORMATTING

4

MODELLING

CE	Coverage	Sequence
12	5000x	[TCTG] ₄ TCA [TCTA] ₈
10	3000x	[TCTG] ₄ TCA [TCTA] ₆

toaSTR algorithm

1

CALLING

2

CLUSTERING

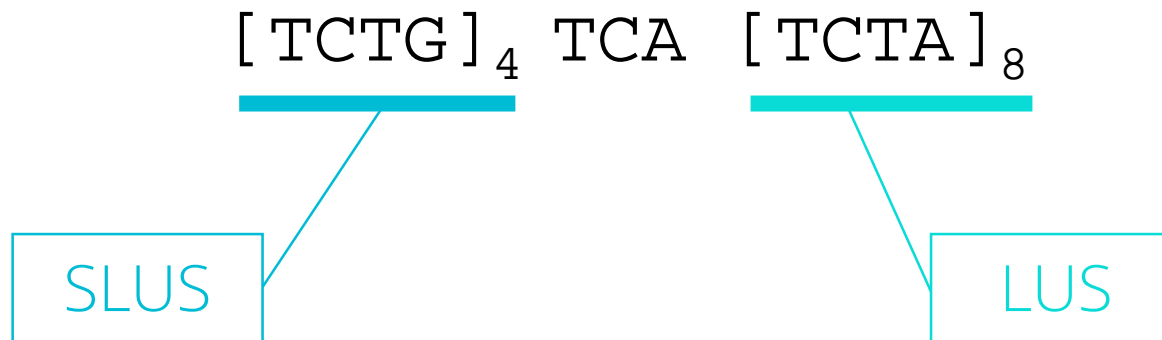
3

FORMATTING

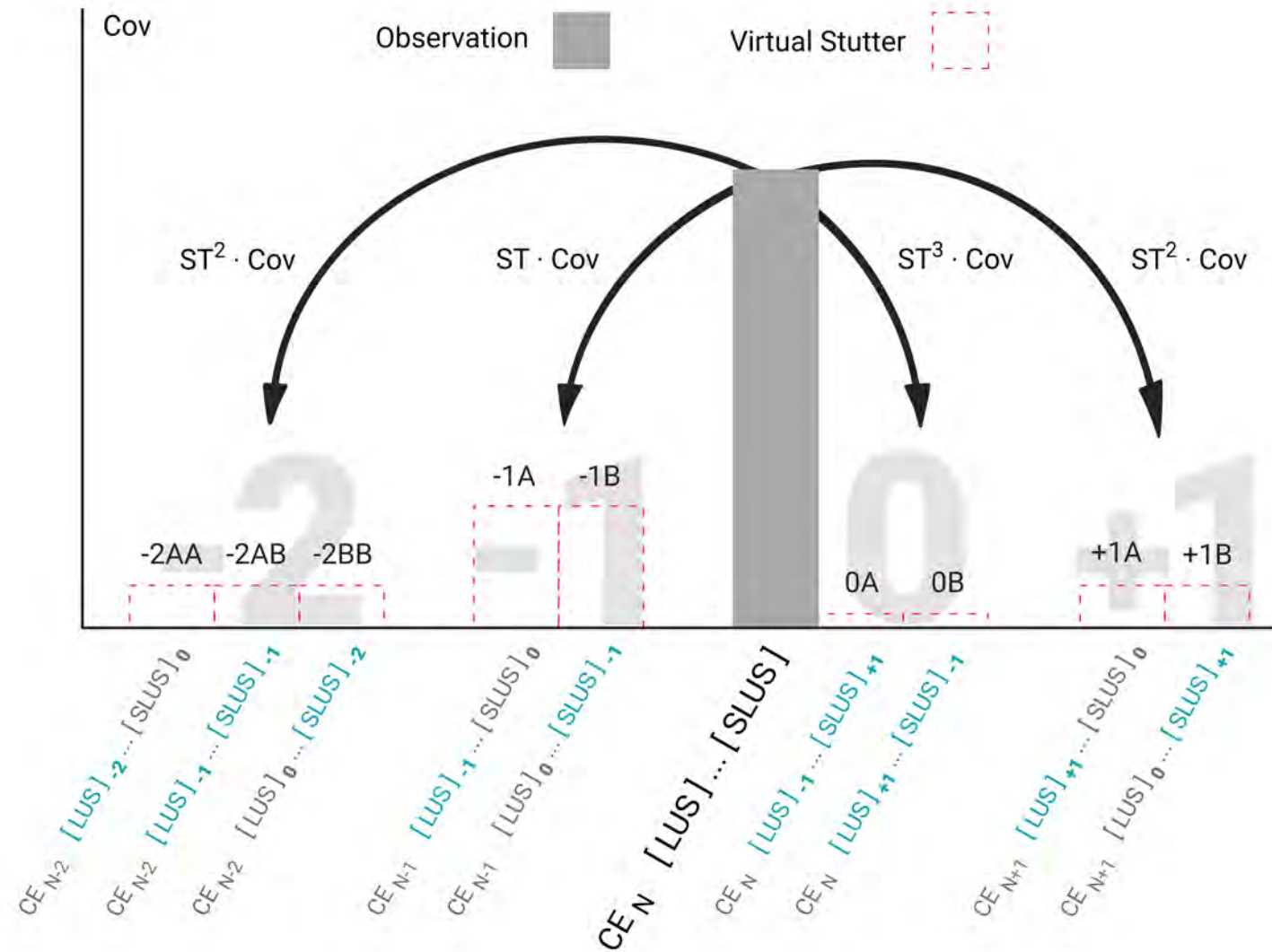
4

MODELLING

- Classification of observations: alleles, stutter, artefact
- Sequence-based stutter prediction
- Assumption: stutter occurs most likely in the LUS and SLUS



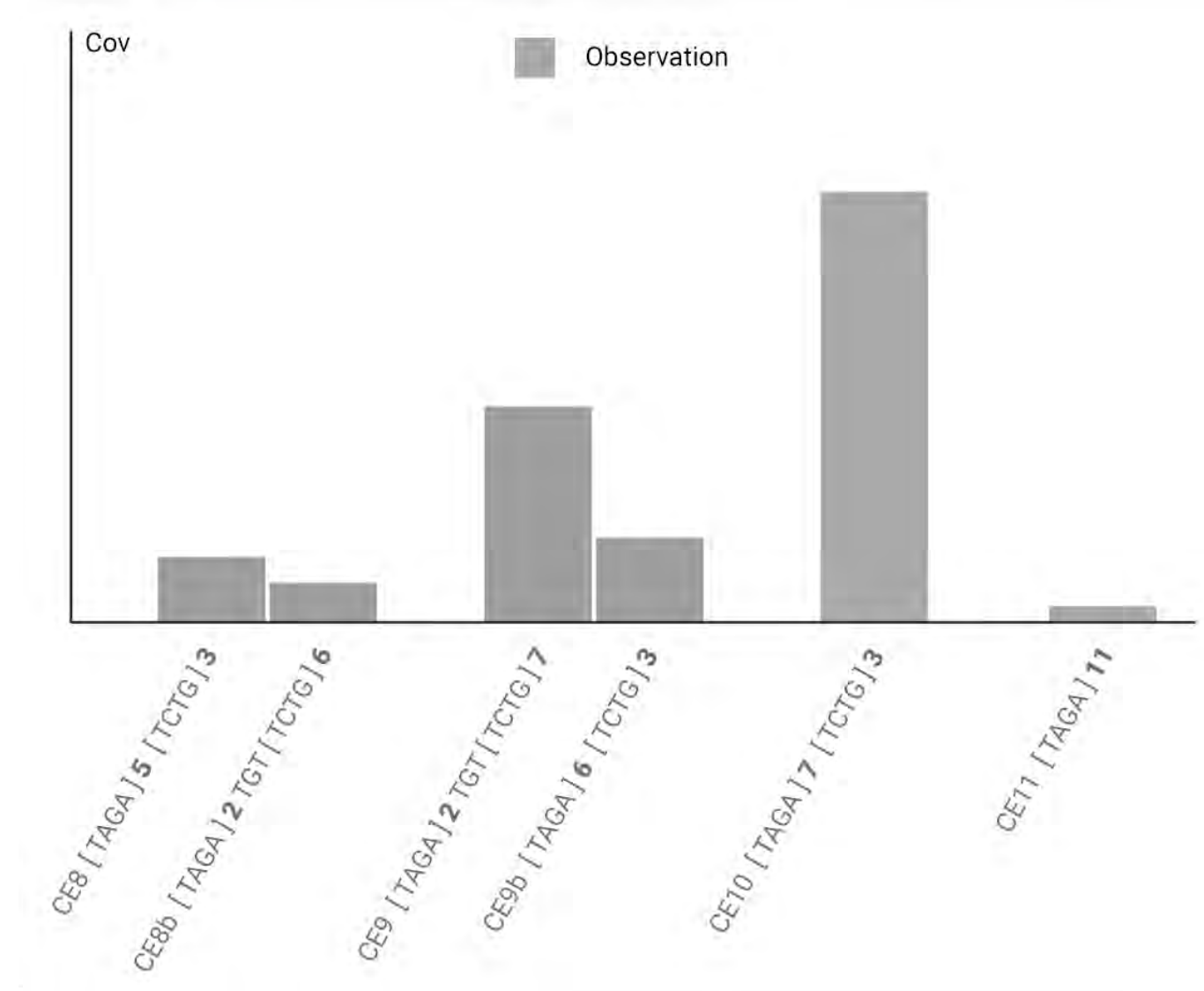
Stutter model



Classification

Stutter modelling enables automatic classification of observations:

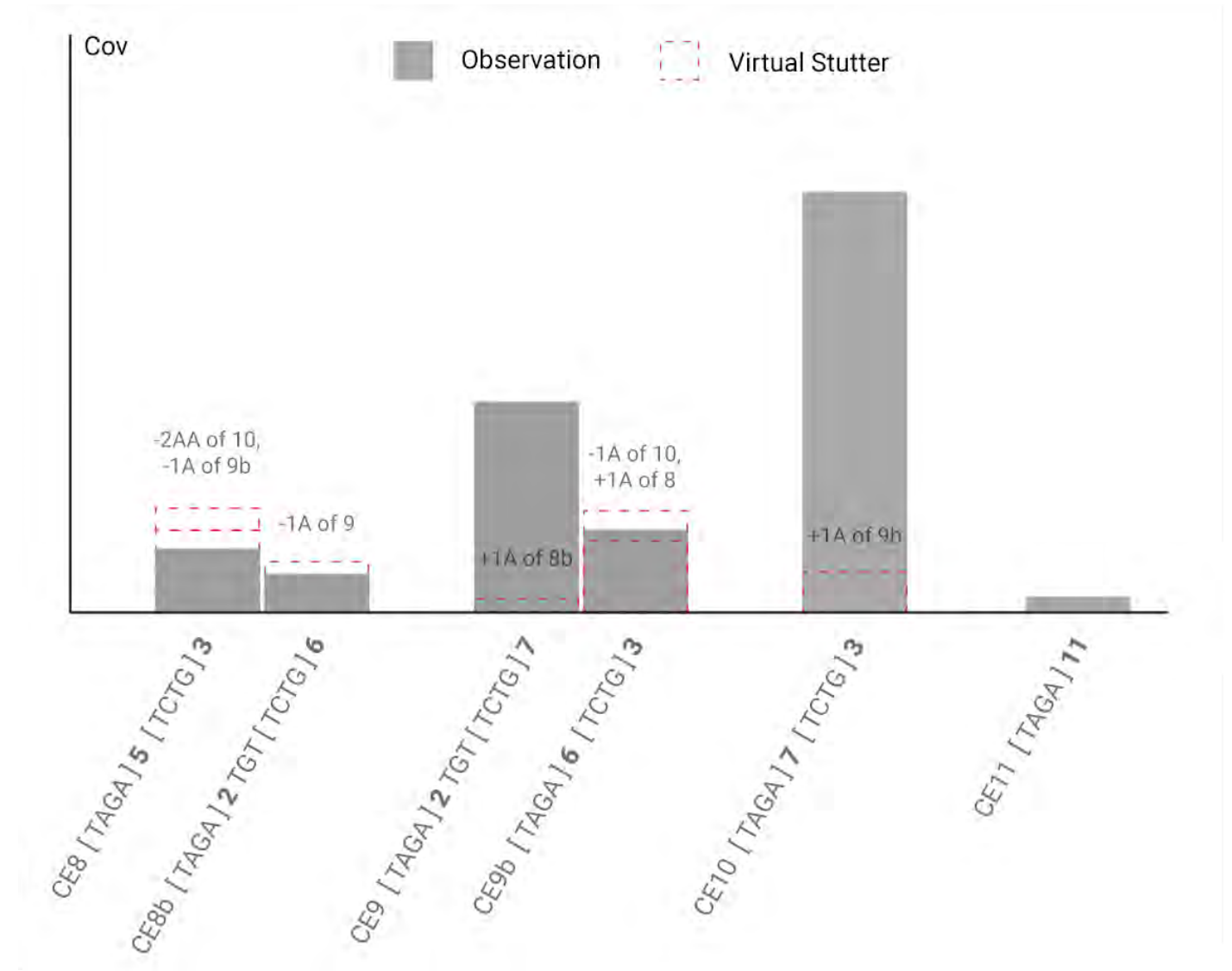
1. Calling of observations



Classification

Stutter modelling enables automatic classification of observations:

1. Calling of observations
2. Calculation of Virtual Stutter

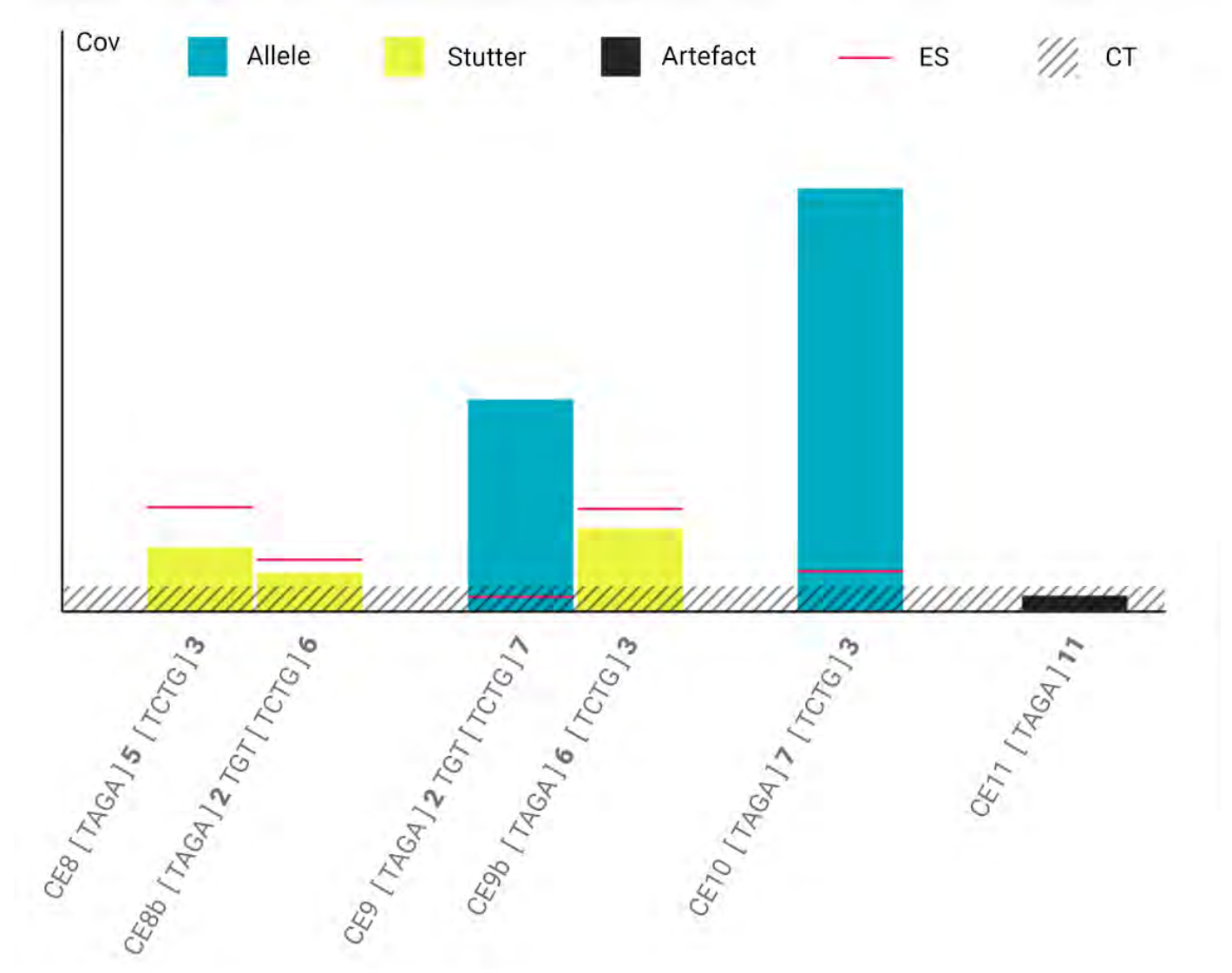


Classification

Stutter modelling enables automatic classification of observations:

1. Calling of observations
2. Calculation of Virtual Stutter
3. Classification

→ Supports the interpretation of mixed samples



Dashboard



Analyses

Start a new analysis, list finished and running analyses, show and edit results

GO



Panels

Create and view custom collections of STR markers and define stutter thresholds

VIEW



Manual


Read the comprehensive documentation of toaSTR features and algorithms

READ

Changelog

All notable changes to this project will be documented in this log.

The format is based on [Keep a Changelog](#) and this project adheres to [Semantic Versioning](#).

The  icon indicates important changes that may affect your results.

1.0.0-beta.12 (2017-10-23)

Panels / **New Panel**

[← back to panels](#)

Panel Name

0/30

▼ Autosomal

D1S1656 %

D2S1338 %

D2S441 %

TPOX %

D3S1358 %

FGA %

CSF1PO %

D5S818 %

D6S1043 %

SE33 %

D7S820 %

D8S1179 %

D10S1248 %

TH01 %

D12S391 %

VWA %

D13S317 %

PENTA E %

D16S539 %

D18S51 %

Panels / **New Panel**

Panel Name

21-plex in-house

Autosomal

D1S1656 15 %

D2S1338 20 %

D3S1358 %

FGA %

D6S1043 %

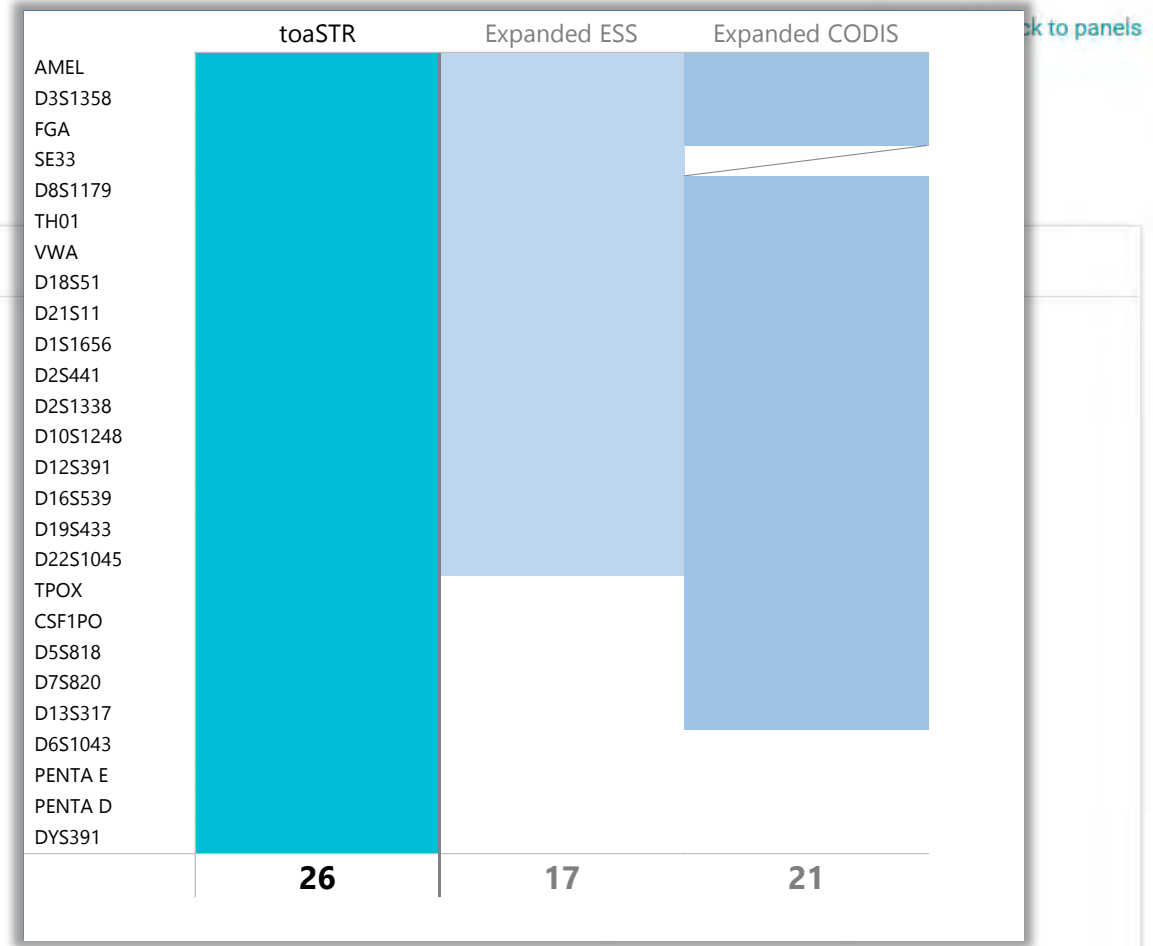
SE33 %

D10S1248 %

TH01 %

D13S317 %

PENTA E %



[Link to panels](#)

FEEDBACK

D16S539 %

D18S51 %

Dashboard



Analyses

Start a new analysis, list finished and running analyses, show and edit results

GO



Panels

Create and view custom collections of STR markers and define stutter thresholds

VIEW



Manual


Read the comprehensive documentation of toaSTR features and algorithms

READ

Changelog

All notable changes to this project will be documented in this log.

The format is based on [Keep a Changelog](#) and this project adheres to [Semantic Versioning](#).

The  icon indicates important changes that may affect your results.

1.0.0-beta.12 (2017-10-23)

Analyses / **New Analysis**

Name ⓘ

Notes ⓘ

Sample type ⓘ

Blood



Constellation ⓘ

This is a single person sample [reference]



Panel ⓘ

- Promega PowerSeq
- Illumina ForenSeq
- Thermo Early Access
- 21-plex in-house

Analyses / **New Analysis**

Name ⓘ

2800M

Notes ⓘ

Reference DNA

Sample type ⓘ

Unknown ▼

Constellation ⓘ

This is a single person sample [reference] ▼

Panel ⓘ

- Promega PowerSeq
- Illumina ForenSeq
- Thermo Early Access
- 21-plex in-house

Unknown

Constellation ⓘ

This is a single person sample [reference]

Panel ⓘ

- Promega PowerSeq
- Illumina ForenSeq
- Thermo Early Access
- 21-plex in-house

Analytical Threshold ⓘ

10 reads (default)

Calling Threshold ⓘ

2 % (default)

DATA FILE

Please refer to the manual for input recommendations.
Maximum file size 500MB. Supported formats:
FASTA (*.fasta, *.fa, *.fa.gz), FASTQ (*.fastq, *.fq, *.fq.gz)

START

INPUT READS ⓘ

250000

CALLED READS ⓘ

71568 (28%)

Quick edit

NAME

2800M

SAMPLE TYPE

Blood ▼



CONSTELLATION

This is a single person sample [ref]

NOTES

SAVE

Analyses / 2800M

 ANALYTICAL THRESHOLD 10 reads CALLING THRESHOLD 2%

Summary

 EXPORT

Promega PowerSeq Auto

SYSTEM	COVERAGE	ALLELES
D1S1656	4031	12 13
D2S1338	3049	22 25
D2S441	3220	10 14
TPOX	3584	11
D3S1358	2992	17 18
FGA	2511	20 23
CSF1PO	3662	12
D5S818	2982	12
D7S820	2514	8 11
D8S1179	3419	14 15
D10S1248	3044	13 15
TH01	3382	6 9.3
D12S391	3330	18 23
VWA	2195	16 19
D13S317	2396	9 11

INPUT READS ⓘ

250000

CALLED READS ⓘ

71568 (28%)

Quick edit

NAME
2800M

SAMPLE TYPE

Blood ▼

CONSTELLATION

This is a single person sample [ref]

NOTES

SAVE

Made with toastR
v1.0.0-beta.14

@LABCON-OWL GmbH

Research use only.
Not for commercial use.

D2S441

Called reads 3228
Stutter threshold 15 %

Results

CE	St	Cov	ES	Seq	Source
9	S	95	220	D2S441[CE9]-Chr02-GRCh38 68011947-68011994 [TCTA]9	-1A of 10, -2AA of 11,
10	A	1457	9	D2S441[CE10]-Chr02-GRCh38 68011947-68011994 [TCTA]10	+1A of 9, -1A of 11,
11	X	46	33	D2S441[CE11]-Chr02-GRCh38 68011947-68011994 [TCTA]11	+1A of 10,
13	S	58	227	D2S441[CE13]-Chr02-GRCh38 68011947-68011994 [TCTA]10 TTTA [TCTA]2	-1A of 14, -2AA of 15,
14	A	1506	6	D2S441[CE14]-Chr02-GRCh38 68011947-68011994 [TCTA]11 TTTA [TCTA]2	+1A of 13, -1A of 15,
14b	X	32	0	D2S441[CE14]-Chr02-GRCh38 68011947-68011994 [TCTA]10 TTTA TTTA [TCTA]2	
15	X	34	34	D2S441[CE15]-Chr02-GRCh38 68011947-68011994 [TCTA]12 TTTA [TCTA]2	+1A of 14,

CE = Capillary electrophoresis name, St = Status, A = Allele, S = Stutter, X = Artefact, (A) = manually excluded, Cov = Coverage, ES = Expected Stutter value, Seq = Comprehensive sequence name.

2800M

Page 4 of 25

EXPORT

D12S391

3330

18 23

VWA

2195

16 19

D13S317

2396

9 11

ALLED READS ⓘ

3330

STUTTER THRESHOLD ⓘ

25 %

AUTOSDMAL

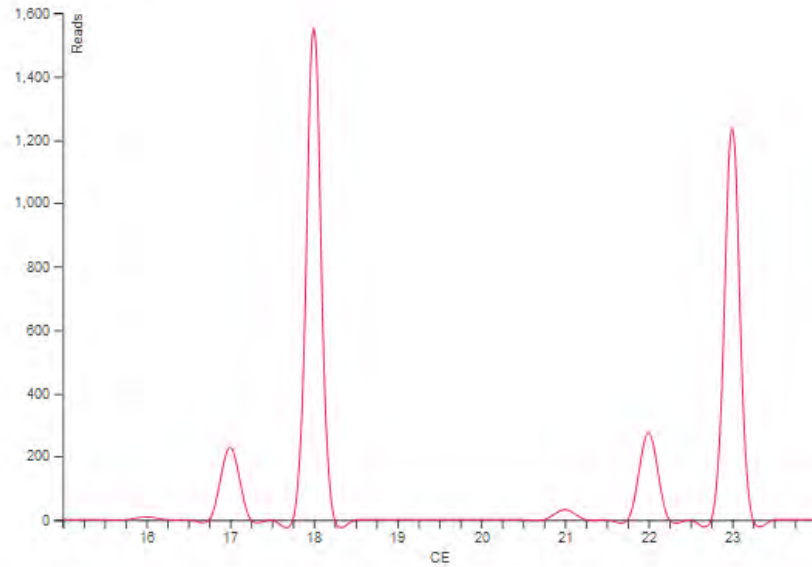
- D1S1656
- D2S1338
- D2S441
- TPOX
- D3S1358
- FGA
- CSF1PO
- D5S818
- D7S820
- D8S1179
- D10S1248
- TH01
- D12S391
- VWA
- D13S317
- PENTA E
- D16S539
- D18S51
- D19S433

2800M / D12S391

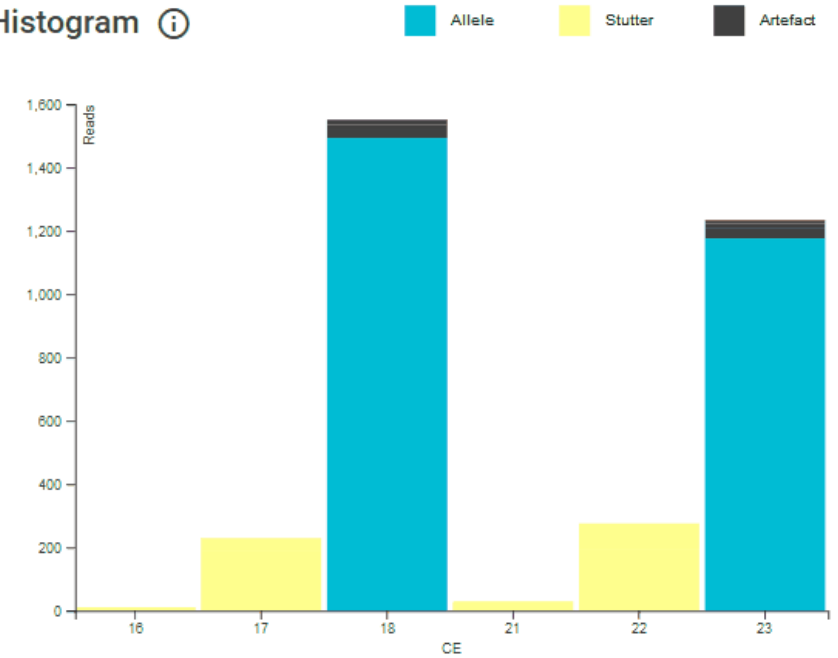
ANALYTICAL THRESHOLD 10 reads

CALLING THRESHOLD 2 % (66 reads)

Electropherogram ⓘ



Histogram ⓘ



Show alleles
 Show stutter
 Show artefacts
 Show source

REPORTED	CE	STATUS	COVERAGE	EXPECTED STUTTER	SEQUENCE
	16	Stutter	11	141	[AGAT]9 [AGAC]6 [AGAT]1
	17	Stutter	191	275	[AGAT]110 [AGAC]16 [AGAT]11

Dashboard



Analyses

Start a new analysis, list finished and running analyses, show and edit results

GO



Panels

Create and view custom collections of STR markers and define stutter thresholds

VIEW



Manual


Read the comprehensive documentation of toaSTR features and algorithms

READ

Changelog

All notable changes to this project will be documented in this log.

The format is based on [Keep a Changelog](#) and this project adheres to [Semantic Versioning](#).

The  icon indicates important changes that may affect your results.

1.0.0-beta.12 (2017-10-23)

Content

Introduction

Workflow recommendations

toaSTR algorithm

Stutter modelling

Classification of alleles, stutter and artefacts

Arranging a marker panel

Overview of analyses

Starting a new analysis

Reading the result report

Exporting results

FAQ

Documentation

Introduction

In recent studies, massively parallel sequencing (MPS) has demonstrated its potential for the forensic analysis of short tandem repeats (STRs). In addition to nominal allele lengths, MPS can discover sequence variation in isoalleles (alleles that are identical by the number of repeats) and thus increase discriminatory power over conventional capillary electrophoresis (CE). However, considering currently available software, data analysis with routine use in mind turns out to be a cumbersome process, especially for laboratories with limited bioinformatical expertise.

We developed the web application toaSTR, a user-friendly tool for STR allele calling in MPS data independent of the instrument platform or the forensic kit used. toaSTR comes up with a clean, intuitive graphical user interface and well-documented parameter settings. Users have the ability to select from a wide range of STR markers to [configure custom marker panels](#). This software supports both commercial and in-house multiplex PCR kits and various library preparation chemistries. Its sequence-based [stutter-modelling algorithm](#) automatically differentiates [biological \(iso-\)alleles from stutter and artefacts](#) to assist the interpretation of mixed samples.

toaSTR features a comprehensive [data visualization](#) with interactive diagrams and an adjustable tabular overview of sequence observations. Results are concordant with CE-based fragment analysis and can be [exported](#) for further analysis in biostatistical software or as an archivable/printable PDF document with sequence description in the ISFG-recommended nomenclature.

Citation:

Ganschow S, Wiegand P, Tiemann C (2017) toaSTR: A web-based forensic tool for the analysis of short tandem repeats in massively parallel sequencing data. <http://dx.doi.org/10.1016/j.fsigs.2017.09.034>

[▲ top](#)

Workflow recommendations

A common MPS genotyping workflow involves a multiplex PCR amplification of STR targets followed by library

Usage statistics

>60

registered users

>2,800

analyses performed

Future directions

Forensic Science International: Genetics 37 (2018) 21–28

Contents lists available at ScienceDirect

 Forensic Science International: Genetics 

journal homepage: www.elsevier.com/locate/fsigen

toaSTR: A web application for forensic STR genotyping by massively parallel sequencing 

Sebastian Ganschow^{a,*}, Janine Silvery^a, Jörn Kalinowski^b, Carsten Tiemann^a

^a LABCON-OWL Analytik, Forschung und Consulting GmbH, Siemensstr. 40, 32105 Bad Salzungen, Germany
^b Center for Biotechnology (CeBITec), Bielefeld University, Sequenz 1, 33615 Bielefeld, Germany

- Refine the stutter model
- Capture flanking variation
- Interact with databases (e.g. NOMAUT)

Acknowledgement



Supported by:



Federal Ministry
for Economic Affairs
and Energy

on the basis of a decision
by the German Bundestag



Ulm University Hospital
Institute for legal medicine
Division forensic genetics



LABCON-OWL
Analytik, Forschung und Consulting GmbH
Bad Salzuflen

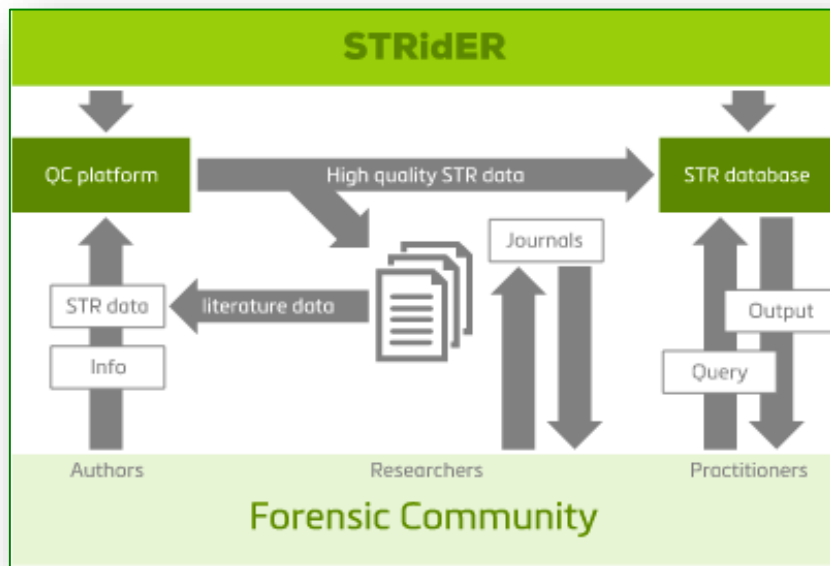
Points of discussion

- Comprehensive **benchmarking** of MPS-STR genotyping tools on standardized data sets
- Standardization of **anchor** sequences (impacts the allelic detection capability):
 - locus specificity (mismatch tolerance)
 - kit compatibility
 - CE concordance
- Guideline for **quality control** of raw sequencing data
 - Quality filtering or quality trimming
 - How to measure quality? Which QC tools and quality thresholds?
 - Implement QM in genotyping software?
- Sequencing mode: single-end, paired-end, **merging** of paired end reads

STR Databasing and Quality control

STRidER

shortened version for circulation



DBS1179								
Allele	AUSTRIA	BELGIUM	BOSNIA AND HERZEGOWINA	CZECH REPUBLIC	DENMARK	FINLAND	FRANCE	GERMANY
	222	206	171	200	200	230	208	
8	1.8018e-2	7.2816e-3	5.8480e-3	7.5000e-3	1.5000e-2	1.7391e-2	2.4039e-2	1.2840
9	1.8018e-2	1.2136e-2	8.7719e-3	5.0000e-3	1.0000e-2	8.6956e-3	9.6154e-3	1.2840
10	9.4595e-2	8.7379e-2	5.8479e-2	5.5000e-2	9.7500e-2	8.2609e-2	8.4135e-2	8.7613
11	1.0135e-1	9.7087e-2	3.2164e-2	1.0000e-1	8.0000e-2	1.3261e-1	8.8942e-2	7.7795
12	1.6216e-1	1.5049e-1	1.8713e-1	1.5250e-1	1.3000e-1	1.3261e-1	1.3462e-1	1.4199
13	2.9054e-1	3.1311e-1	3.4210e-1	3.5000e-1	3.4500e-1	3.5217e-1	3.1490e-1	3.1269
14	1.9144e-1	1.6990e-1	2.1637e-1	2.1250e-1	2.0750e-1	1.8478e-1	2.0433e-1	1.9864
15	1.0360e-1	1.2379e-1	1.1403e-1	9.7500e-2	8.5000e-2	5.6522e-2	1.0336e-1	1.1933
16	1.8018e-2	3.3981e-2	3.2164e-2	2.0000e-2	2.0000e-2	1.5217e-2	3.6058e-2	3.0211
17	2.2522e-3	4.8544e-3	2.9240e-3		5.0000e-3	1.0870e-2		6.0423

Martin Bodner, Walther Parson

STRAND WG meeting, London, 12 April 2019



Outline of presentation

1. (short) history of activities
2. **STRidER** and quality control
3. related activities

Early database considerations



Forensic Science International

Volume 131, Issues 2–3, 28 January 2003, Pages 184–196




A comparison of adjustment methods to test the robustness of an STR DNA database comprised of 24 European populations

Peter Gill ^a, Lindsey Foreman ^b, John S Buckleton ^c, Christopher M Triggs ^d, Heather Allen ^a

An aim of the European Network of Forensic Science Institutes (**ENFSI**) is to produce a **DNA database of second generation multiplex (SGM) STR profiles** that is representative of the resident cosmopolitan populations. To achieve this, data were collected from 24 different populations. All of the data were combined to form one database of 5,700 profiles from which allele proportions were calculated.

ENFSI STR database: STRbASE

- established in 2004
- **high quality** autosomal STR database
- **19** European countries, **16** autosomal STR loci
- allele **frequency tables**, also for download and import into other software
- **query function**: single and batch query
- hosted by 

Carefully enhancing STRbASE

- **STRbASE** became well-established
- **time for update:** more markers, more populations
- **new elements?**

• name changed to **STRidER**

• maintaining the tried-and-tested

STRBase: NIST STR Database

Address <http://www.csl.nist.gov/biotech/strbase/>

Short Tandem Repeat DNA
Internet DataBase

These data are intended to benefit research and application of short tandem repeat DNA markers to human identity testing. The authors are solely responsible for the information herein. [\[Purpose of Database\]](#)

This database has been accessed **30462** times since 10/02/97. (Counter courtesy www.digitals.com - see [disclaimer](#))

Created by [John M. Butler](#) and [Dennis J. Reeder](#) (NIST Biotechnology Division), with invaluable help from Christian Ruitberg and Michael Tung
Site creators' curriculum vitae available using links above.

*Partial support for the design and maintenance of this website is being provided by The National Institute of Justice through the NIST Office of Law Enforcement Stan

- o [STR101: Brief Introduction to STRs](#)
- o [STR Fact Sheets \(observed alleles and PCR product sizes\)](#) **Updated**
- o [Sequence Information \(annotated\)](#) **Updated**
- o [Multiplex STR sets](#) **Updated**
- o [Non-published Variant Allele Reports](#) ♦
- o [FBI/COJIS Core STR Loci](#)
- o [DNA Advisory Board Quality Assurance Standards](#)
- o [NIST Standard Reference Material for PCR-Based Testing](#)
- o [Chromosomal Locations](#)
- o [Mutation Rates for Common Loci](#)

STRBase published in
Nucleic Acids Research!
Ruitberg, C.M., Reeder, D.J., Butler, J.M.
(2001) 29(1): 320-322
[DOWNLOAD PDF](#)

©2002 Academic Press

How will **STRidER** be used

Survey sent out to all ENFSI DNA WG participants (2016)

Use of adjustment methods for genotype probabilities (only complete profiles)

- which formulae to present?
- which are the preferred methods?

27 labs returned answers

Survey on adjustment factors

- vast majority of users **does not use adjustments:**
 - *„not necessary“*
 - calculations using **third party software**
- majority of users prefers to receive **only unadjusted frequencies**
- exception: minor allele frequency **MAF ($5/2n$)**
- **data quality is the most important**

Formulae

Actual matching probability

$P_m = 2p_i p_j$	Heterozygotes
$P_m = p_i^2$	Homocygotes
$P_m = 2p_i - p_i^2$	Single alleles

Balding & Nichols (1994)

$P_m = \frac{2(\Theta + (1 - \Theta)p_i)(\Theta + (1 - \Theta)p_j)}{(1 + \Theta)(1 + 2\Theta)}$	Balding-Nichols heterozygotes
$P_m = \frac{(2\Theta + (1 - \Theta)p_i)(3\Theta + (1 - \Theta)p_i)}{(1 + \Theta)(1 + 2\Theta)}$	Balding-Nichols homocygotes
$P_m = \frac{(2\Theta + (1 - \Theta)p_i)(3\Theta + (1 - \Theta)p_i)}{(1 + \Theta)(1 + 2\Theta)}$	Balding-Nichols single alleles

Balding size bias correction (1995)

$P_m = \frac{2(x_i + 2)(x_j + 2)}{(n + 4)^2}$	heterozygotes
$P_m = \frac{(x_i + 2)^2}{(n + 4)^2}$	homocygotes

Confidence Intervals (NRC-Report 1996)



$Var(\ln(2p_i p_j)) \approx \frac{p_i + p_j - 4p_i p_j}{2N p_i p_j}$	Confidence interval heterozygotes
$Var(\ln(2p_i^2)) \approx \frac{2(1 - p_i)}{N p_i}$	Confidence interval homocygotes
$Var(\ln(2p_i - p_i^2)) \approx \frac{2(1 - p_i)^3}{N p_i (2 - p_i)^2}$	Confidence interval single alleles
$\Gamma = \sqrt{V_1 + V_2 + \dots + V_k}$	
$Upper\ bound = \log^{-1}(\log_{10}(P_m)) + 1.96\Gamma$	

Formulae

Actual matching probability

$P_m = 2p_i p_j$	Heterozygotes
$P_m = p_i^2$	Homocygotes
$P_m = 2p_i - p_i^2$	Single alleles

Data quality is crucial

- high quality data is the most important feature of **STRidER**
- quality control is necessary towards better STR population data
- assistance to authors, reviewers and editors
- enormous experience at  from 



Contents lists available at ScienceDirect

Forensic Science International: Genetics

journal homepage: www.elsevier.com/locate/fsig



**observations indicating
the need for STR data QC**

Inspecting close maternal relatedness: Towards better mtDNA population samples in forensic databases

Martin Bodner^a, Jodi A. Irwin^b, Michael D. Coble^{b,1}, Walther Parson^{a,*}

^aInstitute of Legal Medicine, Innsbruck Medical University, Müllerstr. 44, 6020 Innsbruck, Austria

^bArmed Forces DNA Identification Laboratory, 1413 Research Blvd, Rockville, MD 20850, USA

Gomes et al. *BMC Genomics* (2015) 16:70
DOI 10.1186/s12864-014-1201-x



RESEARCH ARTICLE

Open Access

Human settlement history between Sunda and Sahul: a focus on East Timor (Timor-Leste) and the Pleistocenic mtDNA diversity

Sibylle M Gomes^{1†}, Martin Bodner^{2†}, Luis Souto^{1,3}, Bettina Zimmermann², Gabriela Huber², Christina Strobl², Alexander W Röck², Alessandro Achilli^{4,5}, Anna Olivieri⁴, Antonio Torroni⁴, Francisco Côrte-Real⁶ and Walther Parson^{2,7*}

analysis and detection conditions [100,121-123]. After typing 15 autosomal STR loci and the amelogenin length polymorphism, pedigree construction, and likelihood ratio (LR) calculation using reported STR allele frequencies [7] [correcting the 10.2 allele frequency of D18S51 to 0.0 (L Souto, *pers. comm.*)], no donor pair revealed close maternal relatedness (*i.e.*, mother-child and sibling constellations) applying a cut-off LR of 1,000 [124,125]

Open Access

Erratum

Erratum

Tamyra R. Moretti Ph.D., Bruce Budowle Ph.D., John S. Buckleton Ph.D.

First published: 3 June 2015 [Full publication history](#)

DOI: 10.1111/1556-4029.12806 [View/save citation](#)

Cited by (CrossRef): 3 articles [Check for updates](#) | [Citation tools](#)

Am scores 2

This article corrects:

*Names of commercial manufacturers are provided for identification purposes only, and inclusion does not imply endorsement of the manufacturer or its products and services by the FBI. The views are those of the authors and do not necessarily reflect the official policy or position of the FBI or the US government.

Reference: Budowle B, Moretti TR, Baumstark AL, Defenbaugh DA, Keys KM. Population data on the thirteen CODIS core short tandem repeat loci in African Americans, US Caucasians, Hispanics, Bahamians, Jamaicans, and Trinidadians. J Forensic Sci 1999;44(6):1277-86.



View issue TOC
Volume 60, Issue 4
July 2015
Pages 1114-1116

observations indicating
the need for STR data QC

Re-typing of a widely applied population dataset after 16 years revealed a certain number of **clerical, technical, and data/sample processing errors**

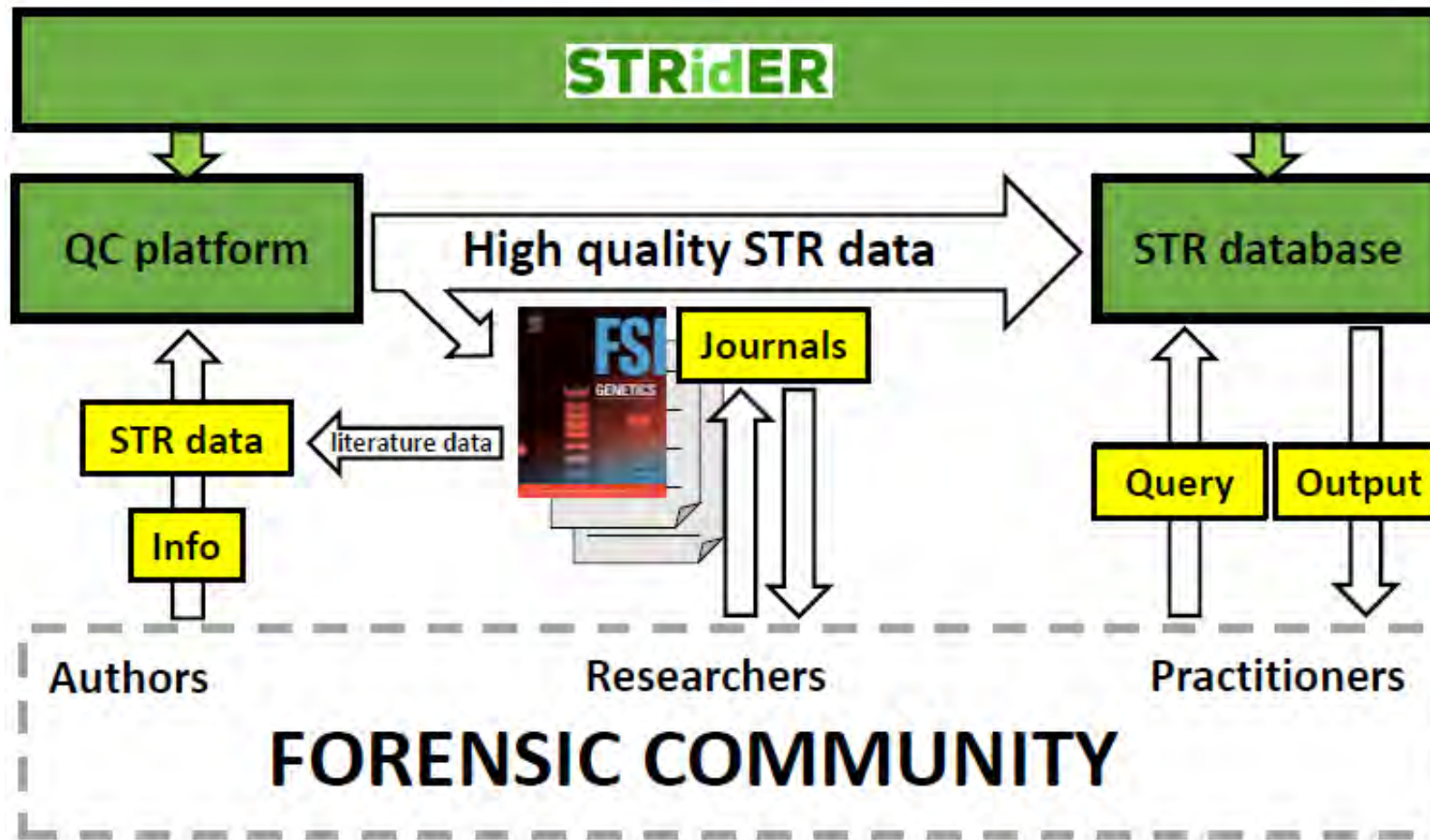
Recommendations of the DNA Commission of the International Society for Forensic Genetics (ISFG) on quality control of autosomal Short Tandem Repeat allele frequency databasing (STRidER)



Martin Bodner^a, Ingo Bastisch^b, John M. Butler^c, Rolf Fimmers^d, Peter Gill^{e,f}, Leonor Gusmão^{g,h,i}, Niels Morling^j, Christopher Phillips^k, Mechthild Prinz^l, Peter M. Schneider^m, Walther Parson^{a,n,*}

Content

- I) Positioning **STRidER** relative to other existing databases (STRbase, ALFRED, popSTR, popAffiliator, ALLST*R); **important element of QC**
- II) Rationale, concept and workflow of **QC** via **STRidER**
- III) **Benefits** to forensic and other scientific community
- IV) Transparency, traceability and protection of data
- V) Outlook: **STR sequence data** in **STRidER** (MPS)



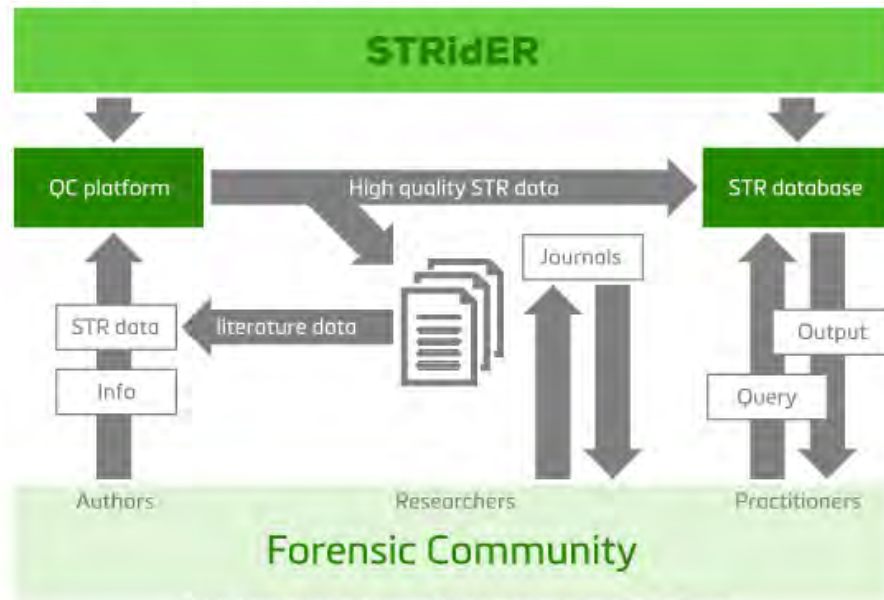
The ENFSI DNA Working Group provides an **updated version** to the previously published 'ENFSI DNA WG STR Population Database'

Welcome to STRidER!

STRidER (STRs for Identity ENFSI Reference Database) is the expanded and enhanced version of the ENFSI STRbase (2004–2016). This curated online high quality STR allele frequency population database enables scientifically reliable **STR genotype probability estimates** and provides **quality control** of autosomal STR data. A suite of software tools has been developed at the Institute of Legal Medicine, Medical University of Innsbruck to scrutinize STR population data and thus increase the quality of datasets to ensure reliable allele frequency estimates. STRidER acts as **frequency database and software platform** for the development of novel tools for STR data QC and other forensic analyses.

STRidER serves the STR community in forensics and beyond in inter-related ways:

- The high-quality autosomal STR allele frequency database can be directly queried
- Allele frequency tables of STR loci from diverse populations can be downloaded and used for third party software
- Centralized STR data quality control is offered prior to publication
- Accepted datasets will become rapidly available online and receive a unique and traceable STRidER accession number
- Allele frequencies and forensic/population genetic parameters are calculated from datasets
- Individual STR genotypes are not accessible on STRidER to comply with privacy regulations



STRidER in the field of forensic STR typing (from Bodner et al. 2016)

The concept of STRidER has been developed together with the DNA Commission of the ISFG and is outlined in Bodner M, Bastisch I, Butler JM, Fimmers R, Gill P, Gusmão L, Marling N, Phillips C, Prinz M, Schneider PM, Parson W (2016) Recommendations of the DNA Commission of the International Society for Forensic Genetics (ISFG) on a quality control of autosomal Short Tandem Repeat allele frequency databasing (STRidER); *Forensic Sci Int Gen* 24:97-102.

The STRidER online platform is work in progress. Additional datasets and features will continuously become available. To receive periodic news and stay updated about STRidER, register here for the STRidER newsletter.

<https://strider.online/>

ongoing addition of
data and functions

STRidER Newsletter enabled
(via HOME tab)

Query

The second version of STRidER holds STR loci defined in the specifications of the ENFSI DNA WG. Additional loci included in commercial kits are not included, as no high quality population data are available. Those loci are dimmed in the input form. [More...](#)

Kit check/uncheck all

- AUSTRIA
- BELGIUM
- BOSNIA AND HERZEGOWINA
- CZECH REPUBLIC
- DENMARK
- FINLAND
- FRANCE
- GERMANY
- GREECE
- HUNGARY
- IRELAND
- MONTENEGRO
- NORWAY
- POLAND
- SLOVAKIA
- SLOVENIA
- SPAIN
- SWEDEN
- SWITZERLAND

D3S1358	VWA	D16S539	CSF1PO	TPOX
<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
Y indel	D8S1179	D21S11	D18S51	DYS391
<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
D2S441	D19S433	TH01	FGA	
<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	
<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	
D22S1045	D5S818	D13S317	D7S820	SE33
<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
D10S1248	D1S1656	D12S391	D2S1338	
<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	
<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	

QUERY

according to questionnaire among ENFSI labs

- only **uncorrected AMP** is calculated
- query profile is not added to database
- **correction factors not offered** any longer
- F alleles no longer allowed
- only correction used is **MAF (5/2n)**

Submit

Clear Form

Query

The CSV file requires *commas (,)* as delimiters and *double quotes (")* as field enclosure characters.
[Download a sample CSV file.](#)

File format: CSV GeneMapper

CSV file: Keine Datei ausgewählt.

- check/uncheck all
- AUSTRIA
- BELGIUM
- BOSNIA AND HERZEGOWINA
- CZECH REPUBLIC
- DENMARK
- FINLAND
- FRANCE
- GERMANY
- GREECE
- HUNGARY
- IRELAND
- MONTENEGRO
- NORWAY
- POLAND
- SLOVAKIA
- SLOVENIA
- SPAIN
- SWEDEN
- SWITZERLAND

BATCH QUERY

Submit

Frequencies

These tables include allele frequencies and number of samples (n) from the most recent database release sorted by marker and country. In these tables, „1” represents all rare alleles shorter than the accepted allele categories. The value „99” represents all rare alleles longer than the accepted categories.

This data can be downloaded as [XML file](#).

VWA

Allele	AUSTRIA	BELGIUM	BOSNIA AND HERZEGOWINA	CZECH REPUBLIC	DENMARK	FINLAND	FRANCE	GERMANY	GREECE	HUNGARY	IRELAND	MONTENEGRO	NORWAY	POLAND	SLOVAKIA	SLOVENIA	SPAIN	SWEDEN	SWITZERLAND	Europe	Entire Database			
11	222	206	171	200	200	230	208	662	208	224	304	200	202	206	247	207	449	424	402	5172	5172			
12								7.5529e-4													9.6674e-5	9.6674e-5		
13			1.1696e-2				2.2659e-3	2.4038e-3	2.2321e-3				2.4753e-3		2.0243e-3	2.4155e-3	6.6815e-3	1.1792e-3			1.8368e-3	1.8368e-3		
14	1.0586e-1	1.0680e-1	1.1111e-1	1.0000e-1	7.0000e-2	1.3043e-1	8.6539e-2	9.7432e-2	9.3750e-2	1.1161e-1	1.1349e-1	1.4500e-1	8.6634e-2	7.7670e-2	1.1943e-1	1.0145e-1	1.1024e-1	9.4340e-2	1.0448e-1	1.0335e-1	1.0335e-1			
15	9.2342e-2	1.2136e-1	1.2573e-1	9.7500e-2	9.7500e-2	5.2174e-2	1.2740e-1	1.0347e-1	7.9327e-2	1.1384e-1	1.0197e-1	9.0000e-2	9.9010e-2	8.4951e-2	1.1943e-1	1.2077e-1	1.2361e-1	8.9623e-2	1.0697e-1	1.0296e-1	1.0296e-1			
16	1.7568e-1	1.9903e-1	2.0468e-1	1.7500e-1	2.6000e-1	1.7609e-1	2.4038e-1	2.2130e-1	1.6827e-1	2.0536e-1	2.1875e-1	1.7500e-1	2.2277e-1	2.2330e-1	1.9231e-1	1.8599e-1	2.4276e-1	2.0991e-1	2.0647e-1	2.0872e-1	2.0872e-1			
17	2.8604e-1	2.7185e-1	2.3977e-1	3.1250e-1	2.3000e-1	2.7174e-1	2.3317e-1	2.5453e-1	3.1731e-1	3.0134e-1	2.7138e-1	2.8750e-1	2.8960e-1	2.7670e-1	2.7530e-1	2.8985e-1	2.7171e-1	2.6533e-1	2.7736e-1	2.7291e-1	2.7291e-1			
18	2.5901e-1	2.0146e-1	2.1053e-1	2.2750e-1	2.4000e-1	2.0435e-1	2.1154e-1	2.2054e-1	2.4279e-1	1.7634e-1	1.9243e-1	2.1250e-1	1.9802e-1	2.4757e-1	2.0445e-1	2.1739e-1	1.7038e-1	2.4174e-1	2.0896e-1	2.1384e-1	2.1384e-1			
19	7.2072e-2	8.0097e-2	9.0643e-2	7.2500e-2	8.2500e-2	1.3696e-1	8.6539e-2	8.6103e-2	7.4519e-2	7.1429e-2	9.3750e-2	7.2500e-2	8.6634e-2	8.0097e-2	7.6923e-2	5.5556e-2	6.1247e-2	7.9009e-2	8.2090e-2	8.0917e-2	8.0917e-2			
20	9.0090e-3	1.9418e-2	5.8480e-3	1.5000e-2	1.7500e-2	2.1739e-2	1.4423e-2	1.2840e-2	1.4423e-2	1.5625e-2	8.2237e-3	1.7500e-2	1.4852e-2	9.7087e-3	1.0122e-2	2.1739e-2	1.3363e-2	1.6509e-2	1.3682e-2	1.4115e-2	1.4115e-2			
21					2.5000e-3	6.5217e-3		7.5529e-4	2.4038e-3	2.2321e-3											4.8309e-3	2.3585e-3	1.0634e-3	1.0634e-3

TH01

Allele	AUSTRIA	BELGIUM	BOSNIA AND HERZEGOWINA	CZECH REPUBLIC	DENMARK	FINLAND	FRANCE	GERMANY	GREECE	HUNGARY	IRELAND	MONTENEGRO	NORWAY	POLAND	SLOVAKIA	SLOVENIA	SPAIN	SWEDEN	SWITZERLAND	Europe	Entire Database		
5	2.2522e-3	2.4272e-3	1.71	2.00	2.00	2.30	2.08	662	208	224	304	200	202	206	247	207	454	425	402	5178	5178		
6	2.0946e-1	2.1359e-1	2.7778e-1	2.3750e-1	2.3750e-1	1.9565e-1	2.4038e-1	2.2659e-1	2.6923e-1	2.2098e-1	2.3191e-1	3.3500e-1	2.2030e-1	2.0631e-1	2.3077e-1	2.3430e-1	2.5330e-1	1.9529e-1	2.1890e-1	2.3165e-1	2.3165e-1		
7	1.2613e-1	1.8689e-1	1.3158e-1	1.6000e-1	1.5250e-1	2.1522e-1	1.5865e-1	1.4804e-1	1.2019e-1	1.6071e-1	2.0230e-1	1.0750e-1	2.3267e-1	1.5291e-1	1.4372e-1	1.3043e-1	1.7181e-1	1.9059e-1	1.7537e-1	1.6348e-1	1.6348e-1		
8	1.1486e-1	1.2136e-1	1.3158e-1	1.0000e-1	1.1250e-1	1.1522e-1	1.2019e-1	1.3897e-1	1.3462e-1	1.1384e-1	9.2105e-2	1.4500e-1	8.6634e-2	8.4951e-2	1.2348e-1	1.2560e-1	1.0573e-1	9.1765e-2	1.2189e-1	1.1529e-1	1.1529e-1		
8.3	2.2522e-3		2.9240e-3					7.5529e-4							2.0243e-3						1.2438e-3	4.8281e-4	4.8281e-4
9	1.6441e-1	1.5049e-1	1.9883e-1	1.7750e-1	1.5250e-1	1.8913e-1	1.8269e-1	1.7523e-1	2.1635e-1	2.1205e-1	1.3651e-1	1.7750e-1	1.3366e-1	2.1359e-1	2.0040e-1	1.9565e-1	1.9053e-1	1.5765e-1	1.7413e-1	1.7748e-1	1.7748e-1		
9.3	3.6712e-1	3.1796e-1	2.4561e-1	3.2250e-1	3.4000e-1	2.7826e-1	2.8365e-1	2.9305e-1	2.3077e-1	2.7679e-1	3.2895e-1	2.2250e-1	3.1436e-1	3.3010e-1	2.9555e-1	3.0193e-1	2.6982e-1	3.5647e-1	2.9353e-1	2.9973e-1	2.9973e-1		
10	1.1261e-2	7.2816e-3	1.1696e-2	2.5000e-3	5.0000e-3	6.5217e-3	1.4423e-2	1.5861e-2	2.8846e-2	1.3393e-2	8.2237e-3	1.2500e-2	9.9010e-3	1.2136e-2	4.0486e-3	1.2077e-2	8.8106e-3	7.0588e-3	1.3682e-2	1.1008e-2	1.1008e-2		
10.3	2.2522e-3																				9.6562e-5	9.6562e-5	

Formulae

Actual matching probability

$$P_m = 2p_i p_j \quad \text{Heterozygotes}$$

$$P_m = p_i^2 \quad \text{Homozygotes}$$

A minimum allele frequency of $5/2n$ [1] is used for calculations.

[1] National Research Council. (1996) The evaluation of forensic DNA evidence. National Academy Press, Washington D.C.

MPS STR Nomenclature: ISFG considerations

Forensic Science International: Genetics 22 (2016) 54–63



Contents lists available at ScienceDirect

Forensic Science International: Genetics

journal homepage: www.elsevier.com/locate/fsig

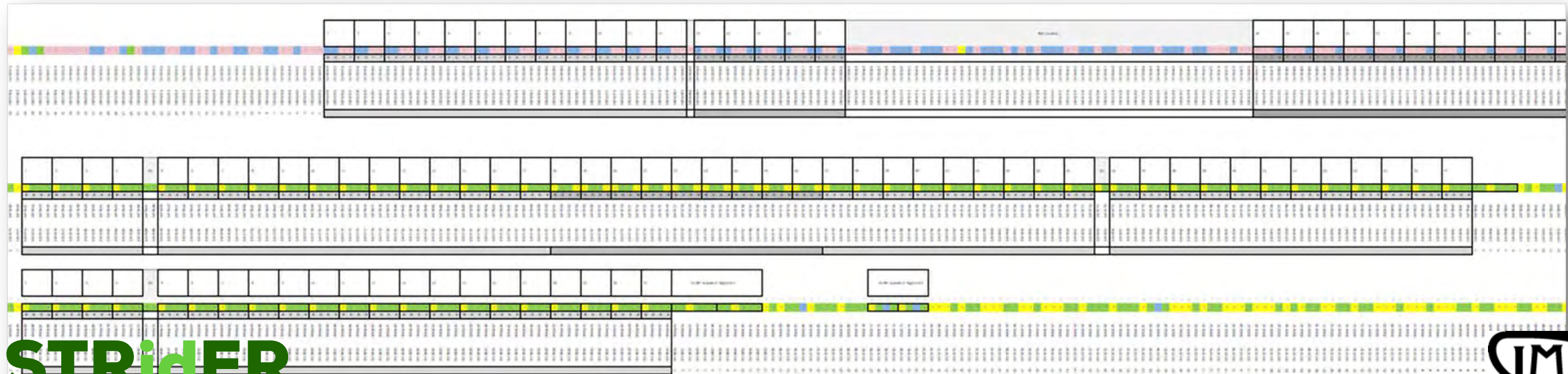


Massively parallel sequencing of forensic STRs: Considerations of the DNA commission of the International Society for Forensic Genetics (ISFG) on minimal nomenclature requirements



Walther Parson^{a,b,*}, David Ballard^c, Bruce Budowle^{d,e}, John M. Butler^f, Katherine B. Gettings^f, Peter Gill^{g,h}, Leonor Gusmão^{i,j,k}, Douglas R. Hares^l, Jodi A. Irwin^l, Jonathan L. King^d, Peter de Knijff^m, Niels Morlingⁿ, Mechthild Prinz^o, Peter M. Schneider^p, Christophe Van Neste^q, Sascha Willuweit^r, Christopher Phillips^s

ESM1 file showing STR sequence structure + flanking region



STRidER



STR Sequence Nomenclature

The 'Forensic STR Sequence Structure' file is an updated set of forensic STR sequences that was originally Parson W, Ballard D, Budowle B, Butler JM, Gettings KB, Gill P, Gusmão L, Hares DR, Irwin JA, King JL, de Kn Willuweit S, Phillips C: **Massively parallel sequencing of forensic STRs: Considerations of the DNA commission (ISFG) on minimal nomenclature requirements.** Forensic Science International Genetics 2016, 22: 54-63 (<http://www.isfg.org/Publication;Parson2016>).

The original file has been expanded, enhanced and revised as described in the publication Phillips C, Gettings KB: **"The devil's in the detail": Release of an expanded, enhanced and dynamically revised forensic STR Sequence Structure file.** The most recent version of this permanently curated and updated Forensic STR sequence structure file can be found [here](#). The updates since the last version are reported in a change log contained in the file. To receive information about the structure file and to stay updated about STRidER, [register here](#) for the STRidER newsletter.

Forensic Science International: Genetics 34 (2018) 162-169

Contents lists available at ScienceDirect

Forensic Science International: Genetics

journal homepage: www.elsevier.com/locate/fsigen

ELSEVIER

FSI GENETICS

Check for updates

"The devil's in the detail": Release of an expanded, enhanced and dynamically revised forensic STR Sequence Guide

C. Phillips^{a,*}, K. Butler Gettings^b, J.L. King^c, D. Ballard^d, M. Bodner^e, L. Borsuk^b, W. Parson^{e,f}

^a Forensic Genetics Unit, Institute of Forensic Sciences, University of Santiago de Compostela, Spain
^b National Institute of Standards and Technology, Biomolecular Measurement Division, Gaithersburg, MD, USA
^c Institute of Applied Genetics, Department of Molecular and Medical Genetics, University of North Texas Health Science Center, Fort Worth, TX, USA
^d King's Forensics, King's College London, Franklin-Wilkins Building, London, UK
^e Institute of Legal Medicine, Medical University of Innsbruck, Innsbruck, Austria
^f Forensic Science Program, The Pennsylvania State University, University Park, PA, USA, USA

Updated Forensic STR Sequence Structure Guide and change log is available on STRidER



Quality Control

STRidER provides quality control of autosomal STR data. STRidER is accepting datasets from diverse worldwide populations and forensically relevant autosomal STR markers that comply with ethical standards. Minimum requirements of journals might apply when datasets are intended for peer reviewed publication. A suite of software tools has been developed to scrutinize STR population data and thus increase the quality of datasets to ensure reliable allele frequency estimates. The board of the International Society of Forensic Genetics (ISFG) and the editors of *Forensic Science International: Genetics* [1] are 15 autosomal STR loci typed in 500 samples (for exceptional populations, the latter number can be smaller, please contact STRidER before submission). Before STR population papers are put forward to the editors for review, the authors are requested to submit the data to STRidER. After positive evaluation, the authors will be contacted with the respective STRidER accession numbers that serve as indicator of successful QC for the editors and reviewers. The necessary steps for submission of CE-based STR data to STRidER are outlined below. Please contact [STRidER](#) in case you want to submit STR sequence data.

Step 1

Prepare your **STR data file** as shown in the example file [that can be downloaded](#) and used as template. It is a tab delimited text file that can be created using standard text software or MS Excel (then, save file under .txt format). The minimum requirements for population datasets for *Forensic Science International: Genetics* [1] are 15 autosomal STR loci typed in 500 samples (for exceptional populations, the latter number can be smaller, please contact STRidER before submission).

The initial lines (identified using the "#" symbol) specify details of the dataset and origin of the samples. Line 1 must contain a description of population(s) reported (e.g., the title of the study), number of samples, geographic origin, and the number of STR loci. Line 2 must indicate the contact author's name with email address. Further text lines marked with "#" can be included for comments or description of the detailed geographic background and the appropriate metapopulation affiliation of the genotypes. Lines below these text lines list the original STR genotypes. Allele nomenclature criteria are applied as described in the "About" tab of this website. The order of loci does not matter. Alleles for the same locus have to be reported in adjacent columns. Loci names must not contain spaces. Report both alleles for homozygous loci. Use "." instead of "," for incomplete alleles, e.g. "9.3" not "9,3". Note that only complete genotypes are accepted. It is imperative that STR genotypes are reported individually and unshuffled using a unique identifier for each genotype in the dataset. The names are necessary for correspondence.

Also prepare an **accompanying STR information file** per population containing additional information on the dataset [as outlined in the example information file](#). This information might be necessary for evaluation of the dataset. Keep raw data files available for any later inquiries. Please also send the allele frequency table and forensic parameters you have calculated from the dataset (no special format required).

Step 2

Submit your files to STRidER by email (see [contact](#)). The genotype data should be submitted as a file containing the following notation: Author_country_number of samples.txt (e.g. Parson_AUT_573.txt), the accompanying file should be named Author_country_number of samples_Info.xls or .xlsx (e.g. Parson_AUT_573_Info.xls). The data will be quality checked as outlined in [2] using in-house software.

Step 3

After STRidER evaluation, communication with respect to individual genotypes may follow. Once your data passed QC you will receive the STRidER accession number(s) for your data together with allele frequencies and forensic/population genetic parameters calculated from the dataset(s). Please provide accession number(s) to the journal editor and cite STRidER [2] in your manuscript.

Step 4

Data that successfully passed QC will be uploaded onto the STRidER database.

References

[1] Gusmão L, Butler JM, Linacre A, Parson W, Roewer L, Schneider PM, Carracedo A (2017) Revised guidelines for the publication of genetic population data; *Forensic Sci Int Gen* 30:160-163

[2] Bodner M, Bastisch I, Butler JM, Fimmers R, Gill P, Gusmão L, Morling N, Phillips C, Prinz M, Schneider PM, Parson W (2016) Recommendations of the DNA Commission of the International Society for Forensic Genetics (ISFG) on quality control of autosomal Short Tandem Repeat allele frequency databasing (STRidER); *Forensic Sci Int Gen* 24:97-102.

QUALITY CONTROL
instructions

FSI:G Guidelines 2017

Forensic Science International: Genetics xxx (2017) xxx–xxx



ELSEVIER

Contents lists available at ScienceDirect

Forensic Science International: Genetics

journal homepage: www.elsevier.com/locate/fsig



**STRidER QC
is now mandatory**

Editorial

Revised guidelines for the publication of genetic population data

Since 2007, when the journal was launched, the number of submissions of manuscripts reporting population genetic data to FSI: Genetics has continuously increased. This type of data is very welcome, considering the importance of having accurate estimates

1. Quality control of population DNA databases

To improve the quality of the submitted to the journal, in 2010 [1]

Gusmao et al., 2017

To also improve the quality of the data generated from autosomal STRs, the ISFG executive board and the editors of FSI: Genetics have now invited STRidER (<http://strider.online>), a publicly available, centrally curated online allele frequency database and quality control platform for autosomal STRs [5], to logistically organize and perform quality control before autosomal STR manuscripts are put forward for review. Upon successful QC, STRidER accession numbers will be assigned to the submitted population data that serve as indicators of successful QC for the editors and reviewers. The necessary steps for submission of autosomal STR genotypes to STRidER are outlined below.

Submission of STR data for QC to STRidER

- instructions on data preparation and submission in publication and on website
- template files on website

Please provide this information in as much detail as possible. It can be helpful for evaluating your STR genotype data.
Please name this file "Author_country_number of samples_info".xls

accompanying information file for STR dataset	Parson_AUT_527.txt
submitting lab	Institute of Legal Medicine, Innsbruck, Austria
contact person name	Walther Parson
contact person e-mail	walther.parson@i-med.ac.at
lab accreditation status if any	ISO 17025
intention for manuscript publication	yes or no
manuscript running title	STR population data for Austria
intended journal	FSI Genetics
informed consent/ethics approval/data generation according to national laws	not applicable / confirmed
type of sample set	tribe / admixed urban population
exclusion/inclusion criteria	including all residents
unrelatedness	four generations, assessed by interviews
geographic origin country/region/city	Austria/Tyrol (province)
metapopulation	unspecified
subpopulations	none
total number of individuals	527
published data from overlapping samples	yes, sample 1 published in XYZ
concordance if applicable	yes
type of specimen	buccal swabs
DNA extraction method/direct amp.	Chelex-100
CE length based alleles/allele sequencing	CE length based
STR typing kit version/ homemade	insert name of commercial kit
allelic ladder	ladder included in kit or other
detection platform	ABI3100
detection chemistry (polymer)	POP5, POP6
data analysis software	insert name of software
data analysis software settings	standard settings
peak detection thresholds	insert rfu thresholds
positive control(s) - pass/fail/none	6 of 6 passed
raw data available	yes (if not: please contact STRidER before submission)
data transfer mode (manual/automated)	manual
suspected null alleles	none
observed discordances from different chemistries	n.a.
additional comments concerning the dataset	none



Quality Control

STRidER provides quality control of autosomal STR data. STRidER is accepting datasets from diverse worldwide populations and forensically relevant autosomal STR markers that comply with ethical standards. Minimum requirements of journals might apply when datasets are intended for peer-reviewed publication. A suite of software tools has been developed to scrutinize STR population data and thus increase the quality of datasets to ensure reliable allele frequency estimates. The board of the International Society of Forensic Genetics (ISFG) and the editors of Forensic Science International: Genetics invited STRidER to logically organize and perform quality control (QC) of autosomal STR population data in the course of manuscript preparations for the journal [1]. Before STR population papers are put forward to the editors for review, the authors are requested to submit the data to STRidER. After positive evaluation, the authors will be contacted with the respective STRidER accession numbers that serve as an indicator of successful QC for the editors and reviewers. The necessary steps for submission of CE-based STR data to STRidER are outlined below. Please contact STRidER in case you want to submit STR sequence data.

Step 1

Prepare your STR data file as shown in the example file that can be downloaded and used as template. It is a tab delimited text file that can be created using standard text software or MS Excel (then, save file under .txt format). The minimum requirements for population datasets for Forensic Science International: Genetics [1] are 15 autosomal STR loci typed in 500 samples (for exceptional populations, the latter number can be smaller, please contact STRidER before submission).

The initial lines (identified using the "#") specify details of the dataset and origin of the samples. Line 1 must contain a description of population(s) reported (e.g., the title of the study), number of samples, geographic origin, and the number of STR loci. Line 2 must indicate the contact author's name with email address. Further text lines marked with "#" can be included for comments or description of the detailed geographic background and the appropriate metapopulation/offspring of the genotypes. Lines below these text lines list the original STR genotypes. Allele nomenclature criteria are applied as described in the "About" tab of this website. The order of loci does not matter. Alleles for the same locus have to be reported in adjacent columns. Locus names must not contain spaces. Report both alleles for homozygous loci. Use "-" instead of "/" for incomplete alleles, e.g. "19.3" not "19.3:". Note that only complete genotypes are accepted. It is imperative that STR genotypes are reported individually and unshuffled using a unique identifier for each genotype in the dataset. The names are necessary for correspondence.

Also prepare an accompanying STR information file per population containing additional information on the dataset as outlined in the example information file. This information might be necessary for evaluation of the dataset. Keep raw data files available for any later inquiries. Please also send the allele frequency table and forensic parameters you have calculated from the dataset (no special format required).

Step 2

Submit your files to STRidER by email (see contact). The genotype data should be submitted as a file containing the following notation: Author_country_number of samples.txt (e.g. Parson_AUT_527.txt), the accompanying file should be named Author_country_number of samples_info.xls or xls (e.g. Parson_AUT_527_info.xls). The data will be quality checked as outlined in [2] using in-house software.

Step 3

After STRidER evaluation, communication with respect to individual genotypes may follow. Once your data passed QC you will receive the STRidER accession number(s) for your data together with allele frequencies and forensic/population genetic parameters calculated from the dataset(s). Please provide accession number(s) to the journal editor and cite STRidER [2] in your manuscript.

Step 4

Data that successfully passed QC will be uploaded onto the STRidER database.

References

[1] Guimó L, Butler JM, Lindors A, Parson W, Roewer L, Schneider PM, Carracedo A (2017) Revised guidelines for the publication of genetic population data. Forensic Sci Int Gen 30:160-163.

[2] Bodnar M, Bostasson I, Butler JM, Fimmers R, Gil P, Guimó L, Morling N, Phillips C, Prinz M, Schneider PM, Parson W (2016) Recommendations of the DNA Commission of the International Society for Forensic Genetics (ISFG) on quality control of autosomal Short Tandem Repeat allele frequency databasing (STRidER). Forensic Sci Int Gen 24:97-102.

```
# 527 unshuffled genotypes from Austria, Tyrol province, at 22 autosomal STR loci using kit XYZ
# submitted by Walther Parson, Institute of Legal Medicine, Medical University of Innsbruck, walther.parson@i-med.ac.at
# admixed urban population, random unrelated sample
# this file should be named "Author_country_number of samples.txt" (e.g., "Parson_AUT_527.txt")
# further lines marked with "#" can be included for comments or description of the detailed geographic background and the appropriate m
Sample ID      AMEL  AMEL  D3S1358  D3S1358  D1S1656  D1S1656  D2S441  D2S441  D10S1248  D10S1248  D13S317  D13S317  PENTA_E
sample1 X      Y      15    18      12      17.3    11.3    14      14      17      12      13      7      23      10      12
sample2 ...    ...    ...    ...    ...    ...    ...    ...    ...    ...    ...    ...    ...    ...    ...
sample3 ...    ...    ...    ...    ...    ...    ...    ...    ...    ...    ...    ...    ...    ...
sample4 ...    ...    ...    ...    ...    ...    ...    ...    ...    ...    ...    ...    ...    ...
...            ...    ...    ...    ...    ...    ...    ...    ...    ...    ...    ...    ...    ...
```



QC of STR data on STRidER

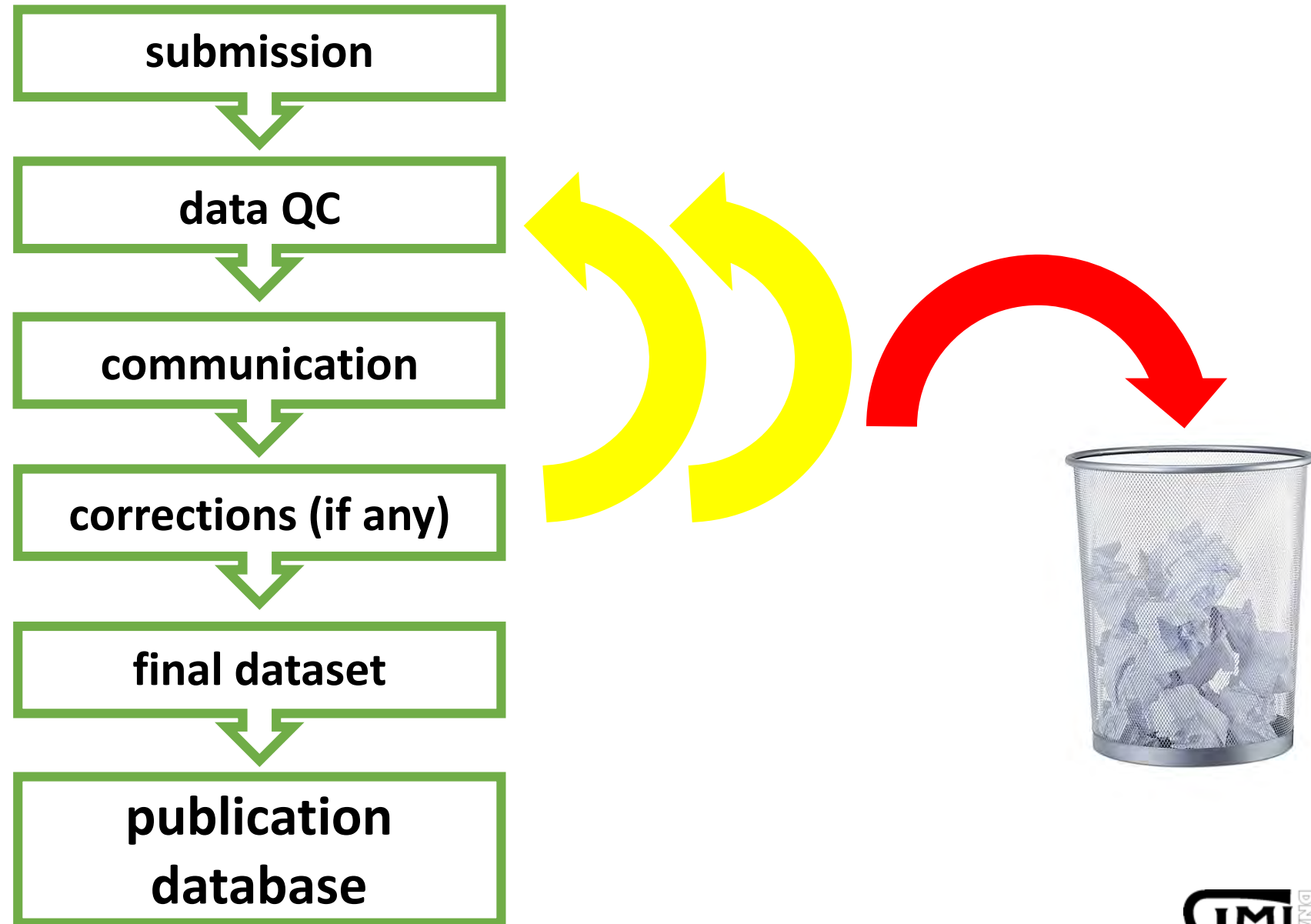
Recommendations of the DNA Commission of the International Society for Forensic Genetics (ISFG) on quality control of autosomal Short Tandem Repeat allele frequency databasing (STRidER)



Martin Bodner^a, Ingo Bastisch^b, John M. Butler^c, Rolf Fimmers^d, Peter Gill^{e,f}, Leonor Gusmão^{g,h,i}, Niels Morling^j, Christopher Phillips^k, Mechthild Prinz^l, Peter M. Schneider^m, Walther Parson^{a,n,*}

- assessment of **non-DNA** information
- assessment of reported **genotypes** and allele frequencies in dataset
- EPGs / **raw data** might be requested for QC: **true variants or errors?**
- **optimized procedure** for the detection of common data idiosyncrasies
- not an independent evaluation of all raw data
- **communication** during entire QC: discussion of all findings

Flowchart of QC on STRidER



Report on QC of STR data on **STRidER**

Information on **QC process**

Detailed **QC report**

- submitted data
- **corrections** made to dataset

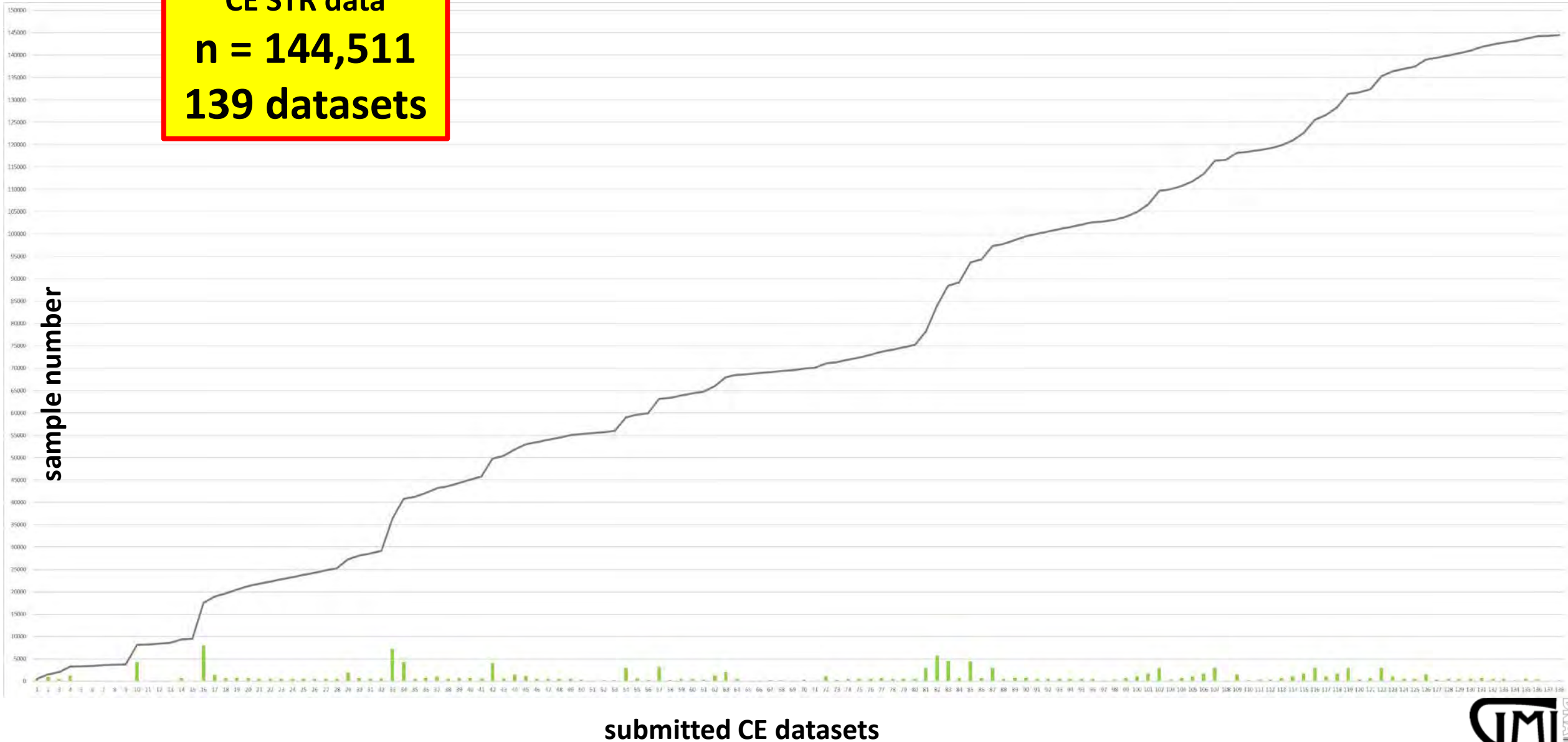
FIGURE REMOVED

STRidER accession number (STR.....) for publication

(corrected) **dataset** + allele frequency table

Submissions to STRidER Aug 2017 – March 2019

CE STR data
n = 144,511
139 datasets



Online submission tool for QC

- **form** to enter required non-genetic information (drop-down/free text)
- upload of **genotype table** in required format
- **initial plausibility tests** performed during submission process
- submission completed: e-mail notification to submitter and **STRidER**
- **external testing** ongoing **dna.bases**

**back to more basic checks
(completeness etc.)**

Related activities



dna.bases

STRidER



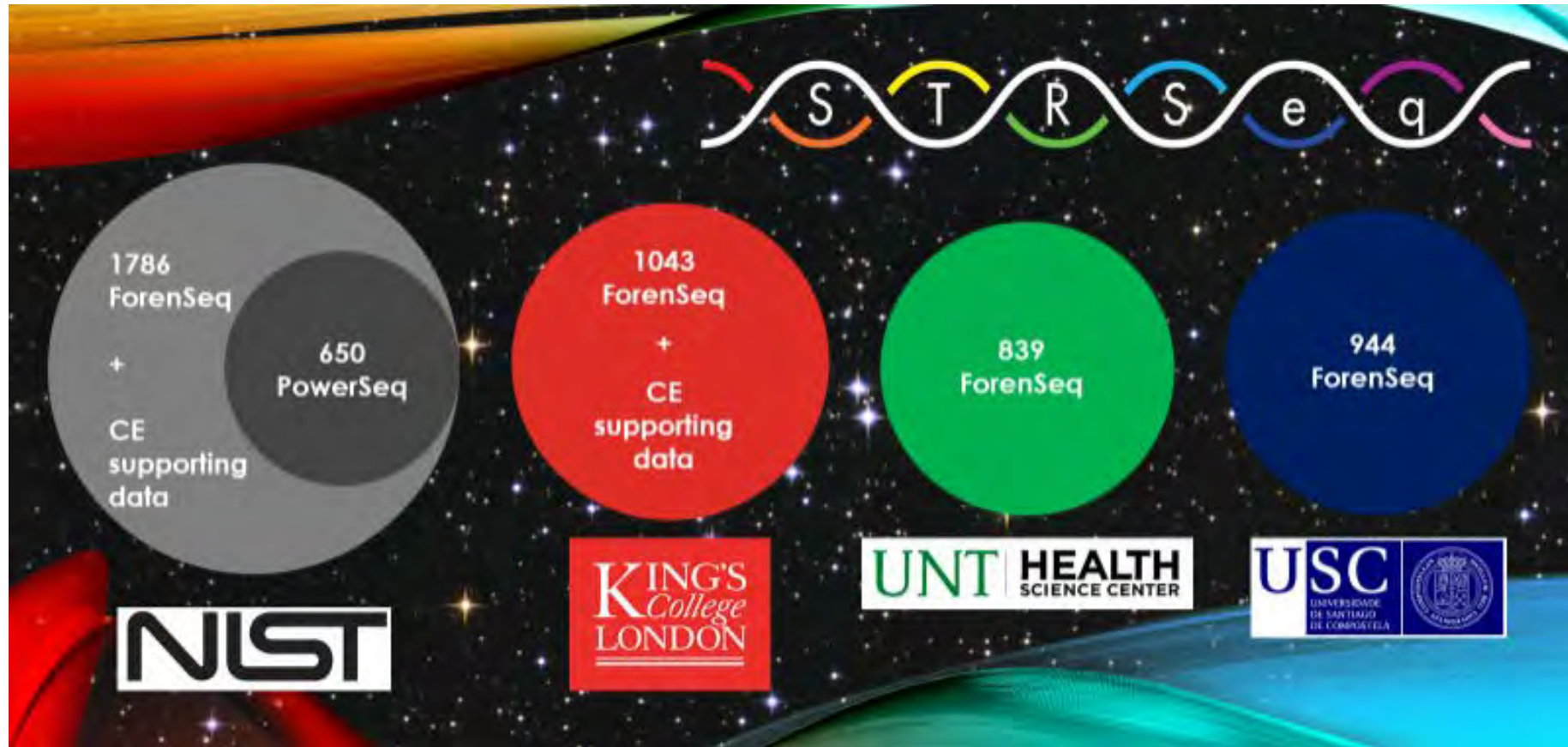
(1) STRSeq NCBI BioProject

Mission: To provide **high-confidence STR allele sequence records** with uniform annotation, facilitating exchange of information across forensic laboratories.

- collaborators with large datasets “seed” the BioProject
- NIST evaluates raw sequence data with agnostic bioinformatic pipeline
- GenBank record for **all unique sequences**
- BioProject searchable by string (BLAST), locus, allele...

STRAND *working group*

align | name | define



www.ncbi.nlm.nih.gov/bioproject/PRJNA380127



ELSEVIER

Contents lists available at [ScienceDirect](#)

Forensic Science International: Genetics

journal homepage: www.elsevier.com/locate/fsigen



Research paper

STRSeq: A catalog of sequence diversity at human identification Short Tandem Repeat loci

Katherine Butler Gettings^{a,*}, Lisa A. Borsuk^a, David Ballard^b, Martin Bodner^c, Bruce Budowle^{d,e}, Laurence Devesse^b, Jonathan King^d, Walther Parson^{c,f}, Christopher Phillips^g, Peter M. Vallone^a



NCBI BioProject-STRseq and **STRidER**
Collaboration in QC and exchange of data

(2) DNASEQEX



DNA-STR Massive Sequencing & International Information Exchange
(HOME/2014/ISFP/AG/LAWX/4000007135)



DNASEqEx - DNA-STR Massive Sequencing & International Information Exchange



UNIVERSITÄTSMEDIZIN BERLIN



2 years (2016-2018)

Project

DNASEQEX

Lutz Roewer Walther Parson Antonio Alonso Sascha Willuweit
 Bruce Budowle

Institutions: Charité Universitätsmedizin Berlin, Medizinische Universität Innsbruck, Instituto Nacional de Toxicología y Ciencias Forenses

Goal: DNASEQEX is an EU-funded project called "DNA-STR Massive Sequencing & International Information Exchange"; includes the validation of a global 23 STR & 27 Y-STR profiling system by Massively Parallel Sequencing (MPS); testing of 50 marker... [+](#)



Objectives

- Promote the implementation of MPS technology for improved STR profiling and international data exchange
→ Inter-laboratory evaluation studies
- Evaluate the impact of STR sequencing on National DNA databases (EU Prüm)
→ Alonso et al. 2017 **FSIG**
- Facilitate and standardize forensic STR sequence allele nomenclature
→ **NOMAUT** - lead Berlin



European survey on forensic applications of massively parallel sequencing

Antonio Alonso

National Institute of Toxicology and Forensic Sciences, Madrid Department, Spain

Petra Müller

Institute of Legal Medicine, Medical University of Innsbruck, Austria

Lutz Roewer, Sascha Willuweit

Institute of Legal Medicine and Forensic Sciences, Charité–Universitätsmedizin Berlin, Germany

Bruce Budowle

Walther Parson

PlumX Metrics

DOI: <https://doi.org/10.1016/j.fsigen.2017.04.017> |



Research paper

Inter-laboratory validation study of the ForenSeq™ DNA Signature Prep Kit

Steffi Köcher^{a,*,1}, Petra Müller^{b,1}, Burkhard Berger^b, Martin Bodner^b, Walther Parson^{b,c}, Lutz Roewer^a, Sascha Willuweit^a, The DNASEqEx Consortium

^a Institute of Legal Medicine and Forensic Sciences, Charité – Universitätsmedizin Berlin, Germany

^b Institute of Legal Medicine, Medical University of Innsbruck, Austria

^c Forensic Science Program, The Pennsylvania State University, PA, USA



Research paper

Systematic evaluation of the early access applied biosystems precision ID Globalfiler mixture ID and Globalfiler NGS STR panels for the ion S5 system

Petra Müller^a, Antonio Alonso^b, Pedro A. Barrio^b, Burkhard Berger^a, Martin Bodner^a, Pablo Martin^b, Walther Parson^{a,c,*}, The DNASEQEX Consortium

^a Institute of Legal Medicine, Medical University of Innsbruck, Innsbruck, Austria

^b National Institute of Toxicology and Forensic Sciences, Madrid Department, Las Rozas de Madrid, Spain

^c Forensic Science Program, The Pennsylvania State University, PA, USA



ELECTROPHORESIS

Review

Current state-of-art of STR sequencing in forensic genetics

Antonio Alonso , Pedro Alberto Barrio, Petra Müller, Steffi Köcher, Burkhard Berger, Pablo Martin, Martin Bodner, Sascha Willuweit, Walther Parson, Lutz Roewer, Bruce Budowle



(3) dna.bases

MONOPOLY 2016 - STEFA - WP G7

Empowering forensic genetic DNA databases for the interpretation of next generation sequencing profiles (**dna.bases**)

STRidER & EmPOP

Jan 2018 - Dec 2019

Sequence alignments

Increase sample size

Increase markers/regions

Further develop QC tools

User-friendly access



STRidER

dna.bases

EMPOP



Objectives of the work package

- 1) Develop a **new database engine concept** that enables event-based query of unaligned nucleotide sequence strings. This is relevant for the determination of DNA profiles with multiple sequence differences that are caused by single mutational events.

Collaborate to use and further develop existing STR sequence alignment tools

Objectives of the work package

- 1) Develop a **new database engine concept** that enables event-based query of unaligned nucleotide sequence strings. This is relevant for the determination of DNA profiles with multiple sequence differences that are caused by single mutational events.
- 2) Extension of the number of **STR markers** and quality-controlled DNA profiles from global populations, including full mtGenome data and **STR sequences** generated by NGS.

Objectives of the work package

- 1) Develop a **new database engine concept** that enables event-based query of unaligned nucleotide sequence strings. This is relevant for the determination of DNA profiles with multiple sequence differences that are caused by single mutational events.
- 2) Extension of the number of STR markers and quality-controlled DNA profiles from global populations, including full mtGenome data and STR sequences generated by NGS.
- 3) Update of existing and development of new tools for **quality control** of relevant genetic data. **Development and deployment of QC tools in database.**

Objectives of the work package

- 1) Develop a **new database engine concept** that enables event-based query of unaligned nucleotide sequence strings. This is relevant for the determination of DNA profiles with multiple sequence differences that are caused by single mutational events.
- 2) Extension of the number of STR markers and quality-controlled DNA profiles from global populations, including full mtGenome data and STR sequences generated by NGS.
- 3) Update of existing and development of new tools for **quality control** of relevant genetic data.
- 4) Enable **user-friendly access** to the databases from diverse platforms including mobile devices and establish links to existing software for data interpretation of DNA evidence including Forensim and LRmix studio.

(4) SeqForSTRs

Sequencing of Forensic STRs

3 years (2017-2020)



Tasks

Population study STRs as data basis to evaluate MPS - STRidER

Concordance between CE and MPS

Validation Study

Cost - benefit study

Feedback to European Laboratories (ENFSI)

Consortium

Federal crime lab **Wiesbaden** (Consortial leader)

State crime labs: Berlin, Bavaria, Rheinland-Pfalz

Legal Medicine labs: Berlin, Cologne, Innsbruck, Ulm



Acknowledgements



Richard Scheithauer, Daniela Niederwieser,
Martin Pircher, Stefan Troger



Monopoly 2010



Monopoly 2014

vXWeb

ISFG Commission on MPS of STRs

ISFG Commission on STRidER



dna.bases

ENFSI laboratories



Peter Gill



Co-funded by the Internal
Security Fund of the
European Union

Thank you to all attendees!

STRAND *working group*

align | name | define

David Ballard
Martin Bodner
Lisa Borsuk
Katherine Gettings
Jonathan King
Walther Parson
Christopher Phillips

