


### Underlying sequence variation within STRs: considerations for nomenclature, storage, searching, and reporting

Peter M. Vallone, Ph.D.  
Leader, Applied Genetics Group

*Workshop: Analyzing and Utilizing Data from Next-Generation Sequencers in the Forensic Genomics Era*  
26<sup>th</sup> annual International Symposium on Human Identification  
October 12, 2015  
Grapevine, TX



### Outline for Today

- STR sequence diversity
- Investigating STR sequence diversity
  - Sequencing of 183 samples for 22 autosomal STRs
  - Informatics
  - Sequence diversity (within the STR motif)
  - Flanking region sequence diversity
- NIST resources for NGS researchers
- Nomenclature: thoughts and discussion


### NGS Activities at NIST - Overview

Samples	Markers	Platforms	Assays
SRM 2391c	STR	MiSeq/FGx	PowerSeq Auto System
NIST Population Samples	miDNA	PGM	ForenSeq Panel
SRM 2392/2392I (mtDNA)	SNP		AmpliSeq Identity SNP Panel
GWU Samples			AmpliSeq Ancestry SNP Panel

### NGS Activities at NIST - Overview

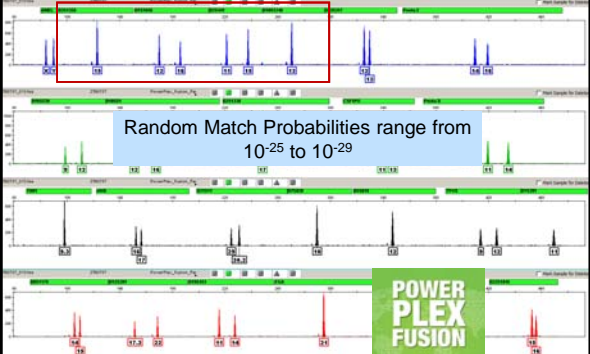
Samples	Markers	Platforms	Assays
SRM 2391c	STR	MiSeq/FGx	PowerSeq Auto System
NIST Population Samples	miDNA	PGM	ForenSeq Panel
SRM 2392/2392I (mtDNA)	SNP		AmpliSeq Identity SNP Panel
GWU Samples			AmpliSeq Ancestry SNP Panel

### NGS has potential for finer resolution of STR amplicons not detectable by CE-length based methods




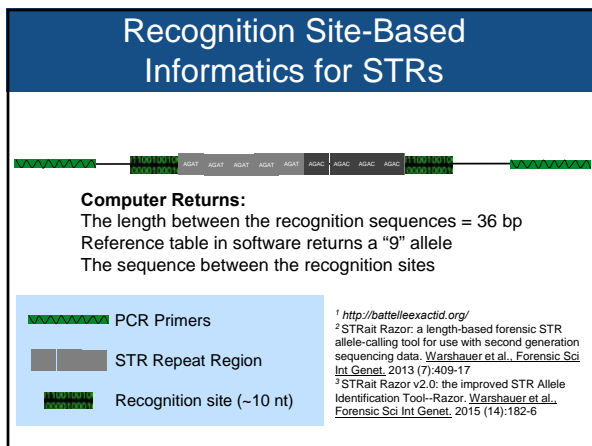
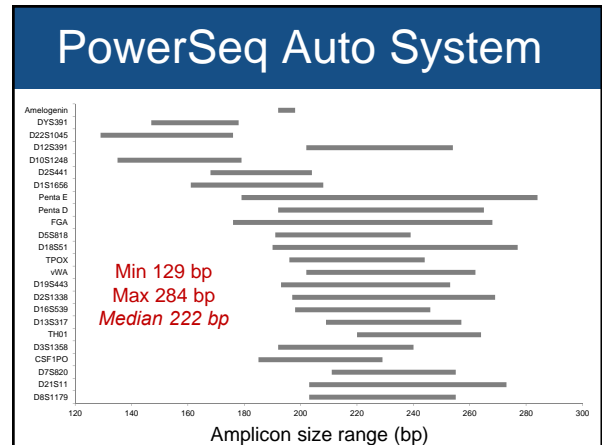
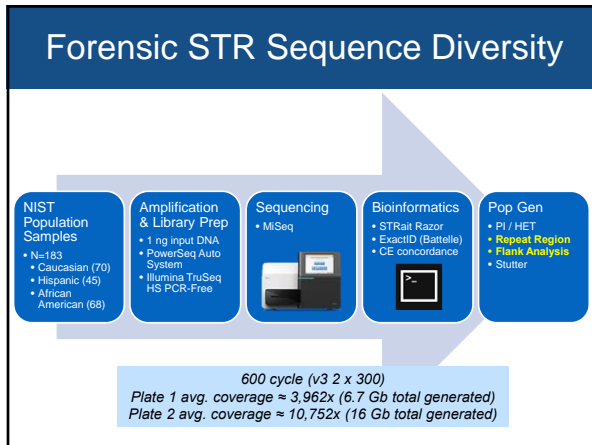
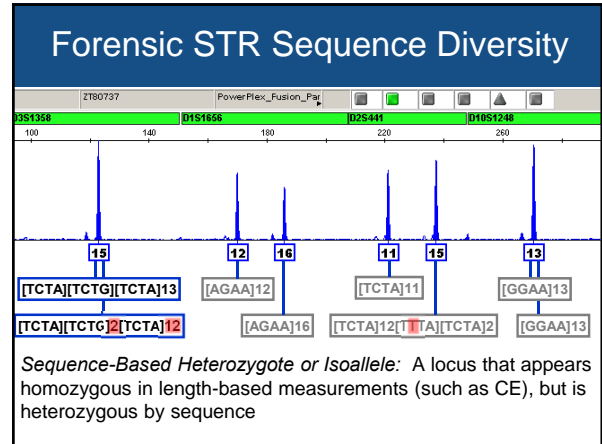
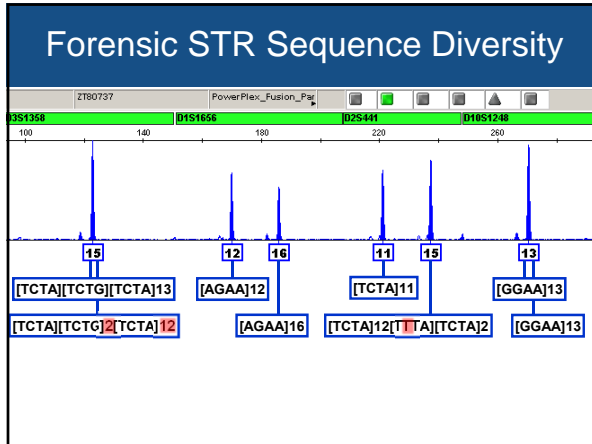
- Additional STR alleles
- Flanking region SNPs and InDels
- *Resolve minor contributor peaks from stutter*
- *Resolve homozygous by length peaks*

### Forensic STR Sequence Diversity



Random Match Probabilities range from  $10^{-25}$  to  $10^{-29}$





### Forensic STR Sequence Diversity

CE (length-based genotype) concordance check results

24 loci x 183 samples = 4392 loci evaluated

➤ 99% concordance with CE data

Why were some discordances observed between the CE and NGS measurements...informatics



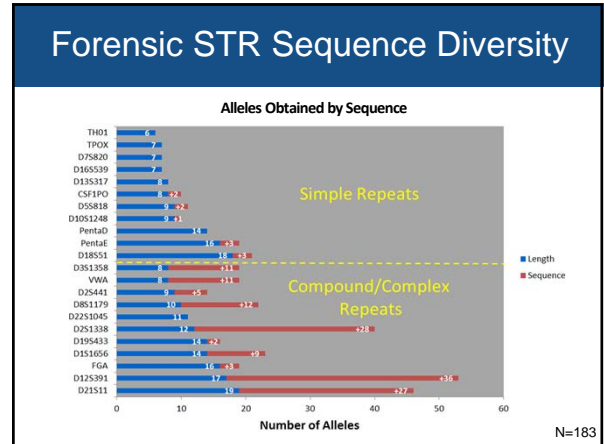
### Forensic STR Sequence Diversity

Additional Alleles by Sequence

CSF1PO		
7	[AGAT]7	AGAT AGAT AGAT AGAT AGAT AGAT AGAT
8	[AGAT]8	AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT
9	[AGAT]9	AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT
10	[AGAT]10	AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT
10	[AGG][AGAT]9	AGG AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT
11	[AGAT]11	AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT
11	[AGAT]3AGG[AGAT]7	AGAT AGAT AGAT AGG AGAT AGAT AGAT AGAT AGAT AGAT
12	[AGAT]12	AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT
13	[AGAT]13	AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT
14	[AGAT]14	AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT

8 alleles by length → 10 alleles by sequence

N=183



### Manuscript submitted to FSIG

Title:  
 Sequence variation of 22 autosomal STR loci detected by next generation sequencing

Authors and Affiliatic  
 Katherine Butler Getti  
 Baker<sup>1</sup>, Brian A. You  
<sup>1</sup>U. S. National Institut

10	[TCTA]10
10	[TCTA]8[TCTG][TCTA]
11	[TCTA]11
11	[TCTA]9[TCTG][TCTA]
11.3	[TCTA]3[TCA][TCTA]8
11.3	[TCTA]4[TCA][TCTA]7
12	[TCTA]10[TCTG][TCTA]
12	[TCTA]12
12.3	[TCTA]4[TCA][TCTA]8
13	[TCTA]10[TTTA][TCTA]2
13	[TCTA]13
14	[TCTA]11[TTTA][TCTA]2
14.3	[TCTA]3[TCA][TCTA]11
15	[TCTA]12[TTTA][TCTA]2

• Repeat region sequences

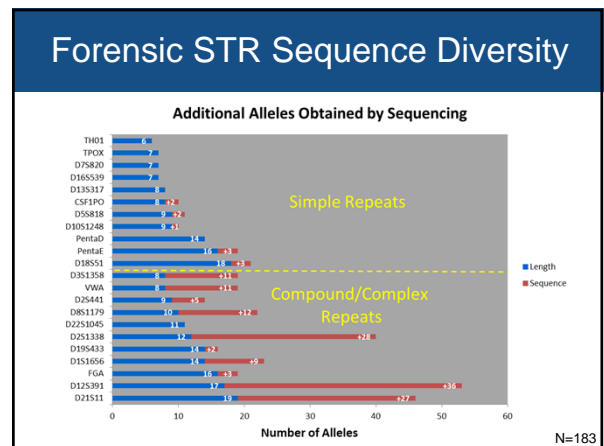
### Flanking Region Variation

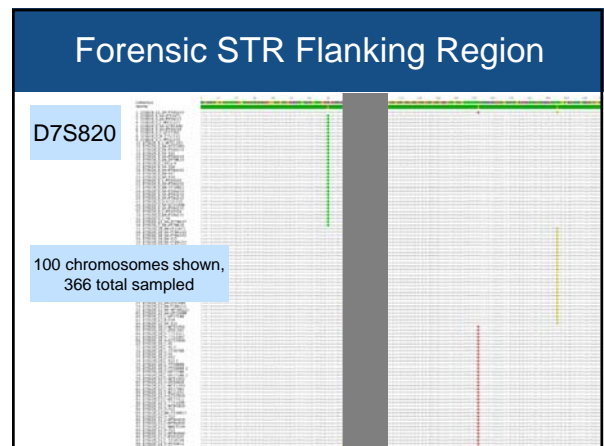
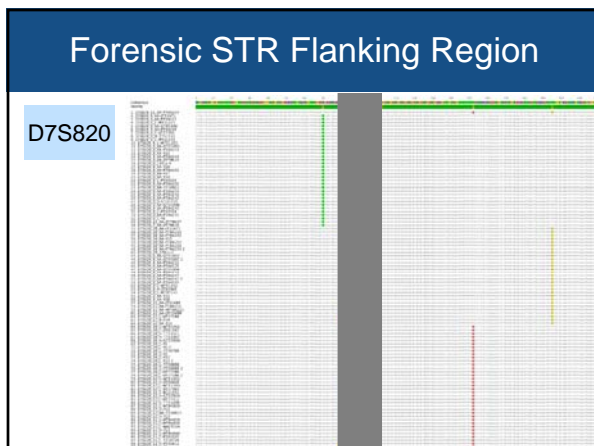
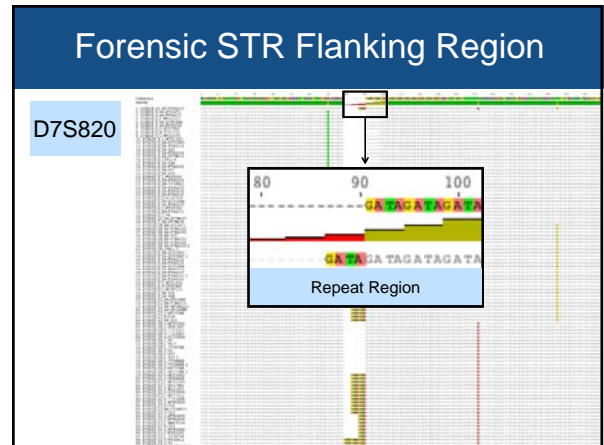
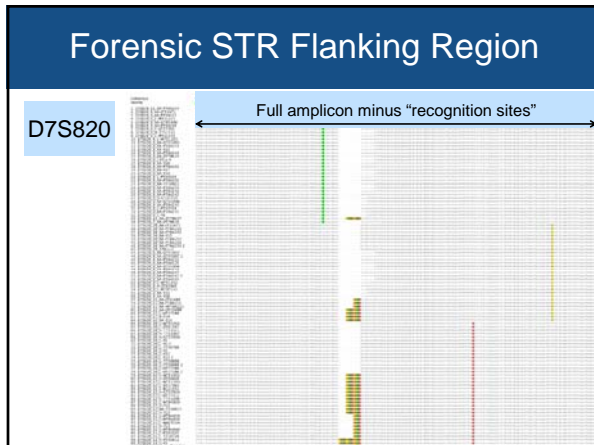
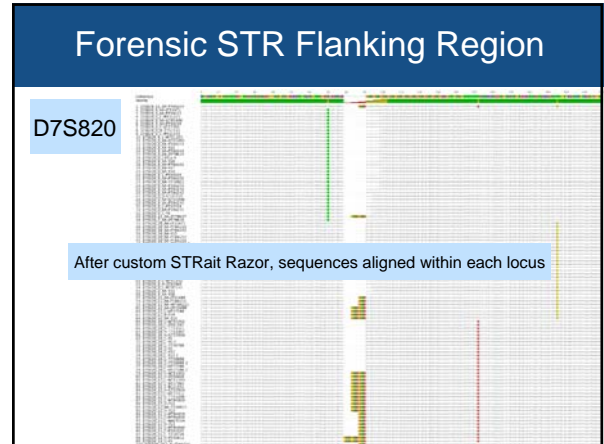
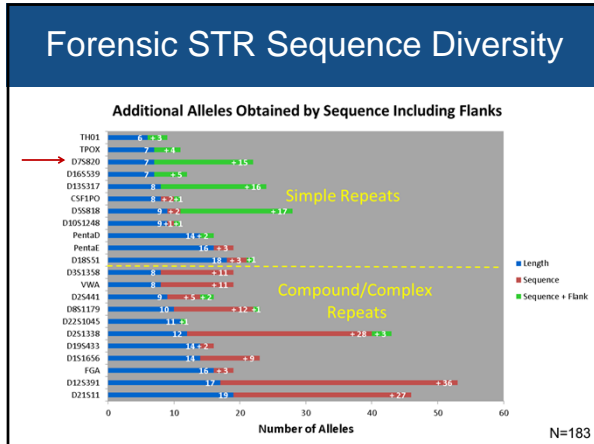
Moving out past the repeat regions...

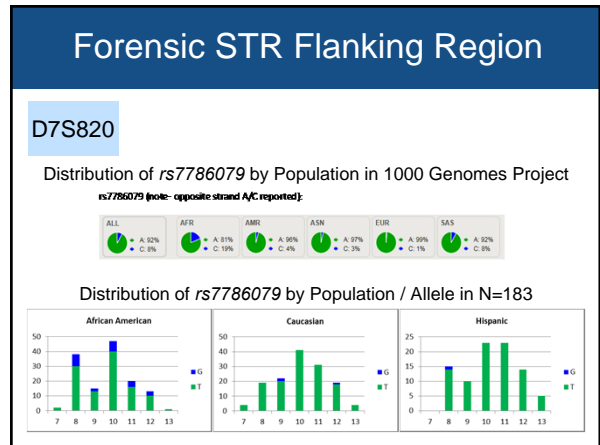
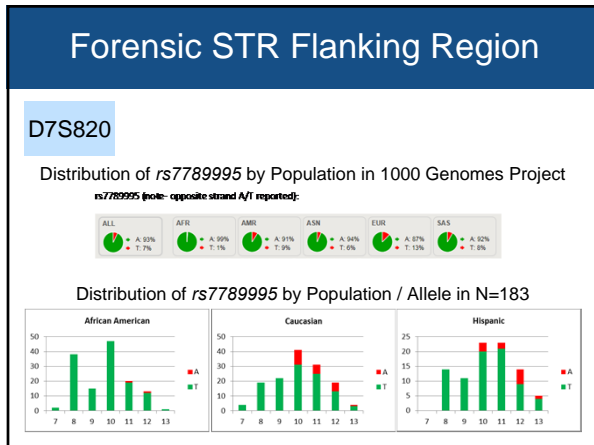
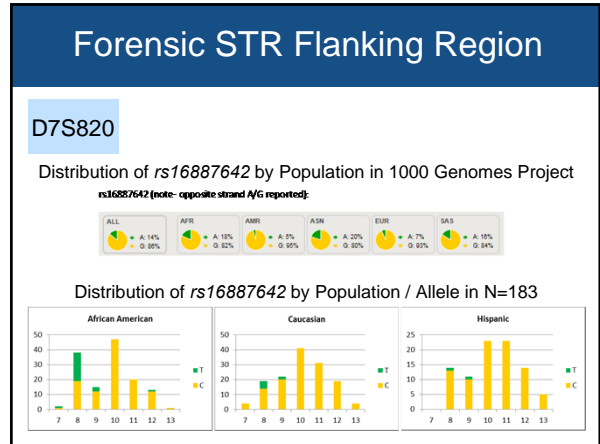
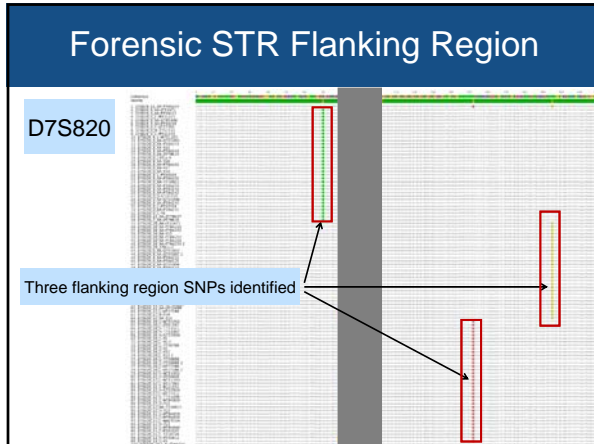
### Recognition Site-Based Informatics for STRs

Moving recognition sites out will capture information within the flanking regions

- PCR Primers
- STR Repeat Region
- Recognition site (~10 nt)





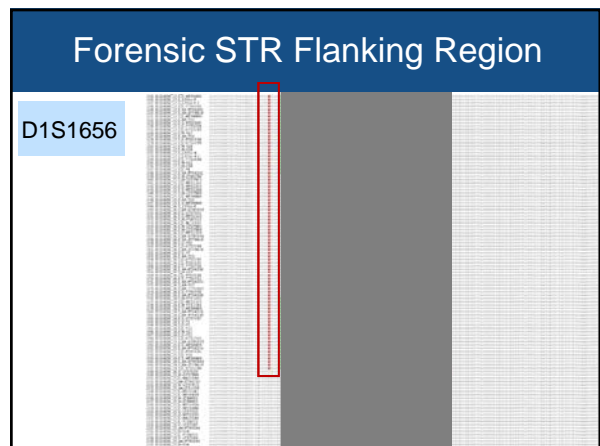


### Forensic STR Flanking Region

D7S820

Allele	Sequence	rs16887642- rs7789995- rs7786079
7	(GATA)7	C-T
8	(GATA)8	C-T
8	(GATA)8	C-T-G
8	(GATA)8	C-T-T
8	(GATA)8	C-T-C
8	(GATA)8	C-T-G
8	(GATA)8	C-T-T
9	(GATA)9	C-T-G
10	(GATA)10	C-T-G
10	(GATA)10	C-T-T
10	(GATA)10	C-T-C
10	(GATA)10	C-T-T (direction)
11	(GATA)11	C-T-G
11	(GATA)11	C-T-T
11	(GATA)11	C-T-G
11	(GATA)11	C-T-T
12	(GATA)12	C-T-G
12	(GATA)12	C-T-T
12	(GATA)12	T-T-T
12	(GATA)12	C-T-G
13	(GATA)13	C-T-T
13	(GATA)13	C-T-T

15 Additional Alleles



### Forensic STR Flanking Region

D1S1656

17 Allele: GTGATG[TAGA]16 [TAGG]  
 16.3 Allele: ATGATG[TAGA]4 TGA [TAGA]11 [TAGG]

In this case the flanking region SNP (G->A) coincides with a .3 genotype  
 Does not provide additional information

### Forensic STR Flanking Region

ARTICLE IN PRESS

Forensic Science International: Genetics Supplement Series xxx (2015) xxx–xxx

Contents lists available at ScienceDirect

Forensic Science International: Genetics Supplement Series

Journal homepage: www.elsevier.com/locate/FSIGSS

The next dimension in STR sequencing: Polymorphisms in flanking regions and their allelic associations

Katherine Butler Gettings<sup>a,\*</sup>, Rachel A. Aponte<sup>b</sup>, Kevin M. Kiesler<sup>c</sup>, Peter M. Vallone<sup>d</sup>

<sup>a</sup>US National Institute of Standards and Technology, 100 Bureau Drive, Gaithersburg, MD 20899-1316, USA  
<sup>b</sup>The George Washington University, Department of Forensic Sciences, 2300 Fordham Road NW, Washington, DC 20007, USA  
<sup>c</sup>US National Institute of Standards and Technology, Special Programs Office, 100 Bureau Drive, Gaithersburg, MD 20899, USA

### Forensic STR Sequence Diversity Comments

- 22 autosomal forensic STR loci were sequenced for a set of 183 samples (U.S. sample groups)
- Sequence variation within the STR region and surrounding flanking regions were characterized
- Loci with **compound and complex** repeat motifs contained the majority of 'additional' information by NGS
- Sequencing repeat region offers significant gains for some loci – (complex/compound)
- Extending analysis to the flanking regions offers additional gains – (simple repeats)

### NIST Support for NGS Research FSI Genetics review article

Forensic Science International: Genetics xxx (2015) xxx–xxx

Contents lists available at ScienceDirect

Forensic Science International: Genetics

Journal homepage: www.elsevier.com/locate/fig

STR allele sequence variation: Current knowledge and future issues

Katherine Butler Gettings<sup>a,\*</sup>, Rachel A. Aponte<sup>b</sup>, Peter M. Vallone<sup>c</sup>, John M. Butler<sup>d</sup>

<sup>a</sup>US National Institute of Standards and Technology, Measurement Standards, 100 Bureau Drive, Gaithersburg, MD 20899, USA  
<sup>b</sup>The George Washington University, Department of Forensic Sciences, 2300 Fordham Road NW, Washington, DC 20007, USA  
<sup>c</sup>US National Institute of Standards and Technology, Special Programs Office, 100 Bureau Drive, Gaithersburg, MD 20899, USA

Updating observed STR sequence variations for 24 autosomal loci on STRBase

### NIST Support for NGS Research

Allele	Repeat Structure	Reference	Platform
8	(TAA)8	Phillips et al. (2010)	Sanger
9	(TAA)9	Phillips et al. (2010)	Sanger
10	(TAA)10	Lareu et al. (1998)	Sanger
11	(TAA)11	Lareu et al. (1998)	Sanger
12	(TAA)12	Lareu et al. (1998)	Sanger
13	(TAA)13	Phillips et al. (2010)	Sanger
14	(TAA)14	Phillips et al. (2010)	Sanger
16	(TAA)16	Gettings et al. (2015)	MiSeq
17	(TAA)17	Phillips et al. (2010)	Sanger
18	(TAA)18	Lareu et al. (1998)	Sanger
19	(TAA)19	Lareu et al. (1998)	Sanger
20	(TAA)20	Lareu et al. (1998)	Sanger
21	(TAA)21	Lareu et al. (1998)	Sanger
22	(TAA)22	Lareu et al. (1998)	Sanger
23	(TAA)23	Lareu et al. (1998)	Sanger
24	(TAA)24	Lareu et al. (1998)	Sanger
25	(TAA)25	Lareu et al. (1998)	Sanger
26	(TAA)26	Lareu et al. (1998)	Sanger
27	(TAA)27	Lareu et al. (1998)	Sanger
28	(TAA)28	Lareu et al. (1998)	Sanger
29	(TAA)29	Lareu et al. (1998)	Sanger
30	(TAA)30	Lareu et al. (1998)	Sanger
31	(TAA)31	Lareu et al. (1998)	Sanger
32	(TAA)32	Lareu et al. (1998)	Sanger
33	(TAA)33	Lareu et al. (1998)	Sanger
34	(TAA)34	Lareu et al. (1998)	Sanger
35	(TAA)35	Lareu et al. (1998)	Sanger
36	(TAA)36	Lareu et al. (1998)	Sanger
37	(TAA)37	Lareu et al. (1998)	Sanger
38	(TAA)38	Lareu et al. (1998)	Sanger
39	(TAA)39	Lareu et al. (1998)	Sanger
40	(TAA)40	Lareu et al. (1998)	Sanger
41	(TAA)41	Lareu et al. (1998)	Sanger
42	(TAA)42	Lareu et al. (1998)	Sanger
43	(TAA)43	Lareu et al. (1998)	Sanger
44	(TAA)44	Lareu et al. (1998)	Sanger
45	(TAA)45	Lareu et al. (1998)	Sanger
46	(TAA)46	Lareu et al. (1998)	Sanger
47	(TAA)47	Lareu et al. (1998)	Sanger
48	(TAA)48	Lareu et al. (1998)	Sanger
49	(TAA)49	Lareu et al. (1998)	Sanger
50	(TAA)50	Lareu et al. (1998)	Sanger

**Excel Workbook**

- Sheet for each STR locus
- Observed alleles, repeat structure, platform
- Broken out by sub-motif
- References
- Not frequency data
- This can be updated as needed

### NIST Support for NGS Research

**Annotations**

- GRCh38 genome with STR repeat region and flanking SNPs identified
- A file for each locus can be downloaded

### NIST Support for NGS Research

Reference SNP	Chromosome	Chromosome Position	Distance from STR repeat	RefSNP alleles	Minor Allele	Minor Allele Frequency	Minor Allele Count
SNP ID	Chr	Position (bp)	Distance (bp)	Forward strand	Reverse strand	Frequency	Count
rs134273070	2	218014007	463	C/T	C	0.0004	152
rs37321055	2	218014076	434	A/G	A	0.004	9
rs142180843	2	218014063	387	A/G	G	0.0122	61
rs12445264	2	218014179	371	A/C	A	0.0014	8
rs12654061	2	218014619	291	C/T	T	0.0122	61
rs148093778	2	218014678	272	C/G	T	0.0004	2
rs16861033	2	218014651	299	C/T	T	0.114	57
rs137961276	2	218014796	152	A/G	A	0.0006	3
rs1736893	2	218014824	126	A/C	A	0.2041	1022
rs1061818	2	218014870	80	A/C	C	0.6468	2178
rs14761823	2	218014912	38	A/G	A	0.0058	29
rs127177707	2	218014919	31	A/G	A	0.0008	4
rs1736895	2	218014925	25	A/C	A	0.3005	1504
rs1322233	2	218014929	21	A/C	A	0.0569	285
rs13920853	2	218014933	17	A/C	A	0.0128	64
rs14809423	2	218014937	13	A/C	A	0.0008	4
Repeat Region	2	218014938-218014939 (bp)	27 repeats			0.013	66
rs1705118	2	218015007		A/G	A	0.0013	66
rs12913952	2	21801510		A/G	A	0.0013	66
rs17361378	2	21801511		A/G	A	0.0006	3
rs162604377	2	21801514		A/G	A	0.0006	3
rs147416182	2	21801512		A/G	A	0.0004	2
rs144450033	2	21801512		A/G	A	0.0004	2
rs166045602	2	21801513		A/G	A	0.0014	8
rs12770379	2	21801513		A/G	A	0.0004	2
rs148094433	2	21801514		A/G	A	0.0006	3
rs185427083	2	21801514		A/G	A	0.001	5

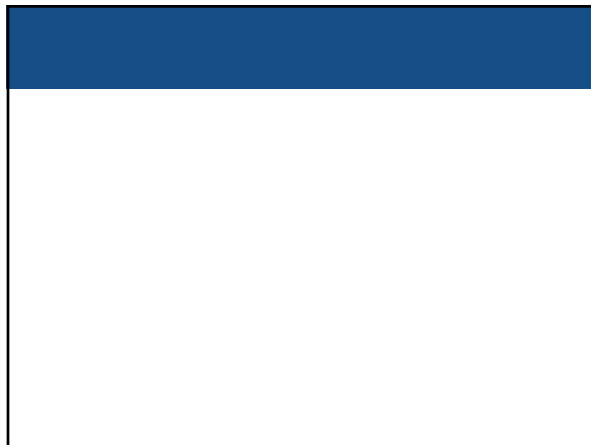
Cataloged SNPs

- SNPs found in 1000 genome data set within 500 bp of STRs
- Useful for NGS primer design and bioinformatics

### NIST Support for NGS Research SRM 2391c

Sequence sequence of the core autosomal and Y STRs motifs included

[https://www-s.nist.gov/srmors/view\\_cert.cfm?srn=2391C](https://www-s.nist.gov/srmors/view_cert.cfm?srn=2391C)



### Considerations for nomenclature, storage, searching, and reporting

How do we communicate and categorize this additional information?

- General thoughts
- Recent paper by Gelardi et al.
- ISFG nomenclature panel discussion

### Adopting sequence-based genotype information What are our concerns?

- Take advantage of the additional information provided by sequencing the STR allele (and flanking regions)
- Maintain back compatibility to existing databases
  - Genotyping by length-based CE methods
- Report and represent data in a consistent manner
- Application of sequence data
  - **Storage, searching, and reporting**

### Nomenclature

- Important so that all labs 'speak the same language'
  - The established nomenclature rules for the core STR loci can be expanded upon if needed
- Designate sequenced alleles in an internationally accepted format
- How will bracket notation fit into our plans?
- Do we need to define amplicon boundaries?
  - Core recognition sites (pros/cons)
- How will we designate sequence that flanks the STR repeat? (SNPs, indels)



### D8S1179 Example

Allele	Repeat Structure
7	[TCTA]7-14
7	[TCTA]7
8	[TCTA]8
9	[TCTA]9
10	[TCTA]10
11	[TCTA]11
12	[TCTA]12
13	[TCTA]13
14	[TCTA]14
<hr/>	
[TCTA]1[TCTG]11[TCTA]10-14	
12	[TCTA]1[TCTG]11[TCTA]10
13	[TCTA]1[TCTG]11[TCTA]11
14	[TCTA]1[TCTG]11[TCTA]12
16	[TCTA]1[TCTG]11[TCTA]14
<hr/>	
[TCTA]2[TCTG]11[TCTA]9-15	
11	[TCTA]2[TCTG]11[TCTA]8
12	[TCTA]2[TCTG]11[TCTA]9
13	[TCTA]2[TCTG]11[TCTA]10
14	[TCTA]2[TCTG]11[TCTA]11
16	[TCTA]2[TCTG]11[TCTA]12
16	[TCTA]2[TCTG]11[TCTA]13
17	[TCTA]2[TCTG]11[TCTA]14
18	[TCTA]2[TCTG]11[TCTA]15
<hr/>	
[TCTA]2[TCTG]12[TCTA]11-15	
16	[TCTA]2[TCTG]12[TCTA]11
16	[TCTA]2[TCTG]12[TCTA]12
17	[TCTA]2[TCTG]12[TCTA]13
19	[TCTA]2[TCTG]12[TCTA]15

Four 'flavors' of the motif are observed (so far...)  
Tempting to just call them A, B, C, D, etc

14 [TCTA]14 (A)

14 [TCTA]1 [TCTG]1 [TCTA]12 (B)

14 [TCTA]2 [TCTG]1 [TCTA]11 (C)

What if/when a new motif and/or SNP is detected within one of the motifs?

Allele	Repeat Structure
10	[TGCC]4-[TTCC]16-18
12	[TGCC]4 [TTCC]16
12	[TGCC]4 [TTCC]18
12	[TGCC]4 [TTCC]14
13	[TGCC]4 [TTCC]19
15	[TGCC]4 [TTCC]11
16	[TGCC]4 [TTCC]12
16	[TGCC]4 [TTCC]10
17	[TGCC]4 [TTCC]13
17	[TGCC]4 [TTCC]11
17	[TGCC]4 [TTCC]12
18	[TGCC]4 [TTCC]15
18	[TGCC]4 [TTCC]11
19	[TGCC]4 [TTCC]13
19	[TGCC]4 [TTCC]12
19	[TGCC]4 [TTCC]11
20	[TGCC]4 [TTCC]14
20	[TGCC]4 [TTCC]13
20	[TGCC]4 [TTCC]12
21	[TGCC]4 [TTCC]14
21	[TGCC]4 [TTCC]13
22	[TGCC]4 [TTCC]18
22	[TGCC]4 [TTCC]14
22	[TGCC]4 [TTCC]13
23	[TGCC]4 [TTCC]14
23	[TGCC]4 [TTCC]13
23	[TGCC]4 [TTCC]12
24	[TGCC]4 [TTCC]15 [TTCC]1 [TTCC]12
<hr/>	
[TGCC]4-[TTCC]15-17 [TTCC]1 [TTCC]12	
19	[TGCC]4 [TTCC]15 [TTCC]1 [TTCC]12
20	[TGCC]4 [TTCC]16 [TTCC]1 [TTCC]12
21	[TGCC]4 [TTCC]17 [TTCC]1 [TTCC]12
22	[TGCC]4 [TTCC]18 [TTCC]1 [TTCC]12
23	[TGCC]4 [TTCC]19 [TTCC]1 [TTCC]12
24	[TGCC]4 [TTCC]20 [TTCC]1 [TTCC]12

The A B designation for D2S1338 would be problematic due to multiple alleles with the same motif

17 [TGCC]4[TTCC]13

17 [TGCC]5[TTCC]12

17 [TGCC]6[TTCC]11

## Recent Paper

Forensic Science International: Genetics 12 (2014) 18-41

Contents lists available at ScienceDirect

Forensic Science International: Genetics

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

Second generation sequencing of three STRs D3S1358, D12S391 and D21S11 in Danes and a new nomenclature for sequenced STR alleles

Chiara Gelardi<sup>a,b,1</sup>, Eszter Rockenbauer<sup>a,1,\*</sup>, Signun Dalsgaard<sup>a</sup>, Claus Bersting<sup>a</sup>, Niels Morling<sup>a</sup>

<sup>a</sup>Section of Forensic Genetics, Department of Forensic Medicine, Faculty of Health and Medical Sciences, University of Copenhagen, Denmark  
<sup>b</sup>Faculty of Mathematical, Physical and Natural Sciences, University of Padova, Padova, Italy

Gelardi nomenclature format:  
Locus | Length | Sequence | Flank Variants

## Recent Paper

**Table 1**  
Novel D3S1358, D12S391 and D21S11 alleles.

CE allele name	SGS allele name	Count	Frequency
14	D3S1358[14][TCTA]1 TCTG[3][TCTA]10]	1	0.0025
14	D3S1358[14][TCTA]1 TCTG[1][TCTA]12]	1	0.0025
20	D3S1358[20][TCTA]1 TCTG[4][TCTA]15]	1	0.0025
18	D12S391[18][AGAT]11] AGAC[7]	3	0.0076
18	D12S391[18][AGAT]12] AGAC[5][AGAT]1]	4	0.0102
19	D12S391[19][AGAT]11] AGAC[7][AGAT]1]	1	0.0025
19	D12S391[19][AGAT]11] AGAC[8]	6	0.0152
20	D12S391[20][AGG]1] AGAT[10][AGAC]9]	2	0.0051
21	D12S391[21][AGG]13] GGAC[1][AGAC]7]	1	0.0025
21	D12S391[21][AGG]11] AGAT[11][AGAC]9]	1	0.0025

## Recent Paper

**Table 1**  
Novel D3S1358, D12S391 and D21S11 alleles.

Gelardi nomenclature example  
Locus | length | sequence | flank variants

D3S1358 [14] TCTA[1]TCTG[3]TCTA[10]

D12S391 [20] AGGT[1]AGAT[10]AGAC[9]

D21S11 [32.2] TCTA[6]TCTG[6]itvsTCTA[11]TA[1]TCTA[1]

*itvs: TCTA[3]TA[1]TCTA[3]TCA[1]TCTA[2]TCCA[1]TA[1] (43 bp)*

21	D12S391[21]AGAT[13] GGAC[1]AGAC[7]	1	0.0025
21	D12S391[21]AGG[11] AGAT[11]AGAC[9]	1	0.0025


## Reporting

- Right now we report: D2S1338 [20,22]
- Is the full sequence string too much in a report?
  - Just the STR motif or include flanking regions (length limits?)
- Should the sequence string be condensed down to something easier to understand in a report?
  - Bracket notation? (example next slide)
- An unique code might be fine for storage and searching, but not intuitive for a report (lookup key)
  - Example: gatagatagatagatagatagatagata = 'TK421'




**Christophe Van Neste, Ghent University  
FLAD: Forensic Loci Allele Database**

- A minimal nomenclature
  - Not reporting the full sequence
  - An identifier can be referenced to the full sequence
  - FL1A12 ↔ AGATAGATAGATAGATAGATAGAT**
  - <https://forensic.ugent.be/FLAD>
- A custom solution
  - Avoiding GenBank refseqs
  - dbSNP rs#
- Identifier stored in FLAD
  - Need QC on data



**Christophe Van Neste, Ghent University  
FLAD: Forensic Loci Allele Database**

- A minimal nomenclature
  - Not reporting the full sequence
  - Article in Press
  - Forensic Loci Allele Database (FLAD): automatically generated, permanent identifiers for sequenced forensic alleles
- A
  - Christophe Van Neste  
Laboratory of Pharmaceutical Biotechnology, Faculty of Pharmaceutical Sciences, Ghent University, Ghent, Belgium
  - dbSNP rs#
- Identifier stored in FLAD
  - Need QC on data



**Christophe Van Neste, Ghent University  
FLAD: Forensic Loci Allele Database**

FA[HEX] → FA1A, FA23 (FA: Forensic Allele + random hex number)

↓ Adding focus information

FL[D]A[HEX] → FL3A1A, FL1A23 (FLA: Forensic Locus ... Allele ...)

↓ Adding primerset information

FL[D(.D)]A[HEX] → FL3.1A1A, FL3.2A3C

↓ Adding allele annotation

FL[D](X|A|P)[HEX] → FL3X1A, FL1P23

(X = unvalidated allele, P = validated, but rare allele, A = common population allele)

**Christophe Van Neste, Ghent University  
FLAD: Forensic Loci Allele Database**

FA[HEX]	Loci	Genotype
↓	Amelogenin [FL1]	X[A01], Y[A03]
↓	D3S1358 [FL3]	16[P0A]
↓	THO1 [FL2]	6[A3B], 9.3[A20]
↓	D21S11 [FL11]	30[A05], 30.2[AFF]
↓	D18S51 [FL13]	15[A10]
↓	D8S1179 [FL6]	14[A5D]
↓	TPOX [FL5]	8[A10], 12[A02]
↓	FGA [FL8]	23[A41]

(X = unvalidated allele, P = validated, but rare allele, A = common population allele)

**Kristiaan van der Gaag, Rick de Leeuw & Peter de Knijff,  
Department of Human Genetics Leiden University  
Medical Center**

**Fully descriptive and relative to a reference genome**

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

D13S317[CE12]-Chr13-GRCh37-g.82.722.160:82.722.223-TATC[13]AAATC[1]ATCT[3]-g.x.136G>A  
Locus and CE allele

D13S317[CE12]-Chr13-GRCh37-g.82.722.160:82.722.223-TATC[13]AAATC[1]ATCT[3]-g.x.136G>A  
CHR and human genome reference version

D13S317[CE12]-Chr13-GRCh37-g.82.722.160:82.722.223-TATC[13]AAATC[1]ATCT[3]-g.x.136G>A  
STR repeat region coordinates of reference allele

D13S317[CE12]-Chr13-GRCh37-g.82.722.160:82.722.223-TATC[13]AAATC[1]ATCT[3]-g.x.136G>A  
Full description of STR motif

D13S317[CE12]-Chr13-GRCh37-g.82.722.160:82.722.223-TATC[13]AAATC[1]ATCT[3]-g.x.136G>A  
Genome coordinate of derived SNP

**Sascha Willuweit**  
Department of Forensic Genetics, Institute of Legal Medicine and Forensic Sciences,  
Charité – Universitätsmedizin, Berlin, Germany

- Reporting full sequences is not feasible (yet)
- Long lists of sequences are correct, but not practical
- Need legacy compatibility (with CE data)
- A Nomenclature Authority (NA) would assign a code/type to a new allele
  - Centralized server
  - Reliable, secure
  - ISFG organized/maintained

**Sascha Willuweit**  
 Department of Forensic Genetics, Institute of Legal Medicine and Forensic Sciences,  
 Charité – Universitätsmedizin, Berlin, Germany

- Submit sequenced allele to NA
  - If already observed (at least twice) a fixed value is returned
  - If not, a temporary identifier is assigned until confirmed
- Example, submit a sequence of a 14 allele
  - **Already observed**, returned the value '14\*f'
  - **New allele**, assigned '14\*68'
  - **When observed again** – this is converted into '14\*g' (fixed value)

**Chris Phillips**  
 Univ of Santiago de Compostela

Is a 'synthesis' possible and can we simplify sequences still further?

- Define a sequence only by its **differences from the reference genome**
- Upload raw sequences as, say, Excel lists which are then processed by the NA

Region of interest	Region of primer	Region of primer 2	Region of primer 3
D13S317 CE 10,11 (RR only)	D13S317 CE 10,11 (amplicon)	D13S317 CE 11,13 (extra repeats)	
A1	A1	A1	
13:82722200-82722203 —	13:82722196 G	A2	
13:82722204 T	13:82722200-82722203 —	13:82722203 +1-8 TATC	
A2	13:82722204 T	13:82722203 +2 G	
	A2		

**Thoughts**

- Pros/cons of using a **code** versus **full descriptor**
  - How difficult is it to deal with full sequence strings?
- References to genome builds can/will change
  - How to deal with this issue?
- Commercial kits/software will have to be informed
- QC of sequence submitted to nomenclature assigning websites (if they exist) will be important
  - See EMPOP

*Lots of aspects to think about---stay tuned*

**Acknowledgements**

**NIST**  
 Katherine Gettings  
 Kevin Kiesler  
 Nate Olson  
 Jo Lynne Harenza  
 Mike Coble  
 Becky Hill  
 Margaret Kline

**Student Interns**  
 Rachel Aponte (GWU)  
 Harish Swaminathan (Rutgers)  
 Anna Blendermann (MC)

**Battelle**  
 Seth Faith (now @NCSU)  
 Rich Guerrieri  
 Brian Young

**Promega**  
 Doug Storts  
 Jay Patel

**ISFG Nomenclature Panel Members**  
 Peter de Knijff  
 Sascha Willuweit  
 Christophe Van Neste  
 Chris Phillips

Contact Information  
 peter.vallone@nist.gov

NIST Disclaimer: Certain commercial equipment, instruments and materials are identified in order to specify experimental procedures as completely as possible. In no case does such identification imply a recommendation or it imply that any of the materials, instruments or equipment identified are necessarily the best available for the purpose. Funding FBI: DNA as a Biometric