

**Forensic Genetic NGS:
Sequence Diversity of STRs**

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

NGS/MPS Workshop
26th Congress of the International Society for Forensic Genetics
August 31, 2015
Krakow, Poland

NIST
National Institute of
Standards and Technology
U.S. Department of Commerce

Disclaimer

We will mention commercial products and information, but we are in no way attempting to endorse any specific products.

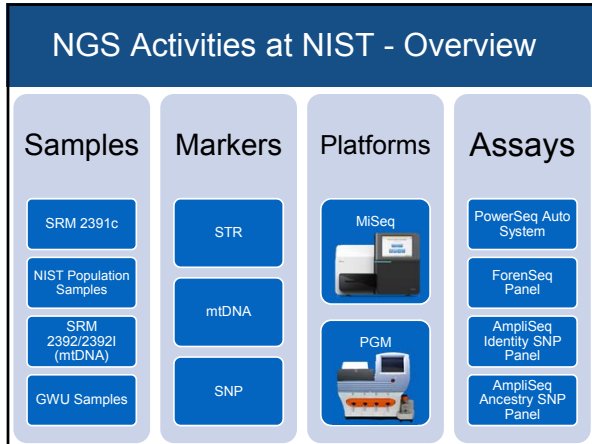
NIST Disclaimer: Certain commercial equipment, instruments and materials are identified in order to specify experimental procedures as completely as possible. In no case does such identification imply a recommendation or it imply that any of the materials, instruments or equipment identified are necessarily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

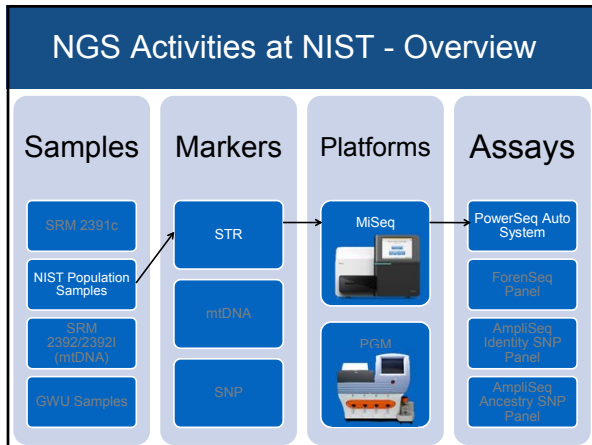
Information presented does not necessarily represent the official position of the National Institute of Standards and Technology or the U.S. Department of Justice.

Our group receives or has received funding from the FBI Laboratory and the National Institute of Justice.

Outline for Today


- STR sequence diversity
- Investigating STR sequence diversity at NIST
 - Informatics
 - Sequence and genotype diversity
 - Flanking region diversity
 - Stutter
- NIST resources for NGS researchers

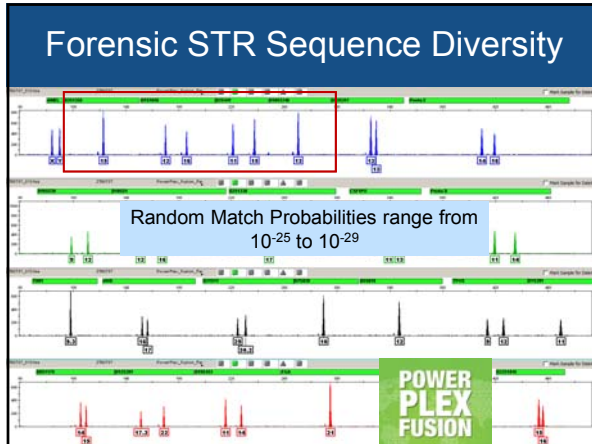


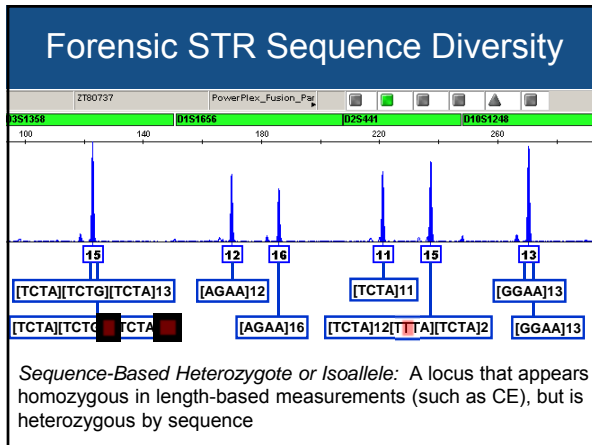


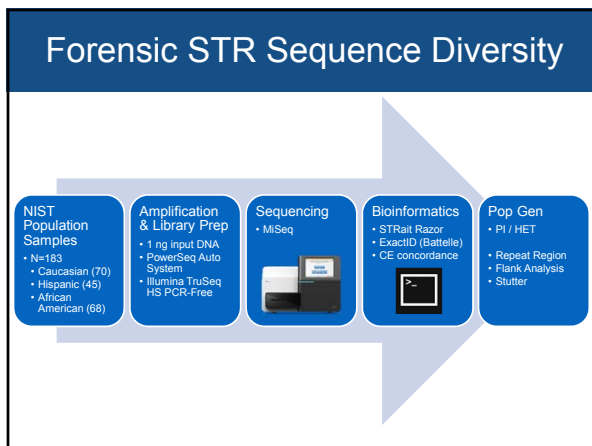
NGS has potential for finer resolution of STR amplicons not detectable by CE-length based methods

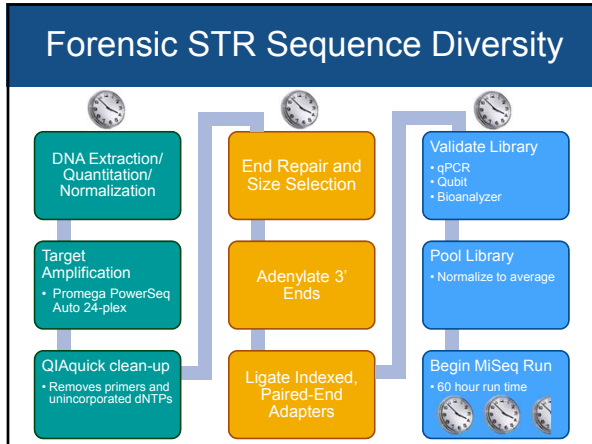
- Additional STR alleles
- Flanking region SNPs and InDels
- *Resolve homozygous by length peaks*
- *Resolve minor contributor peaks from stutter*

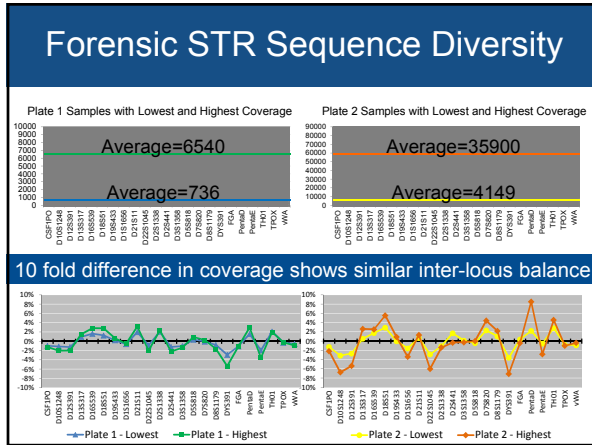


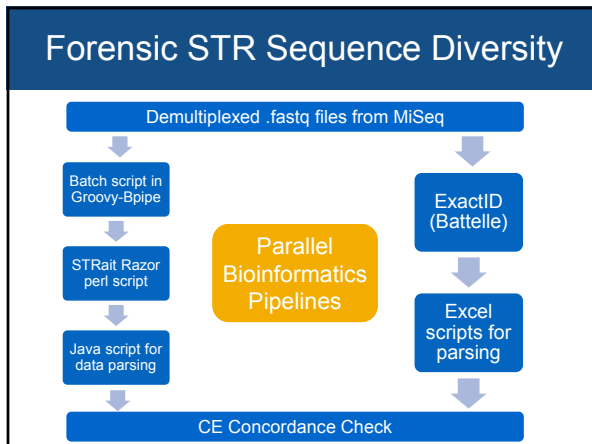












Recognition Site-Based Informatics for STRs

Computer Returns:
 The length between the recognition sequences = 36 bp
 Reference table in software returns a "9" allele
 The sequence between the recognition sites

TATC PCR Primers

(TATC)_n STR Repeat Region

TATC Recognition site (~10 nt)

¹ <http://batelleeexacid.org/>

² STRait Razor: a length-based forensic STR allele-calling tool for use with second generation sequencing data. [Warshauer et al., Forensic Sci Int Genet, 2013 \(7\):409-17](#)

³ STRait Razor v2.0: the improved STR Allele Identification Tool-Razor. [Warshauer et al., Forensic Sci Int Genet, 2015 \(14\):182-6](#)

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

Obs 5/183

Sequence [TATC]₁₀_[TATC]₁₁

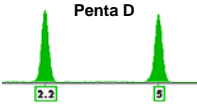
← Repeat Region NGS Recognition Region 4 bp Deletion CE Primer Binding Site →

```
TATC TATC TATC AATCAATCATCTATCTATCTTTCTGTC----TTTGGGCTGCCTATGGCTCAA
TATC TATC TATC AATCAATCATCTATCTATCTTTCTGTCGTCTTTTGGGCTGCCTATGGCTCAA
```

Flanking region InDel: Bioinformatic pipelines may reduce the region used for genotyping, resulting in deletions not being "counted" as they would via CE

Forensic STR Sequence Diversity

Obs 15/183



Penta D

Sequence [AAAGA]₅

13 bp deletion in Recognition Region

← CE Primer Binding Site Recognition Region Repeat Region →

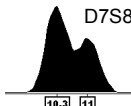
```

TAGGTTCACAGAGCAAGACACCATCTCAAG-----AAAGA AAAGA AAAGA AAAGA AA
TAGGTGACAGAGCAAGACACCATCTCAAGAAAGAAAAAAAAAG AAAGA AAAGA AAAGA AA
    
```

Bioinformatic Null Allele: A true allele that is present within the raw sequence data but is not detected by the bioinformatic pipeline

Forensic STR Sequence Diversity

Obs 1/183



D7S820

NGS Sequence [GATA]₁₁

Sequence [GATA]₁₁

REPEAT REGION RECOGNITION SITE 1 bp DELETION

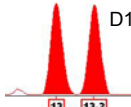
```

GATA GATA GATA GATA GATA GACAGATTGATAGTTTTTTTATCTACTAAATAGTCTATAGT
GATA GATA GATA GATA GATA GACAGATTGATAGTTTTTTTATCTACTAAATAGTCTATAGT
    
```

Flanking region InDel: Bioinformatic pipelines may reduce the region used for genotyping, resulting in deletions not being "counted" as they would via CE.

Forensic STR Sequence Diversity

Obs 3/183



D19S433

NGS Sequence [AAGG]₁₃

Sequence [AAGG]₁₃

2 bp deletion in Recognition Region

REPEAT REGION RECOGNITION REGION

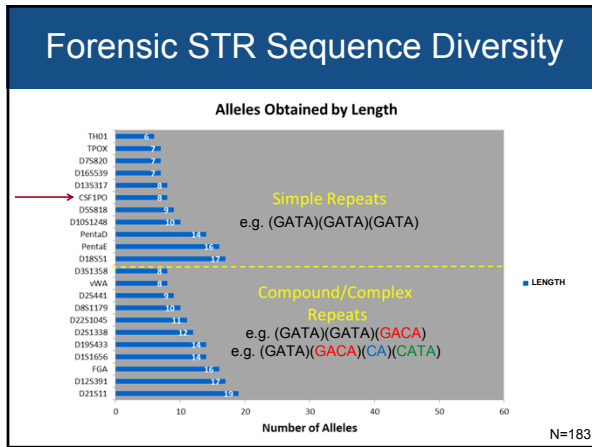
```

AAGG AAGG AAGG AAGG AAGG AAGG AAGG AGAGAGGAGAAAGAGAGAAGATTTTTATTCGGGT
AAGG AAGG AAGG AAGG AAGG AAGG AAGG AGAG--GAAGAAAGAGAGAAGATTTTTATTCGGGT
    
```

Bioinformatic Null Allele: A true allele that is present within the raw sequence data but is not detected by the bioinformatic pipeline

Comments

- This is not a criticism of the informatics tools/methods
 - After discovery the listed discordance issues are easily addressed
- No issues with sequencing the amplicons
- How do we want to capture STR allele information?
 - Defined by recognition sites adjacent to the flanking region?
 - Defined by PCR primers?
 - Some set of core recognition sites? Maintain back compatibility with existing technology (CE) enable comparisons between labs



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	[AGGT][AGAT]9	AGGT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT
11	[AGAT]11	AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT
11	[AGAT]9AGGT[AGAT]7	AGAT AGAT AGAT AGAT 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

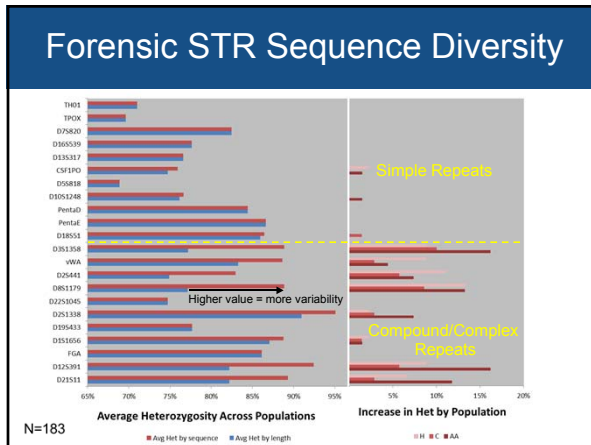
N=183

Forensic STR Sequence Diversity

Heterozygosity

heterozygotes observed
of loci tested

Indicates genetic variability at a locus



Forensic STR Sequence Diversity

Probability of Identity

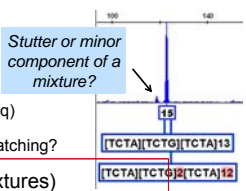
Sum of each genotype frequency² at each locus

$$= \sum_{i=1}^n x_i^2$$

Probability that two unrelated individuals
selected at random
will have the same genotype at a locus

Isoalleles

- Within an individual
 - Sequence based heterozygote
 - Increased heterozygosity ($p^2 \rightarrow 2pq$)
 - Marginal benefits for one-to-one matching?
- Between individuals (DNA mixtures)
 - Resolve overlapping alleles
 - 15 → '15 from person A' and '15 from person B'
 - Resolve stutter artifacts versus minor components



Manuscript submitted to FSIG

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

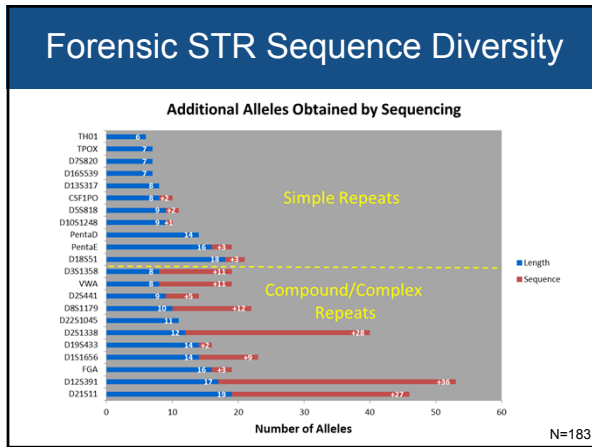
Authors and Affiliations:
Katherine Butler Gettings^{1*}, Kevin M. Kiesler¹, Seth A. Faith², Elizabeth Montano³, Christine H. Baker³, Brian A. Young³, Richard A. Guerrieri³, and Peter M. Vallone¹

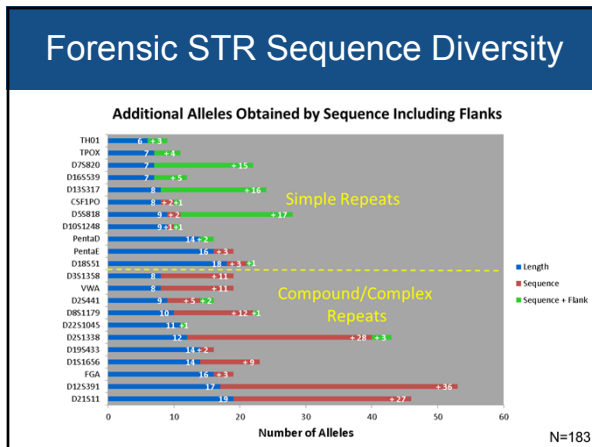
Flanking Region Variation

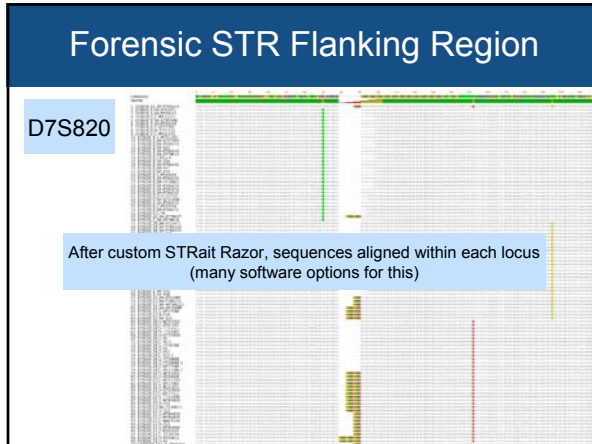
Recognition Site-Based Informatics for STRs

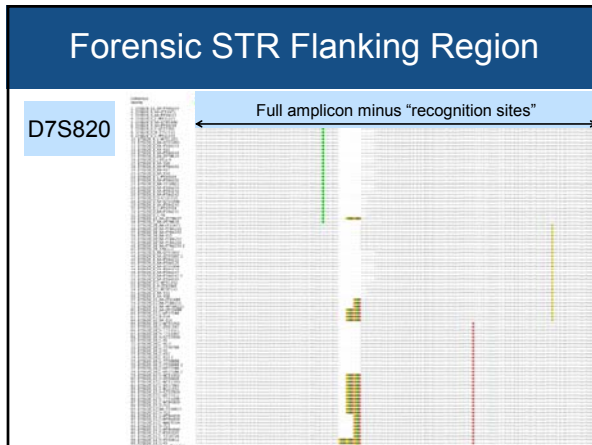
Moving recognition sites out will capture information within the flanking regions

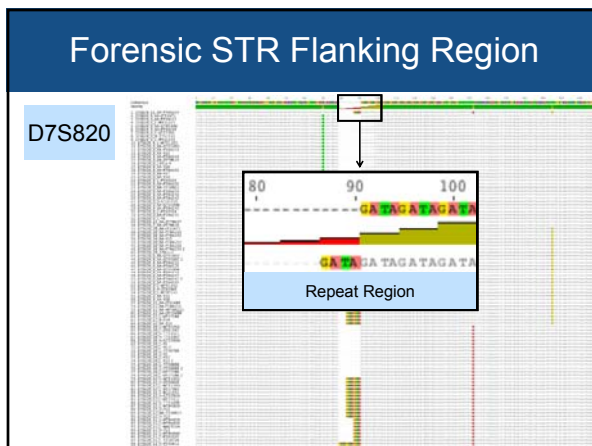
- PCR Primers
- STR Repeat Region
- Recognition site (~10 nt)

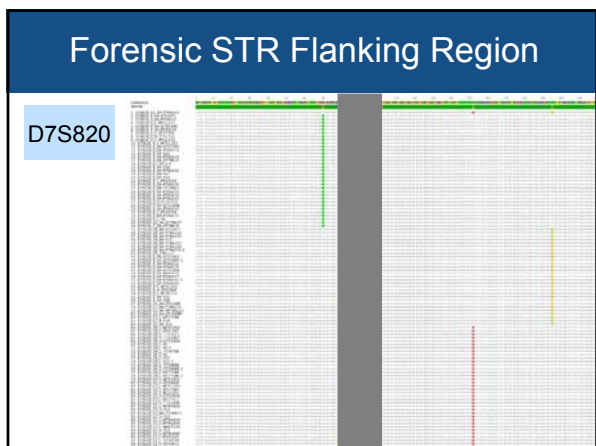


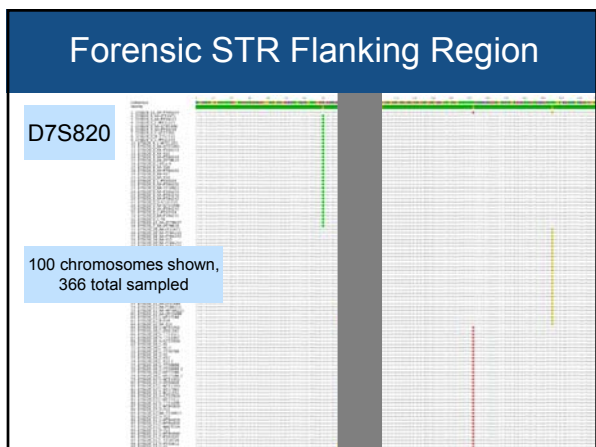


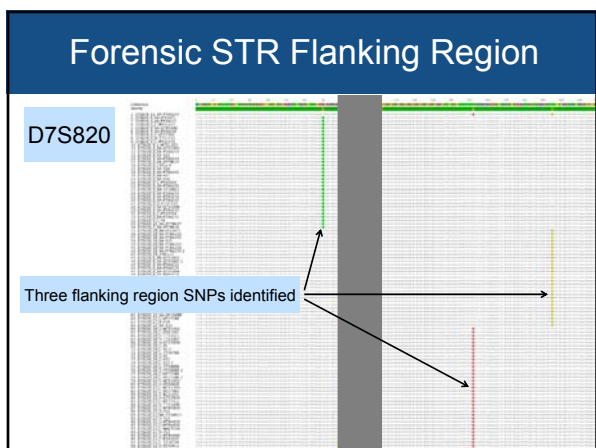


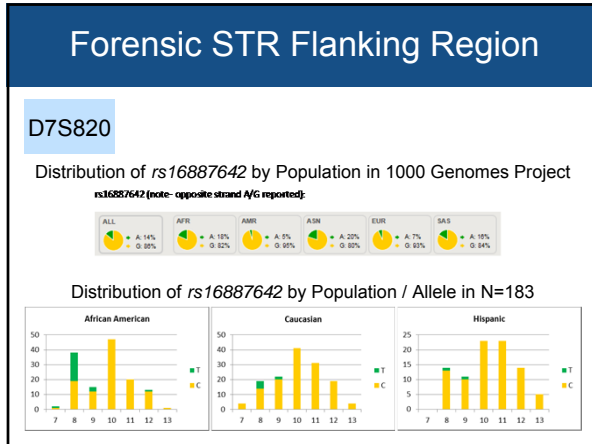


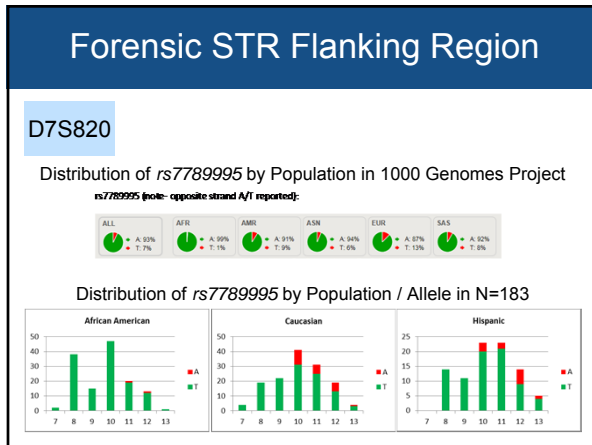


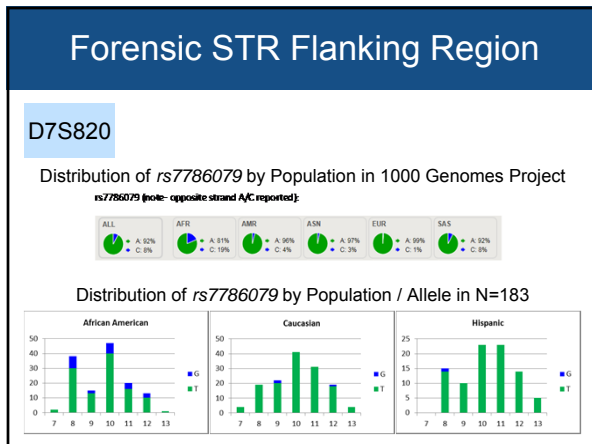








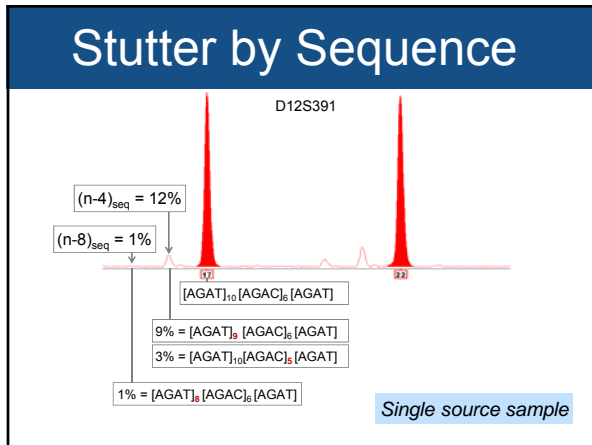


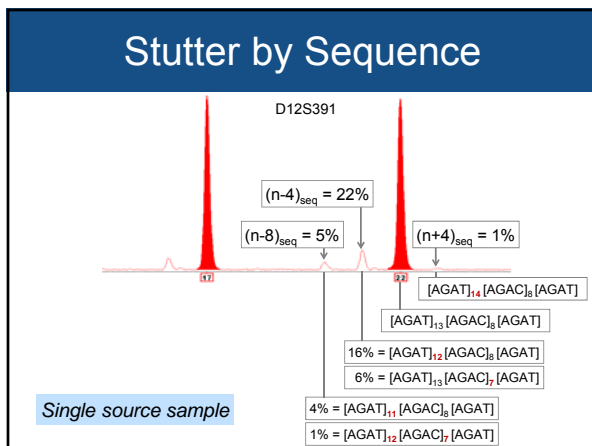


**Flanking region manuscript in preparation
 & ISFG oral presentation**

Wednesday

11.00-11.30	Coffee & Posters
11.30-11.45	David Ballard Complex mixture interpretation using massively parallel sequencing
11.45-12.00	Kristiaan van der Gaag Massive parallel sequencing of autosomal short tandem repeats and SNPs, the next level of forensic mixture analysis
* 12.00-12.15	Katherine Gettings The next dimension in STR sequencing: polymorphisms in flanking regions and their allelic associations
12.15-13.00	Panel discussion Walther Parson, Chris Phillips, Katherine Gettings, Peter de Knijff, David Ballard, Bruce Budowle, Niels Morling Sequencing based STR nomenclature
13.00-14.30	Lunch





Signal and Noise

- Sources of noise
 - Biological and sequencing chemistry
- Even the 'majority peaks' (alleles) will contain some sequence variants
- Interpretable for single source samples
- BUT we need to understand S/N in order to use NGS to resolve mixtures
 - Is the 'other' sequence from another contributor or just noise?

Poster #100

SEQUENCE-BASED ANALYSIS OF STUTTER AT STR LOCI: CHARACTERIZATION AND UTILITY (P)

R.A. Aponte¹, K.B. Gettings², D.L. Diewer², M.D. Coble², and P.M. Vallone²

¹ Department of Forensic Sciences, The George Washington University, Washington, DC, U.S.

² National Institute of Standards and Technology, Gaithersburg, MD, U.S.
katherine.gettings@nist.gov

Comments on mixtures

What is needed to understand the true benefits of sequencing to mixtures?

- Validate analytical thresholds for NGS work
- Assess sensitivity of NGS methods to detect minor components
 - Still PCR front end – will we observe better than 10:1?
- In practice - how often will *resolvable* overlapping alleles be observed in a mixture?
- Need the allele frequencies for 'new' alleles
 - What size population databases are needed? Greater than 200?
- Incorporate NGS sequence data into probabilistic genotyping software (STRMix, True Allele, etc)
 - What are the gains in stats (Log LR)? Improved contributor ratio estimates?

More loci – better loci?

NIST Support for Forensic NGS Research FSI Genetics review article

Forensic Science International: Genetics
 Contents lists available at ScienceDirect
 ELSEVIER
 journal homepage: www.elsevier.com/locate/fsig

STR allele sequence variation: Current knowledge and future issues
 Katherine Butler Gettings^{a,*}, Rachel A. Aponte^b, Peter M. Vallone^a, John M. Butler^c

**Updating observed STR sequence variations for
 24 autosomal loci on STRBase**

NIST Support for NGS Research

Allele	Repeat Structure	Reference	Platform
8	(TAA)9-14(TTT)	Phillips et al. (2010)	Sanger
9	(TAA)10	Phillips et al. (2010)	Sanger
10	(TAA)11	Laird et al. (1998)	Sanger
11	(TAA)12	Laird et al. (1998)	Sanger
12	(TAA)13	Laird et al. (1998)	Sanger
13	(TAA)14	Phillips et al. (2010)	Sanger
14	(TAA)15	Phillips et al. (2010)	Sanger
15	(TAA)16	Gettings et al. (2015)	Illumina
16	(TAA)17	Phillips et al. (2010)	Sanger
17	(TAA)18	Laird et al. (1998)	Sanger
18	(TAA)19	Laird et al. (1998)	Sanger
19	(TAA)20	Laird et al. (1998)	Sanger
20	(TAA)21	Laird et al. (1998)	Sanger
21	(TAA)22	Laird et al. (1998)	Sanger
22	(TAA)23	Laird et al. (1998)	Sanger
23	(TAA)24	Laird et al. (1998)	Sanger
24	(TAA)25	Laird et al. (1998)	Sanger
25	(TAA)26	Laird et al. (1998)	Sanger
26	(TAA)27	Laird et al. (1998)	Sanger
27	(TAA)28	Laird et al. (1998)	Sanger
28	(TAA)29	Laird et al. (1998)	Sanger
29	(TAA)30	Laird et al. (1998)	Sanger
30	(TAA)31	Laird et al. (1998)	Sanger
31	(TAA)32	Laird et al. (1998)	Sanger
32	(TAA)33	Laird et al. (1998)	Sanger
33	(TAA)34	Laird et al. (1998)	Sanger
34	(TAA)35	Laird et al. (1998)	Sanger
35	(TAA)36	Laird et al. (1998)	Sanger
36	(TAA)37	Laird et al. (1998)	Sanger
37	(TAA)38	Laird et al. (1998)	Sanger
38	(TAA)39	Laird et al. (1998)	Sanger
39	(TAA)40	Laird et al. (1998)	Sanger
40	(TAA)41	Laird et al. (1998)	Sanger
41	(TAA)42	Laird et al. (1998)	Sanger
42	(TAA)43	Laird et al. (1998)	Sanger
43	(TAA)44	Laird et al. (1998)	Sanger
44	(TAA)45	Laird et al. (1998)	Sanger
45	(TAA)46	Laird et al. (1998)	Sanger
46	(TAA)47	Laird et al. (1998)	Sanger
47	(TAA)48	Laird et al. (1998)	Sanger
48	(TAA)49	Laird et al. (1998)	Sanger
49	(TAA)50	Laird et al. (1998)	Sanger
50	(TAA)51	Laird et al. (1998)	Sanger
51	(TAA)52	Laird et al. (1998)	Sanger
52	(TAA)53	Laird et al. (1998)	Sanger
53	(TAA)54	Laird et al. (1998)	Sanger
54	(TAA)55	Laird et al. (1998)	Sanger
55	(TAA)56	Laird et al. (1998)	Sanger
56	(TAA)57	Laird et al. (1998)	Sanger
57	(TAA)58	Laird et al. (1998)	Sanger
58	(TAA)59	Laird et al. (1998)	Sanger
59	(TAA)60	Laird et al. (1998)	Sanger
60	(TAA)61	Laird et al. (1998)	Sanger
61	(TAA)62	Laird et al. (1998)	Sanger
62	(TAA)63	Laird et al. (1998)	Sanger
63	(TAA)64	Laird et al. (1998)	Sanger
64	(TAA)65	Laird et al. (1998)	Sanger
65	(TAA)66	Laird et al. (1998)	Sanger
66	(TAA)67	Laird et al. (1998)	Sanger
67	(TAA)68	Laird et al. (1998)	Sanger
68	(TAA)69	Laird et al. (1998)	Sanger
69	(TAA)70	Laird et al. (1998)	Sanger
70	(TAA)71	Laird et al. (1998)	Sanger
71	(TAA)72	Laird et al. (1998)	Sanger
72	(TAA)73	Laird et al. (1998)	Sanger
73	(TAA)74	Laird et al. (1998)	Sanger
74	(TAA)75	Laird et al. (1998)	Sanger
75	(TAA)76	Laird et al. (1998)	Sanger
76	(TAA)77	Laird et al. (1998)	Sanger
77	(TAA)78	Laird et al. (1998)	Sanger
78	(TAA)79	Laird et al. (1998)	Sanger
79	(TAA)80	Laird et al. (1998)	Sanger
80	(TAA)81	Laird et al. (1998)	Sanger
81	(TAA)82	Laird et al. (1998)	Sanger
82	(TAA)83	Laird et al. (1998)	Sanger
83	(TAA)84	Laird et al. (1998)	Sanger
84	(TAA)85	Laird et al. (1998)	Sanger
85	(TAA)86	Laird et al. (1998)	Sanger
86	(TAA)87	Laird et al. (1998)	Sanger
87	(TAA)88	Laird et al. (1998)	Sanger
88	(TAA)89	Laird et al. (1998)	Sanger
89	(TAA)90	Laird et al. (1998)	Sanger
90	(TAA)91	Laird et al. (1998)	Sanger
91	(TAA)92	Laird et al. (1998)	Sanger
92	(TAA)93	Laird et al. (1998)	Sanger
93	(TAA)94	Laird et al. (1998)	Sanger
94	(TAA)95	Laird et al. (1998)	Sanger
95	(TAA)96	Laird et al. (1998)	Sanger
96	(TAA)97	Laird et al. (1998)	Sanger
97	(TAA)98	Laird et al. (1998)	Sanger
98	(TAA)99	Laird et al. (1998)	Sanger
99	(TAA)100	Laird 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
