

Sequencing STRs in forensic mixtures: Current perspective on the benefits and challenges

Green Mountain DNA Conference
Burlington, VT
August 2, 2016

Katherine Gettings PhD
Applied Genetics Group

Dr. Pete Vallone
Lisa Borsuk
Kevin Kiesler
Becky (Hill) Steffen

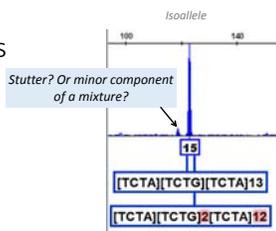
Sequencing STRs in forensic mixtures: Current perspective on the benefits and challenges

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.

STR Sequence Applications

- One to one matching?
 - new core loci = higher CE stats
 - Partial profiles
 - Kinship
- Degraded samples?
 - Smaller PCR amplicons relative to CE
- Mixtures
 - Resolve alleles identical by length, but differ by sequence
 - Separate stutter from low level contributors (based on sequence)



Sampling of forensic NGS kits

- Illumina ForenSeq (FGx)
 - 27 aSTR + 24 YSTR + 7 XSTR + Amel
 - 94 HID-SNPs + 56 ancestry SNPs + 22 pheno SNPs
 - Multiplex B – all markers
 - Multiplex A – STRs and HID SNPs
- Promega PowerSeq kits (Auto, Y, Mito)
 - 22 aSTR + DYS391 + Amel
 - 22 aSTR + 21 YSTR
- Thermo Fisher Precision ID (STR, SNP, Mixtures & Mito)
 - Mixture panel: STR, MH, IISNP
 - 29 aSTR + DYS391 + Amel
 - 124 Identity SNPs
 - 165 Ancestry SNPs
 - Mito control region or whole genome

Run on the FGW or SS

Run on the MiSeq/FGx



STR Sequence Gains

- PowerSeq Auto (MiSeq)
- 183 samples (3 US pops)
- STRait Razor (& ExactID)



- ForenSeq (FGx)
- Currently working sequencing NIST population samples (>1000)
- Sequencing nearly complete
- Bioinformatic checks
- CE data backfilling
- The goal is to provide the allele frequencies of the sequenced STRs

A little bit of informatics

Recognition Site-Based Informatics for STRs



Software returns:
The length between the recognition sequences (= 24)
A reference table returns a "6" allele **and the sequence between the recognition sites**

- PCR primers
- Recognition site (~10 nt)
- STR repeat region

¹STRait Razor: a length-based forensic STR allele-calling tool for use with second generation sequencing data. [Worshauer et al. Forensic Sci Int Genet. 2013 \(7\):409-17](#)
²STRait Razor v2.0: the improved STR allele identification Tool-Razor. [Worshauer et al. Forensic Sci Int Genet. 2015 \(14\):182-6](#)
³<http://atellforensictd.org/>

Recognition Site-Based Informatics for STRs



Moving the recognition sites out further:
Captures the flanking region SNPs and indels
Still returns the allele length and **more sequence between the recognition sites**

- PCR primers
- Recognition site (~10 nt)
- STR repeat region

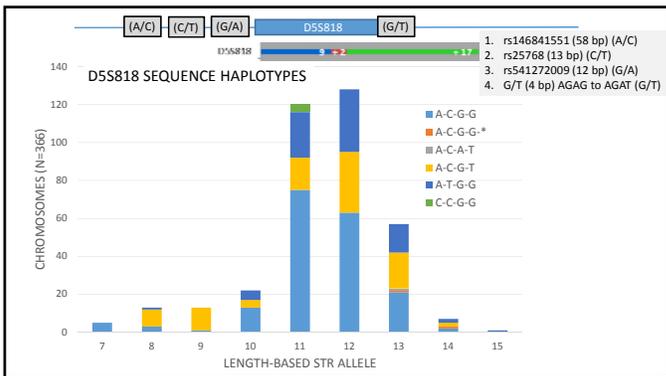
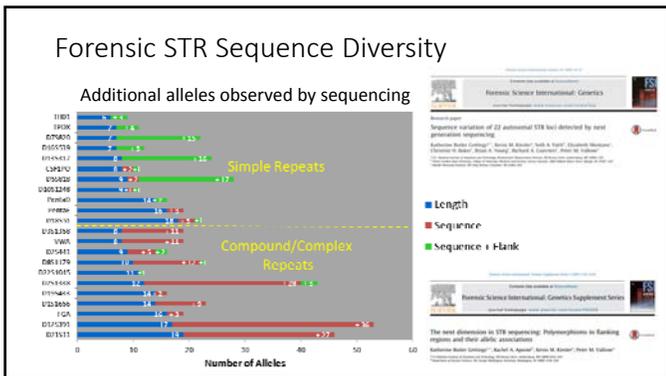
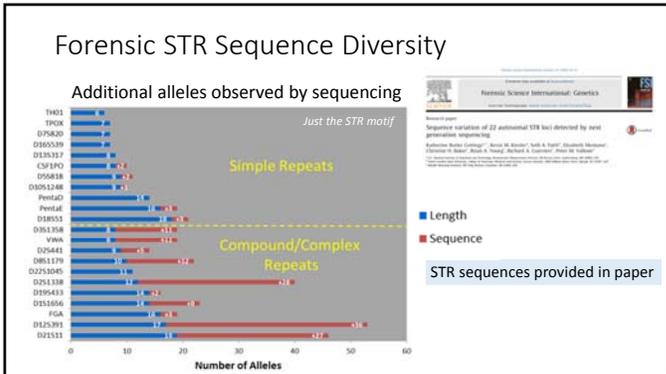
STR sequence bracketing (custom scripts)

- Makes it easier for humans to 'read' the STR sequence
- Check known list of sequence orientation
 - reverse complement sequence if needed
- Select known motif(s) specific for locus
 - **D2S1338 – GGAA, GGCA**
- 'Move' through sequence identifying repeats
 - **GGAAGGAAGGAAGGAAGGAAGGAAGGAAGGAAGGAAGGAAGGACAGGCAGGCAGGCAGGCA**
 - **[GGAA]11 GGACAGGCAGGCAGGCAGGCAGGCA**
 - **[GGAA]11 [GGCA]6**
- Additional trimming for sequences with known flanking sequence present
- Additional processing for loci which require specific handling – **special rules and exceptions**



Data analysis

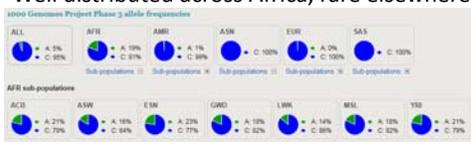
- FASTQ files were analyzed by STRaitRazor
 - Allele not present in lookup table
 - SNP or InDel caused a 'bioinformatic null'
 - Presence of an InDel resulted in a length difference from CE
- STRAitRazor data was parsed with custom scripts
- Allele calling: majority coverage and/or het balance (0.4)
 - Data compared to CE calls; discordant allele calls addressed
- Allele sequences were counted for each locus



Category	Locus	rs Number	Minor Allele Frequency		
			African American	Caucasian	Hispanic
Population / Allele Specific SNPs	TPOX	rs13422969	0.135	-	0.011
		rs115644759	0.022	-	-
		rs149212737	-	0.007	-

TPOX: **rs13422969** and two additional rare SNPs

- Well distributed across Africa, rare elsewhere

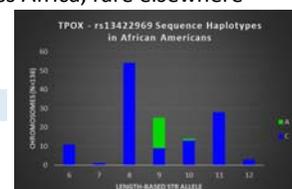


Category	Locus	rs Number	Minor Allele Frequency		
			African American	Caucasian	Hispanic
Population / Allele Specific SNPs	TPOX	rs13422969	0.135	-	0.011
		rs115644759	0.022	-	-
		rs149212737	-	0.007	-

TPOX: **rs13422969** and two additional rare SNPs

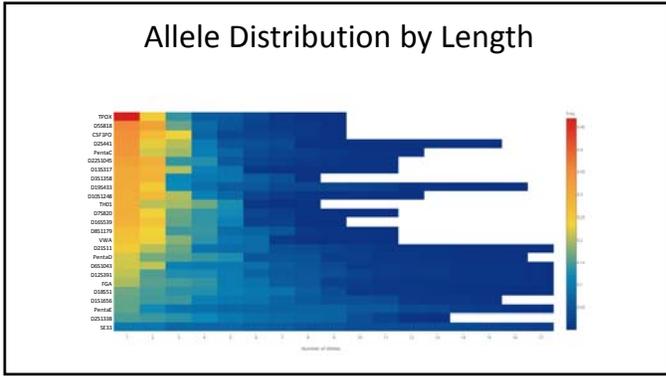
- Well distributed across Africa, rare elsewhere

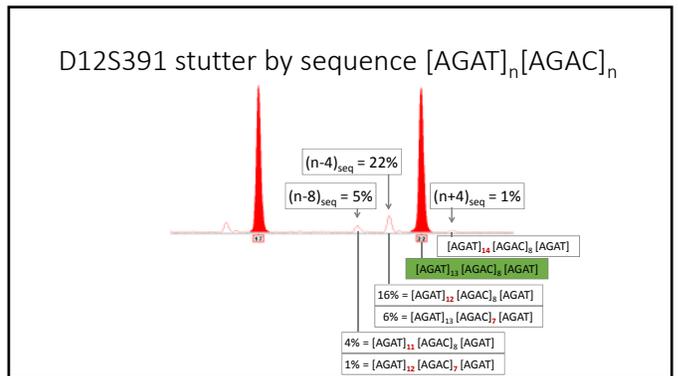
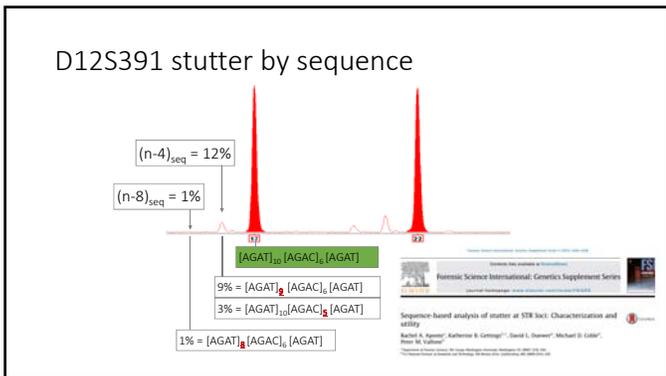
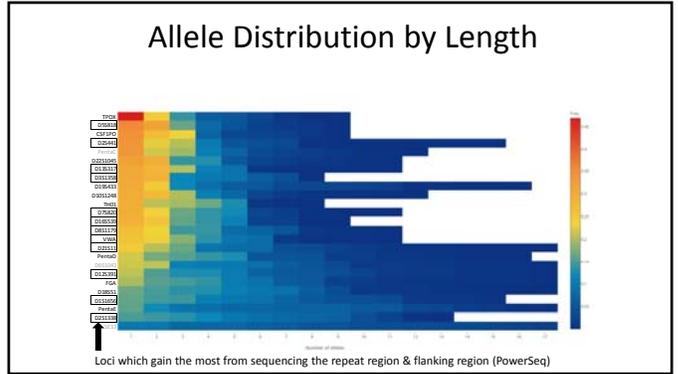
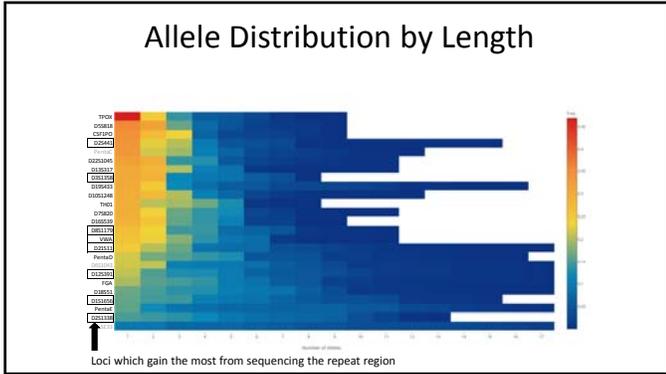
Associated with "9" allele



Flanking region SNP categories from paper

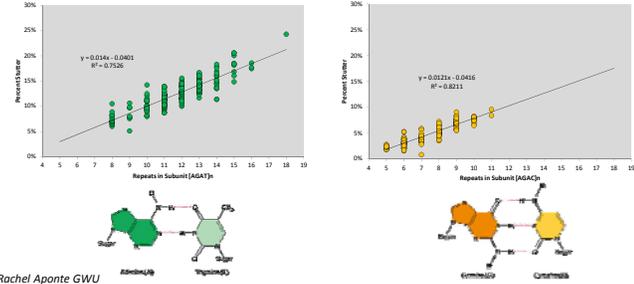
- Multiple polymorphisms in haplotype (D5S818 example)
- "Old" single polymorphisms
- Population or allele specific polymorphisms (TPOX example)
- Polymorphisms associated with STR sequence variants
- Rare polymorphisms (< 5%)
- No polymorphisms



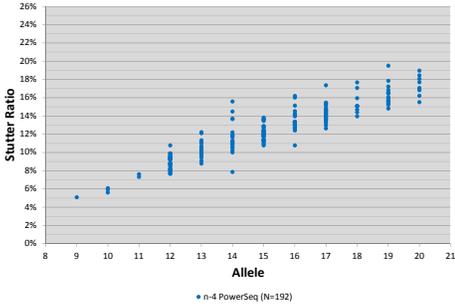
D12S391 stutter by sequence

[AGAT]_n trends ~3% higher stutter than [AGAC]_n

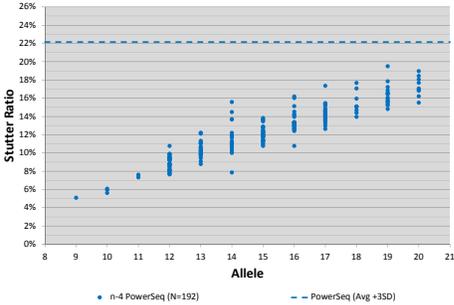


Rachel Aponte GWU

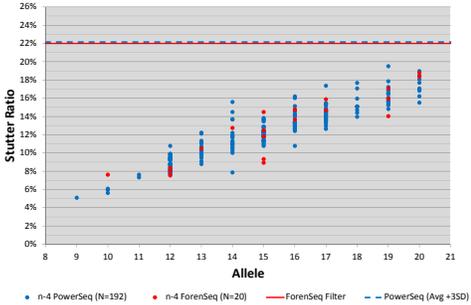
D18S51 n-4 Stutter by Kit

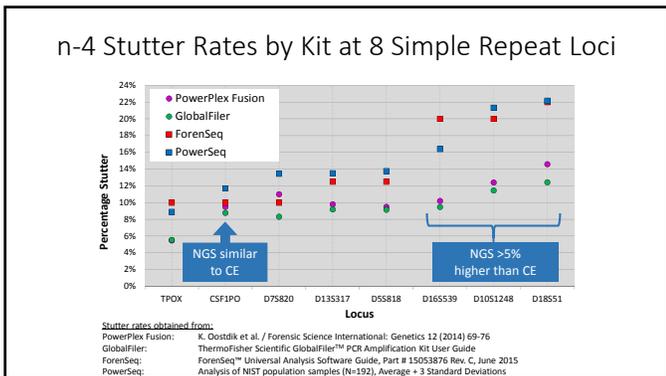
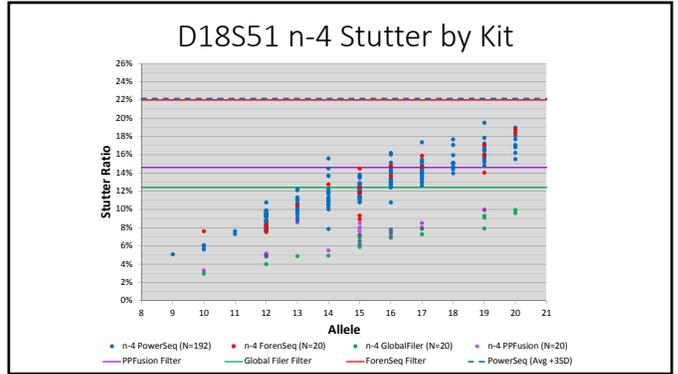
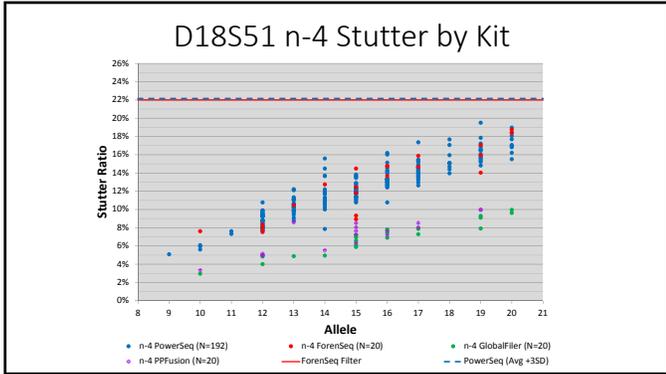


D18S51 n-4 Stutter by Kit



D18S51 n-4 Stutter by Kit





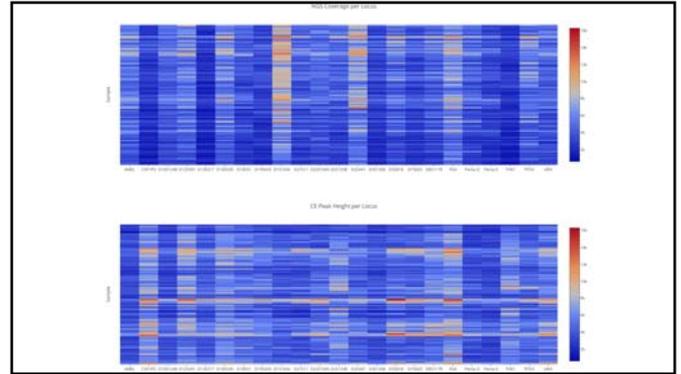
NGS Mixture Study

Are mixture ratios by NGS the same as mixture ratios by CE?

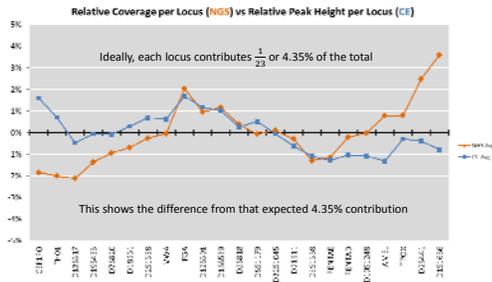
	CE		NGS
Loci	PowerPlex Fusion + PowerPlex Y23		PowerSeq Auto + Y
Input DNA	0.5 ng each		0.5 ng total
Amp Parameters	30 cycles		30 cycles, same as PPF
Everything Else	3500xL		TruSeq PCR free Library Prep, MiSeq v3

NGS Mixture Study

	1	2	3	4	5	6	7	8	9	10	11	12
A	A	B	M	1:1 AB								
B	1:1 AB	1:1 AM	1:1 BM	1:1 BA								
C	3:1 AB	3:1 AB	3:1 AB	1:1 AM	1:1 AM	1:1 AM	1:1 BM	1:1 BM	1:1 BM	1:1 AB	1:1 AB	1:1 AB
D	3:1 BA	3:1 BA	3:1 BA	1:1 MA	1:1 MA	1:1 MA	1:1 MB	1:1 MB	1:1 MB	1:1 AB	1:1 AB	1:1 AB
E	3:1 AM	3:1 AM	3:1 AM	1:1 BM	1:1 BM	1:1 BM	1:1 AB					
F	3:1 MA	3:1 MA	3:1 MA	1:1 MB	1:1 MB	1:1 MB	1:1 AB					
G	3:1 BM	3:1 BM	3:1 BM	1:1 AB								
H	3:1 MB	3:1 MB	3:1 MB	1:1 BA	1:1 BA	1:1 BA	1:1 AB					



NGS Mixture Study



NGS Mixture Study

Are mixture ratios by NGS the same as mixture ratios by CE?

1:1:1 Mixture

3 alleles "10"
6 alleles total = 0.5



CSF1PO	Expected	Rep 1	Rep 2	Rep 3	Average
10	0.500	0.457	0.543	0.501	0.501
11	0.167	0.248	0.185	0.172	0.202
12	0.333	0.294	0.271	0.327	0.298

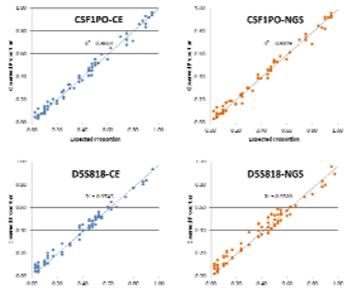
$$\frac{3490}{7630} = 0.457$$

D5S818	A	B	M
	11,12	12,13	12,12

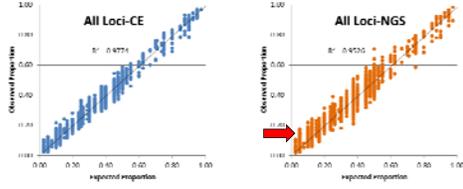
D5S818	Expected	Rep 1	Rep 2	Rep 3	Average
11	0.167	0.217	0.278	0.285	0.257
12	0.667	0.605	0.659	0.533	0.599
13	0.167	0.177	0.163	0.182	0.174

- CE: no alleles below 75 RFU
- NGS: no alleles below 75X coverage
- Average of 3 replicates

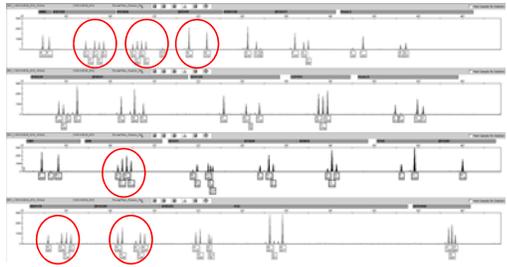
Expected vs Observed Mixture Contributions



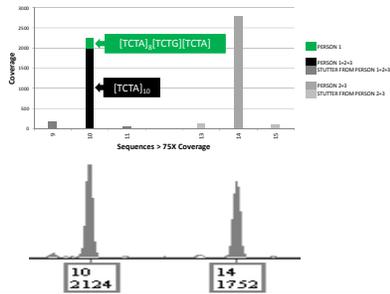
Expected vs Observed Mixture Contributions

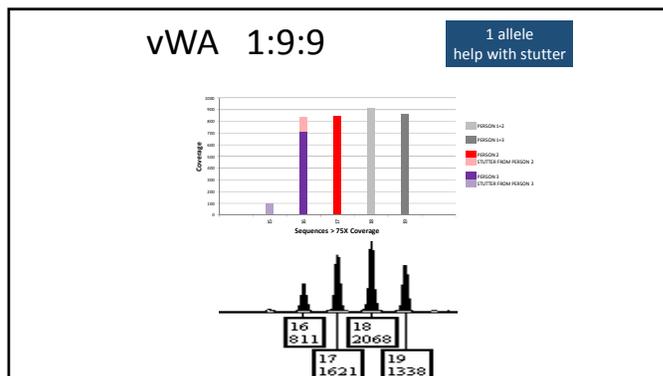
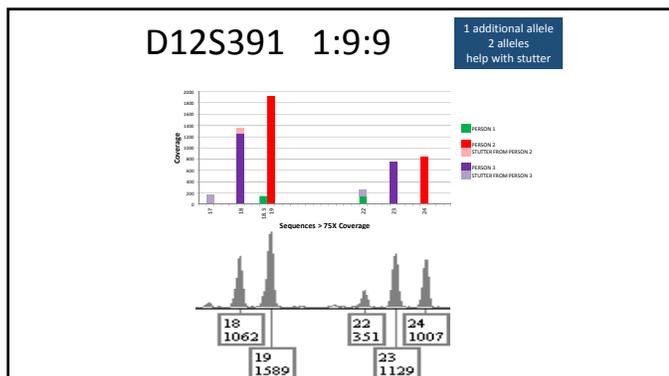
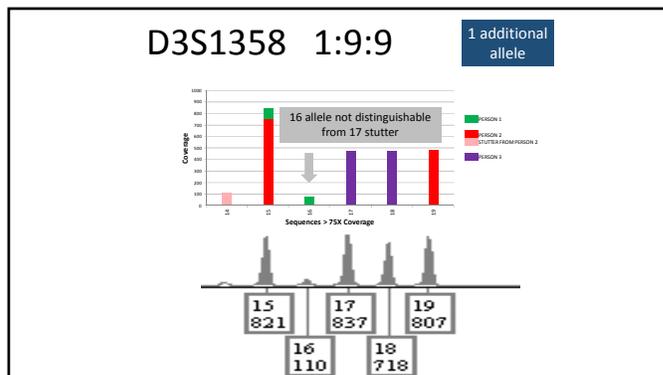
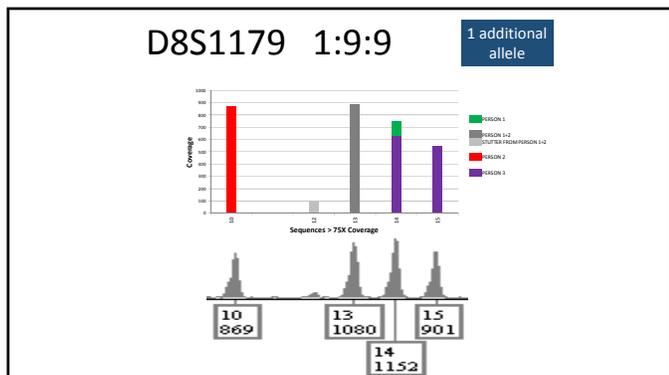


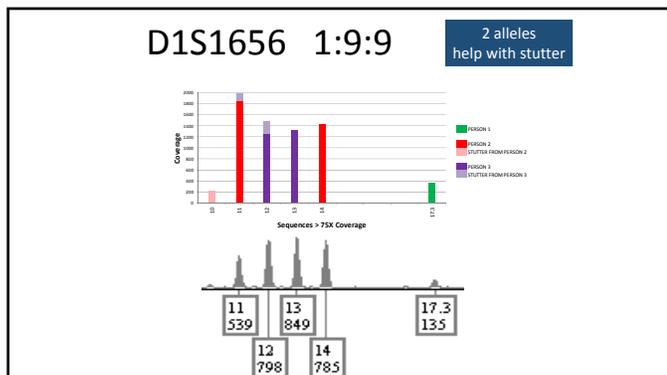
NGS Mixture Study 1:9:9



D2S441 1:9:9







Summary 1:9:9

Locus	Additional Alleles	Help with Stutter
D2S441	1	
D8S1179	1	
D3S1358	1	
D12S391	1	1
vWA	1	1
D1S1656	2	2

NGS profile contains four additional alleles and improved stutter attribution for four alleles

NGS Implications for Mixtures

Conclusions

- Sequencing forensic STR loci can uncover underlying sequence variation in the repeat and flanking regions
- This will increase allelic diversity, thus increasing the ability to discriminate among individuals in a mixture
- Additionally, sequence specific stutter ratios may improve mixture models

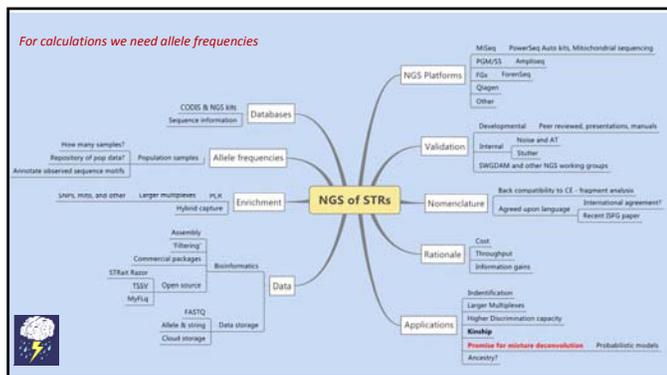
NGS Implications for Mixtures

Conclusions

Prior to implementation:

- Sequence-based allele frequency databases
- Characterization of any NGS-specific effects on
 - peak height ratios
 - minor components/stutter
- Probabilistic genotyping software amenable to sequence data (and sequence-based stutter!)

← assay and locus specific



On November 8-9, 2016, NIST will host a two-day symposium showcasing NIST's Forensic Science Center of Excellence's work and NIST's latest research addressing forensic science. Attendees will learn how NIST's world-class laboratories and staff support many forensic science disciplines.

Agenda:

- November 8, 2016: The latest research in the areas of ballistics, biometrics, DNA, drug analysis, trace, and statistics will be highlighted.
- November 9, 2016: The Forensic Science Center of Excellence: Center for Statistics and Applications in Forensic Science (CSAFE) will have the opportunity to showcase its work in the application of probabilistics to pattern evidence and digital evidence.

Acknowledgments

- Rachel Aponte (GWU – now at Bode)
- Funding
 - NIST - Forensic DNA
 - FBI - DNA as a Biometric

Dr. Pete Vallone
 Lisa Borsuk
 Kevin Kiesler
 Becky (Hill) Steffen
 Dr. Mike Coble

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 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.

katherine.gettings@nist.gov