

# Sequencing STRs: Variation and Nomenclature

Forensics@NIST  
December 3, 2014

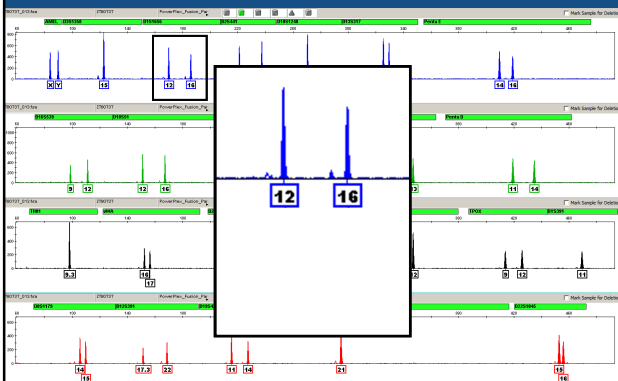
Katherine Butler Gettings, Ph.D.  
Applied Genetics Group  
Biomolecular Measurement Division



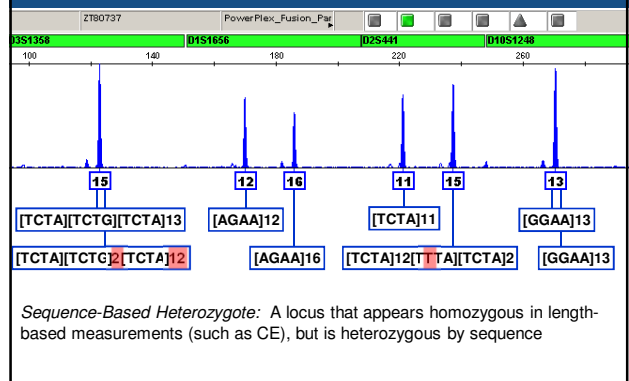
## Disclaimer

Points of view in this document are those of the authors and do not necessarily represent the official position or policies of the U.S. Department of Commerce. Certain commercial equipment, instruments, and materials are identified in order to specify experimental procedures as completely as possible. In no case does such identification imply a recommendation or endorsement by NIST, nor does it imply that any of the materials, instruments, or equipment identified are necessarily the best available for the purpose.

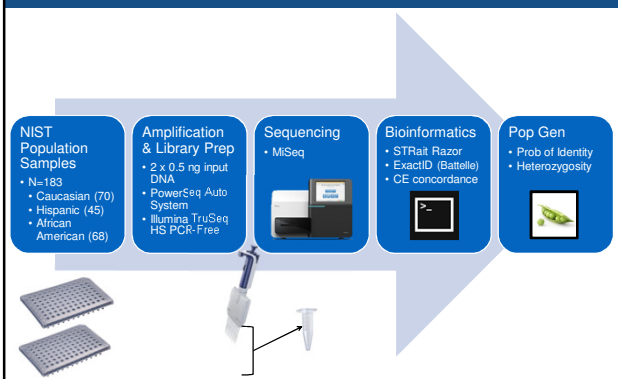
## Forensic STR Sequence Diversity



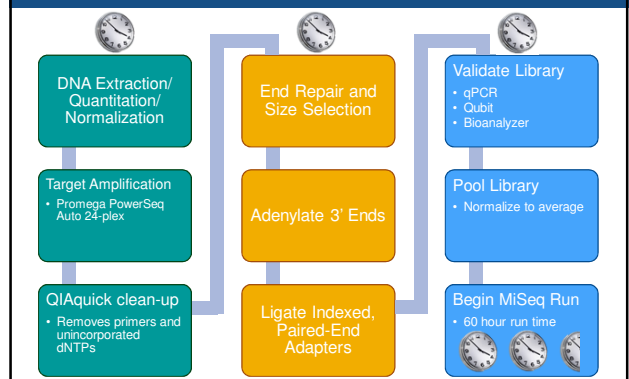
## Forensic STR Sequence Diversity



## Forensic STR Sequence Diversity



## Forensic STR Sequence Diversity





## Forensic STR Sequence Diversity

### Additional Alleles by Sequence

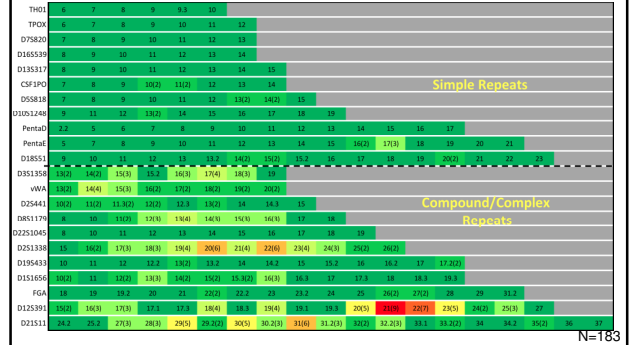
CSF1PO		
7	[AGAT]7	AGAT AGAT AGAT AGAT AGAT AGAT AGAT
8	[AGAT]8	AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT
9	[AGAT]9	AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT
10	[AGAT]10	AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT
10	[AGAT]10	AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT
11	[AGAT]11	AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT
11	[AGAT]11	AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT AGAT
11	[AGAT]3AGGT[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

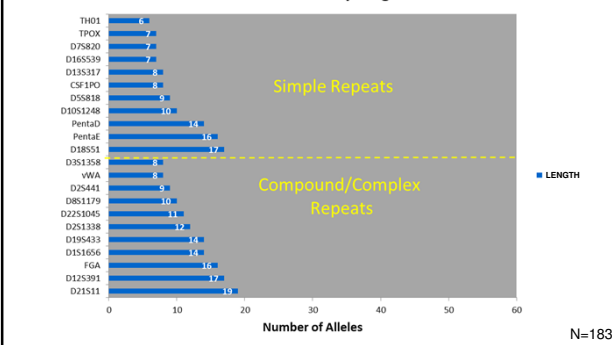
### Additional Alleles Obtained by Sequencing



N=183

## Forensic STR Sequence Diversity

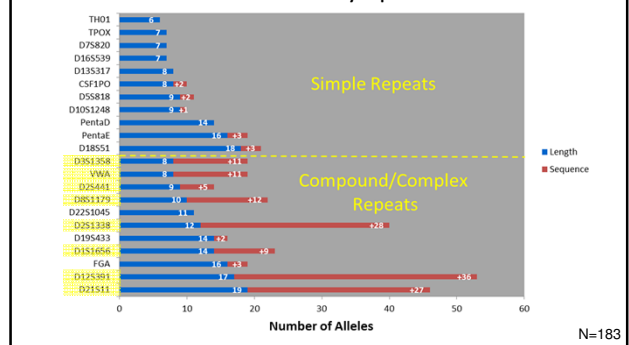
### Alleles Obtained by Length



N=183

## Forensic STR Sequence Diversity

### Alleles Obtained by Sequence



N=183

## Forensic STR Sequence Diversity

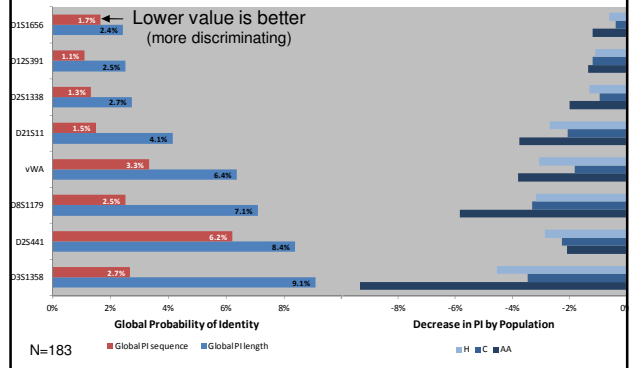
### Probability of Identity

Sum of each genotype frequency<sup>2</sup> 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

## Forensic STR Sequence Diversity



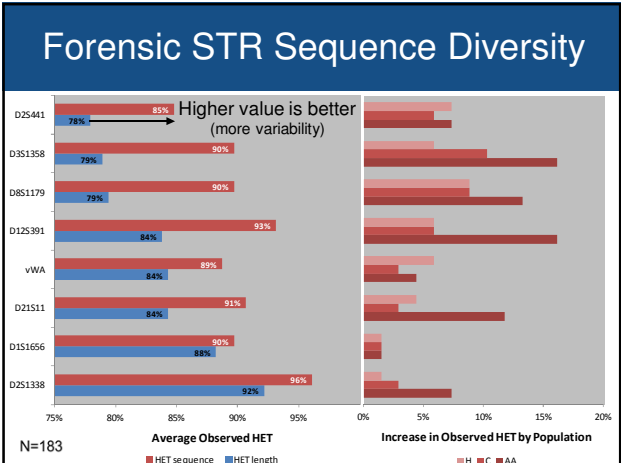
N=183

## Forensic STR Sequence Diversity

**Heterozygosity**

$$\frac{\# \text{ heterozygotes observed}}{\# \text{ of loci tested}}$$

Indicates genetic variability at a locus



## Forensic STR Sequence Diversity

**Conclusions**

- Sequencing forensic STR loci in a HTP manner is possible (automation is needed)
- Bioinformatic tools are in their infancy, testing across platforms and pipelines is important
- At some loci, sequencing will offer significant gains ("core set" for mixture analysis)
- Extending analysis to the flanking regions will increase effective number of alleles
- Infrastructure such as **nomenclature guidelines** and allele frequency databases are needed prior to implementation

## Forensic STR Sequence Nomenclature

Options for representing sequence data, and possible applications:

- Complete Sequence String** entire string of generated sequence  

```

TGACTATGGAGTTATTTTAAAGGTTAATATATATAAAGGGTATGATAGAACACTTGTTCATAGTTTAGAAGCACTAAC
GATAGATAGATAGATAGATAGATAGATAGATAGATAGATAGACAGATTGATAGTTTTTTTTTATCTCAGTAAATAGCTATAGTA
AACATTTAATTACCAATATTTGGTGCAATTCGTC
            
```
- Bracketed sequence**
  - repetitive elements enclosed in brackets and a numeric representation of the repeat length  
 [AATG]<sub>6</sub> A-TG [AATG]<sub>3</sub> = TH01 9.3 allele
  - polymorphisms (SNPs or InDels) in flanking regions identified by "rs" number
- Unique Identifier**
  - 13d rs206432C where 13 = repeat length, d = sequence version, rs number = flank polymorphism
  - @j\*5 = computer-generated code applied to each unique sequence string within a defined region

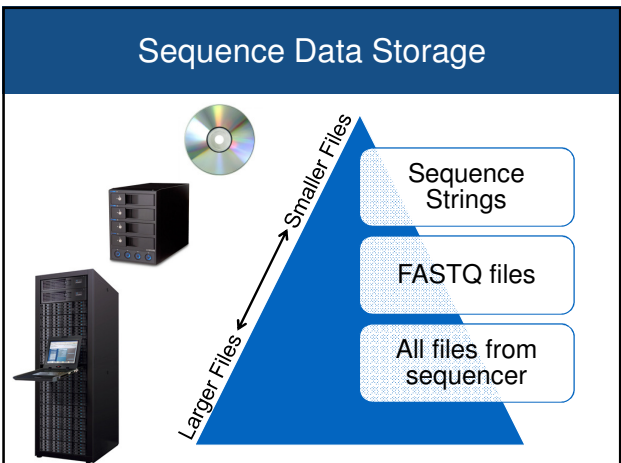
## Forensic STR Sequence Nomenclature

**Reporting/Manual Comparisons**

- Meaningful unique identifier (e.g. 13d) may be helpful for quick comparisons
- Bracketed sequence is intuitive and may help in explaining results to investigator
- Complete sequence could be appended to report

**Database Searching**

- Database searching must be unambiguous and computationally inexpensive (i.e. fast)
- Two most likely possibilities are unique identifier and complete sequence string



## Acknowledgements

### NIST

Peter Vallone  
Erica Butts  
Mike Coble  
David Duewer  
Jo Lynne Harenza  
Becky Hill  
Kevin Kiesler  
Margaret Kline  
Nate Olson  
Harish Swaminathan

Promega  
Doug Storts  
Jay Patel

Battelle  
Seth Faith (NCSU)  
Rich Guerrieri  
Brian Young

### Funding

FBI: DNA as a Biometric

Contact Information  
katherine.gettings@nist.gov