## NGS and Implications for Mixture Analysis

**Michael Coble, Ph.D.**
**Katherine Butler Gettings, Ph.D.**
Research Biologists
Applied Genetics Group

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

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

Updated slides:
http://www.cstl.nist.gov/biotech/strbase/pub_pres/
ISHI_NGS_Workshop_2015_Coble-Gettings.pdf

---

## NGS Implications for Mixtures
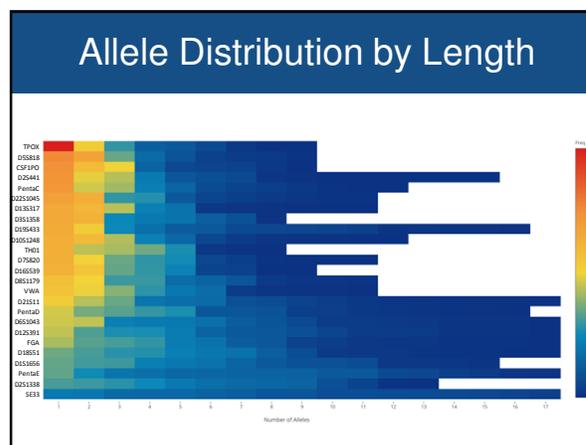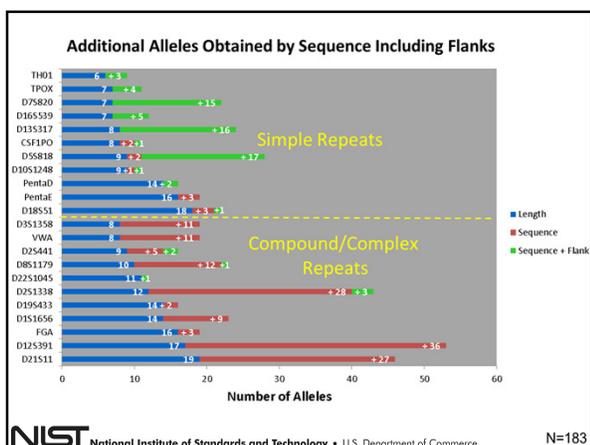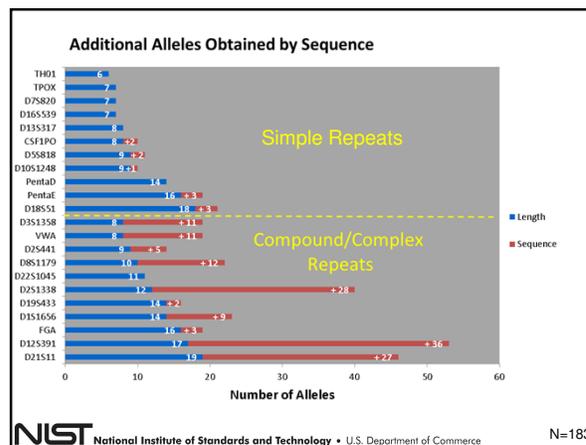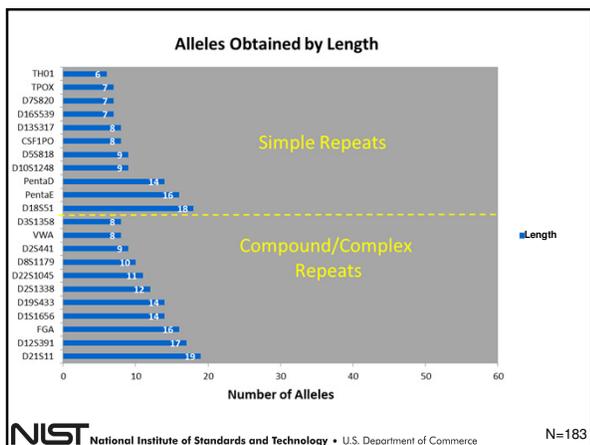### Questions

Questions of Utility
- Which STR loci have the most overlapping alleles?
- Which STR loci are the most likely to be aided by sequence?
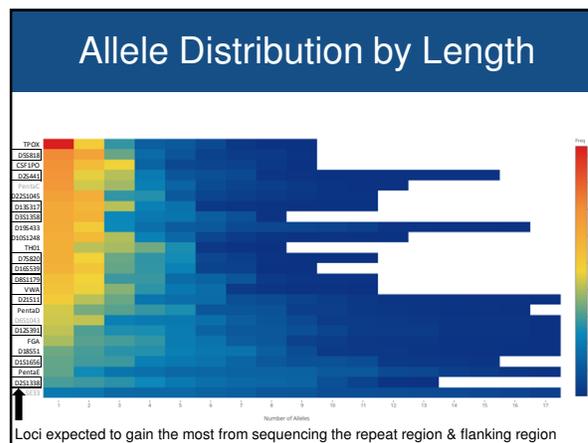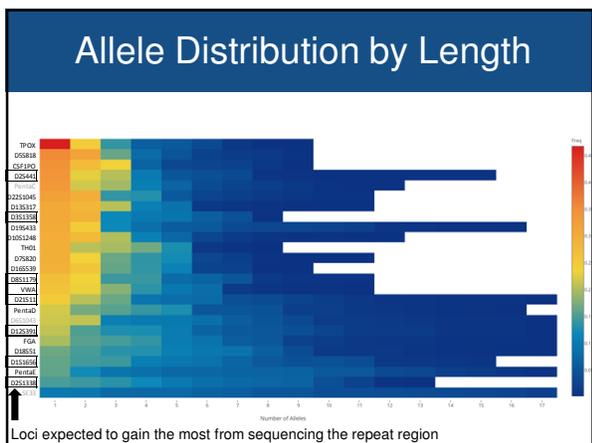  - repeat region vs flank

Validation Questions - General
- What are the appropriate analytical and stochastic thresholds for NGS data?
- Are PHR in NGS consistent with CE?
  - this can be greatly affected by library preparation size selection steps
- How many individuals do we need in sequence allele frequency databases?
  - How will we handle the increased population specificity in repeat region sequences and flanking SNPs?
- How will NGS data affect the interpretation of stutter artifacts?

Validation Questions – Mixture Specific
- How will NGS affect the determination of number of contributors?
  - Probabilistic software is making this moot for CE data
- Are mixture ratios in NGS consistent with CE?

---



Alleles Obtained by Length



Additional Alleles Obtained by Sequence



Additional Alleles Obtained by Sequence Including Flanks



Allele Distribution by Length

---

## Allele Distribution by Length



Loci expected to gain the most from sequencing the repeat region

## Allele Distribution by Length



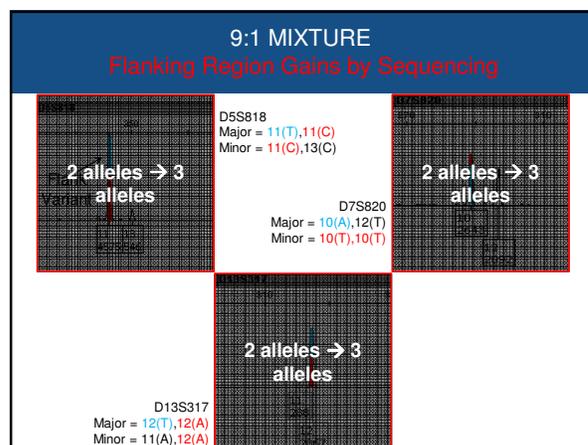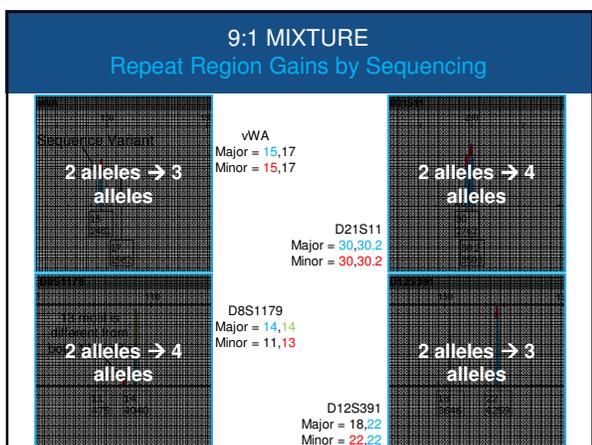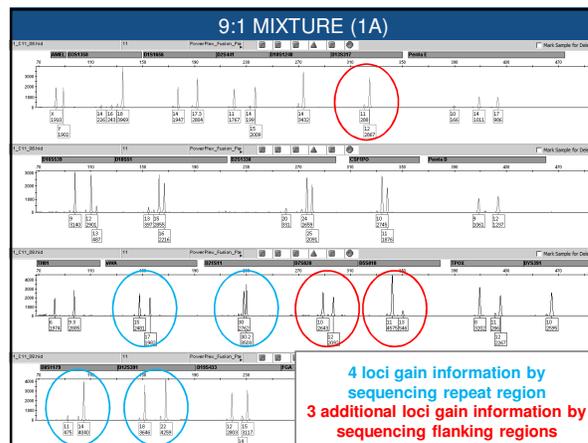Loci expected to gain the most from sequencing the repeat region & flanking region

## NGS of STR Mixtures – Proof of Concept
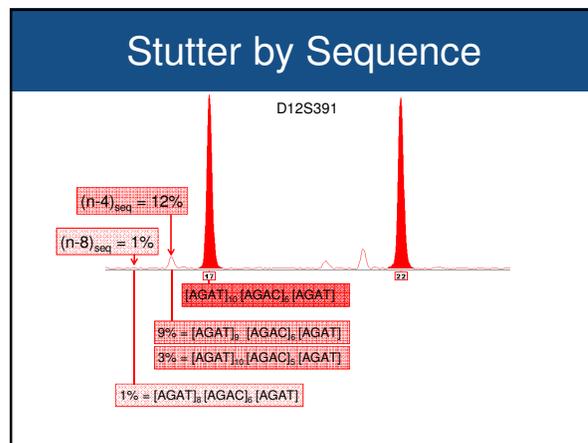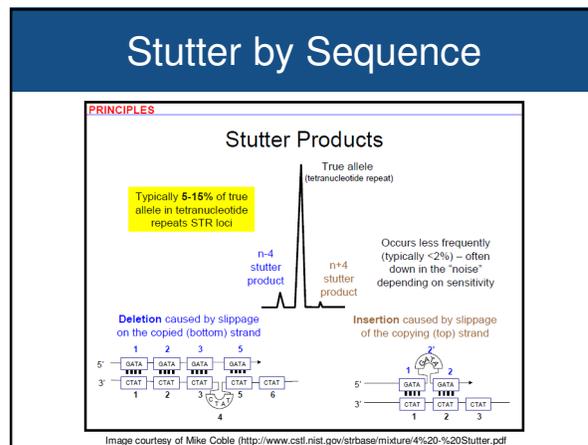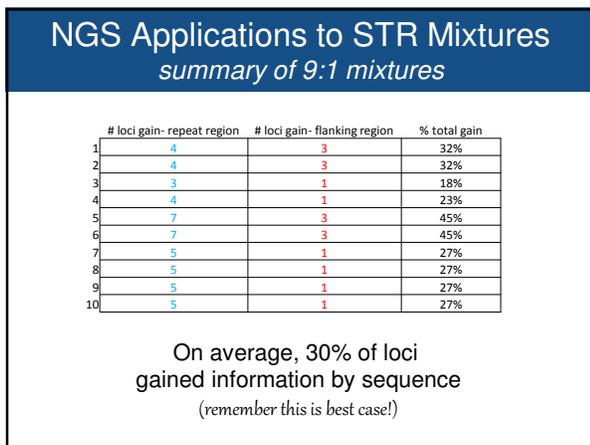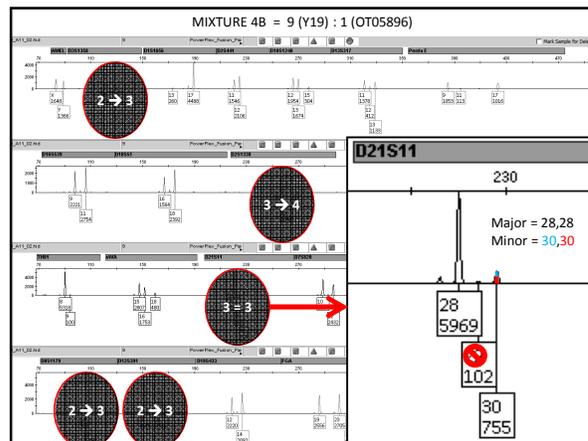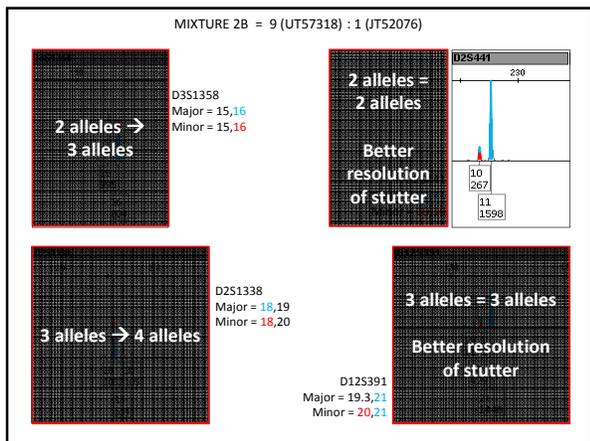
Created 10 difficult 2-person mixtures by CE
- 9:1 ratio
- Maximal overlapping alleles
- Very few loci have four distinguishable alleles by CE

Inferred sequences from NGS data

"Best case" help from NGS in 2-person mixtures

Does not include help with stutter

## 9:1 MIXTURE (1A)



**4 loci gain information by sequencing repeat region**
**3 additional loci gain information by sequencing flanking regions**

## 9:1 MIXTURE
### Repeat Region Gains by Sequencing



vWA
Major = 15,17
Minor = 15,17

D21S11
Major = 30,30.2
Minor = 30,30.2

D8S1179
Major = 14,14
Minor = 11,13

D12S391
Major = 18,22
Minor = 22,22

2 alleles → 3 alleles
2 alleles → 4 alleles
2 alleles → 4 alleles
2 alleles → 3 alleles

## 9:1 MIXTURE
### Flanking Region Gains by Sequencing



D5S818
Major = 11(T),11(C)
Minor = 11(C),13(C)

D7S820
Major = 10(A),12(T)
Minor = 10(T),10(T)

D13S317
Major = 12(T),12(A)
Minor = 11(A),12(A)

2 alleles → 3 alleles
2 alleles → 3 alleles
2 alleles → 3 alleles

http://www.cstl.nist.gov/biotech/strbase/pub_pres/ISHI_NGS_Workshop_2015_Coble-Gettings.pdf

MIXTURE 2B = 9 (UT57318) : 1 (JT52076)

2 alleles → 3 alleles

D3S1358
Major = 15,16
Minor = 15,16

2 alleles = 2 alleles

Better resolution of stutter

D2S441

3 alleles → 4 alleles

D2S1338
Major = 18,19
Minor = 18,20

3 alleles = 3 alleles

Better resolution of stutter

D12S391
Major = 19.3,21
Minor = 20,21



MIXTURE 4B = 9 (Y19) : 1 (OT05896)

2 → 3

3 → 4

3 = 3

2 → 3     2 → 3

D21S11
230
Major = 28,28
Minor = 30,30

## NGS Applications to STR Mixtures
### *summary of 9:1 mixtures*

| | # loci gain- repeat region | # loci gain- flanking region | % total gain |
|---|---|---|---|
| 1 | 4 | 3 | 32% |
| 2 | 4 | 3 | 32% |
| 3 | 3 | 1 | 18% |
| 4 | 4 | 1 | 23% |
| 5 | 7 | 3 | 45% |
| 6 | 7 | 3 | 45% |
| 7 | 5 | 1 | 27% |
| 8 | 5 | 1 | 27% |
| 9 | 5 | 1 | 27% |
| 10 | 5 | 1 | 27% |

On average, 30% of loci
gained information by sequence
*(remember this is best case!)*

## Stutter by Sequence



Image courtesy of Mike Coble (http://www.cstl.nist.gov/strbase/mixture/4%20-%20Stutter.pdf

## Stutter by Sequence



Sequence analysis and characterization of stutter
products at the tetranucleotide repeat locus vWA
P. Sean Walsh*,*, Nicola J. Fildes* and Rebecca Reynolds
*Nucleic Acids Research, 1996, Vol. 24, No. 14*

## Stutter by Sequence



D12S391

$(n-4)_{seq}$ = 12%

$(n-8)_{seq}$ = 1%

[AGAT]$_{10}$ [AGAC]$_6$ [AGAT]

9% = [AGAT]$_9$ [AGAC]$_6$ [AGAT]

3% = [AGAT]$_{10}$ [AGAC]$_5$ [AGAT]

1% = [AGAT]$_8$ [AGAC]$_6$ [AGAT]

## Stutter by Sequence

D12S391

$(n-4)_{seq} = 22\%$

$(n-8)_{seq} = 5\%$

$(n+4)_{seq} = 1\%$

17     22

$[AGAT]_{14} [AGAC]_6 [AGAT]$

$[AGAT]_{13} [AGAC]_6 [AGAT]$

$16\% = [AGAT]_{12} [AGAC]_6 [AGAT]$

$6\% = [AGAT]_{13} [AGAC]_7 [AGAT]$

$4\% = [AGAT]_{11} [AGAC]_6 [AGAT]$

$1\% = [AGAT]_{12} [AGAC]_7 [AGAT]$

## Stutter by Sequence

**D8S1179**

| Allele | Repeat Structure |
|---|---|
| | **[TCTA]10–14** |
| 10 | [TCTA]10 |
| 11 | [TCTA]11 |
| 12 | [TCTA]12 |
| 13 | [TCTA]13 |
| 14 | [TCTA]14 |

## Stutter by Sequence

D8S1179 : Stutter Ratio vs Allele Repeat Number

- [TCTA]n

## Stutter by Sequence

**D8S1179**

| Allele | Repeat Structure |
|---|---|
| | **[TCTA]10–14** |
| 10 | [TCTA]10 |
| 11 | [TCTA]11 |
| 12 | [TCTA]12 |
| 13 | [TCTA]13 |
| 14 | [TCTA]14 |
| | **[TCTA][TCTG][TCTA]10–14** |
| 12 | [TCTA][TCTG][TCTA]10 |
| 13 | [TCTA][TCTG][TCTA]11 |
| 14 | [TCTA][TCTG][TCTA]12 |
| 16 | [TCTA][TCTG][TCTA]14 |

## Stutter by Sequence

D8S1179 : Stutter Ratio vs Allele Repeat Number

- [TCTA]n
- [TCTA][TCTG][TCTA]n

## Stutter by Sequence

**D8S1179**

| Allele | Repeat Structure |
|---|---|
| | **[TCTA]10–14** |
| 10 | [TCTA]10 |
| 11 | [TCTA]11 |
| 12 | [TCTA]12 |
| 13 | [TCTA]13 |
| 14 | [TCTA]14 |
| | **[TCTA][TCTG][TCTA]10–14** |
| 12 | [TCTA][TCTG][TCTA]10 |
| 13 | [TCTA][TCTG][TCTA]11 |
| 14 | [TCTA][TCTG][TCTA]12 |
| 16 | [TCTA][TCTG][TCTA]14 |
| | **[TCTA][TCTA][TCTG][TCTA]8–15** |
| 11 | [TCTA][TCTA][TCTG][TCTA]8 |
| 12 | [TCTA][TCTA][TCTG][TCTA]9 |
| 13 | [TCTA][TCTA][TCTG][TCTA]10 |
| 14 | [TCTA][TCTA][TCTG][TCTA]11 |
| 15 | [TCTA][TCTA][TCTG][TCTA]12 |
| 16 | [TCTA][TCTA][TCTG][TCTA]13 |
| 17 | [TCTA][TCTA][TCTG][TCTA]14 |
| 18 | [TCTA][TCTA][TCTG][TCTA]15 |

http://www.cstl.nist.gov/biotech/strbase/pub_pres/ISHI_NGS_Workshop_2015_Coble-Gettings.pdf

Stutter by Sequence



Stutter by Sequence



Stutter by Sequence



Stutter by Sequence



Stutter by Sequence – D8S1179



Stutter by Sequence – D2S441

http://www.cstl.nist.gov/biotech/strbase/pub_
pres/ISHI_NGS_Workshop_2015_Coble-
Gettings.pdf

## Stutter by Sequence – D2S441



## 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





## NGS Mixture Study



## NGS Mixture Study

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

1:1:1 Mixture



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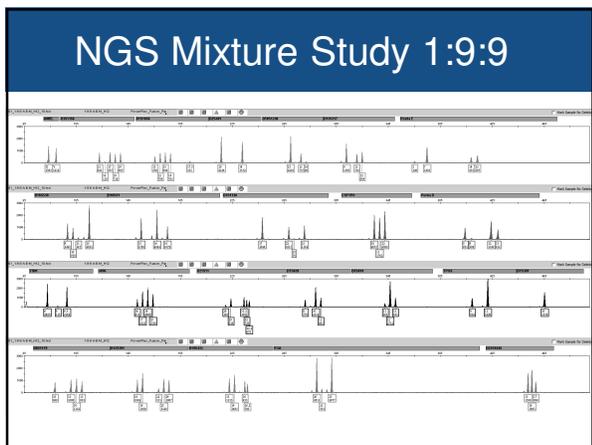- Average of 3 replicates

http://www.cstl.nist.gov/biotech/strbase/pub_
pres/ISHI_NGS_Workshop_2015_Coble-
Gettings.pdf

Expected vs Observed Mixture Contributions


Expected vs Observed Mixture Contributions


NGS Mixture Study 1:9:9


D2S441   1:9:9   1 additional allele


D8S1179   1:9:9   1 additional allele


D3S1358   1:9:9   1 additional allele

http://www.cstl.nist.gov/biotech/strbase/pub_pres/ISHI_NGS_Workshop_2015_Coble-Gettings.pdf

## D12S391 1:9:9

1 additional allele
2 alleles
help with stutter



## vWA 1:9:9

1 allele
help with stutter



## D1S1656 1:9:9

2 alleles
help with stutter



## Summary 1:9:9

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

*The gain is difficult to quantify*

Prior to implementation:

- Sequence-based allele frequency databases
- Characterization of peak height ratios and stutter by NGS (assay and locus specific)
- Probabilistic genotyping software amenable to sequence data (and sequence-based stutter!)

http://www.cstl.nist.gov/biotech/strbase/pub_
pres/ISHI_NGS_Workshop_2015_Coble-
Gettings.pdf

Acknowledgements

**NIST**
Pete Vallone
Kevin Kiesler
Lisa Borsuk
Becky Hill
Margaret Kline

**Student Interns**
Rachel Aponte (GWU)
Harish Swaminathan
(Rutgers)

**Promega**
Doug Storts
Jay Patel

Updated slides:
http://www.cstl.nist.gov/biotech/strbase/pub_pres/
ISHI_NGS_Workshop_2015_Coble-Gettings.pdf

**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**

Contact Information
katherine.gettings@nist.gov
michael.coble@nist.gov

http://www.cstl.nist.gov/biotech/strbase/pub_pres/ISHI_NGS_Workshop_2015_Coble-Gettings.pdf