**DNA Interpretation Workshop 2**

# Probabilistic Genotyping

Michael D Coble, PhD
U.S. National Institute of Standards and Technology (NIST)

http://www.cstl.nist.gov/strbase/training.htm

**ISFG Pre-Conference Workshop**
**Melbourne, Australia**
**September 2-3, 2013**

ISFG

NIST

---

## NIST and NIJ Disclaimer

---

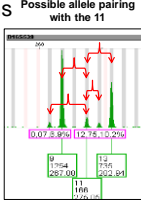Whatever way uncertainty is approached, probability is the *only* sound way to think about it.

-Dennis Lindley

---

## Do You Have Uncertainty in Your Data?

- **If allele dropout is a possibility** (e.g., in a partial profile), then there is uncertainty in whether or not an allele is present in the sample…and therefore what genotype combinations are possible

**Possible allele pairing with the 11**

- **If different allele combinations are possible** in a mixture, then there is uncertainty in the genotype combinations that are possible…

---

## Uncertainty and Probability

- "Contrary to what many people think, **uncertainty is present throughout any scientific procedure**."
    – Dennis V. Lindley, in his foreword to Aitken & Taroni (2004)
    *Statistics and the Evaluation of Evidence for Forensic Scientists, Second Edition*

- "It is now recognized that **the only tool for handling uncertainty is probability**."
    – Dennis V. Lindley, in his foreword to Aitken & Taroni (2004)
    *Statistics and the Evaluation of Evidence for Forensic Scientists, Second Edition*

---

## "On the Threshold of a Dilemma"

- Gill and Buckleton (2010)
- Although most labs use thresholds of some description, this philosophy has always been problematic because there is an inherent illogicality which we call the falling off the cliff effect.

JOURNAL OF **FORENSIC SCIENCES**

**Commentary on:** Budowle B, Onorato AJ, Callaghan TF, Della Manna A, Gross AM, Guerrieri RA, Luttman JC, McClure DL. Mixture interpretation: defining the relevant features for guide-lines for the assessment of mixed DNA profiles in forensic casework. J Forensic Sci 2009;54(4):810–21.

*J Forensic Sci*, January 2010, Vol. 55, No. 1
doi: 10.1111/j.1556-4029.2009.01257.x
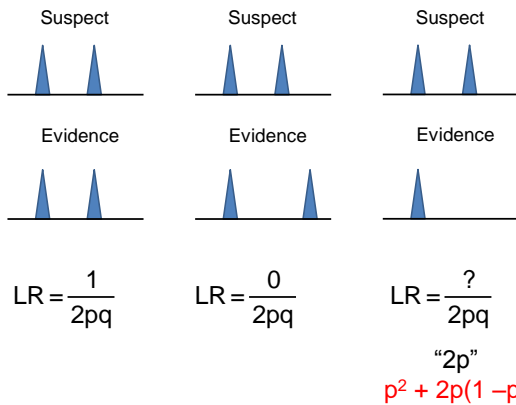Available online at: interscience.wiley.com

## "Falling off the Cliff Effect"

- If T = an arbitrary level (e.g., 150 rfu), an allele of 149 rfu is subject to a different set of guidelines compared with one that is 150 rfu even though they differ by just 1 rfu (Fig. 1).



Gill and Buckleton *JFS* 55: 265-268 (2010)

## Gill and Buckleton *JFS* **55:** 265-268 (2010)

- "The purpose of the ISFG DNA commission document was to provide a way forward to demonstrate the use of ***probabilistic models to circumvent the requirement for a threshold*** and to safeguard the legitimate interests of defendants."



$$LR = \frac{1}{2pq} \qquad LR = \frac{0}{2pq} \qquad LR = \frac{?}{2pq}$$

"2p"
$p^2 + 2p(1-p)$

## What should we do with discordant data?

- Continue to use RMNE (CPI, CPE) (not optimal)
- Use the Binary LR with 2p (not optimal)
- Semi-continuous methods with a LR (Drop models)

## Some Drop Model Examples

- LR mix (Haned and Gill)
- Balding (likeLTD - R program)
- FST (NYOCME, Mitchell *et al.*)
- Kelly *et al.* (University of Auckland, ESR)
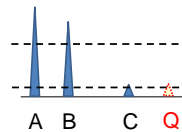- Lab Retriever (Lohmueller, Rudin and Inman)

## Semi-continuous methods

- Use a Pr(DO) and LRs
- Speed of analysis – "relatively fast"

- The methods do not make full use of data - only the alleles present.

## What should we do with discordant data?

- Continue to use RMNE (CPI, CPE) (not optimal)
- Use the Binary LR with 2p (not optimal)
- Semi-continuous methods with a LR (Drop models)
- Fully continuous methods with LR

## Continuous Models

- Mathematical modeling of "molecular biology" of the profile (mix ratio, PHR (Hb), stutter, etc…) to find optimal genotypes, giving **WEIGHT** to the results.

Probable Genotypes
AC – 40%
BC – 25%
CC – 20%
CQ – 15%

A  B   C   Q

## Some Continuous Model Examples

- TrueAllele (Cybergenetics)
- STRmix (ESR [NZ] and Australian collaboration)
- Cowell et al. (FSI-G (2011) **5:**202-209)

Weights are determined by performing simulations of the data (Markov Chain Monte Carlo - MCMC)

JOURNAL OF **FORENSIC SCIENCES**

PAPER

CRIMINALISTICS

*Mark W. Perlin,[1] M.D., Ph.D.; Matthew M. Legler,[1] B.S.; Cara E. Spencer,[1] M.S.; Jessica L. Smith,[1] M.S.; William P. Allan,[1] M.S.; Jamie L. Belrose,[2] M.S.; and Barry W. Duceman,[3] Ph.D.*

Validating TrueAllele® DNA Mixture Interpretation[*,†]

- Quantitative computer interpretation using Markov Chain Monte Carlo testing
- Models peak uncertainty and infers possible genotypes
- Results are presented as a Combined LR
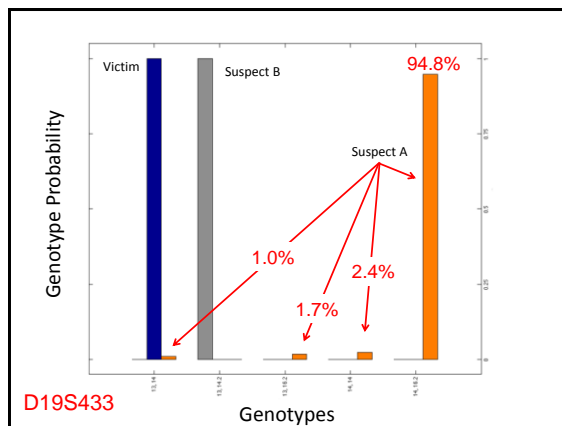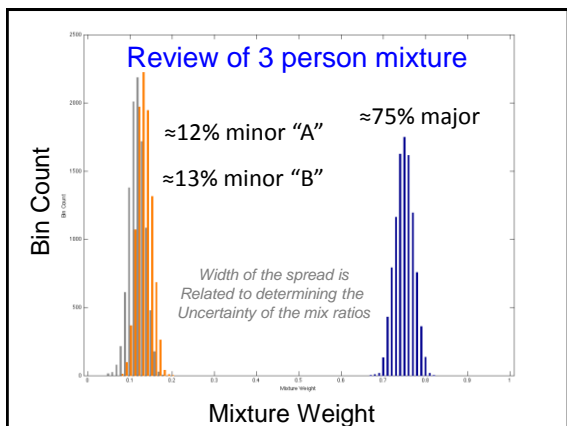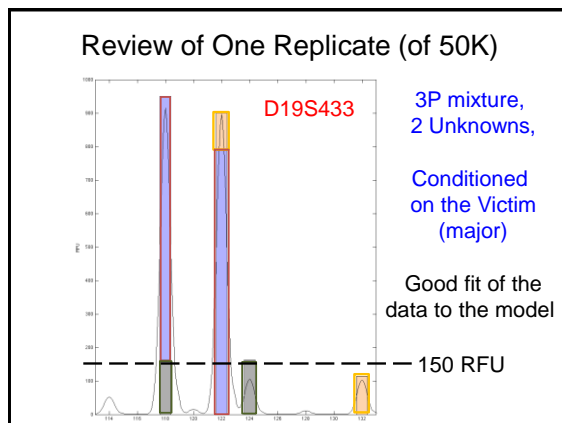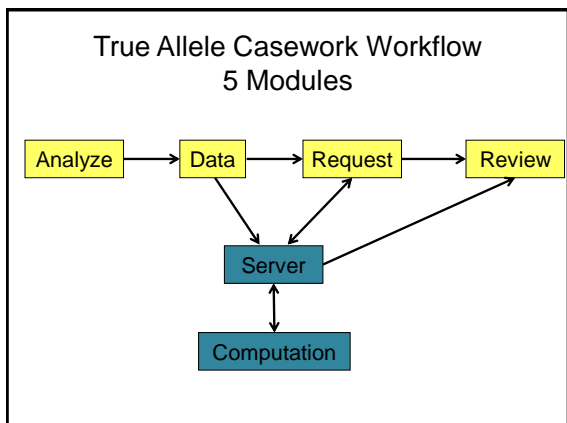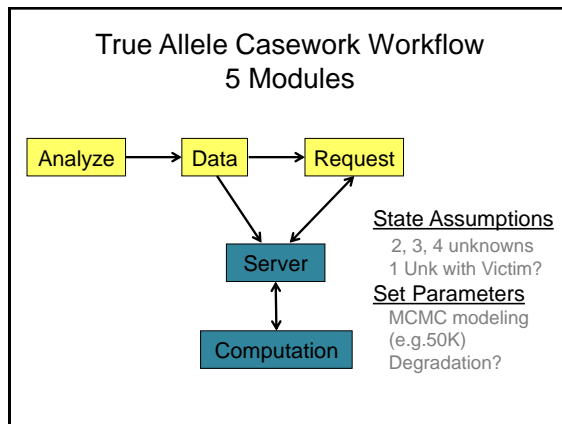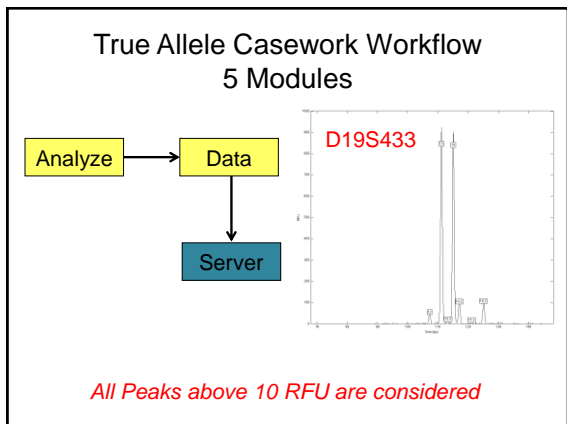
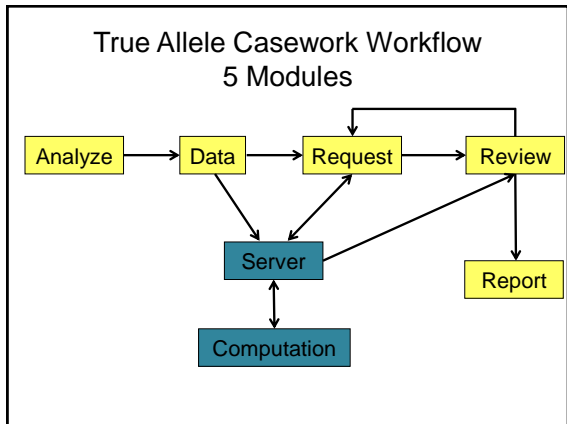## True Allele Software (Cybergenetics)

- Software runs on a Linux Server with a Mac interface.

## True Allele Casework Workflow
### 5 Modules

Analyze

.fsa files imported
Size Standard check
Allelic Ladder check
Alleles are called

## True Allele Casework Workflow
### 5 Modules

Analyze → Data
Data → Server

D19S433

*All Peaks above 10 RFU are considered*

---

## True Allele Casework Workflow
### 5 Modules

Analyze → Data → Request
Data → Server
Request → Server
Server ↔ Computation

<u>State Assumptions</u>
2, 3, 4 unknowns
1 Unk with Victim?
<u>Set Parameters</u>
MCMC modeling
(e.g.50K)
Degradation?

---

## True Allele Casework Workflow
### 5 Modules

Analyze → Data → Request → Review
Data → Server
Request → Server
Server → Review
Server ↔ Computation

---

## Review of One Replicate (of 50K)

D19S433

3P mixture,
2 Unknowns,

Conditioned
on the Victim
(major)

Good fit of the
data to the model

— 150 RFU

---

## Review of 3 person mixture

≈12% minor "A"

≈13% minor "B"

≈75% major

*Width of the spread is
Related to determining the
Uncertainty of the mix ratios*

Bin Count

Mixture Weight

---

Genotype Probability

Victim   Suspect B   94.8%

Suspect A

1.0%   2.4%

1.7%

D19S433   Genotypes

## True Allele Casework Workflow
## 5 Modules



## Combined LR = 5.6 Quintillion

| locus | allele pair x | Likelihood l(x) | Genotype Probability Distribution Questioned q(x) | Reference r(x) | Suspect s(x) | Weighted Likelihood Numerator l(x)*s(x) | Denominator l(x)*r(x) | Likelihood Ratio LR | log(LR) |
|---|---|---|---|---|---|---|---|---|---|
| CSF1PO | 11, 12 | 0.686 | 0.778 | 0.1448 | 1 | 0.68615 | 0.1292 | 5.31 | 0.725 |
| D13S317 | 9, 12 | 1 | 1 | 0.0291 | 1 | 0.99952 | 0.02913 | 34.301 | 1.535 |
| D16S539 | 9, 11 | 0.985 | 0.995 | 0.1238 | 1 | 0.98451 | 0.12188 | 8.036 | 0.905 |
| D18S51 | 13, 17 | 0.999 | 1 | 0.0154 | 1 | 0.99915 | 0.01543 | 64.677 | 1.811 |
| D19S433 | 14, 16.2 | 0.967 | 0.948 | 0.012 | 1 | 0.96715 | 0.01222 | 79.143 | 1.898 |
| D21S11 | 28, 30 | 0.968 | 0.98 | 0.0872 | 1 | 0.96809 | 0.08648 | 11.194 | 1.049 |
| D2S1338 | 23, 24 | 0.998 | 1 | 0.0179 | 1 | 0.99831 | 0.01787 | 55.866 | 1.747 |
| D3S1358 | 15, 17 | 0.988 | 0.994 | 0.1224 | 1 | 0.98759 | 0.12084 | 8.14 | 0.911 |
| D5S818 | 11, 11 | 0.451 | 0.394 | 0.0537 | 1 | 0.45103 | 0.07309 | 6.17 | 0.79 |
| D7S820 | 11, 12 | 0.984 | 0.978 | 0.0356 | 1 | 0.98383 | 0.03617 | 27.198 | 1.435 |
| D8S1179 | 13, 14 | 0.203 | 0.9 | 0.1293 | 1 | 0.20267 | 0.02993 | 6.771 | 0.831 |
| FGA | 21, 25 | 0.32 | 0.356 | 0.028 | 1 | 0.31986 | 0.01906 | 16.783 | 1.225 |
| TH01 | 7, 7 | 0.887 | 0.985 | 0.1739 | 1 | 0.88661 | 0.15588 | 5.687 | 0.755 |
| TPOX | 8, 8 | 1 | 1 | 0.1375 | 1 | 1 | 0.13746 | 7.275 | 0.862 |
| vWA | 15, 20 | 0.998 | 0.996 | 0.0057 | 1 | 0.99808 | 0.00569 | 174.834 | 2.243 |

## Results

- Results are expressed as logLR values

$$LR = 1{,}000{,}000 = 10^6$$

$$\log(LR) = \log 10^6$$

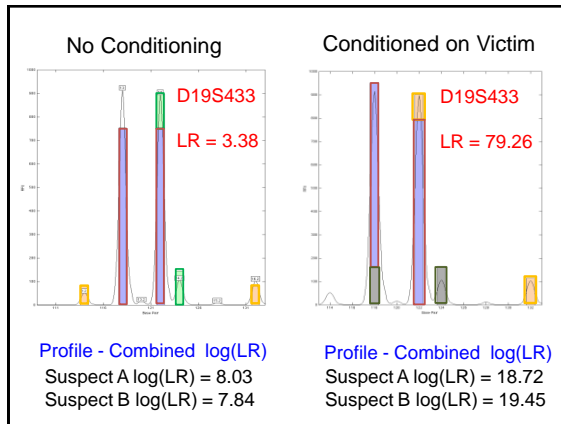$$\log(LR) = 6 * \log 10 \,(1)$$

$$\log(LR) = 6$$

## Review of One Replicate (of 50K)



D19S433

150 RFU

3P mixture,

3 Unknowns

Poor fit of the data to the model

## No Conditioning (3 Unknowns)



D19S433

Major contributor ≈ 75%
(13, 14)
Pr = 1

Genotype Probability

Genotypes

## No Conditioning (3 Unknowns)



D19S433

8.1%

ns for the two
ibutors

Genotype Probability

---

**Slide 1 (top left):**



Suspect "A" Genotype

39 probable genotypes

D19S433

---

**Slide 2 (top right):**

No Conditioning | Conditioned on Victim

D19S433
LR = 3.38

D19S433
LR = 79.26

Profile - Combined log(LR)
Suspect A log(LR) = 8.03
Suspect B log(LR) = 7.84

Profile - Combined log(LR)
Suspect A log(LR) = 18.72
Suspect B log(LR) = 19.45

---

**Slide 3 (middle left):**

## Exploring the Capabilities

- **Degree of Allele Sharing**

- **Mixture Ratios**

- DNA Quantity

---

**Slide 4 (middle right):**

## Mixture Data Set

- Mixtures of pristine male and female DNA amplified at a total concentration of 1.0 ng/μL using Identifiler (standard conditions).
- Mixture ratios ranged from 90:10, 80:20, 70:30 60:40, 50:50, 40:60, 30:70, 20:80, and 10:90
- Each sample was amplified twice.

---

**Slide 5 (bottom left):**

## Mixture Data Set

- Three different combinations:



"Low" Sharing

4 alleles – 10 loci
3 alleles – 5 loci
2 alleles – 0 loci
1 allele – 0 loci

"Medium" Sharing

4 alleles – 3 loci
3 alleles – 8 loci
2 alleles – 4 loci
1 allele – 0 loci

"High" Sharing

4 alleles – 0 loci
3 alleles – 6 loci
2 alleles – 8 loci
1 allele – 1 loci

*Virtual MixtureMaker - http://www.cstl.nist.gov/strbase/software.htm*

---

**Slide 6 (bottom right):**

## Match Score in Duplicate Runs

RMP

Match Rarity (log(LR))

10:90  20:80  30:70  50:50  60:40  70:30  80:20  90:10

Minor Component

Major Component

"Easy" for Deconvolution

---

6

Match Score in Duplicate Runs
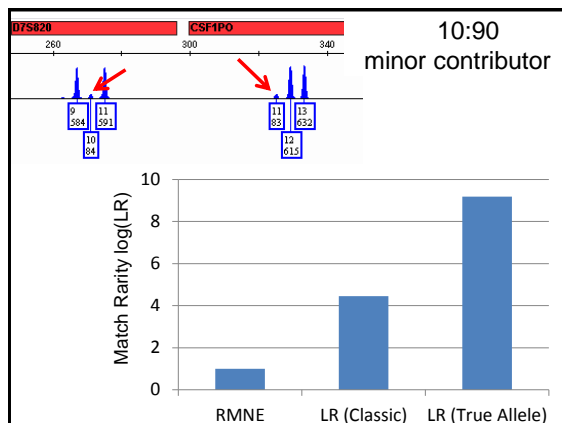


Match Score in Duplicate Runs
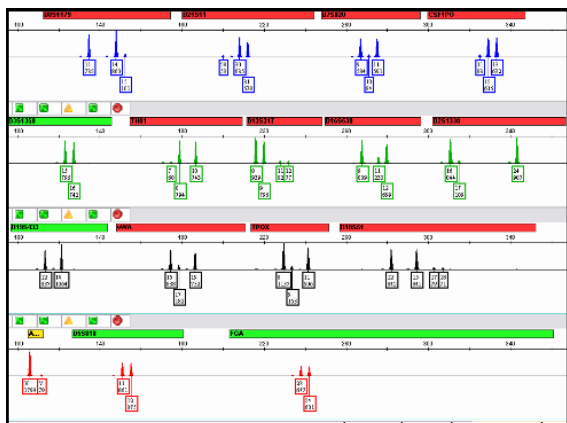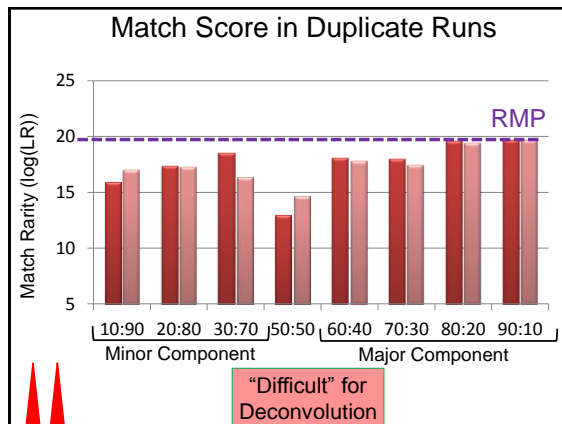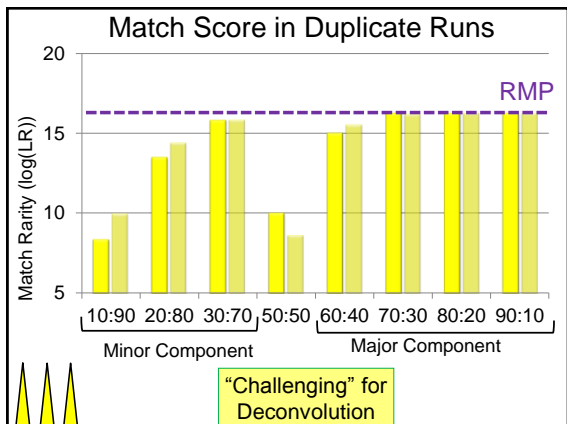




10:90 minor contributor

## Exploring the Capabilities

- Degree of Allele Sharing

- Mixture Ratios

- **DNA Quantity**

D8S1179    "True Genotypes"

A = 13, 16

B = 11, 13

C = 14, 15

3 person Mixture – No Conditioning
Major Contributor ≈ 83 pg input DNA
2 Minor Contributors ≈ 21 pg input DNA



"True Genotypes"

A = 13,16    A = 13,16

B = 11,13    B = 11,13

C = 14,15

C = 12,14



Contributor B (green)
(16%)

Contributor A
(66%)

Contributor C (blue)
(18%)



Genotype Probabilities

A = 13,16

B = 11,13

C = 14,15



Contributor B (gray)
(16%)

Contributor A
(66%)

Contributor C (blue)
(18%)

Conditioned on the Victim



The Power of Conditioning

Victim    Suspect A

C = 14,15

## The Power of Conditioning

|  | LR (no conditioning, 3unk) |
|---|---|
| Contributor A | 1.21 Quintillion |
| Contributor B (victim) | 1.43 Million |
| Contributor C | 9.16 Thousand |

|  | LR (conditioned on victim + 2unk) |
|---|---|
| Contributor A | 1.32 Quintillion |
| Contributor B (victim) | 2.19 Million |
| Contributor C | 59.8 Thousand |

↑

Ranged from 1.13 to 800K

## Summary

- True Allele utilizes probabilistic genotyping and makes better use of the data than the RMNE approach.

- However, the software is computer intensive. On our 4 processor system, it can take 12-16 hours to run up to four 3-person mixture samples.

## Summary

- **Allele Sharing:** Stacking of alleles due to sharing creates more uncertainty.

- **Mixture Ratio:** With "distance" between the two contributors, there is greater certainty. Generally, True Allele performs better than RMNE and the classic LR with low level contributors.

## Summary

- **DNA Quantity:** Generally, with high DNA signal, replicates runs on True Allele are very reproducible.
- However, with low DNA signal, higher levels of uncertainty are observed (as expected).
- There is a need to determine an appropriate threshold for an inclusion log(LR).
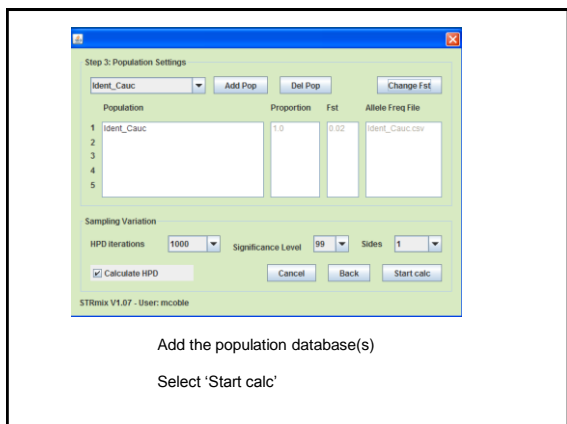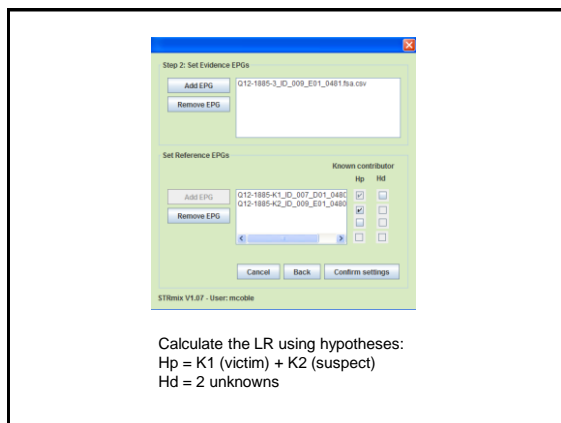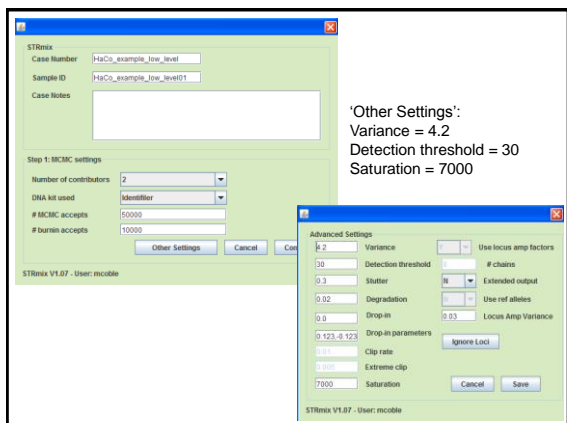
## STRmix

http://strmix.com/

## Challenging Mixture



D19S433

CPI = 1 in 1.7

Michael Donley
Dr. Roger Kahn
Harris Co. (TX) IFS



STRmix

Start Mixture Analysis     Settings

LR from Previous Analysis     Model Maker

Search Database     Exit

About

STRmix V1.07 - User: Coble_F...

---



'Other Settings':
Variance = 4.2
Detection threshold = 30
Saturation = 7000

---



Calculate the LR using hypotheses:
Hp = K1 (victim) + K2 (suspect)
Hd = 2 unknowns

---



Add the population database(s)

Select 'Start calc'

---

Mixture Proportions
Contributor 1 - 87%
Contributor 2 - 13%

D19S433
[14,14] [12,12]     0.2155689382469614
[14,14] (-1,)12]    0.0573843980749224
[14,14] [12,13]     0.1815558506514194
[14,14] [12,14]  ←  0.5454908130266968

Locus 10(D19S433): Pr(E|Hp) = 0.54549, Pr(E|Hd) = 0.00704. LR = 77.51793

LR total = 1.73E16

## Fully continuous methods

- Use a Pr(DO) and LRs
- Speed of analysis – can vary

- Attempts to use all of the data

## Mixture 1.1 – 3P





TrueAllele

19%

37.3%

43.7%

| Evidence | Contributor | Weight | 1 | 10 | 11 |
|----------|-------------|--------|-----|------|------|
| Mixt_1.1_3unk | 1 | 0.4367 | 5.7252 | 18.8886 | -19.1664 |
| Mixt_1.1_3unk | 2 | 0.3734 | 10.8278 | 13.4535 | -18.1592 |
| Mixt_1.1_3unk | 3 | 0.19 | -19.6287 | -19.4984 | 12.7172 |

logLR = 12.7

logLR = 13.4 – 18.8

logLR = 5.7 – 10.8





STRmix
Case Number    mix1.1
Sample ID      mix1.1
Case Notes

Step 1: MCMC settings
Number of contributors    3
DNA kit used              Identifiler
# MCMC accepts            50000
# burnin accepts          10000
                          Other Settings   Cancel   Confirm

STRmix V1.07 - User: mcobie

$H_P$ = Suspect 11 + Suspect 01 + Suspect 10 are in the mixture



$H_D$ = 3 unknown, unrelated individuals are in the mixture







12

Cannot change settings (no MCMC)



TrueAllele
logLR = 5.7 – 10.8

Two Person Mixture

4.1 – low level

D7

Major – 9,12

Minor – 11, 13



Major – 9,12

Minor – 11, 13

| 9 | 11 | 0.223 |
| 11 | 13 | 0.174 |
| 12 | 13 | 0.173 |
| 11 | 12 | 0.171 |
| 9 | 12 | 0.143 |
| 9 | 13 | 0.052 |
| 12 | 12 | 0.051 |
| 9 | 9 | 0.009 |

### TrueAllele Results

| Mixt_4.1_2unk | 0.8709 | 19.5021 | -26.5579 |
| Mixt_4.1_2unk | 0.1291 | -24.5606 | 19.0005 |

86% Major
14% Minor





Ref 12 – 4.95 x $10^{17}$  Ref 16 – 2.52 x $10^{16}$

| Mixt_4.1_2unk | 0.8709 | 19.5021 | -26.5579 |
| _Mixt_4.1_2unk | 0.1291 | -24.5606 | 19.0005 |

### Acknowledgements

- ISFG
- John Buckleton and co-presenters
- Cybergenetics
- Catherine Grgicak and Robin Cotton (Boston U.)
- Charlotte Word (Charlotte Word Consulting)