

DNA Identification Information: Science and Statistics

Mark W Perlin, PhD, MD, PhD
Cybergenetics, Pittsburgh, PA

June, 2010



Cybergenetics © 2003-2010

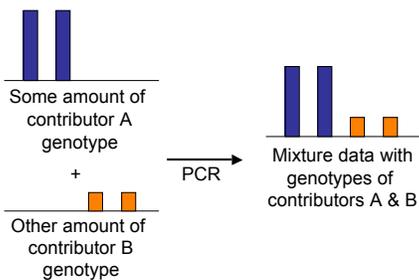
Objective DNA Identification

Given uncertain STR data d

(1) Infer questioned genotype Q

(2) Match with suspect genotype S
relative to random genotype R
to form likelihood ratio LR

DNA Mixture Data



Data Model: Likelihood

$$\mathbf{d}_l \sim N_+(\mu_l, \Sigma_l)$$

$$\mu_l = m_l \cdot \sum_{k=1}^K w_{k,l} \cdot \mathbf{g}_{k,l}$$

$$\mathbf{w}_l \sim N_{[0,1]^{K-1}}(\mathbf{w}, \psi^2 \cdot \mathbf{I})$$

$$\Sigma_l = \sigma^2 \cdot V_l + \tau^2$$

Prior Probability

$$\mathbf{g}_{k,l} \sim \begin{cases} f_i^2, & i = j \\ 2f_i f_j, & i \neq j \end{cases}$$

$$\mathbf{w} \sim \text{Dir}(\mathbf{1})$$

$$m_l \sim N_+(5000, 5000^2)$$

$$\sigma^{-2} \sim \text{Gam}(10, 20)$$

$$\tau^{-2} \sim \text{Gam}(10, 500)$$

$$\psi^{-2} \sim \text{Gam}(1/2, 1/200)$$

Genotype Inference

$$\Pr\{Q = x | d_{1,1}, d_{1,2}, \dots, d_{l,2}, \dots\} \propto \Pr\{Q = x\} \cdot \prod_{j=1}^J \Pr\{d_j | Q = x, \dots\}$$

$$\Pr\{W = w | d_1, d_2, \dots, d_j, \dots\} \propto \Pr\{W = w\} \cdot \prod_{j=1}^J \Pr\{d_j | W = w, \dots\}$$

$$\Pr\{\sigma^2 = s^2 | d_1, d_2, \dots, d_j, \dots\} \propto \Pr\{\sigma^2 = s^2\} \cdot \prod_{j=1}^J \Pr\{d_j | \sigma^2 = s^2, \dots\}$$

$$\Pr\{\tau^2 = t^2 | d_1, d_2, \dots, d_j, \dots\} \propto \Pr\{\tau^2 = t^2\} \cdot \prod_{j=1}^J \Pr\{d_j | \tau^2 = t^2, \dots\}$$

Information Gain (LR)

identification hypothesis:
the suspect contributed to the evidence

$$\text{information gain (likelihood ratio)} = \frac{\text{Odds(hypothesis | data)}}{\text{Odds(hypothesis)}}$$

↑ data
after
before

Additive information units: log(LR)
Order of magnitude, powers of ten

(apply Bayes theorem)

$$\begin{aligned} \frac{O(H|d_Q, d_R, d_S)}{O(H)} &= \frac{\Pr\{H|d_Q, d_R, d_S\} / \Pr\{\bar{H}|d_Q, d_R, d_S\}}{\Pr\{H\} / \Pr\{\bar{H}\}} \\ &= \frac{\Pr\{H|d_Q, d_R, d_S\} / \Pr\{H\}}{\Pr\{\bar{H}|d_Q, d_R, d_S\} / \Pr\{\bar{H}\}} \\ &= \frac{\Pr\{d_Q|H, d_R, d_S\} / \Pr\{d_Q\}}{\Pr\{d_Q|\bar{H}, d_R, d_S\} / \Pr\{d_Q\}} \\ &= \frac{\Pr\{d_Q|H, d_R, d_S\}}{\Pr\{d_Q|\bar{H}, d_R, d_S\}} \end{aligned}$$

Likelihood Ratio (LR)

$$LR = \frac{\Pr\{d_Q|H, d_R, d_S\}}{\Pr\{d_Q|\bar{H}, d_R, d_S\}}$$

(extend conversation)

$$\begin{aligned} \frac{\Pr\{d_Q|H,d_R,d_S\}}{\Pr\{d_Q|\bar{H},d_R,d_S\}} &= \frac{\sum_{x \in G} \Pr\{d_Q|H,d_R,d_S,Q=x\} \cdot \Pr\{Q=x|H,d_R,d_S\}}{\sum_{x \in G} \Pr\{d_Q|\bar{H},d_R,d_S,Q=x\} \cdot \Pr\{Q=x|\bar{H},d_R,d_S\}} \\ &= \frac{\sum_{x \in G} \Pr\{d_Q|d_R,d_S,Q=x\} \cdot \Pr\{S=x|d_R,d_S\}}{\sum_{x \in G} \Pr\{d_Q|d_R,d_S,Q=x\} \cdot \Pr\{R=x|d_R,d_S\}} \\ &= \frac{\sum_{x \in G} \Pr\{Q=x\} \cdot \Pr\{S=x|d_S\}}{\sum_{x \in G} \Pr\{Q=x\} \cdot \Pr\{R=x|d_R\}} \\ &= \frac{\sum_{x \in G} \lambda_Q(x) \cdot s(x)}{\sum_{x \in G} \lambda_Q(x) \cdot r(x)} \end{aligned}$$

Genotype-Weighted Likelihood

$$LR = \frac{\sum_{x \in G} \lambda_Q(x) \cdot s(x)}{\sum_{x \in G} \lambda_Q(x) \cdot r(x)}$$

(apply Bayes theorem)

$$\begin{aligned} LR &= \frac{\sum_{x \in G} \lambda_Q(x) \cdot s(x)}{\sum_{x \in G} \lambda_Q(x) \cdot r(x)} \\ &= \sum_{x \in G} \frac{\lambda_Q(x)}{\sum_{y \in G} \lambda_Q(y) \cdot \pi_Q(y)} \cdot s(x) \\ &= \sum_{x \in G} \frac{q(x)}{\pi_Q(x)} \cdot s(x) \\ &= \sum_{x \in G} \frac{q(x) \cdot s(x)}{r(x)} \end{aligned}$$

Genotype Probability Gain

$$LR = \sum_{x \in G} \frac{q(x) \cdot s(x)}{r(x)}$$

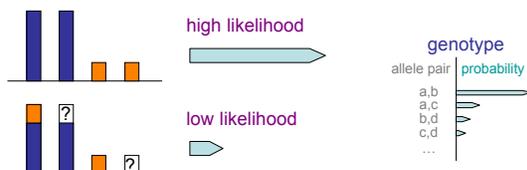
Summary

- Other modeling variables (stutter, ...)
- LR can also consider coancestry
- Infer genotypes, then match them
- Objective: inference never sees suspect
- Goal: preserve identification information
 - + Bayesian modeling for genotype
 - + Bayesian information gain (LR)

Quantitative Mixture Interpretation

Step 1: infer genotype

- consider every possible allele pair
- compare pattern with DNA data
- Rule: *better fit's more likely it*



Information Gain (LR)

Step 2: match genotypes

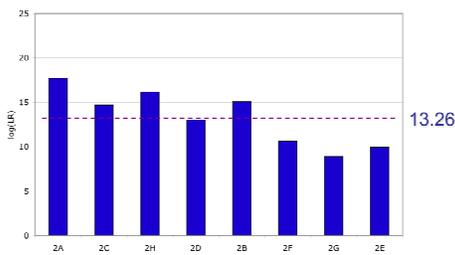
At the suspect's genotype allele pair,
what is the locus *information gain*?

$$\text{information gain (likelihood ratio)} = \frac{\text{Prob(allele pair | data)}}{\text{Prob(allele pair)}}$$

after
↑ data
before
(population)

Computer objectivity:
(Step 1) infer evidence genotype from data
(Step 2) compare genotype with suspect

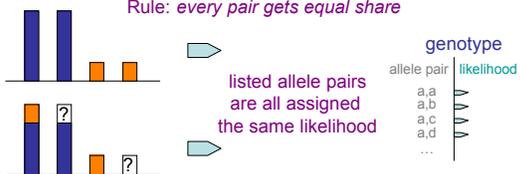
Efficacy (2 unknown)



Qualitative Manual Review

Step 1: infer genotype

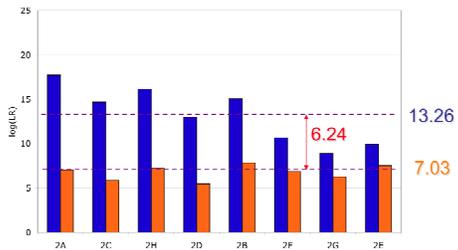
Rule: every pair gets equal share



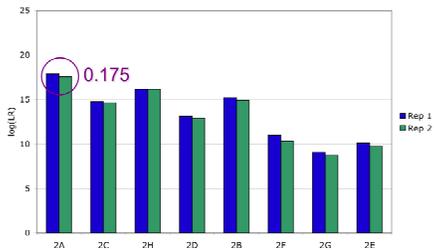
Step 2: match genotype

lower probability means lower information gain (LR)

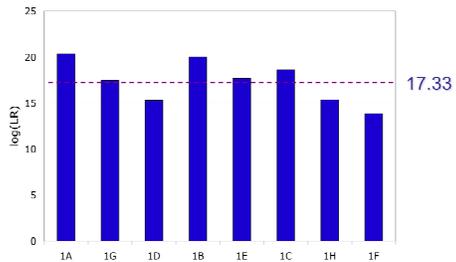
Improvement



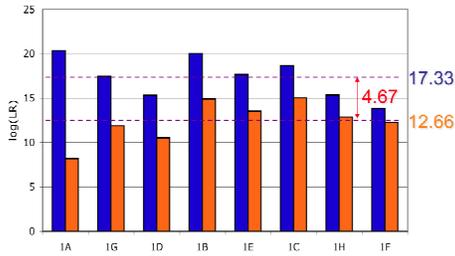
Reproducibility



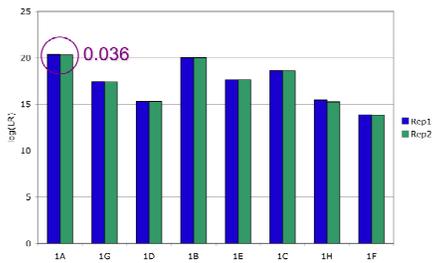
Efficacy (1 unknown)



Improvement



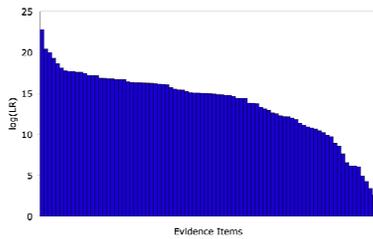
Reproducibility



Comparison

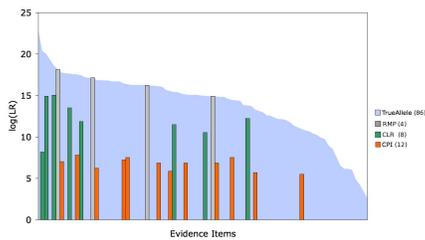
interpretation method	two unknown (without victim)	one unknown (with victim)
quantitative computer	13.26 (0.175) (ten trillion)	17.33 (0.036) (hundred quadrillion)
qualitative human	7.03 (ten million)	12.66 (fifty trillion)
improvement	6.24 (one million)	4.67 (fifty thousand)

TrueAllele Information: 86 Match Stats (100%)



Preserves all the identification information

Human Review Information: 24 Match Stats (28%)



Preserves 20% of the identification information

Summary

- **information gain** (LR) is a universal DNA metric
- **efficacy**: computer extracts useful information
- **improvement**: computer mixture interpretation is more informative than human review
with victim 50,000x - without victim 1,000,000x
- **reproducibility**: tenths of a log(LR) unit
- **objectivity**: "parallel unmasking", infer then match
- **productivity**: lab gives statistic for 1 of 3 items
- **utility**: science, investigation and evidence

Commonwealth vs. Foley

Apr 2006: Blairsville Dentist John Yelenic murdered

Nov 2007: Trooper Kevin Foley charged with crime



Feb 2008: Defense questions 13,000 DNA match score

DNA Evidence

- DNA from under victim's fingernails (Q83)
- two contributors to DNA mixture
- 93.3% victim & 6.7% unknown
- 1,000 pg DNA in 25 ul
- STR analysis with ProfilerPlus®, Cofiler®
- know victim contributor genotype (K53)
- TrueAllele® computer interpretation (using genotype addition method)
infer unknown contributor genotype
- only after having inferred unknown, compare with suspect genotype (K2)

Three DNA Match Statistics

<u>Score</u>	<u>Method</u>
13 thousand	inclusion
23 million	subtraction
189 billion	addition

- Why are there different match results?
- How do mixture interpretation methods differ?
- What should we present in court?

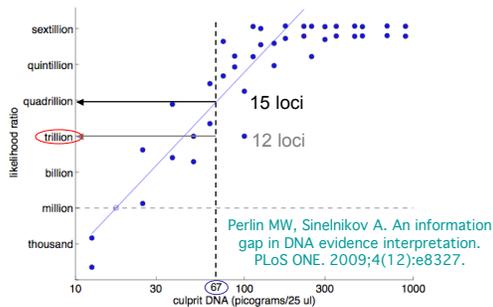
Different Interpretation Methods

Data Used	inclusion	subtraction	addition
victim profile	NO	YES	YES
quantitative data	NO	NO	YES

Frye: General Acceptance in the Relevant Community

- Quantitative STR Peak Information
- Genotype Probability Distributions
- Computer Interpretation of STR Data
- Statistical Modeling and Computation
- Likelihood Ratio Literature
- Mixture Interpretation Admissibility
- Computer Systems for Quantitative DNA Mixture Deconvolution
- TrueAllele Casework Publications

Expected Result

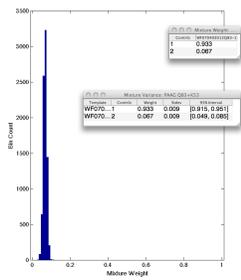


Expert Testimony

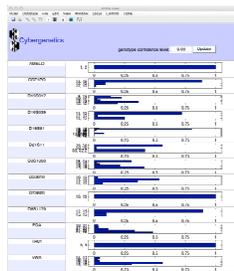
Dr. Perlin explained to the jury why these apparently different results were expected by DNA science. "The less informative methods ignored some of the data," said Dr. Perlin, "while the TrueAllele computation considered all of the available DNA data."

"A scientist may look at the same slide using the naked eye, a magnifying glass, or a microscope," analogized Dr. Perlin. "A computer that considers all the data is a more powerful DNA microscope."

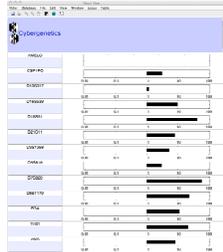
Mixture Weight



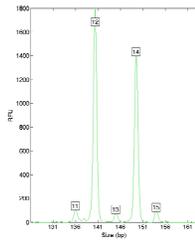
Inferred Genotype



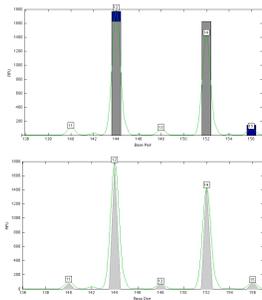
log(LR) Match Information



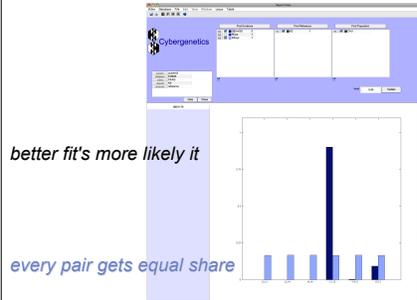
Locus D8S1179 Data



Explain D8S1179 Genotype



Likelihood Comparison



Generate Report

DE3 + K53 contributor: 2 vs. K2 (CAU)

File Signature Statement Summary Calculation

The LR calculation assumes one unknown contributor in the evidence with one known contributor reference relative to a Caucasian human population having a coancestry coefficient of 0.01.

The joint LR is approximately 22.1 billion.

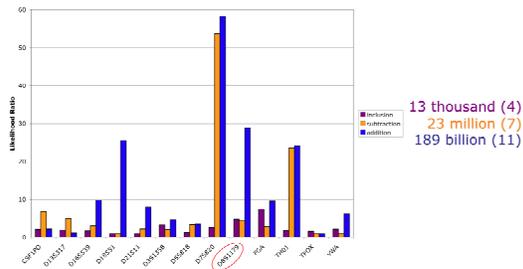
The log(LR) information is 10.54.

Locus	Allele pair	Q	R	S	LR	log(LR)
CSF1PO	12, 13	0.091	0.0518	1	1.755	0.244
D13S317	8, 11	0.136	0.0683	1	1.990	0.299
D16S539	11, 13	0.722	0.0928	1	7.775	0.891
D18S51	12, 13	0.803	0.0554	1	22.683	1.356
D21S11	29, 30	0.561	0.0877	1	6.388	0.805
D3S1358	15, 18	0.213	0.0839	1	2.538	0.405
D5S818	12, 13	0.358	0.1077	1	3.324	0.522
D7S820	10, 13	1	0.0226	1	44.188	1.645
D8S1179	12, 15	0.895	0.0365	1	24.525	1.390
FGA	21, 24	0.483	0.0514	1	9.388	0.973
TH01	8, 9	1	0.0450	1	22.201	1.346
VWA	17, 18	0.562	0.1199	1	4.689	0.671

Locus information gain is genotype probability ratio:
LR = after/before

Joint information is the sum of the locus information

More Data In, More Information Out



Case Observations

- objective review never saw suspect
- easy to testify about in court
- understandable to judge and jury
- have precedent: admitted, testified
- preserve match information in data
- rapid response to attorney
- multiple match scores presented
all information to the triers of fact –
nothing was withheld from the jury
this should be standard practice

The screenshot shows the top portion of a news article on the Indiana Gazette website. The header includes the site name 'Indiana Gazette.com', the date '08 17 09', and the page number '41'. A navigation menu lists categories like HOME, SUBSCRIBERS, MARKETPLACE, NEWS, OBITUARIES, SPORTS, BUSINESS, MULTIMEDIA, and FYI JOURNALS. The article title is 'Jury convicts trooper of dentist slaying', with a sub-headline 'One Verdict'. The text below the title states that an Indiana County Court jury convicted state trooper Kevin Foley of first degree murder in the April 13, 2006, slaying of dentist John Yelenic. A quote from prosecutor Anthony Krastek is provided, along with a link to a DNA Investigator Newsletter.

Indiana Gazette.com 08 17 09 41
IN PRINT DAILY. ONLINE ALWAYS.
HOME SUBSCRIBERS MARKETPLACE NEWS OBITUARIES SPORTS BUSINESS MULTIMEDIA O
FYI JOURNALS
ARCHIVES > NEWS One Verdict
Jury convicts trooper of dentist slaying
Published: Thursday, March 19, 2009 12:46 AM EDT
An Indiana County Court jury this evening convicted state trooper Kevin Foley of first degree murder in the April 13, 2006, slaying death of Blairsville dentist John Yelenic.
"John Yelenic provided the most eloquent and poignant evidence in this case," said the prosecutor, senior deputy attorney general Anthony Krastek. "He managed to reach out and scratch his assailant," capturing the murderer's DNA under his fingernails.
The DNA Investigator Newsletter. Same Data, More Information - Murder, Match and DNA, Cybergenetics, 2009. www.cybgen.com/information/newsletters/CybgenNews1.pdf
