

# The Science of DNA Search

The Changing Face of DNA:  
The Science, Law and Ethics of Familial Searches

Duquesne University  
May, 2011  
Pittsburgh, PA

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



Cybergenetics © 2003-2011

---

---

---

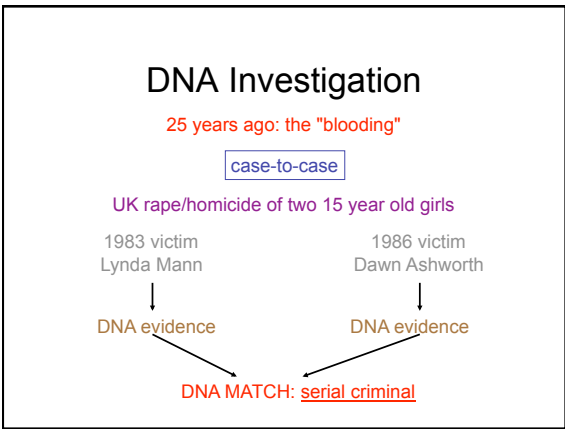
---

---

---

---

---




---

---

---

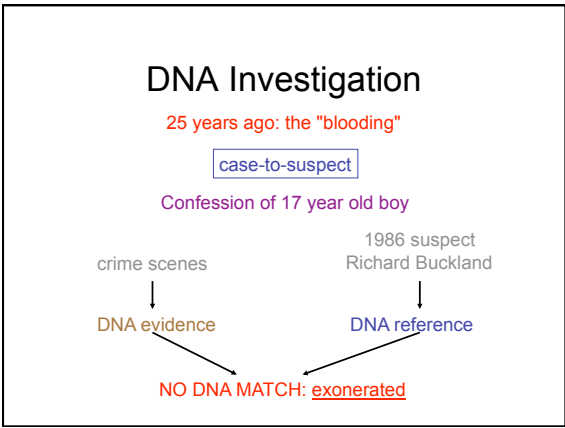
---

---

---

---

---




---

---

---

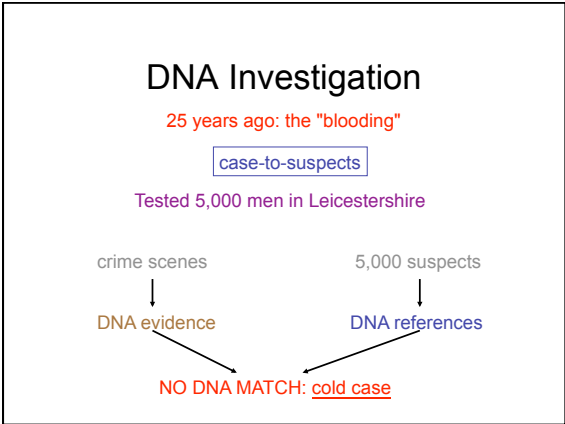
---

---

---

---

---




---

---

---

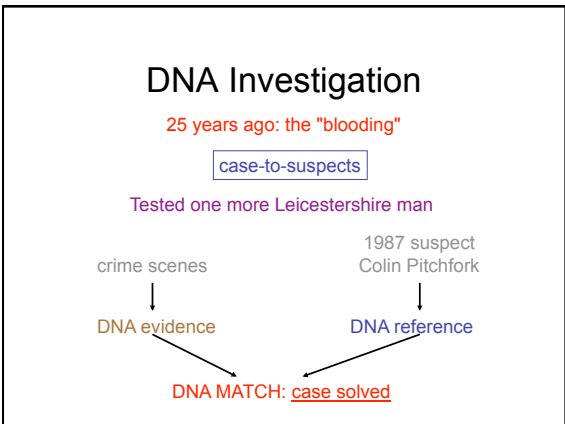
---

---

---

---

---




---

---

---

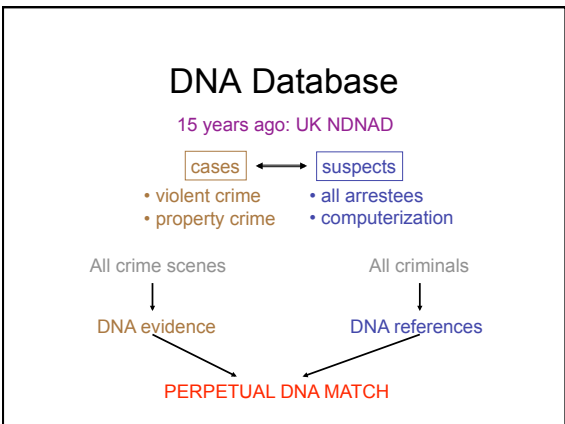
---

---

---

---

---




---

---

---

---

---

---

---

---

## Vision: Preventing Crime

Continual crime scene surveillance by DNA

2004. Ray Wickenheiser: Complete DNA databasing  
(sexual assault, property crime, convicted offender)  
would prevent 100,000 stranger rapes

The UK did exactly this,  
with over 50% DNA database hit rate,  
stopping criminal careers early on at property crime

Protect the public through DNA databases  
by stopping criminals before they strike again

---

---

---

---

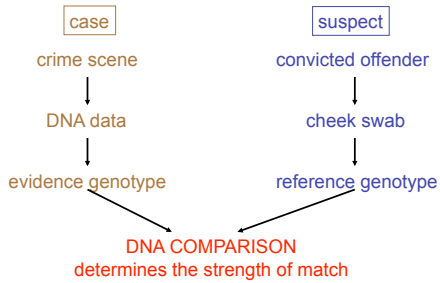
---

---

---

---

## How DNA Matching Works



---

---

---

---

---

---

---

---

## Simple DNA Evidence

evidence



reality

---

---

---

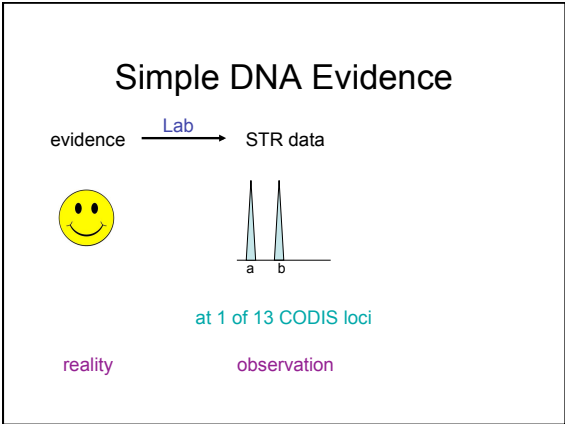
---

---

---

---

---




---

---

---

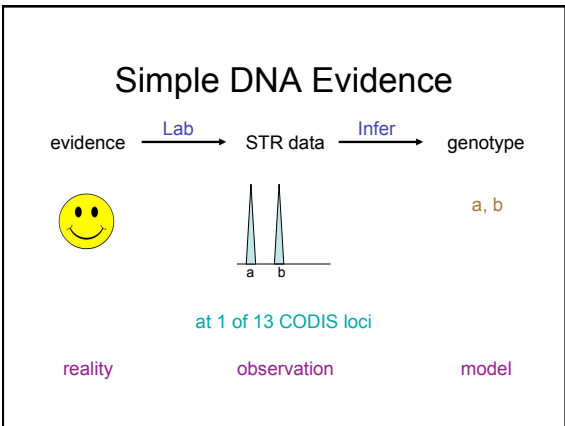
---

---

---

---

---




---

---

---

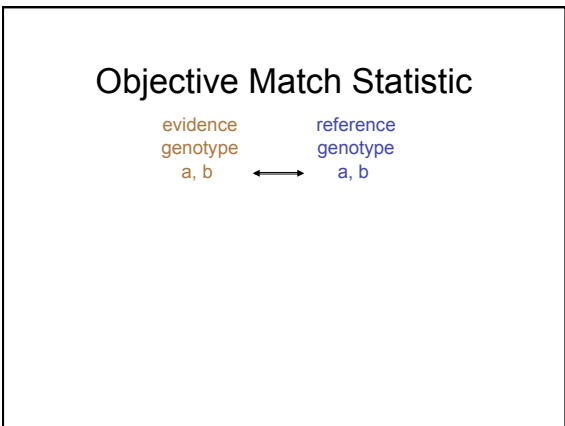
---

---

---

---

---




---

---

---

---

---

---

---

---

## Objective Match Statistic

evidence genotype a, b ↔ reference genotype a, b

likelihood ratio  $LR = \frac{\text{Pr}(\text{match})}{\text{Pr}(\text{coincidence})} = \frac{100\%}{5\%} = 20$

---

---

---

---

---

---

---

## Objective Match Statistic

evidence genotype a, b ↔ reference genotype a, b

likelihood ratio  $LR = \frac{\text{Pr}(\text{match})}{\text{Pr}(\text{coincidence})} = \frac{100\%}{5\%} = 20$

information measure  $\log(LR) = 1.30$  (since  $10^{1.30} = 20$ , in powers of 10)

---

---

---

---

---

---

---

## Objective Match Statistic

evidence genotype a, b ↔ reference genotype a, b

likelihood ratio  $LR = \frac{\text{Pr}(\text{match})}{\text{Pr}(\text{coincidence})} = \frac{100\%}{5\%} = 20$

information measure  $\log(LR) = 1.30$  (since  $10^{1.30} = 20$ , in powers of 10)

total information  $13 \times \log(LR) = 13 \times 1.30 = 16.9$  17 zeros is 100 quadrillion

---

---

---

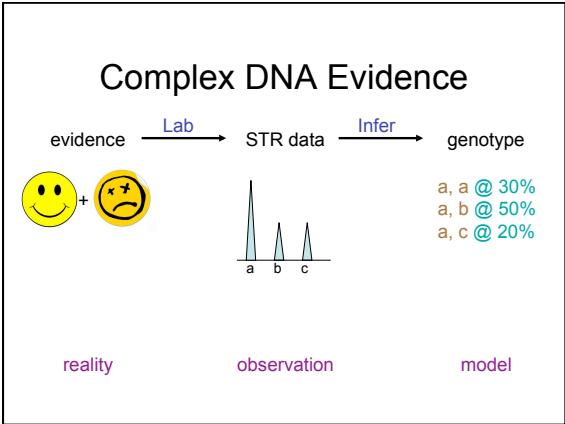
---

---

---

---






---

---

---

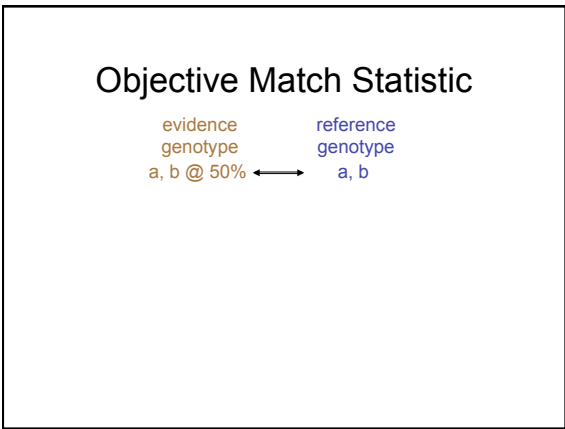
---

---

---

---

---




---

---

---

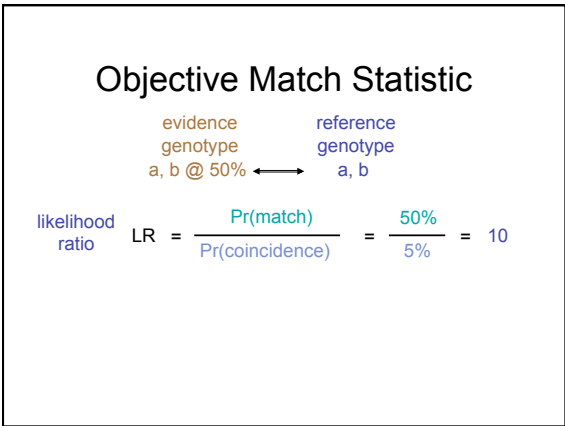
---

---

---

---

---




---

---

---

---

---

---

---

---

## Objective Match Statistic

evidence genotype a, b @ 50% ↔ reference genotype a, b

likelihood ratio  $LR = \frac{\text{Pr}(\text{match})}{\text{Pr}(\text{coincidence})} = \frac{50\%}{5\%} = 10$

information measure  $\log(LR) = 1$  (since  $10^1 = 10$ , in powers of 10)

---

---

---

---

---

---

---

---

## Objective Match Statistic

evidence genotype a, b @ 50% ↔ reference genotype a, b

likelihood ratio  $LR = \frac{\text{Pr}(\text{match})}{\text{Pr}(\text{coincidence})} = \frac{50\%}{5\%} = 10$

information measure  $\log(LR) = 1$  (since  $10^1 = 10$ , in powers of 10)

total information  $13 \times \log(LR) = 13 \times 1 = 13$  13 zeros is 10 trillion

---

---

---

---

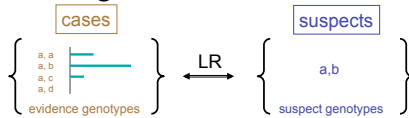
---

---

---

---

## Investigative DNA Database



- genotype probability representation
- fully preserves DNA identification information
- enables LR calculation with every match
  - connect crimes to criminals
  - disaster victim identification (WTC)
  - find missing people
  - automatic familial search
  - combat terrorism through DNA

---

---

---

---

---

---

---

---



## Down A Different Path

15 years ago: United States CODIS  
a second generation DNA database

New feature: represents mixture evidence

The FBI used *simplified* mixture interpretation

+ easier for crime labs to do and testify about

- loses considerable identification information

---

---

---

---

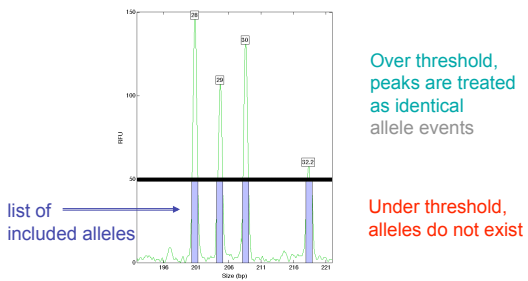
---

---

---

---

## Qualitative Thresholds



---

---

---

---

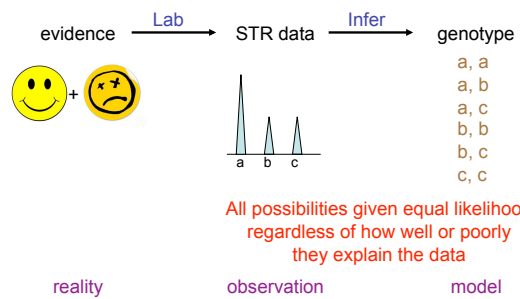
---

---

---

---

## Simplified Mixture Interpretation



---

---

---

---

---

---

---

---

### Reduced Match Statistic

diffuses probability evidence genotype reference genotype  
a, b @ 15% ↔ a, b

likelihood ratio  $LR = \frac{\text{Pr}(\text{match})}{\text{Pr}(\text{coincidence})} = \frac{15\%}{5\%} = 3$

information measure  $\log(LR) = 0.48$  (since  $10^{0.48} = 3$ , in powers of 10)

total information  $13 \times \log(LR) = 13 \times 0.48 = 6.2$  6 zeros is one million

---

---

---

---

---

---

---

---

The DNA Investigator™ Newsletter, 2009  
Same Data, More Information – Murder, Match and DNA

### Fingernail: 7% Mixture

Commonwealth v. Foley

ScoreMethod  
13 thousand threshold  
23 million use victim  
189 billion quantitative

- probability modeling preserves information
- peak threshold discards information

---

---

---

---

---

---

---

---

MW Perlin, MM Legler, CE Spencer, JL Smith, WP Allan, JL Belrose, BW Duceman. Validating TrueAllele DNA mixture interpretation. Journal of Forensic Sciences, 2011.

### Preserve vs. Discard

Category	Quantitative Interpretation (log(LR))	Threshold method (log(LR))
2A	~18	~7
2C	~15	~6
2H	~16	~7
2D	~13	~5
2B	~15	~8
2F	13.26	6.24
2G	~10	~6
2E	~10	~7

---

---

---

---

---

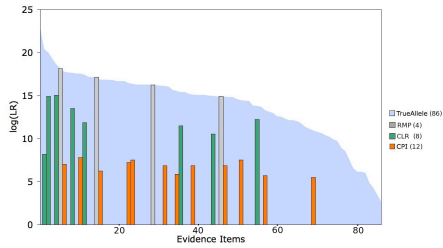
---

---

---

Perlin MW, Ducean BW. Profiles in productivity: greater yield at lower cost with computer DNA interpretation. Twentieth International Symposium on the Forensic Sciences of the Australian and New Zealand Forensic Science Society, Sydney, Australia. 2010.

## Preserve vs. Discard



- quantitative interpretation **preserves** information - every time
- peak threshold **discards** information - 70% of the time

---

---

---

---

---

---

---

---

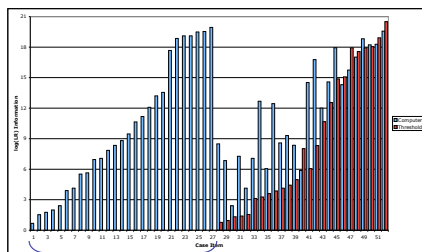
---

---

---

---

## Preserve vs. Discard



Computer: **informative**  
Threshold: **inconclusive**

---

---

---

---

---

---

---

---

---

---

---

---

## CODIS: Second Generation

less informative  
evidence genotype

allele pairs

- a, a
- a, b
- a, c
- b, b
- b, c
- c, c

simplified genotype  
representation

allele list

- a, b, c

- identification information lost
- sacrifices probability accuracy
- probative evidence discarded
- match *trillions* reduced to *millions*

---

---

---

---

---

---

---

---

---

---

---

---

## Truth or Consequences

Exchanging best science for practical procedures

"High stringency" search. With high thresholds, lose sensitivity - can't find criminals

"Low stringency" search. Overly broad, so lose specificity - finds hundreds of CODIS matches

For example, with 100,000 convicted offenders, finding 100 database hits means a 1/1000 hit rate

Many dead end leads, and wasted police work

---

---

---

---

---

---

---

---

## Everyone is a Suspect

Isn't it "fair" to put everyone on the DNA database?  
Why store only criminal genotypes?

With high CODIS specificity (low false positives), OK

But the CODIS mixture representation reduces *trillion*-fold evidence specificity down to *million*-fold

Each evidence item would falsely implicate *hundreds* of innocent people, soon making all Americans suspects

Weak science can have poor consequences

---

---

---

---

---

---

---

---

## Familial Search: A Good Consequence

Serendipity from low stringency searches

CODIS "low stringency": genotypes share at least one allele



---

---

---

---

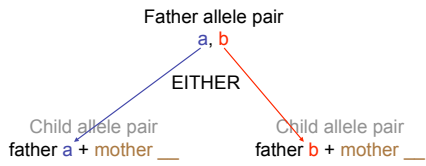
---

---

---

---

## Family Genetics



Parent and child share at least one allele  
(siblings and other kin also share alleles, but fewer)

So "low stringency" match can be used for familial search

---

---

---

---

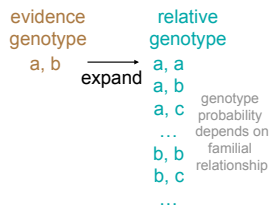
---

---

---

---

## How the Science Works




---

---

---

---

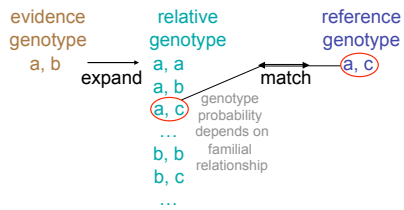
---

---

---

---

## How the Science Works




---

---

---

---

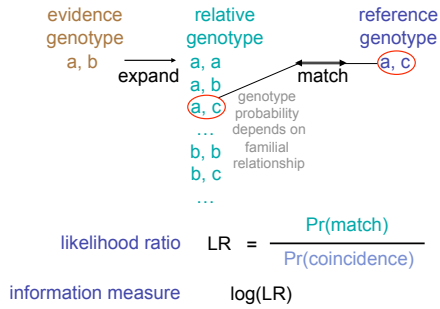
---

---

---

---

## How the Science Works




---

---

---

---

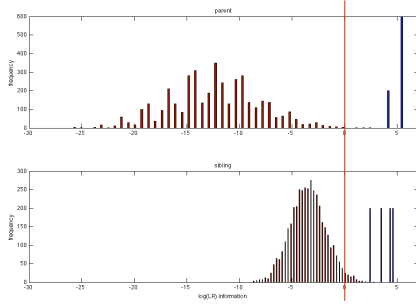
---

---

---

---

## Results: Evidence Mode




---

---

---

---

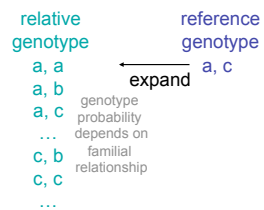
---

---

---

---

## Continuous Mode




---

---

---

---

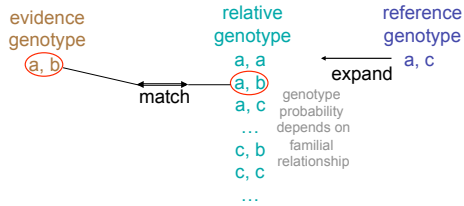
---

---

---

---

## Continuous Mode




---

---

---

---

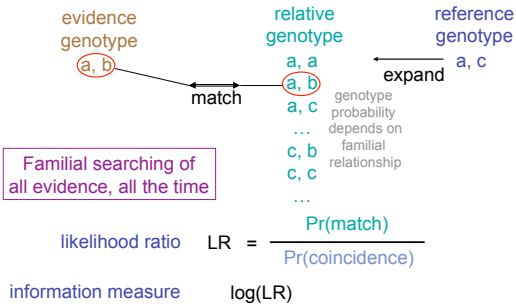
---

---

---

---

## Continuous Mode




---

---

---

---

---

---

---

---

## Mixtures in Familial Search

- Most DNA evidence items are *mixtures*
- CODIS representation has weak information, greatly increasing *familial search* false hits

But probabilistic genotypes overcome these obstacles

- informative computer *interpretation* (SWGAM 2010 guidelines, paragraph 3.2.2)
- information preserving DNA *database* (ANSI/NIST-ITL standard, fields 18.020 & 18.021)

So the full DNA evidence information is preserved

---

---

---

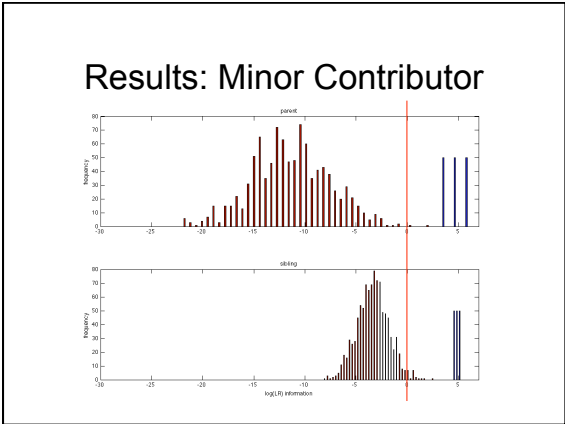
---

---

---

---

---




---

---

---

---

---

---

---

---

### Familial Search Cost/Benefit

most agree that familial search is worthwhile in catching criminals and preventing further crime

- mainly nonmixture evidence with CODIS
- only dozens of hits made in ten years
- cost is \$20,000 per search (Y-STR confirm)
- yield is 10% hits per search

Hit cost is \$200,000 (\$20,000 / 10%)

---

---

---

---

---

---

---

---

### Preserving Information Cost/Benefit

better use of the DNA data *we already have* can catch more criminals and prevent more crime

- all evidence types (including mixtures)
- no confirmation testing required (non-familial)
- computer system costs less than one familial hit
- better information doubles evidence yield
- can produce thousands of hits every year
- expands familial search capability to mixtures

More information  
More hits  
Lower cost

---

---

---

---

---

---

---

---



## Science & Technology

- past: DNA evidence and DNA databases  
1st generation, single allele pair
- present: familial search and other  
innovative DNA paths to public safety  
2nd generation, lists of alleles
- future: more informative DNA methods  
using computers to preserve evidence  
3rd generation, probabilistic genotypes

DNA identification is an information science

---

---

---

---

---

---

---

## Learning More

The science of DNA search

[www.cybgen.com/information](http://www.cybgen.com/information)

- Newsletters  
gentle introduction to ideas
- Courses  
for lawyers and scientists
- Presentations  
handouts, movies, transcripts
- Publications  
abstracts, manuscripts

---

---

---

---

---

---

---