# Computing Without Exposing Data
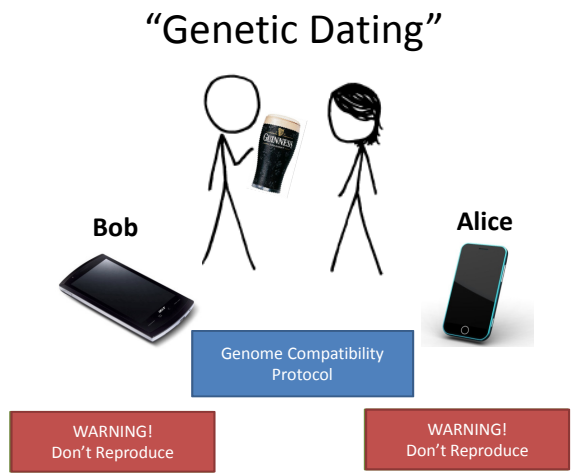
Radford University
CSAT STEM Club
17 March 2011

David Evans
University of Virginia
http://www.cs.virginia.edu/evans
http://www.MightBeEvil.com

1

---

## "Genetic Dating"

Bob

Alice

Genome Compatibility Protocol

WARNING!
Don't Reproduce

WARNING!
Don't Reproduce

2

---

GenePartner™
Love is no coincidence

23andMe

ScientificMatch.com
"The Science of Love"

TheScientist   News   Current Issue   Archive   Sur

SHARE

2 comments
Comment on this news story
By Kerry Grens
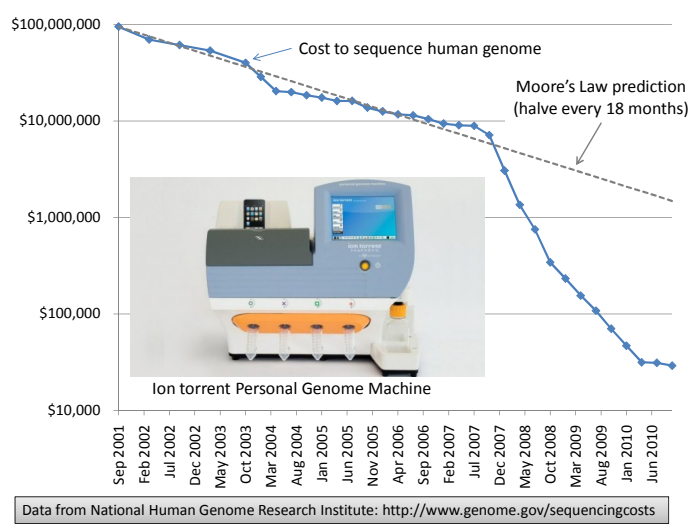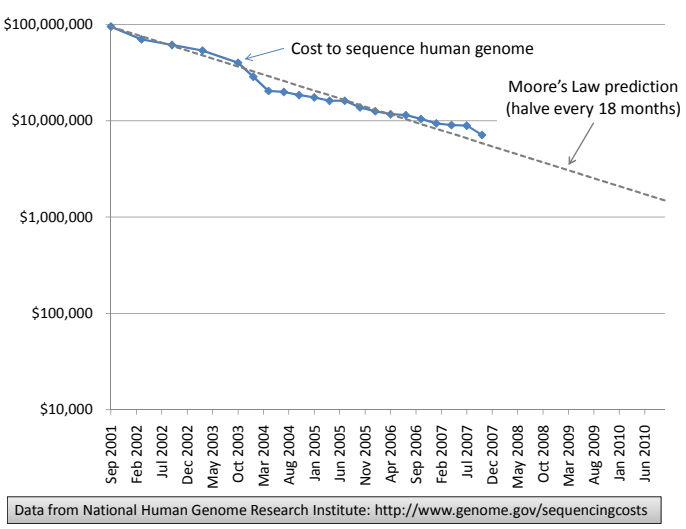
### Forget mistletoe - what about DNA?

A new dating service matches singles using major histocompatibility complex genes

3

---

## Genome Sequencing

1990: Human Genome Project starts, estimate $3B to sequence one genome ($0.50/base)

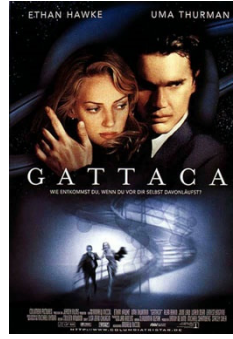2000: Human Genome Project declared complete, cost ~$300M

Whitehead Institute, MIT

4

---

Cost to sequence human genome

Moore's Law prediction (halve every 18 months)

Data from National Human Genome Research Institute: http://www.genome.gov/sequencingcosts

---

Cost to sequence human genome

Moore's Law prediction (halve every 18 months)

Ion torrent Personal Genome Machine

Data from National Human Genome Research Institute: http://www.genome.gov/sequencingcosts

| Year | r... | | Reported sequencing consumables cost | Estimated cost per 40-fold coverage |
|---|---|---|---|---|
| | | | $10,000,000 | $57,000,000 |
| | | | $1,000,000 | $5,700,000 |
| | | | $250,000 | $330,000 |
| | | | $48,000 | $69,000 |
| 2009 | t... | | $8,005 | $3,700 |
| 2009 | t... | | $3,451 | $2,200 |
| 2009 | t... | | $1,726 | $1,500 |

*Human Gen... ...ing DNA Nanoarrays.* Radoje Drmanac, An... ...L. Burns, Bahram G. Kermani, Paolo Carnevali, Igo... ...Andres Fernandez, Bryan Staker, Krishna P. Par... ..., Ryan Cedeno, Linsu Chen, Dan Chernikoff, Alex... ...t, Coleen R. Hacker, Robert Hartlage, Brian Hauser, S... ...lvin Kong, Tom Landers, Catherine Le, Jia Liu, Celeste ... ...Helena Perazich, Kimberly Perry, Brock A. Peters, Joe Peterson, Charit L. Pethiyagoda, Kaliprasad Pothuraju, Claudia Richter, Abraham M. Rosenbaum, Shaunak Roy, Jay Shafto, Uladzislau Sharanhovich, Karen W. Shannon, Conrad G. Sheppy, Michel Sun, Joseph V. Thakuria, Anne Tran, Dylan Vu, Alexander Wait Zaranek, Xiaodi Wu, Snezana Drmanac, Arnold R. Oliphant, William C. Banyai, Bruce Martin, Dennis G. Ballinger, George M. Church, Clifford A. Reid. *Science*, **January 2010.**

George Church (Personal Genome Project) Richmond Forum, Saturday

---

## Dystopia









Personalized Medicine

8

---

## Secure Two-Party Computation



Bob's Genome: ACTG…
Markers (~1000): [0,1, …, 0]

**Bob**

Alice's Genome: ACTG…
Markers (~1000): [0, 0, …, 1]

**Alice**

$$x = f(g_A, g_B)$$

Can Alice and Bob compute a function of their private data, without exposing anything about their data besides the result?

9

---

## Secure Function Evaluation

**Alice (circuit generator)**  **Bob (circuit evaluator)**

Agree on

$\text{Picks } a \in \{0,1\}^s$  $f(a,b) \to x$  $\text{Picks } b \in \{0,1\}^t$
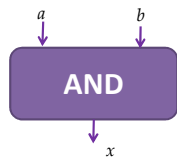
Garbled Circuit Protocol

Outputs $x = f(a,b)$
without revealing $a$
to Bob or $b$ to Alice.

Andrew Yao, 1982/1986
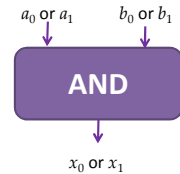
---

## Computing with Lookup Tables

| Inputs | | Output |
|---|---|---|
| $a$ | $b$ | $x$ |
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

$a$     $b$

**AND**

$x$

---

## Computing with Meaningless Values?

| Inputs | | Output |
|---|---|---|
| $a$ | $b$ | $x$ |
| $a_0$ | $b_0$ | $x_0$ |
| $a_0$ | $b_1$ | $x_0$ |
| $a_1$ | $b_0$ | $x_0$ |
| $a_1$ | $b_1$ | $x_1$ |

$a_i, b_i, x_i$ are **random** values, chosen by the **circuit generator** but **meaningless** to the **circuit evaluator**.

$a_0$ or $a_1$     $b_0$ or $b_1$

**AND**

$x_0$ or $x_1$

# Computing with Garbled Tables

| Inputs | | Output |
|--------|--------|--------|
| $a$ | $b$ | $x$ |
| $a_0$ | $b_0$ | $Enc_{a_0,b_0}(x_0)$ |
| $a_0$ | $b_1$ | $Enc_{a_0,b_1}(x_0)$ |
| $a_1$ | $b_0$ | $Enc_{a_1,b_0}(x_0)$ |
| $a_1$ | $b_1$ | $Enc_{a_1,b_1}(x_1)$ |

Bob can only decrypt **one** of these!

$a_i$, $b_i$, $x_i$ are **random** values, chosen by the **circuit generator** but **meaningless** to the **circuit evaluator**.

$a_0$ or $a_1$    $b_0$ or $b_1$

**AND**

$x_0$ or $x_1$

**Garbled And Gate**
$Enc_{a_0,b_1}(x_0)$
$Enc_{a_1,b_1}(x_1)$
$Enc_{a_1,b_0}(x_0)$
$Enc_{a_0,b_0}(x_0)$

---

# Garbled Circuit Protocol

**Alice (circuit generator)**    **Bob (circuit evaluator)**

Creates random keys: $a_0, a_1, b_0, b_1, x_0, x_1$

**And Gate**
$Enc_{a_0,b_1}(x_0)$
$Enc_{a_1,b_1}(x_1)$
$Enc_{a_1,b_0}(x_0)$
$Enc_{a_0,b_0}(x_0)$

Sends $a_i$ to Bob based on her input value    $a_0$

How does the Bob learn his own input wires?

---

# Primitive: **Oblivious Transfer**

**Alice**    **Bob**

Knows $b_0, b_1$    Picks $i \in \{0,1\}$
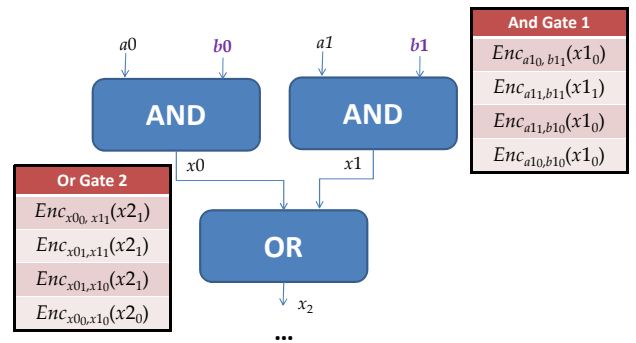
Oblivious Transfer Protocol

Learns nothing    Learns $b_i$ (only)

**Oblivious:** Alice doesn't learn which secret Bob obtains
**Transfer:** Bob learns one of Alice's secrets

Rabin, 1981; Even, Goldreich, and Lempel, 1985; many subsequent papers

---

# Chaining Garbled Circuits

$a0$    $b0$    $a1$    $b1$

**AND**    **AND**

**And Gate 1**
$Enc_{a1_0,b1_1}(x1_0)$
$Enc_{a1_1,b1_1}(x1_1)$
$Enc_{a1_1,b1_0}(x1_0)$
$Enc_{a1_0,b1_0}(x1_0)$

$x0$    $x1$

**Or Gate 2**
$Enc_{x0_0,x1_1}(x2_1)$
$Enc_{x0_1,x1_1}(x2_1)$
$Enc_{x0_1,x1_0}(x2_1)$
$Enc_{x0_0,x1_0}(x2_0)$

**OR**

$x_2$

...

We can do **any** computation privately this way!

---

# Building Computing Systems

$Enc_{x0_0,x1_1}(x2_1)$
$Enc_{x0_1,x1_1}(x2_1)$
$Enc_{x0_1,x1_0}(x2_1)$
$Enc_{x0_0,x1_0}(x2_0)$

| Digital Electronic Circuits | Garbled Circuits |
|---|---|
| Operate on **known data** | Operate on **encrypted wire labels** |
| One-bit logical operation requires moving a few electrons a few nanometers (many Billions per second) | One-bit logical operation requires performing (up to) 4 encryption operations (~100,000 gates per second) |
| Reuse is great! | Reuse is not allowed! |
| All basic operations have similar cost | Some logical operations nearly free (XOR) |

---

# Offspring Immune System Test

Major Histocompatibility Complex (MHC)

Alice's MHC genes:  [0 1 1 0 ... 0 1 1]
Bob's MHC genes:    [1 1 1 0 ... 1 0 1]

Diversity is key to good immune systems!

**Goal:** count number of indices where A[i] ≠ B[i]

## XOR
### Every Cryptographer's Favorite Function

$\oplus$

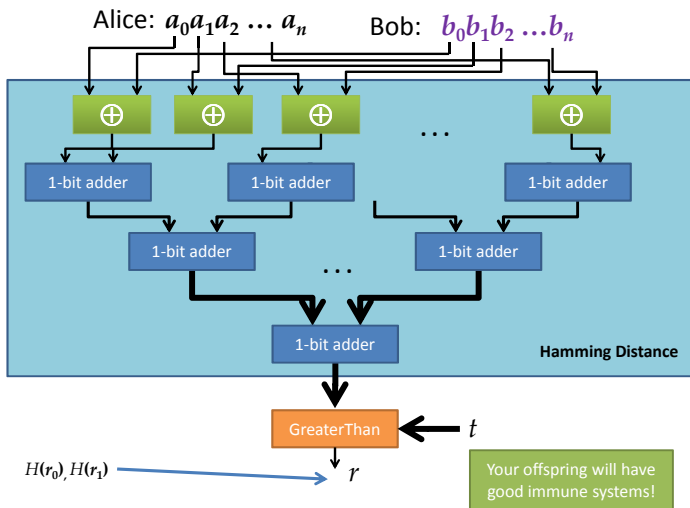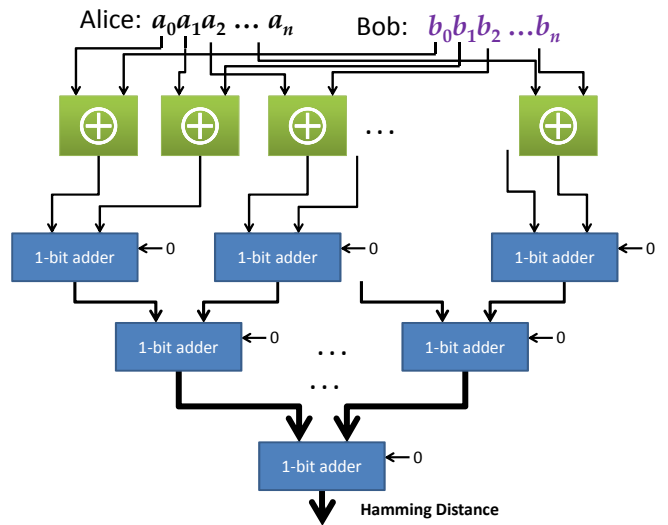| XOR | | |
|---|---|---|
| $a$ | $b$ | $x$ |
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

$a \oplus a = 0$

$a \oplus r$ is uniformly random

$a \oplus r \oplus r = a$

Can compute $a \oplus b$ on garbled inputs without and encryptions ("free").

19

---

Alice: $a_0 a_1 a_2 \ldots a_n$     Bob: $b_0 b_1 b_2 \ldots b_n$



Hamming Distance

---

Alice: $a_0 a_1 a_2 \ldots a_n$     Bob: $b_0 b_1 b_2 \ldots b_n$



Hamming Distance

GreaterThan  $t$

$H(r_0), H(r_1)$  $r$

Your offspring will have good immune systems!

---

## Heterozygous Recessive Risk

**Alice**



carrier

cystic fibrosis

Alice's Heterozygous Recessive genes: { 5283423, 1425236, 839523, … }
Bob's Heterozygous Recessive genes:   { 5823527, 839523, 169325, … }
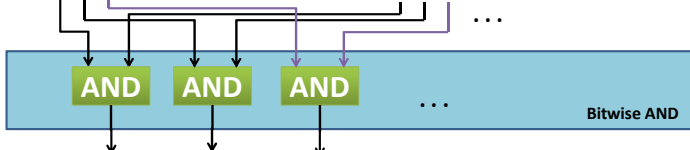
**Goal:** find the intersection of A and B

22

---

## Bit Vector Intersection

Alice's Recessive genes:
{ 5283423, 1425236, 839523, … }

Bob's Recessive genes:
{ 5823527, 839523, 169325, … }

[ PAH, PKU, CF, … ]

[ 0, 0, 1, 0, 0, 0, 1, 0, 1, 1, 0]     [ 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0]
...

AND  AND  AND   ...

Bitwise AND

23

---

## Scaling

What if there are millions of possible diseases?

Length of bit vector:

number of possible values

($2^L$ where $L$ is number of bits for each value)

Other private set intersection problems:
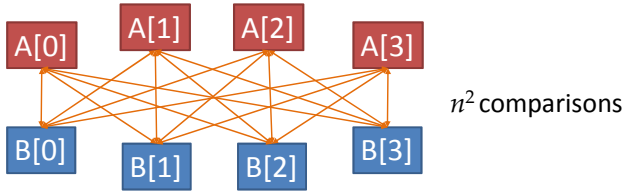   Do Alice and Bob have any friends in common?
   Data mining problems: combine medical records across hospitals
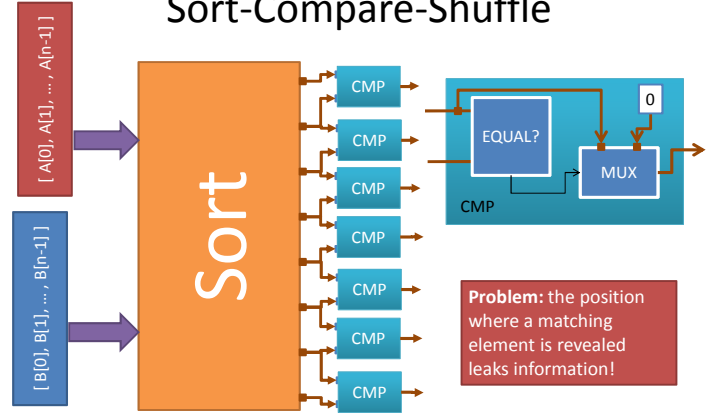   Two companies want to do joint marketing to common customers

24

## Pairwise Comparison

for i in range(0, n-1):
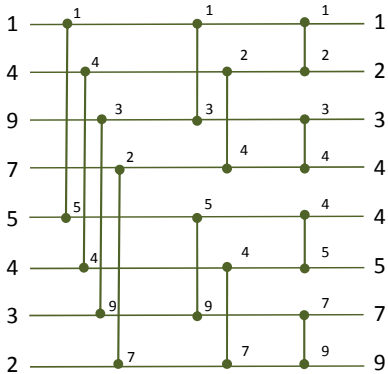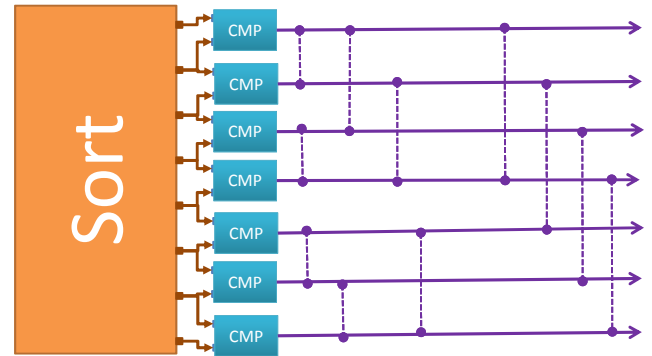 for j in range(0, n-1):
  if A[i] = B[j] output A[i]

A[0] A[1] A[2] A[3]

B[0] B[1] B[2] B[3]

$n^2$ comparisons

---

## Sort-Compare-Shuffle

[ A[0], A[1], … , A[n-1] ]

[ B[0], B[1], … , B[n-1] ]

Sort

CMP
CMP
CMP
CMP
CMP
CMP
CMP

EQUAL?
MUX
0
CMP

**Problem:** the position where a matching element is revealed leaks information!

---

## Bitonic Sorting

1 4 9 7 5 4 3 2

1 2 3 4 4 5 7 9

---

## (Imperfect) Shuffling

Sort

CMP
CMP
CMP
CMP
CMP
CMP
CMP

---

## Some Results

| Problem | Best Previous Result | Our Result | Speedup |
|---------|---------------------|------------|---------|
| **Hamming Distance** (Face Recognition, Genetic Dating) – two 900-bit vectors | 213s [SCiFI, 2010] | 0.051s | 4176 |
| **Levenshtein Distance** (genome, text comparison) – two 200-character inputs | 534s [Jha+, 2008] | 18.4s | 29 |
| **Smith-Waterman** (genome alignment) – two 60-nucleotide sequences | [Not Implementable] | 447s | - |
| **AES Encryption** | 3.3s [Henecka, 2010] | 0.2s | 16.5 |
| **Fingerprint Matching** (1024-entry database, 640x8bit vectors) | ~83s [Barni, 2010] | 18s | 4.6 |

Scalable: 1 Billion gates evaluated at ~100,000 gates/second on laptop

---

**Yan Huang**
(UVa Computer Science PhD Student)

**Jonathan Katz**
(University of Maryland)

**Aaron Mackey**
(UVa Center for Public Health Genomics)

## Introduction to
# Computing

Explorations in Language, Logic, and Machines
Spring 2010

## www.computingbook.org

David Evans
University of Virginia

Shameless
Plug #1

Shameless
Plug #2

Rice Hall, Computer Science Building
Opening August 2011

David Evans
evans@cs.virginia.edu
http://www.cs.virginia.edu/evans