

Privacy-preserving Information Sharing: Tools and Applications (Volume 2)

Emiliano De Cristofaro

University College London (UCL)

<https://emilianodc.com>

Prologue

Privacy-Enhancing Technologies (PETs):

Increase privacy of users, groups, and/or organizations

PETs often respond to privacy threats

Protect personally identifiable information

Support anonymous communications

Privacy-respecting data processing

Another angle: privacy as an enabler

Actively enabling scenarios otherwise impossible w/o clear privacy guarantees

Sharing Information w/ Privacy

Needed when parties with limited mutual trust willing or required to share information

Only the required minimum amount of information should be disclosed in the process

Private Set Intersection?

DHS (Terrorist Watch List) and **Airline** (Passenger List)

Find out whether any suspect is on a given flight

IRS (Tax Evaders) and **Swiss Bank** (Customers)

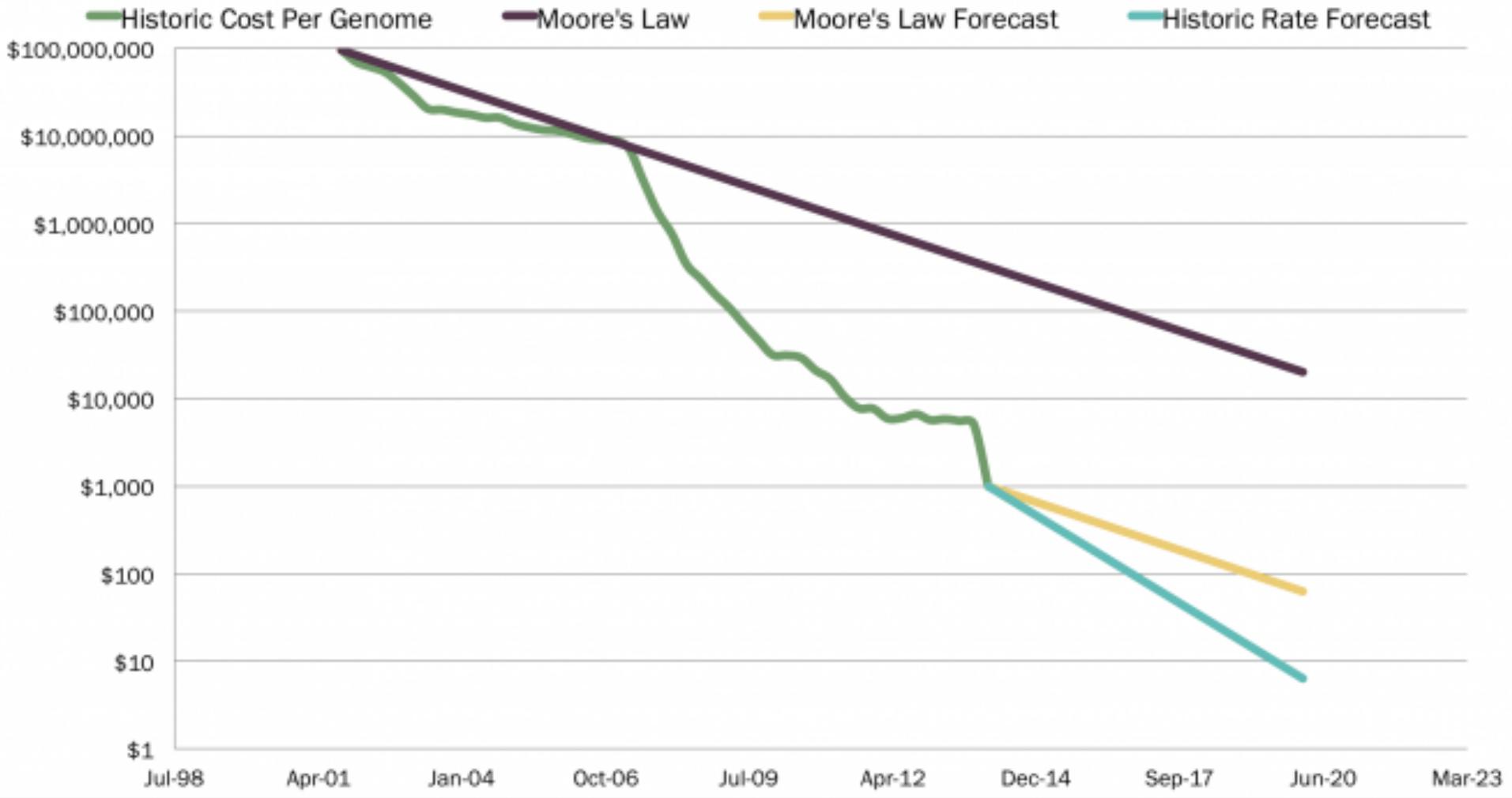
Discover if tax evaders have accounts at foreign banks

Hoag Hospital (Patients) and **SSA** (Social Security DB)

Patients with fake Social Security Number

Genomics

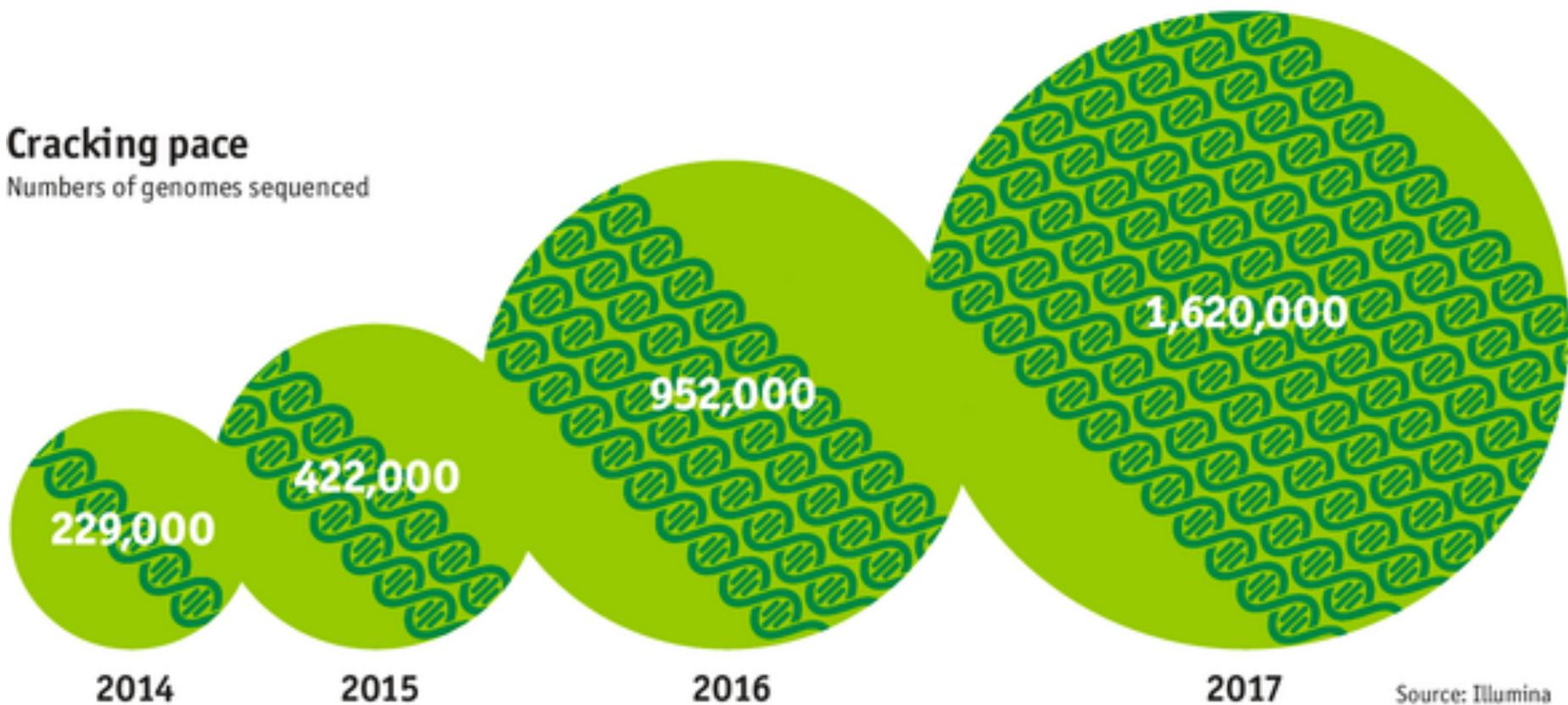
Cost Declines of Genome Sequencing



From: The Economist

Cracking pace

Numbers of genomes sequenced



Source: Illumina

1/05/2011 @ 4:57PM | 30,076 views

The First Child Saved By DNA Sequencing

+ Comment Now + Follow Comments



In Treatment for Leukemia, Glimpses of the Future



LETTER

doi:10.1038/nature13394

Genome sequencing identifies major causes of severe intellectual disability

Christian Gilissen^{1*}, Jayne Y. Hehir-Kwa^{1*}, Djie Tjwan Thung¹, Maartje van de Vorst¹, Bregje W. M. van Bon¹, Marjolein H. Willemsen¹, Michael Kwint¹, Irene M. Janssen¹, Alexander Hoischen¹, Annette Schenck¹, Richard Leach², Robert Klein², Rick Tearle², Tan Bo^{1,3}, Rolph Pfundt¹, Helger G. Yntema¹, Bert B. A. de Vries¹, Tjitske Kleefstra¹, Han G. Brunner^{1,4*}, Lisenka E. L. M. Vissers^{1*} & Joris A. Veltman^{1,4*}

TIME

THE ANGELINA EFFECT

Angelina Jolie's double mastectomy puts genetic testing in the spotlight. What her choice reveals about calculating risk, cost and peace of mind

BY JEFFREY KLUGER & ALICE PARK

Genetic Risk Factors (11) ?

| REPORT | RESULT |
|---|------------------------------|
| Alpha-1 Antitrypsin Deficiency | Variant Absent; Typical Risk |
| Alzheimer's Disease (APOE Variants) | ε4 Variant Absent |
| Early-Onset Primary Dystonia (DYT1-TOR1A-Related) | Variant Absent; Typical Risk |
| Factor XI Deficiency | Variant Absent; Typical Risk |
| Familial Hypercholesterolemia Type B (APOB-Related) | Variant Absent; Typical Risk |

[See all 11 genetic risk factors...](#)

Traits (41) ?

| REPORT | RESULT |
|---|----------------|
| Alcohol Flush Reaction | Does Not Flush |
| Bitter Taste Perception | Can Taste |
| Blond Hair | 28% Chance |
| Earwax Type | Wet |
| Eye Color | Likely Brown |

[See all 41 traits...](#)

Inherited Conditions (43) ?

| REPORT | RESULT |
|--|-----------------|
| Beta Thalassemia | Variant Present |
| ARSACS | Variant Absent |
| Agenesis of the Corpus Callosum with Peripheral Neuropathy (ACCPN) | Variant Absent |
| Autosomal Recessive Polycystic Kidney Disease | Variant Absent |
| Bloom's Syndrome | Variant Absent |

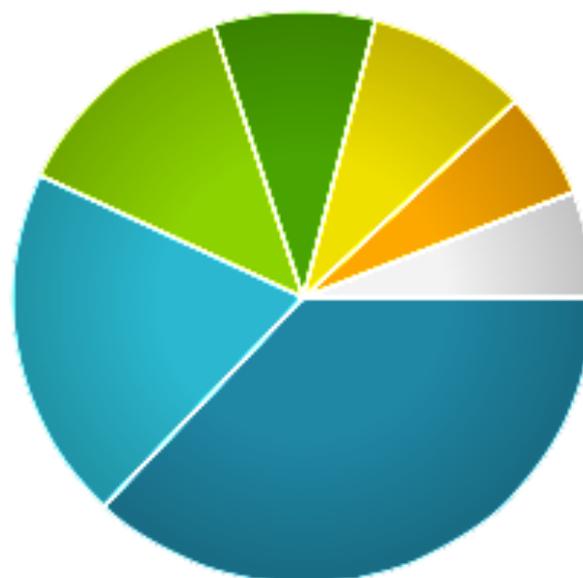
[See all 43 carrier status...](#)

Drug Response (12) ?

| REPORT | RESULT |
|--|-----------------|
| Proton Pump Inhibitor (PPI) Metabolism (CYP2C19-related) | Rapid |
| Warfarin (Coumadin®) Sensitivity | Increased |
| Phenytoin Sensitivity (Epilepsy Drug) | Increased |
| Sulfonylurea Metabolism | Greatly reduced |
| Abacavir Hypersensitivity | Typical |

[See all 12 drug response...](#)

Genetic Ethnicity



| | | |
|---|------------------------------|------------|
|  | Southern European | 37% |
|  | West African | 20% |
|  | British Isles | 13% |
|  | Native South American | 9% |
|  | Finnish/Volga-Ural | 9% |
|  | Eastern European | 6% |
|  | Uncertain | 6% |

List View Map View Surname View

search matches Show: both sides Sort: relationship 25 per page 1 - 25 of 424

| | | | | |
|---|---------------|---|--|--|
|  | Male | You | | UPDATE YOUR PROFILE |
|  | Female | 2nd to 3rd Cousin 1.68% shared, 5 segments | J2a2 | Send an Introduction |
|  | Female | 3rd to 4th Cousin 1.30% shared, 3 segments | United States Alsace-Lorraine (Strasbourg), Fr... Paternal Senape 5 more U5b2 | Public Match Send a Message |
|  | Male | 3rd to 4th Cousin 1.03% shared, 2 segments | H13a1a R1b1b2 | Send an Introduction |
|  | Female | 3rd to 5th Cousin 0.45% shared, 2 segments | H7 | Send an Introduction |
|  | Female | 3rd to 5th Cousin 0.42% shared, 2 segments | H1 | Send an Introduction |
|  | Male | 3rd to 5th Cousin 0.40% shared, 2 segments | United States Reno, Nevada San Diego, California Tucker Littlefield Warga 4 more H1c G2a | Public Match Send a Message |
|  | Male | 3rd to 5th Cousin 0.37% shared, 2 segments | United States fathers father prince Edward isla... R1b1b2a1a K1a1b | Public Match Send a Message |
|  | Male, b. 1978 | 3rd to 6th Cousin 0.40% shared, 1 segment | United States New Jersey Utah California Northern Europe U3b1 T | Send an Introduction |

Privacy Researcher's Perspective

Treasure trove of **sensitive** information

Ethnic heritage, predisposition to diseases

Genome = the ultimate **identifier**

Hard to anonymize / de-identify

Sensitivity is **perpetual**

Cannot be “revoked”

Leaking one's genome \approx leaking relatives' genome

Secure Genomics?

Privacy:

Individuals remain in control of their genome

Allow doctors/clinicians/labs to run genomic tests, while disclosing the required minimum amount of information, i.e.:

(1) Individuals don't disclose their entire genome

(2) Testing facilities keep test specifics (“secret sauce”) confidential

[BBDGT11]: Secure genomics via PSI

Most personalized medicine tests in < 1 second

Works on Android too

Genetic Paternity Test

A Strawman Approach for Paternity Test:

On average, ~99.5% of any two human genomes are identical

Parents and children have even more similar genomes

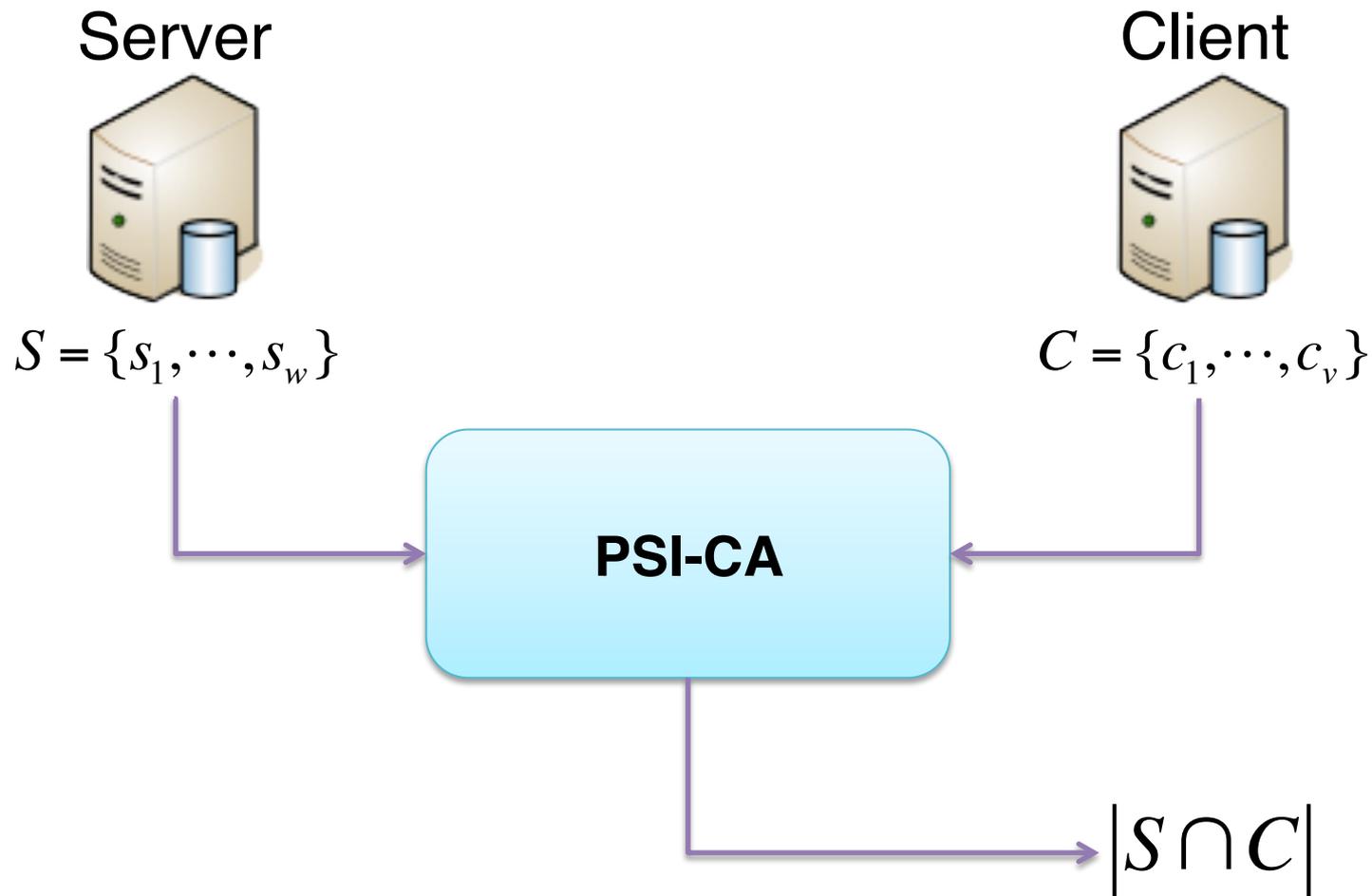
Compare candidate's genome with that of the alleged child:

Test positive if percentage of matching nucleotides is $> 99.5 + \tau$

First-Attempt Privacy-Preserving Protocol:

Use an appropriate secure two-party protocol for the comparison

Private Set Intersection Cardinality (PSI-CA)



Genetic Paternity Test

A Strawman Approach for Paternity Test:

On average, ~99.5% of any two human genomes are identical

Parents and children have even more similar genomes

Compare candidate's genome with that of the alleged child:

Test positive if percentage of matching nucleotides is $> 99.5 + \tau$

First-Attempt Privacy-Preserving Protocol:

Use an appropriate secure two-party protocol for the comparison

PROs: High-accuracy and error resilience

CONs: Performance not promising (3 billion symbols in input)

In our experiments, computation takes a few days

Genetic Paternity Test

Wait a minute!

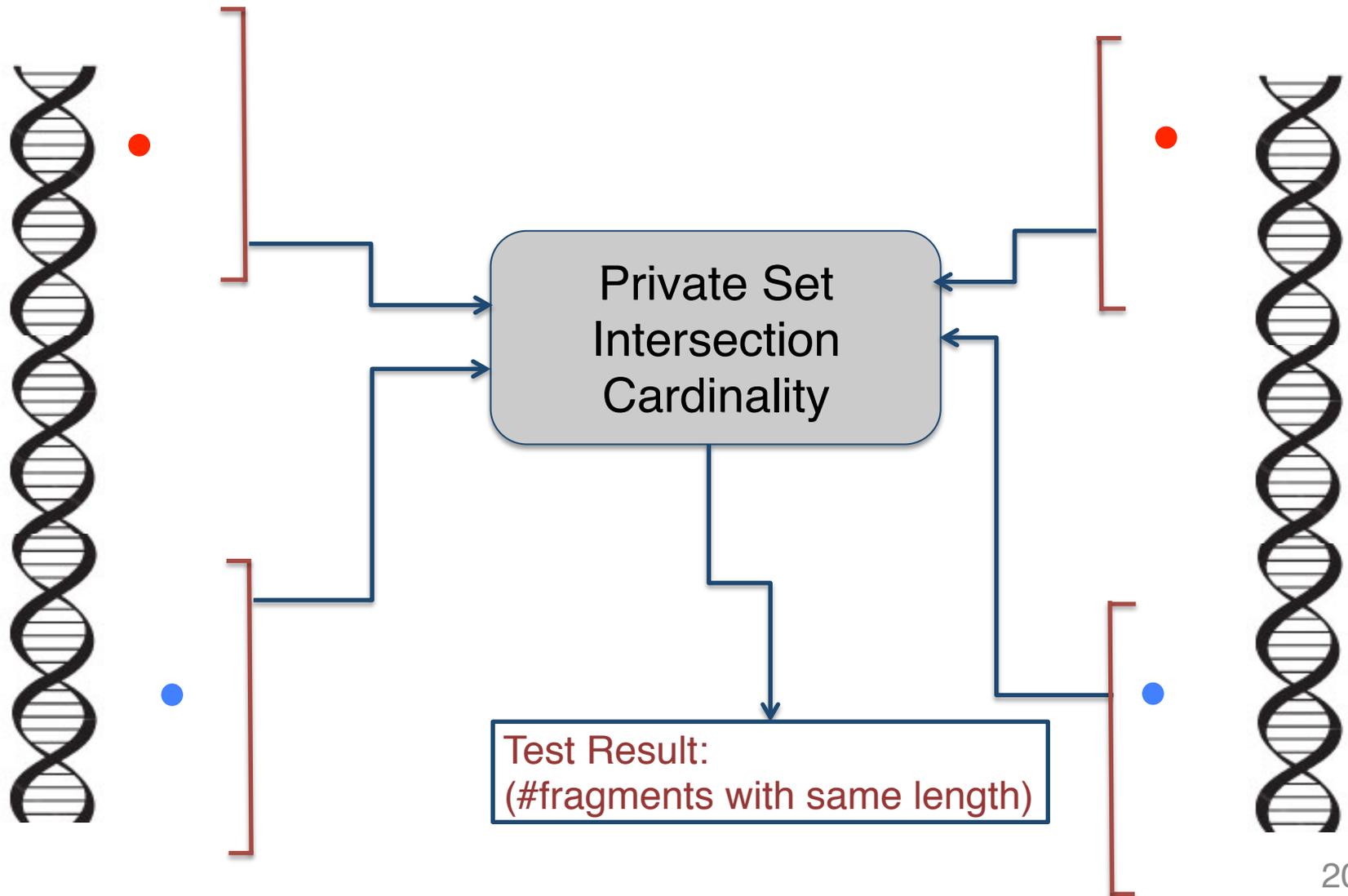
~99.5% of any two human genomes are identical

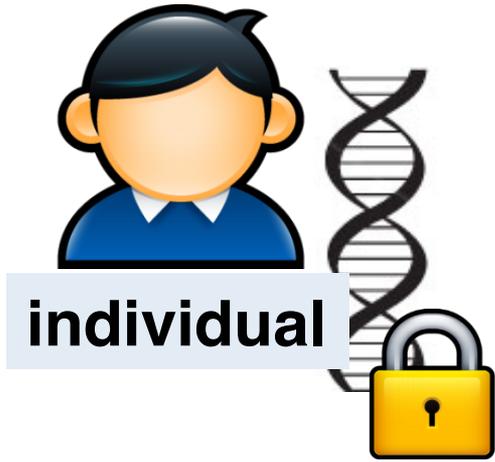
Why don't we compare *only* the remaining 0.5%?

We can compare by counting how many

But... We don't know (yet) where *exactly* this 0.5% occur!

Private RFLP-based Paternity Test





individual

genome

- Private Set Intersection (PSI)
- Authorized PSI
- Private Pattern Matching
- Homomorphic Encryption
- Garbled Circuits
- [...]



**doctor
or lab**

test specifics

**Secure
Function
Evaluation**

- Paternity/Ancestry Testing
- Testing of SNPs/Markers
- Compatibility Testing
- Disease Predisposition [...]

test result

test result

**Output reveals nothing beyond
test result**

Personalized Medicine (PM)

Drugs designed for patients' genetic features

Associating drugs with a unique genetic fingerprint

Max effectiveness for patients with matching genome

Test drug's "genetic fingerprint" against patient's genome

Examples:

tpmt gene – relevant to leukemia

(1) G->C mutation in pos. 238 of gene's c-DNA, or (2) G->A mutation in pos. 460 and one A->G is pos. 419 cause the *tpmt* disorder (relevant for leukemia patients)

hla-B gene – relevant to HIV treatment

One G->T mutation (known as *hla-B*5701* allelic variant) is associated with extreme sensitivity to abacavir (HIV drug)

Privacy-preserving PM Testing (P³MT)

Challenges:

Patients may refuse to unconditionally release their genomes

Or may be sued by their relatives...

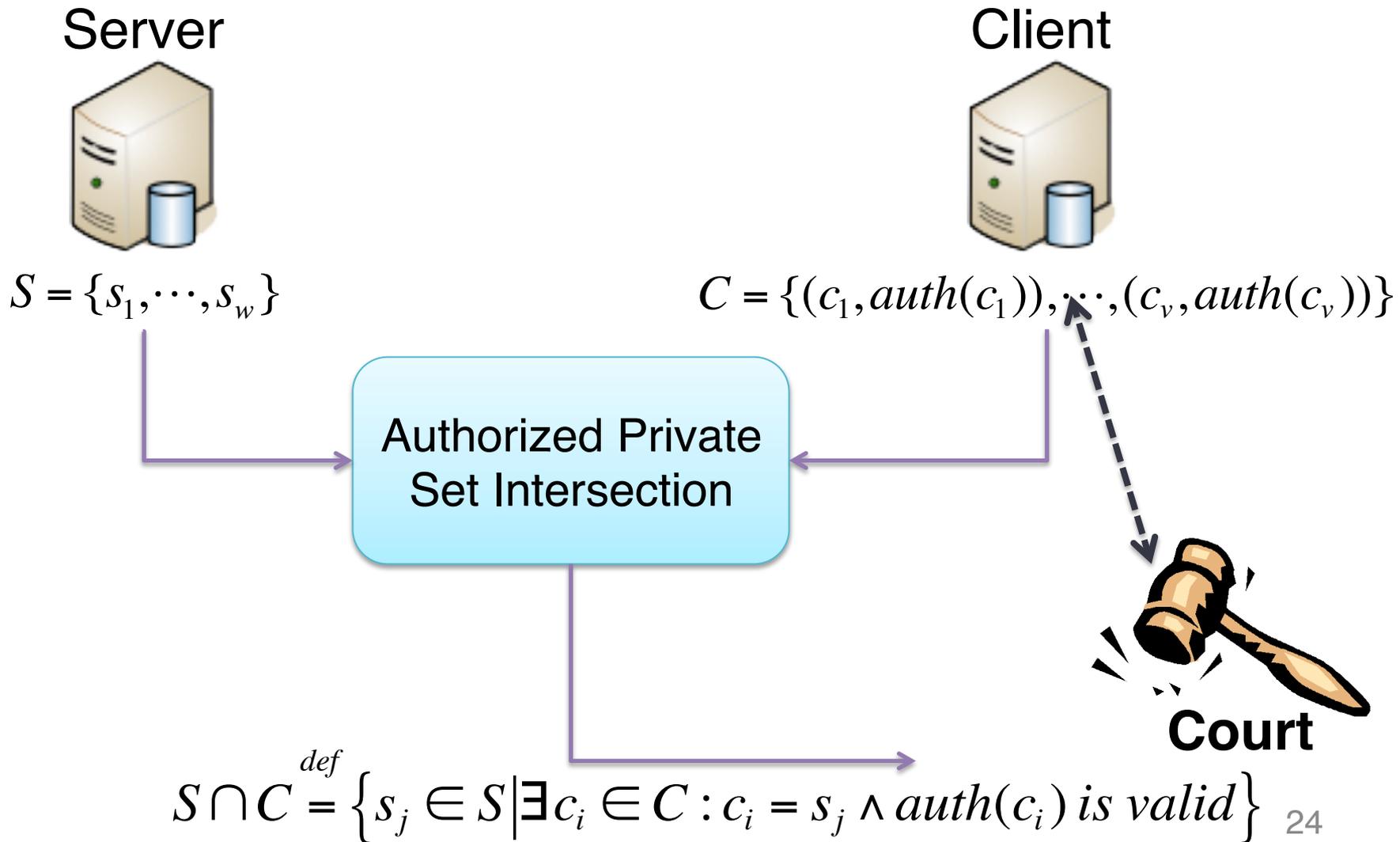
DNA fingerprint corresponding to a drug may be proprietary:

✓ **We need privacy-protecting fingerprint matching**

But we also need to enable FDA approval on the drug/fingerprint

✓ **We reduce P³MT to Authorized Private Set Intersection (APSI)**

Authorized Private Set Intersection (APSI)



Reducing P³MT to APSI

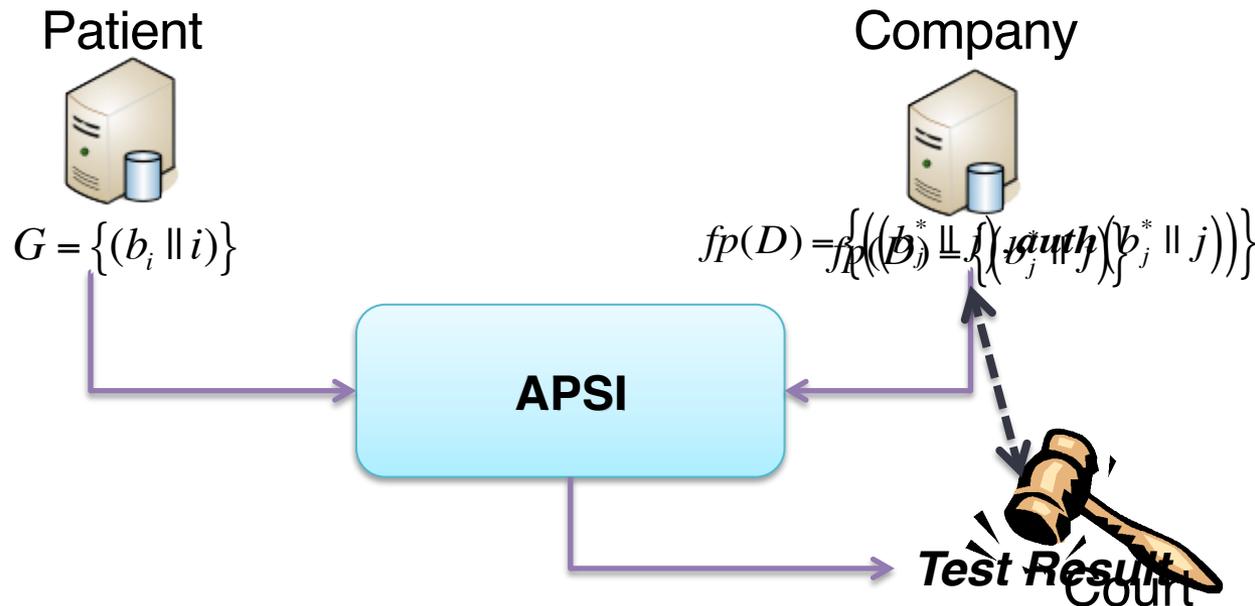
Intuition:

FDA = Court, Pharma = *Client*, Patient = *Server*

Patient's private input set: $G = \{(b_i \parallel i) \mid b_i \in \{A, C, G, T\}\}_{i=1}^{3 \cdot 10^9}$

Pharmaceutical company's input set: $fp(D) = \{(b_j^* \parallel j)\}$

Each item in $fp(D)$ needs to be authorized by FDA



P³MT – Performance Evaluation

Pre-Computation

Patient's pre-processing of the genome: a few days

Optimization:

Patient applies reference-based compression techniques

Input all *differences* with “reference” genome (0.5%)

Online Computation

Depend (linearly) on fingerprint size – typically a few nucleotides, <1s for most tests

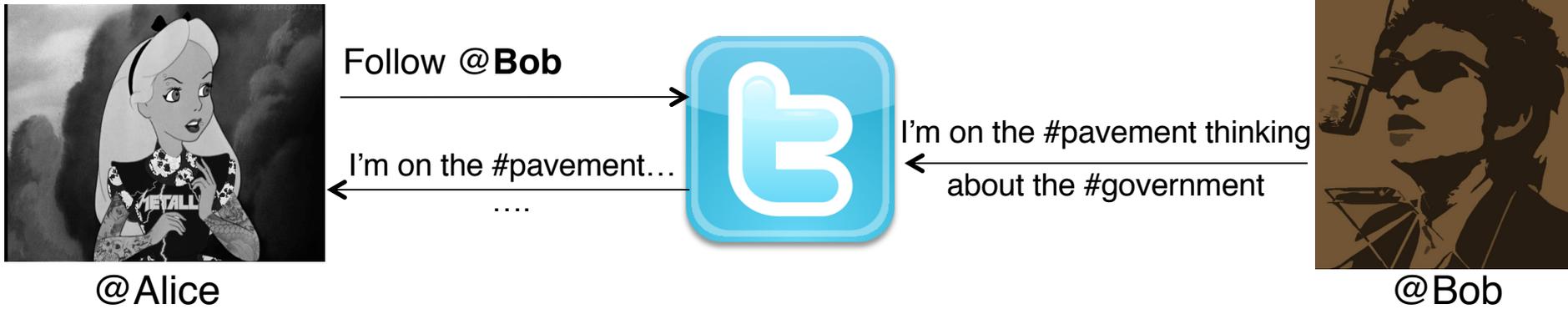
Communication

Depends on the size of encrypted genome (about 4GB)

Open Problems?

Micro-blogging

@Alice and @Bob – Twitter edition



There might be no mutual knowledge/trust between Alice and Bob

Follow requests are approved by default (opt-out)

Tweets are public by default

Streamed into www.twitter.com/public_timeline, available through API

But Bob can restrict his tweets to followers

All public tweets are searchable by hashtag

#Privacy and Twitter



Twitter.com is “trusted” to

- Get all tweets

- Enforce coarse-grained access control (follower-only)

- Monitor relations between users

Privacy and Twitter

- Targeted advertisement, PII collected and shared with third parties

- Trending topics, real-time “news”

I don't care about #privacy on @Twitter... but

Remember @Wikileaks? Snowden?

Our proposal: Hummingbird

Follow by hashtag:

E.g., @Alice follows @Bob only on hashtag #privacy

Tweeter (@Bob)

Learns who follows him but not which hashtags have been subscribed to

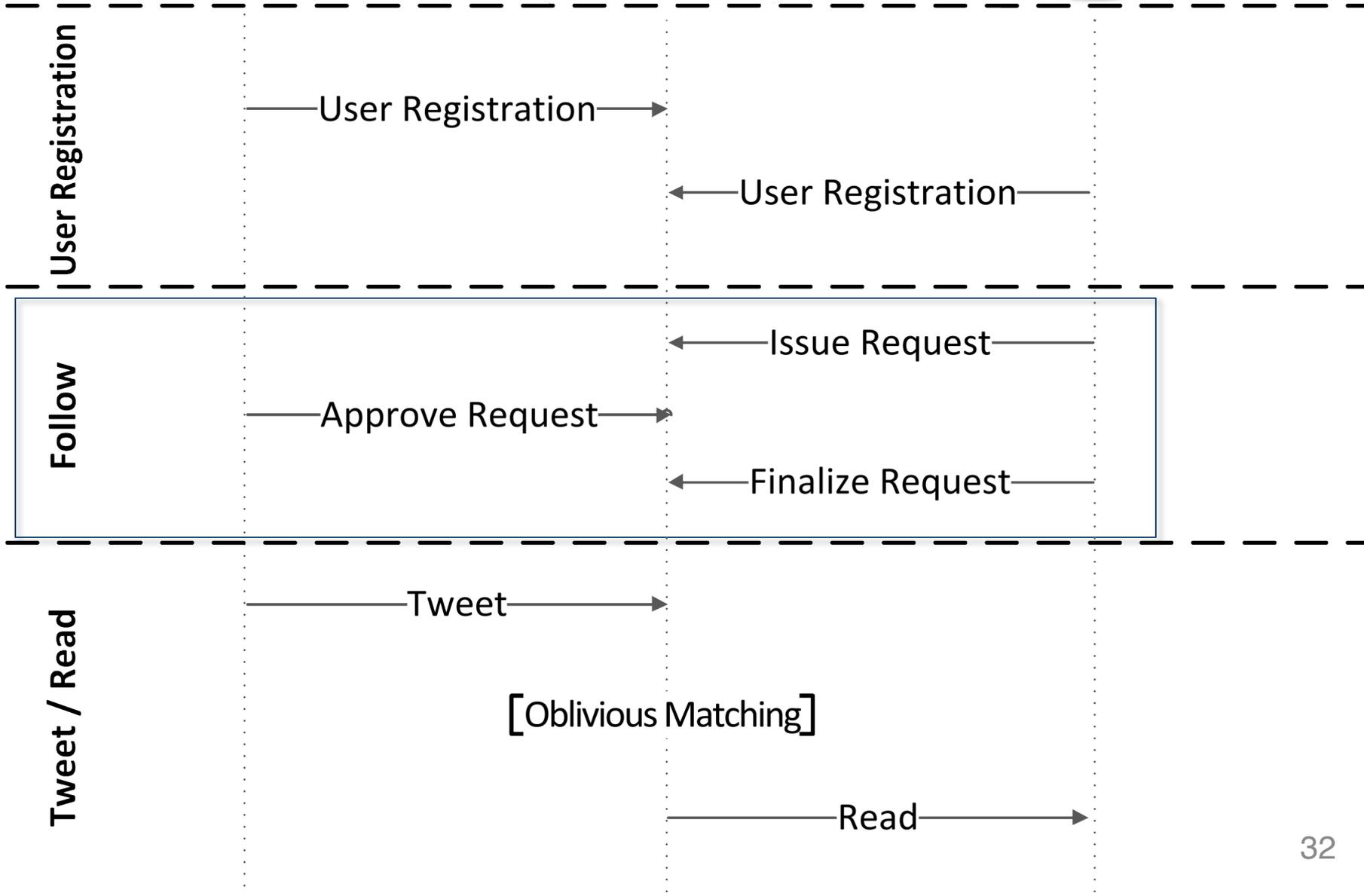
Follower (@Alice)

Learns nothing beyond her own subscriptions

Server (HS)

Doesn't learn tweets' content or hashtags of in
(But can scale to million of tweets/users)





HS

Issue Request

$$(N_b, e_b)$$

$$(Alice, Bob, \mu)$$

Alice (ht)

$$r \in \mathbb{Z}_{N_b}$$

$$\mu = H(ht) \cdot r^{e_b}$$

HS

Approve

$$(Alice, \mu)$$

$$(\mu')$$

Bob

$$\mu' = \mu^{d_b}$$

HS

Finalize Request

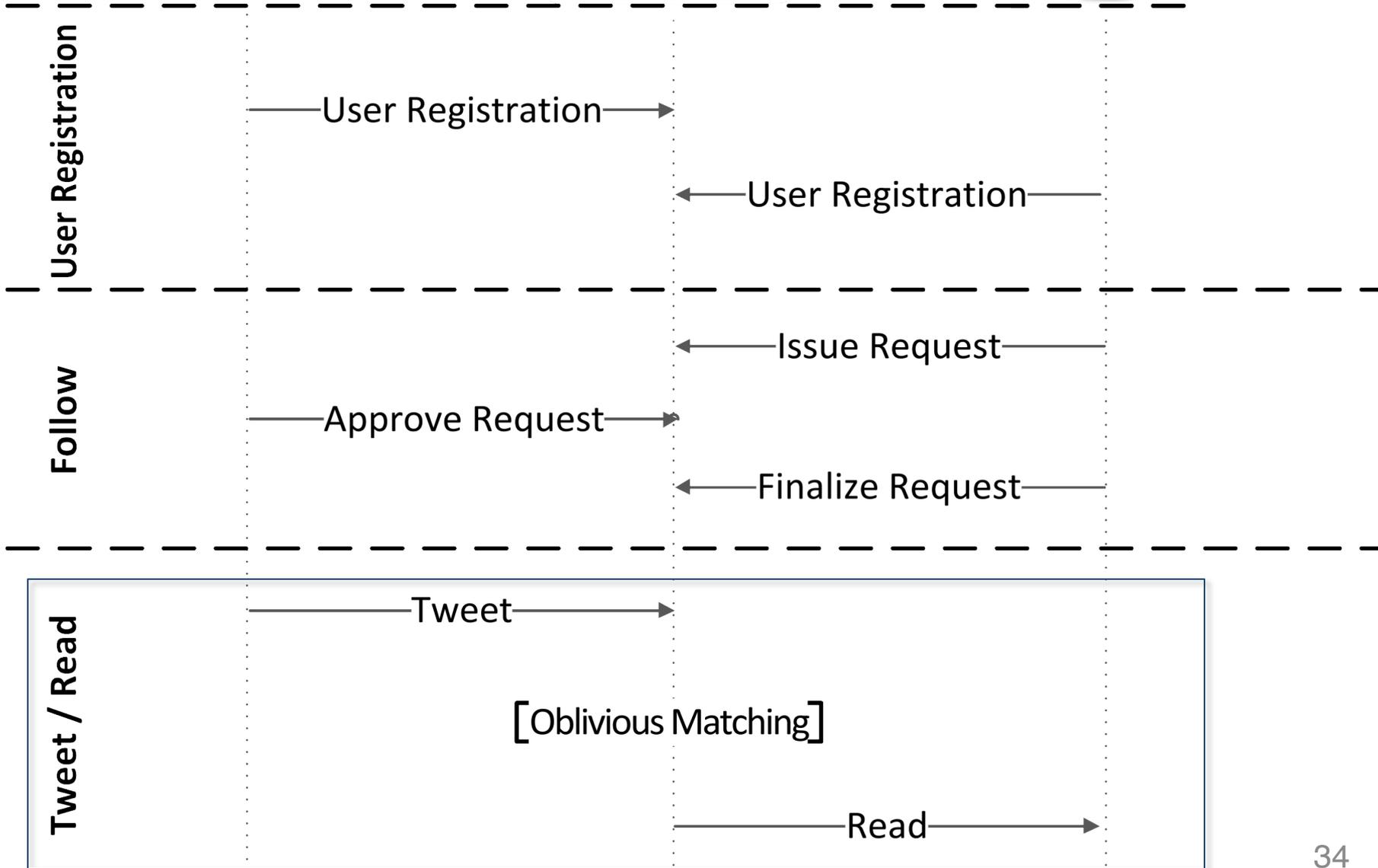
$$(Bob, \mu')$$

$$(Alice, Bob, t)$$

Alice (ht)

$$\delta = \mu' / r$$

$$t = H'(\delta)$$



HS

Tweet

Bob(d_b, M, ht^*)

$$\begin{aligned}\delta &= H(ht^*)^{d_b} \\ t^* &= H'(\delta) \quad k^* = H''(\delta) \\ ct^* &= Enc_{k^*}(M)\end{aligned}$$

(t^*, ct^*)

HS

For all (U, V, t) s.t. $V = \text{'Bob'}$ and $t = t^*$:
Store and mark (Bob, t^*, ct^*) for delivering (t^*, ct^*) to Alice

Oblivious Matching

HS

Read

(Bob, t^*, ct^*)

Alice(δ, t)

$$\begin{aligned}k &= H''(\delta) \\ M &= Dec_k(ct^*)\end{aligned}$$

Overhead

Follow protocol: Alice wants to follow Bob on #privacy

Bob's computation: 1 CRT-RSA signature (<1ms) per hashtag

Alice's computation: 2 mod multiplications per hashtag

Communication: 2 RSA group elements (<1KB)

Tweet: Bob tweets "I'm at #fosad!"

Computation: 1 CRT-RSA signature (<1ms) per hashtag, 1 AES enc

Communication: 1 hash output (160-bit)

Read

Computation: 1 AES decryption

Communication: 1 hash output (160-bit)

Server

No crypto!

Overhead: matching of PRF outputs, 160-bit

Can do efficiently, just like for cleartexts



Collecting Statistics Privately?

Collaboratively Train Machine Learning Models, Privately?

Why are statistics important?

Examples:

1. Recommender systems for online streaming services
2. Statistics about mass transport movements
3. Traffic statistics for the Tor Network

How about privacy?

Private Recommendations

BBC keeps 500-1000 free programs on iPlayer

No account, no tracking, no ads

Still, BBC wants to collect statistics, offer recommendations to its users

E.g., you have watched Dr Who, maybe you'll like Sherlock Homes too!

Item-KNN Recommendation

Predict favorite items for users based on their own ratings and those of “similar” users

Consider N users, M TV programs and binary ratings (viewed/not viewed)

Build a co-views matrix C , where C_{ab} is the number of views for the pair of programs (a,b)

Compute the **Similarity Matrix**

$$\{Sim\}_{ab} = \frac{C_{ab}}{\sqrt{C_a \cdot C_b}}$$

Identify K-Neighbours (**KNN**) based on matrix

Privacy-Preserving Aggregation

Goal: aggregator collects matrix, s.t.

Can only learn aggregate counts (e.g., 237 users have watched both a and b)

Not who has watched what

Use additively homomorphic encryption?

$$\text{Enc}_{PK}(a) * \text{Enc}_{PK}(b) = \text{Enc}_{PK}(a+b)$$

How can I use it to collect statistics?

Keys summing up to zero

Users U_1, U_2, \dots, U_N , each has k_1, k_2, \dots, k_N s.t.

$$k_1 + k_2 + \dots + k_N = 0$$

Now how can I use this?

User \mathcal{U}_i ($i \in [1, N]$)

Tally

$$x_i \in_r \mathbb{G}, y_i := g^{x_i} \bmod q \xrightarrow{y_i}$$

$$k_{i_\ell} := \sum_{j \neq i} \mathbf{H}(y_j^{x_i} \parallel \ell \parallel s) \cdot (-1)^{i > j} \bmod 2^{32} \xleftarrow{\{y_j\}_{j \in [1, N]}}$$

$$b_{i_\ell} := X_{i_\ell} + k_{i_\ell} \bmod 2^{32} \xrightarrow{\{b_{i_\ell}\}_{\ell=1}^L} \text{Fault recovery (if needed)}$$

$$\xleftarrow{\mathcal{U}^{on}}$$

$$k'_{i_\ell} := \sum_{\substack{j \neq i, \\ j \notin \mathcal{U}^{on}}} \mathbf{H}(y_j^{x_i} \parallel \ell \parallel s) \cdot (-1)^{i > j} \bmod 2^{32} \xrightarrow{\{k'_{i_\ell}\}_{\ell=1}^L} C'_\ell := \left(\sum_{i \in \mathcal{U}^{on}} b_{i_\ell} - \sum_{i \in \mathcal{U}^{on}} k'_{i_\ell} \right) \bmod 2^{32}$$

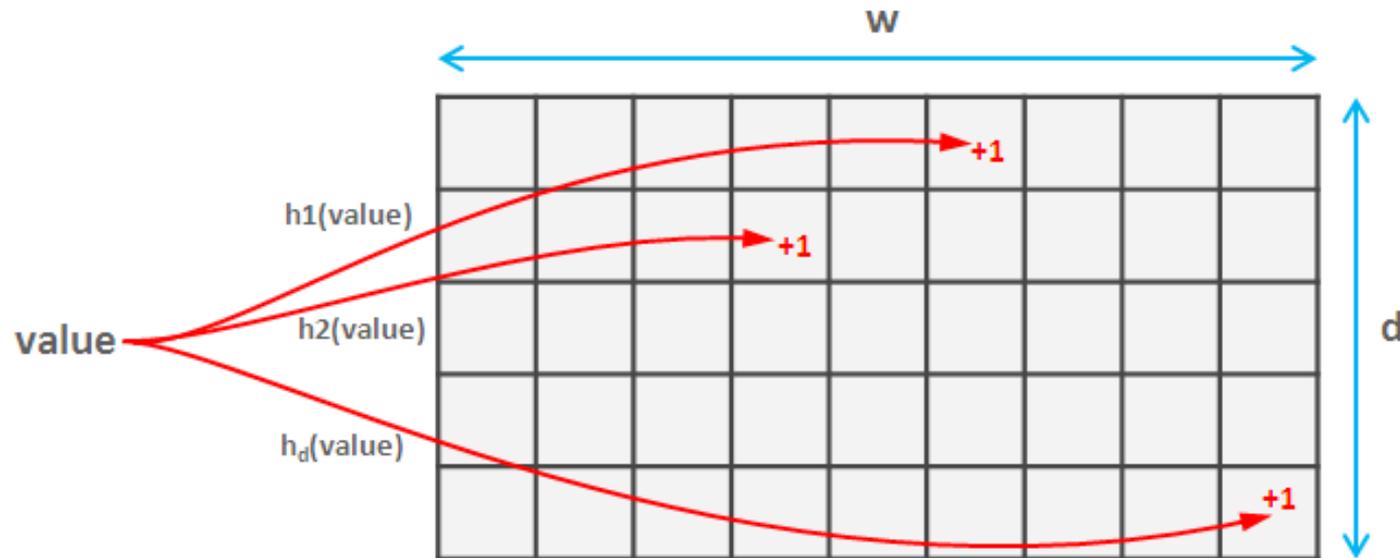
Is this efficient?

Preliminaries: Count-Min Sketch

An estimate of an item's frequency in a stream

Mapping a stream of values (of length T) into a matrix of size $O(\log T)$

The sum of two sketches results in the sketch of the union of the two data streams



Security & Implementation

Security

In the honest-but-curious model under the CDH assumption

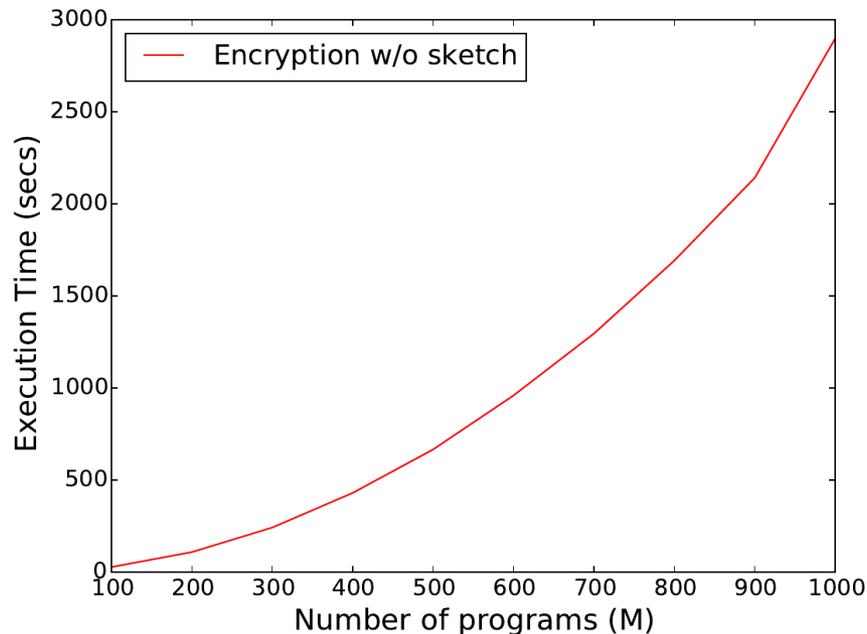
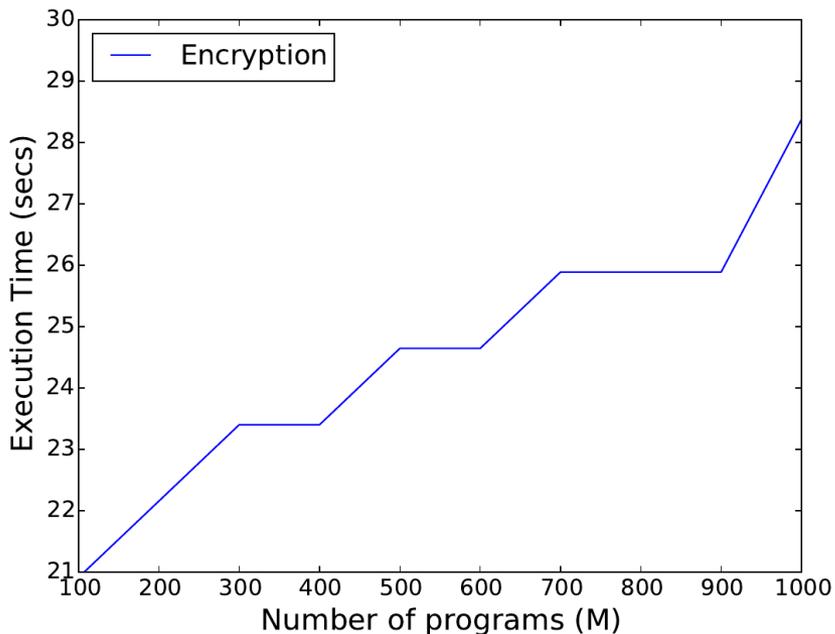
Prototype implementation:

Tally as a Node.js web server

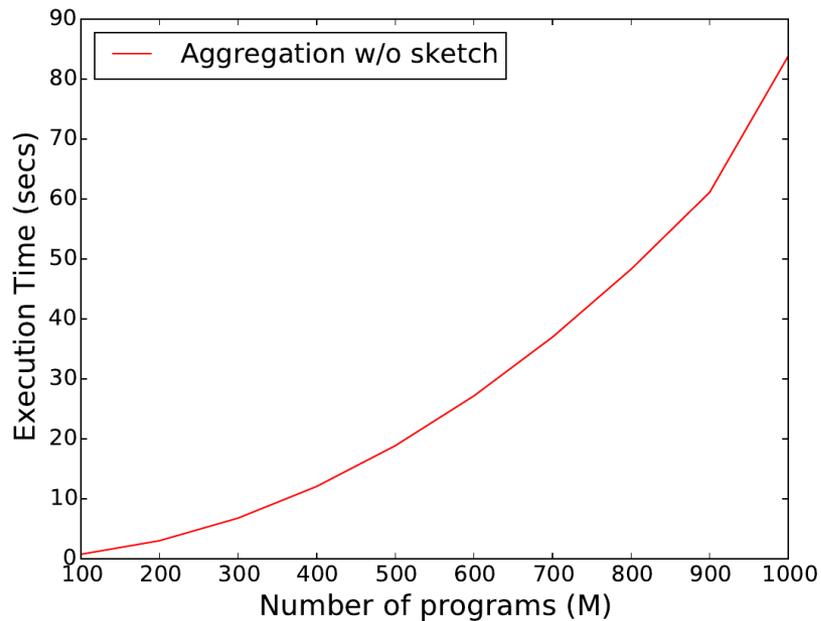
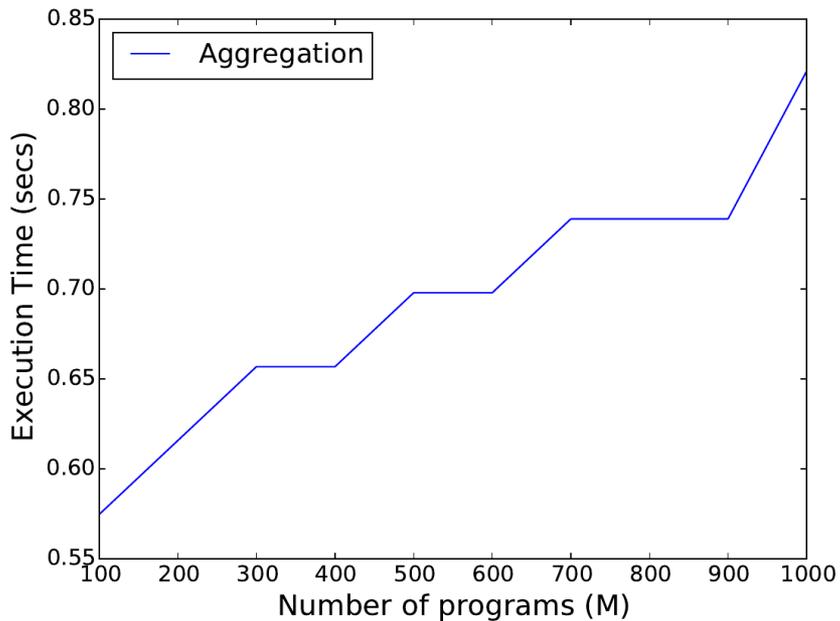
Users run in the browser or as a mobile cross-platform application (Apache Cordova)

Transparency, ease of use, ease of deployment

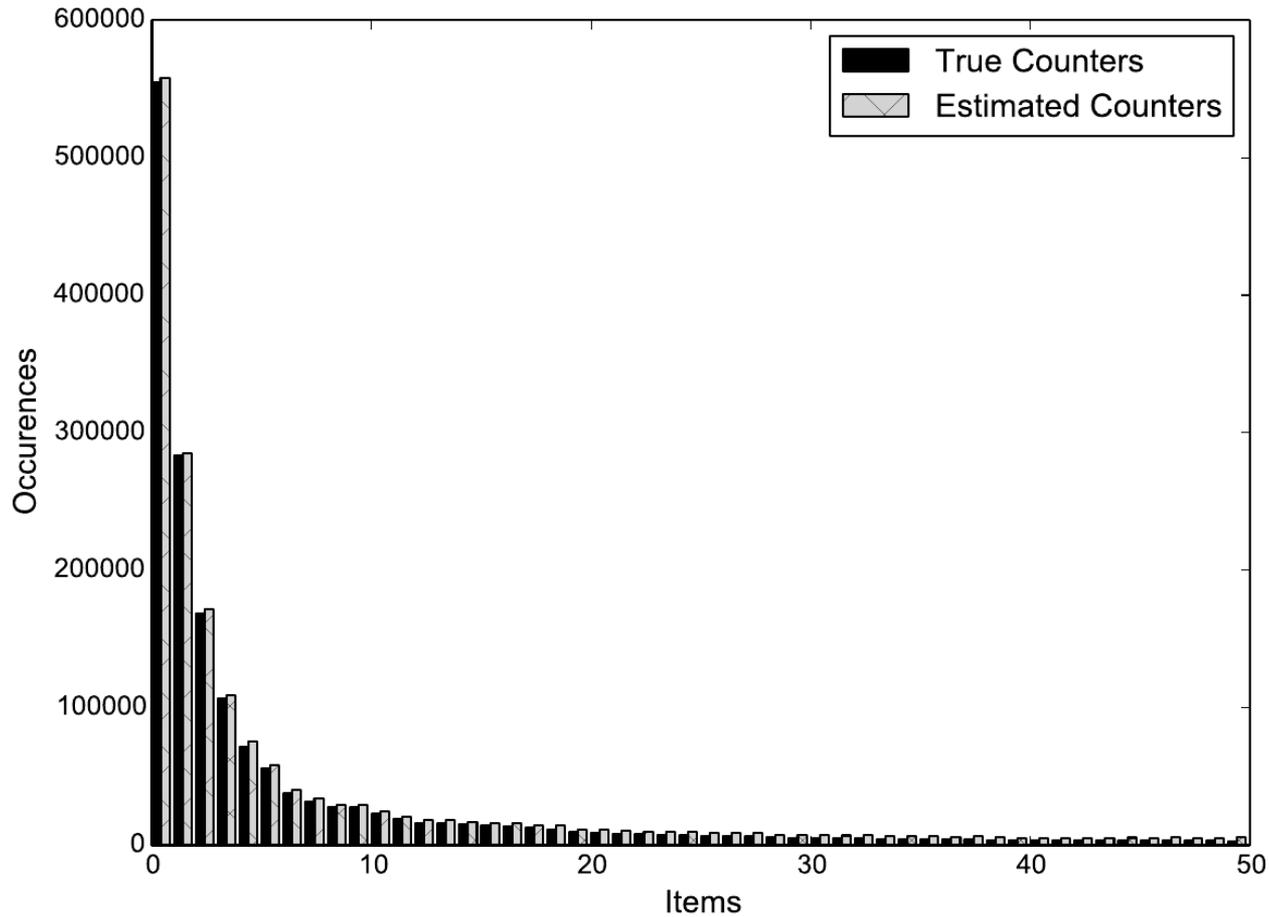
User side



Server side



Accuracy



Tor Hidden Services

Aggregate statistics about the number of hidden service descriptors from multiple HSDirs

Median statistics to ensure robustness

Problem: Computation of statistics from collected data can potentially de-anonymize individual Tor users or hidden services

Private Tor Statistics?

We rely on:

- A set of authorities

- A homomorphic public-key scheme (AH-ECC)

- Count-Sketch (a variant of CMS)

Setup phase

- Each authority generates their public and private key

- A group public key is computed

Private Tor Statistics?

Each HSDir (router) builds a Count-Sketch, inserts its values, encrypts it, sends it to a set of authorities

The authorities:

- Add the encrypted sketches element-wise to generate one sketch characterizing the overall network traffic

- Execute a divide and conquer algorithm on this sketch to estimate the median

How we do it (1/2)

The range of the possible values is known

On each iteration, the range is halved and the sum of all the elements on each half is computed

Depending on which half the median falls in, the range is updated and again halved

Process stops once the range is a single element

How we do it (2/2)

Output privacy:

Volume of reported values within each step is leaked

Provide *differential privacy* by adding Laplacian noise to each intermediate value

Evaluating

Experimental setup:

1200 samples from a mixture distribution

Range of values in $[0, 1000]$

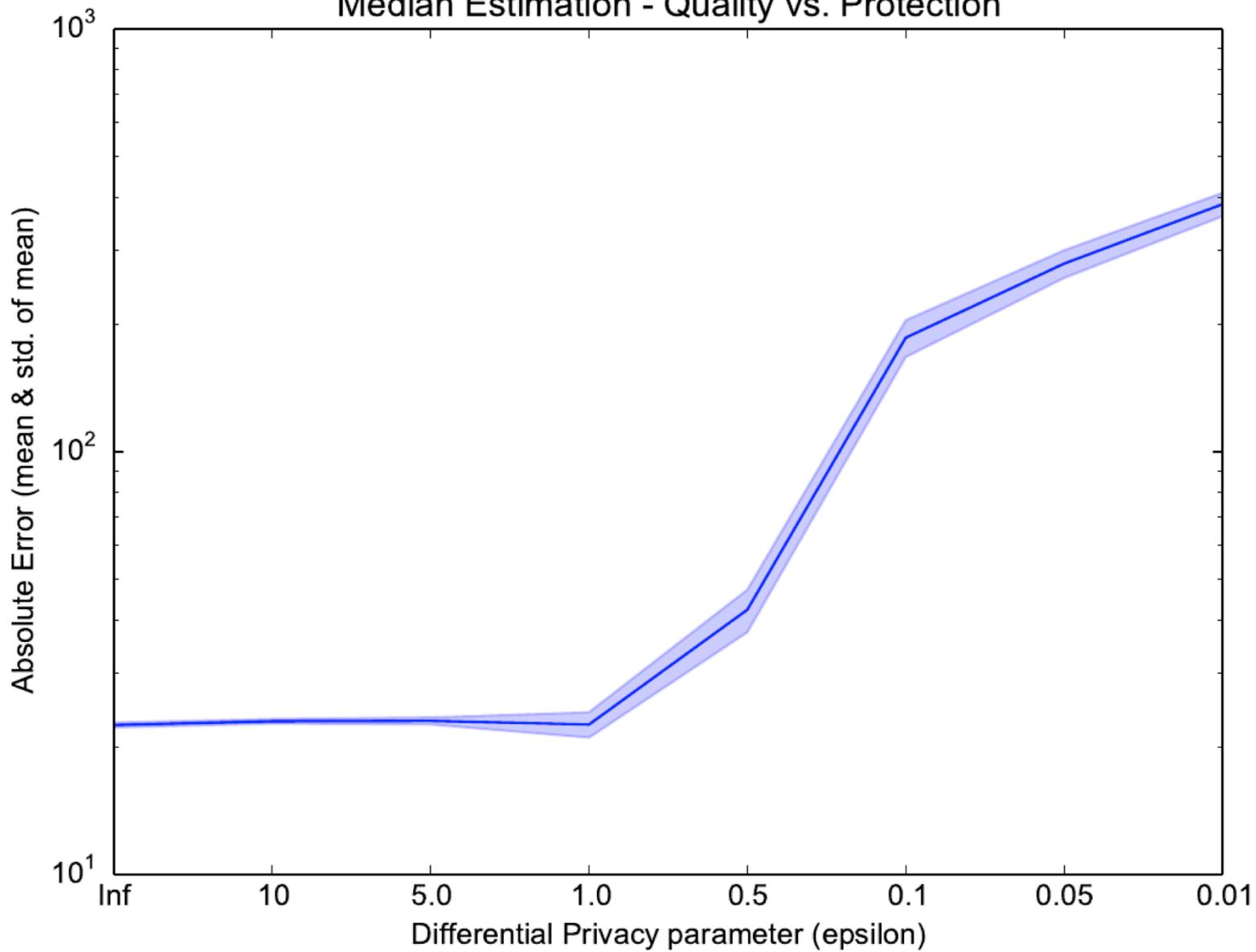
Performance evaluation:

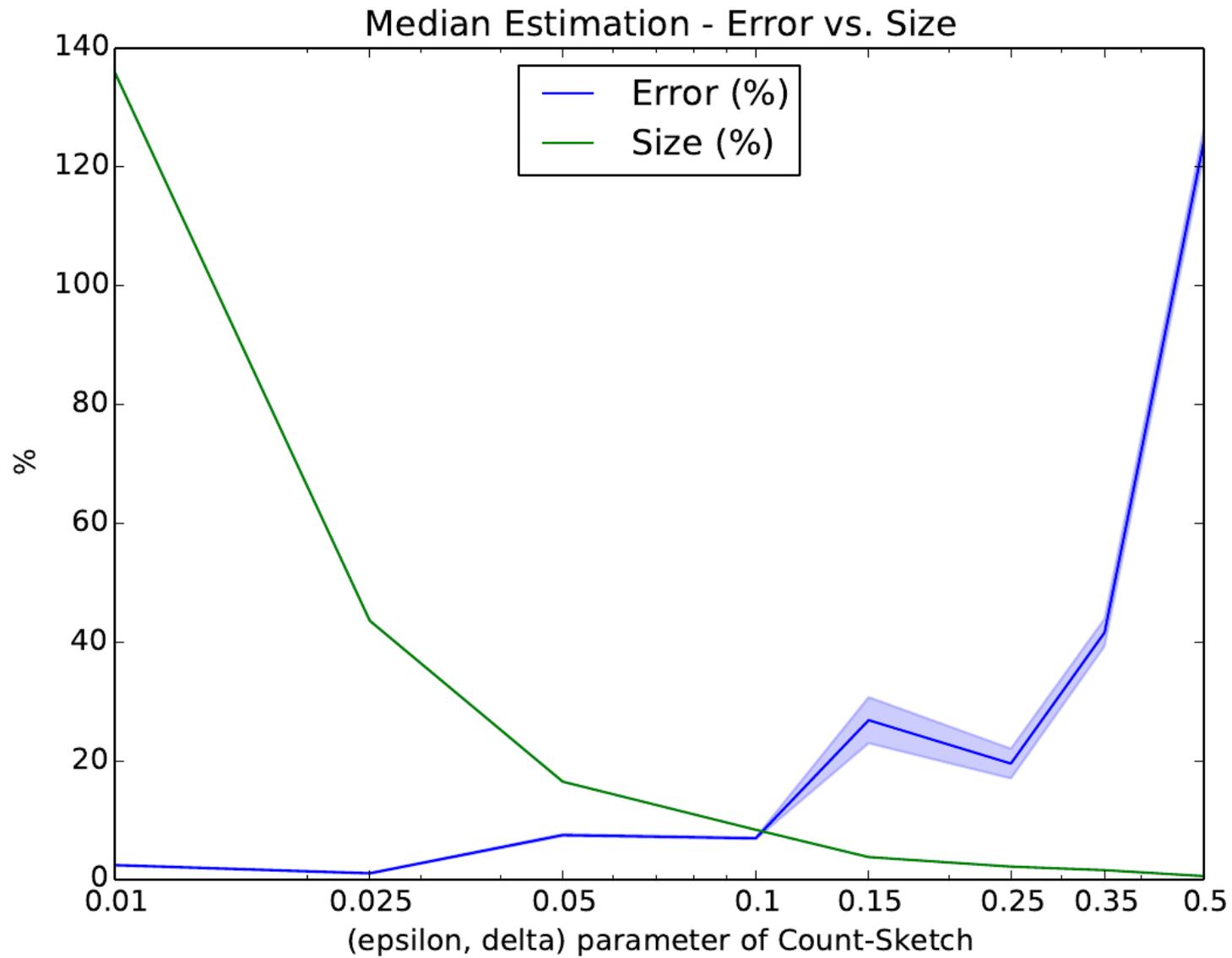
Python implementation (*petlib*)

1 ms to encrypt a sketch (of size 165) for each HSDir and

1.5 sec to aggregate 1200 sketches

Median Estimation - Quality vs. Protection





Collaborative Threat Mitigation

Collaborative Anomaly Detection

Anomaly detection is hard

Suspicious activities deliberately mimic normal behavior

But, malevolent actors often use same resources

Wouldn't it be better if organizations collaborated?

It's a w

“It is the policy of the United States Government to increase the volume, timelines, and quality of cyber threat information shared with U.S. private sector entities so that these entities may better protect and defend themselves against cyber attacks.”

**Barack Obama
2013 State of the Union Address**

Problems with Collaborations

Trust

Will others leak my data?

Legal Liability

Will I be sued for sharing customer data?

Will others find me negligent?

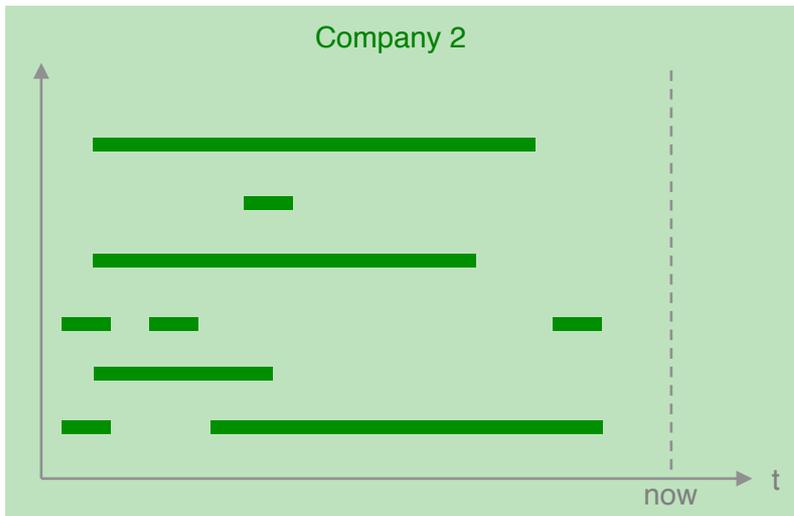
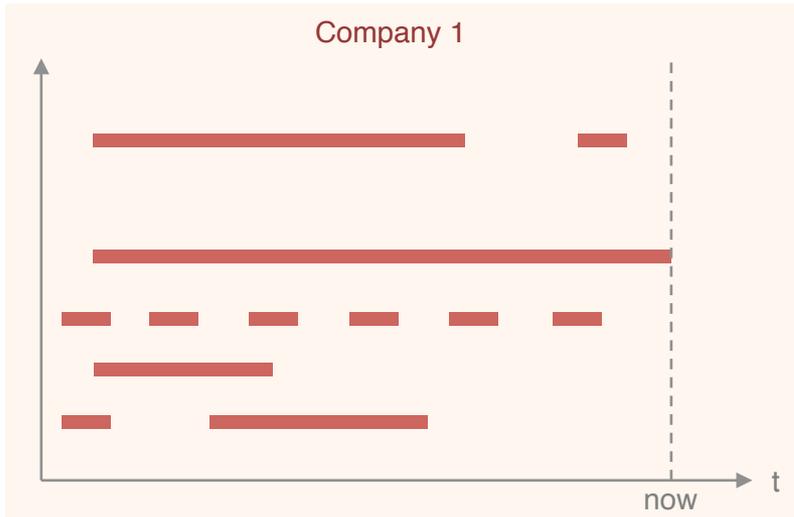
Competitive concerns

Will my competitors outperform me?

Shared data quality

Will data be reliable?

Solution Intuition [FDB15]



Sharing
Information
w/ Privacy

**Securely assess
the benefits of
sharing**

**Securely assess
the risks of
sharing**

Better Analytics

1. Estimate Benefits

What are good **indicators** of the fact that sharing will be beneficial?

- Many attackers in common?
- Many similar attacks in common?
- Many correlated attacks in common?

2. Select Partners

How do I **choose** who to collaborate with?

- Collaborate with the top-k?
- Collaborate if benefit above threshold?
- Hybrid?

3. Merge

Once we partnered up, **what** do we share?

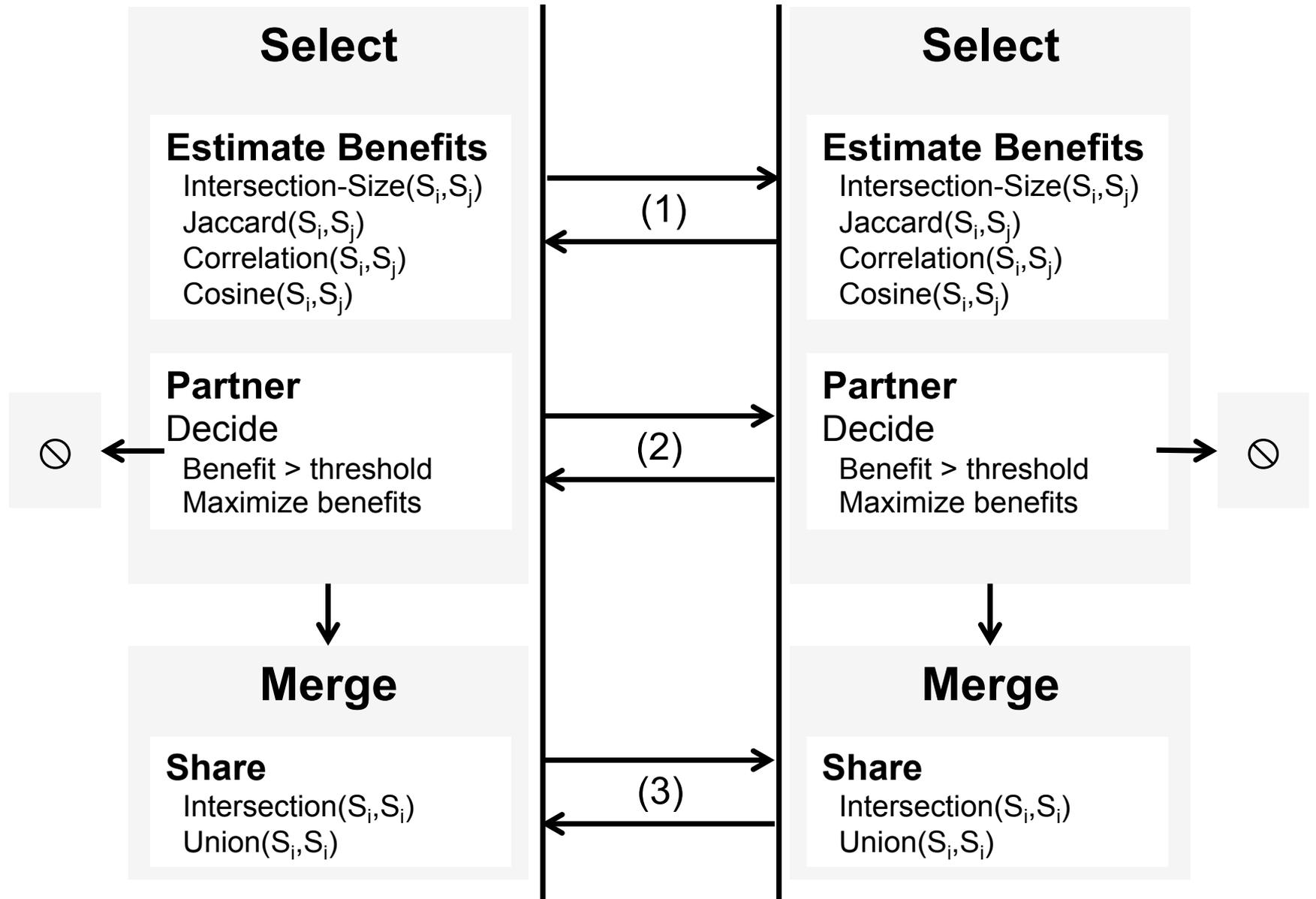
- Everything?
- Just what we have in common?
- Just information about attacks or also metadata?

System Model

Network of n entities $\{V_i\}$ (for $i=1, \dots, n$)

Each V_i holds a dataset S_i of suspicious events

E.g., events in the form $\langle IP, time, port \rangle$ as observed by a firewall or an IDS



Privacy-preserving benefit estimation

| Metric | Operation | Private Protocol |
|--------------------------|---|---|
| Intersection-Size | $ S_i \cap S_j $ | Private Set Intersection Cardinality (PSI-CA) |
| Jaccard | $\frac{ S_i \cap S_j }{ S_i \cup S_j }$ | Private Jaccard Similarity (PJS) |
| Pearson | $\sum_{l=1}^N \frac{(s_{il} - \mu_i)(s_{jl} - \mu_j)}{N\sigma_i\sigma_j}$ | Garbled Circuits (2PC) |
| Cosine | $\frac{\vec{S}_i \vec{S}_j}{\ \vec{S}_i\ \ \vec{S}_j\ }$ | Private Cosine Similarity (PCS) |

Privacy-preserving data sharing

| Metric | Operation | Private Protocol |
|--|---|--|
| Intersection | $ S_i \cap S_j $ | Private Set Intersection (PSI) |
| Intersection with Associated Data | $\{\langle \text{IP}, \text{time}, \text{port} \rangle \mid \text{IP} \in S_i \cap S_j\}$ | Private Set Intersection w/ Data Transfer (PSI-DT) |
| Union with Associated Data | $\{\langle \text{IP}, \text{time}, \text{port} \rangle \mid \text{IP} \in S_i \cup S_j\}$ | - |

The Road Ahead...

This slide is intentionally left blank