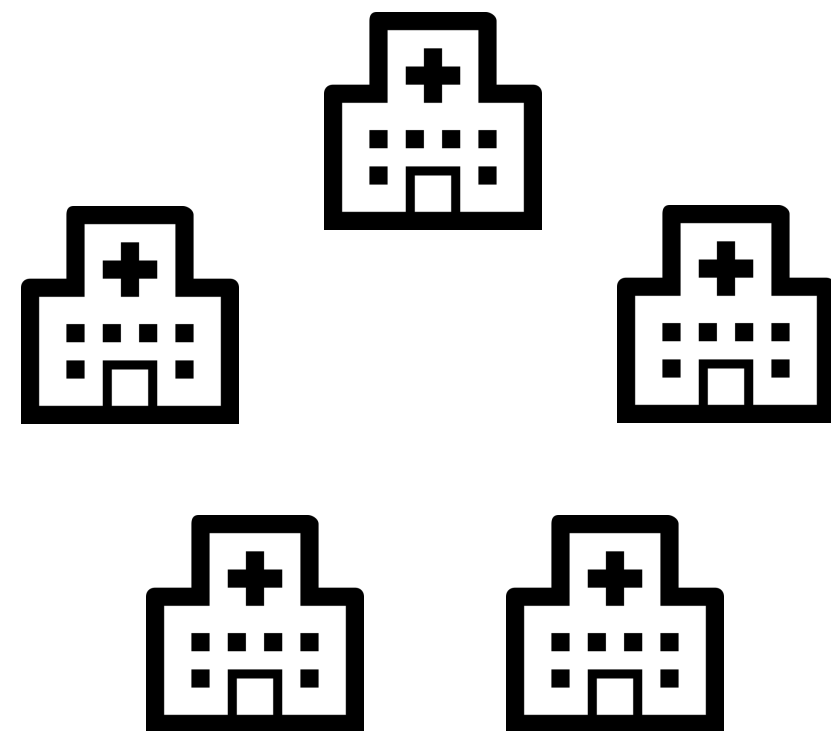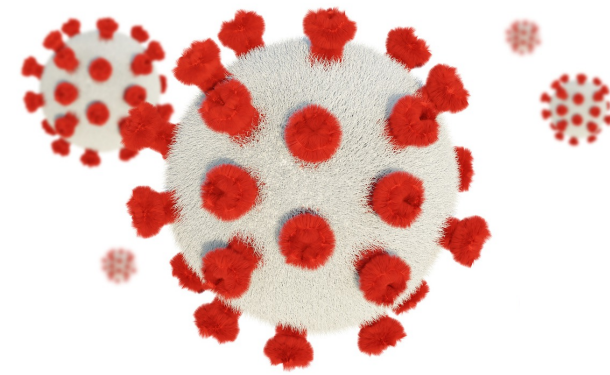# SECRECY: SECURE COLLABORATIVE ANALYTICS IN UNTRUSTED CLOUDS

John Liagouris, Vasiliki Kalavri, Muhammad Faisal, Mayank Varia

Boston University
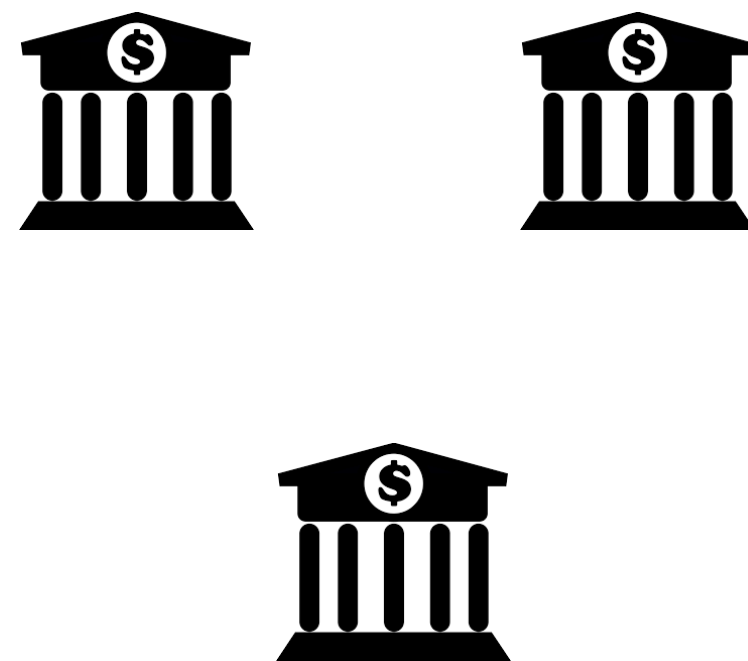
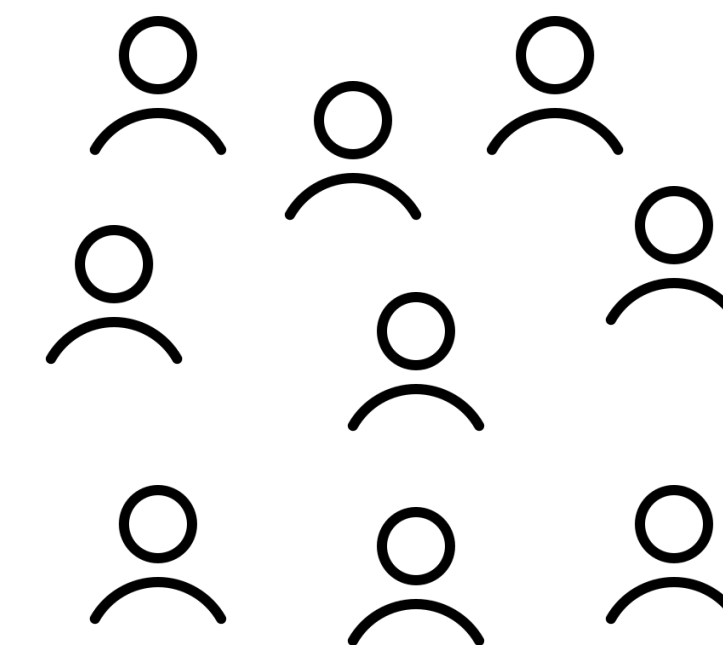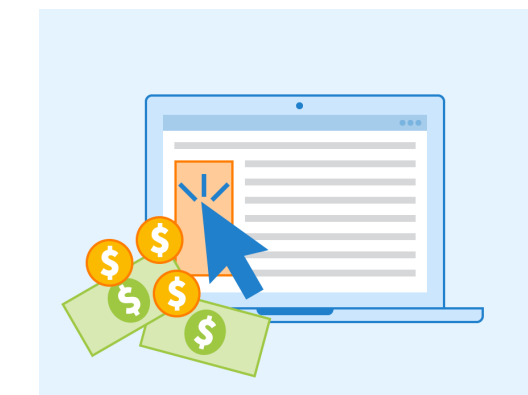# MOTIVATION: SECURE COLLABORATIVE ANALYTICS

*Medical
Studies*



*Market
Analyses*



*Privacy-preserving
advertising*



*Healthcare providers*

*Credit score agencies*

*Web users*

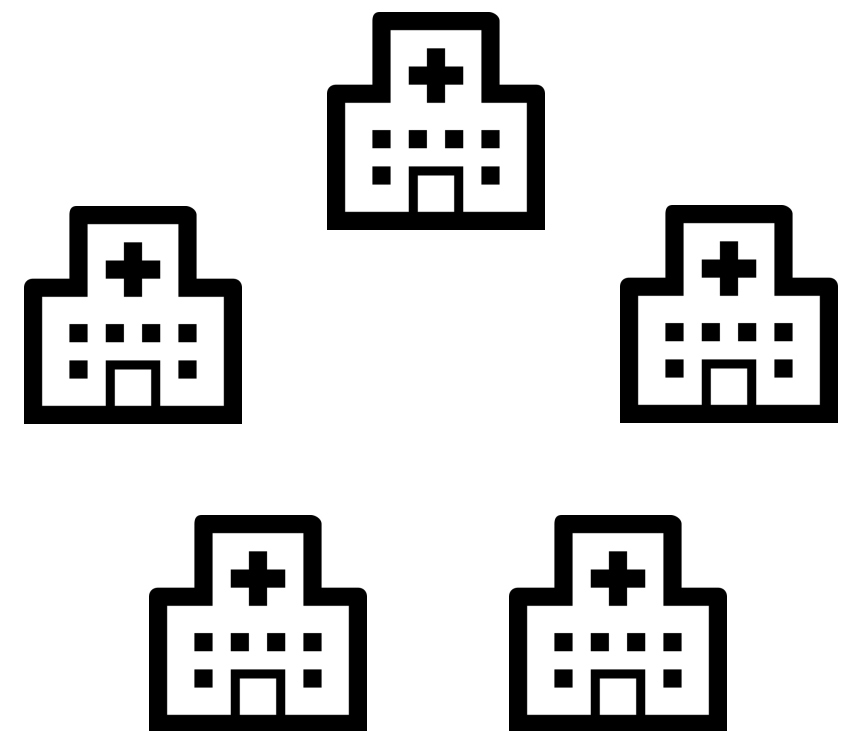# MOTIVATION: SECURE COLLABORATIVE ANALYTICS
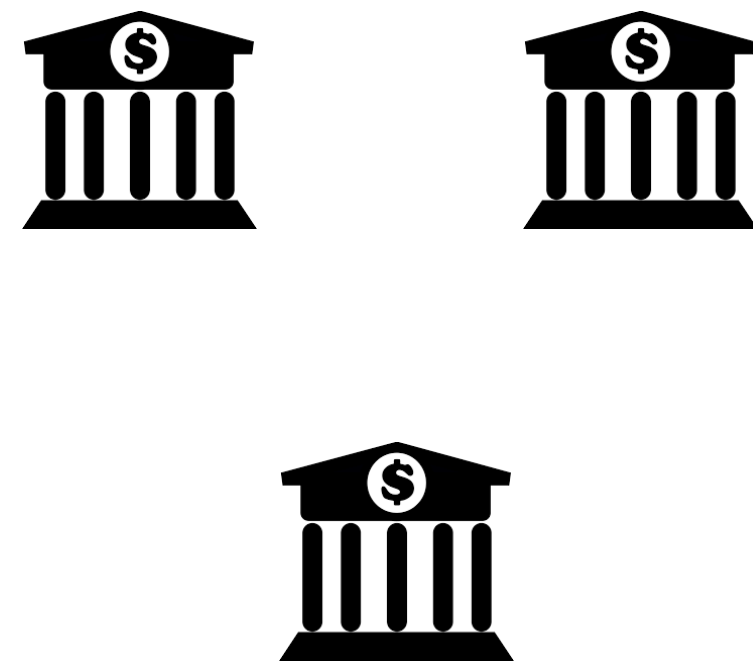
*Medical Studies*

*Market Analyses*

*Privacy-preserving advertising*

*Healthcare providers*

*Credit score agencies*

*Web users*

**Requirements:**

- **No information leakage to untrusted entities**

- **No reliance on trusted resources**

- **Relational analytics**

- **Practical performance**

# APPROACHES TO SECURE COLLABORATIVE ANALYTICS

*Fully Homomorphic Encryption (FHE)*

*Secure Multi-Party Computation (MPC)*

*Trusted Execution Environments (TEEs)*



Computation on encrypted data

Collective computation on encoded data

Computation on plaintext data inside the TEE (e.g. Intel's SGX)

*Security via homomorphic encryption (very high computational cost)*

*Security via decentralized trust (high communication cost)*

*Security via physically protected HW (prone to side-channel attacks)*

# CHALLENGE: HOW TO REDUCE THE MPC COST?

*"Running the query entirely under MPC […] **fails to scale beyond 3,000 total records**…"*

*"Computing a function f on millions of client inputs […] could potentially take an **astronomical amount of time** in a full MPC."*

*"The primary source of the slowdown arises from their join operators that have **hundreds of input tuples**…"*



Aggregation — Runtime [sec] vs Total input records [$\log_{10}$]
- Insecure (Spark)
- Secure (SM)
- Secure (Obliv-C)

Join — Runtime [sec] vs Total input records [$\log_{10}$]
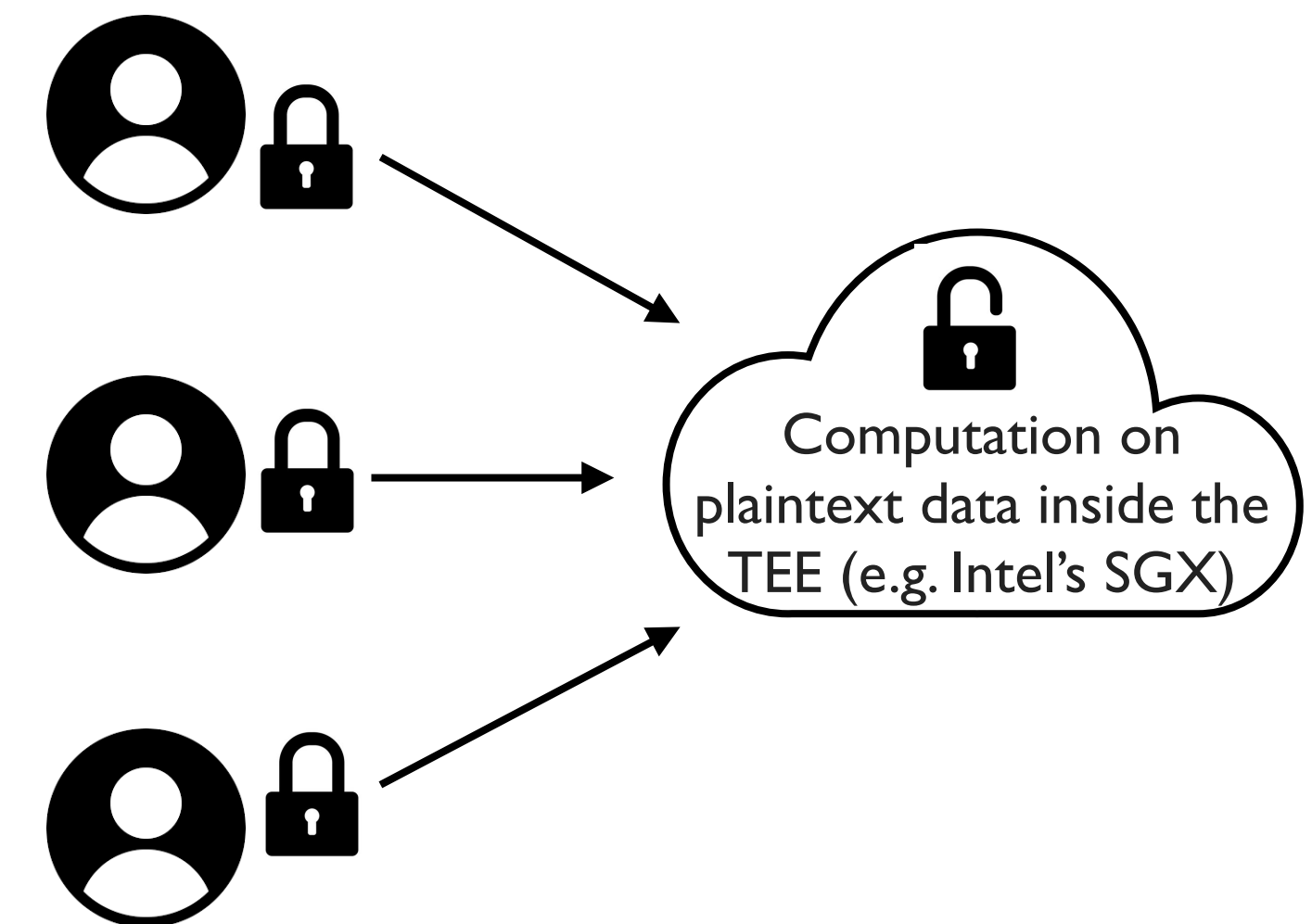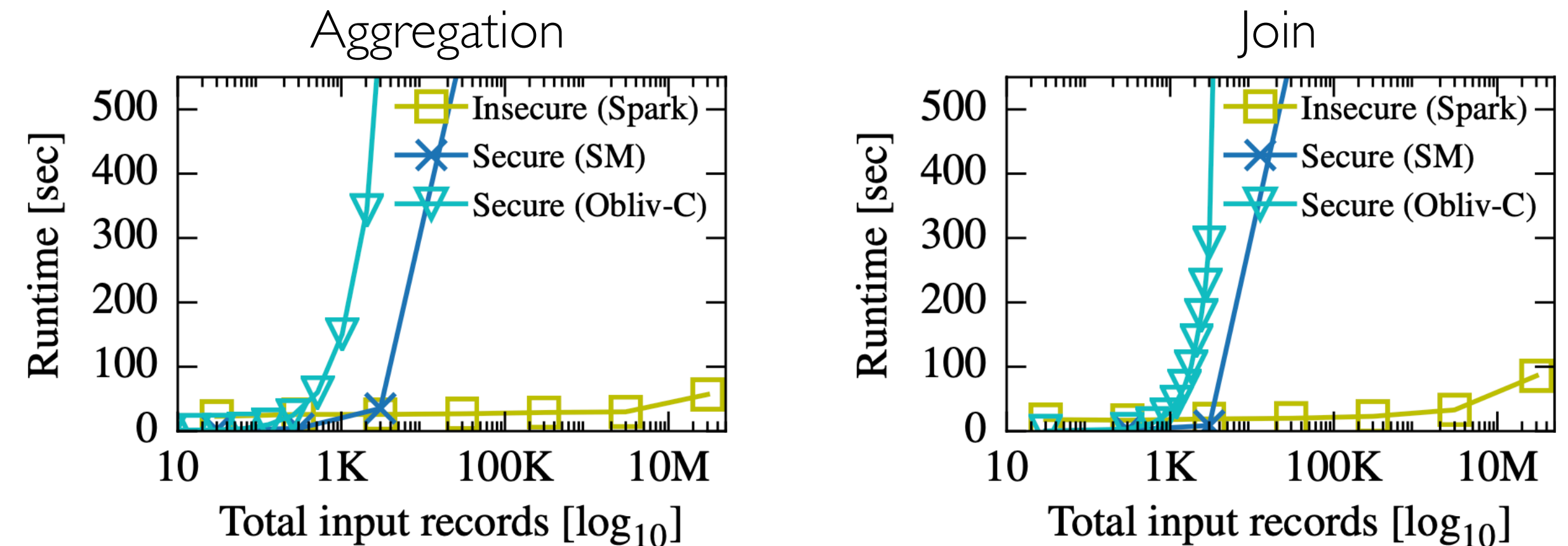- Insecure (Spark)
- Secure (SM)
- Secure (Obliv-C)

| Plaintext | Secure | Slowdown |
|---|---|---|
| 158 | 253,894 | 1,609X |
| 165 | 159,145 | 967X |
| 193 | 8,195,317 | 43,337X |

[1] N. Volgushev, M. Schwarzkopf, B. Getchell, M. Varia, A. Lapets, and A. Bestavros. *Conclave: secure multi-party computation on big data*. EuroSys, 2019.

[2] J. Bater, G. Elliott, C. Eggen, S. Goel, A. N. Kho, and J. Rogers. *SMCQL: secure querying for federated databases*. PVLDB, 10(6):673–684, 2017.

[3] H. Corrigan-Gibbs and D. Boneh. *Prio: Private, Robust, and Scalable Computation of Aggregate Statistics*, NSDI, 2017.

# PRIOR WORK ON RELATIONAL MPC



Peer-to-peer MPC

Data owners act as computing parties using trusted resources

- Data owners may not have domain expertise or private infrastructure

- MPC does not scale well with the number of data owners

# OUR FOCUS: OPTIMIZE MPC IN THE CLOUD

Peer-to-peer MPC

Outsourced MPC

Data owners act as computing parties using trusted resources

Data owners outsource secret shares of their data to untrusted third parties

- Data owners may not have domain expertise or private infrastructure

- MPC does not scale well with the number of data owners

+ Data owners can use untrusted cloud resources on demand

+ A small number of third parties can support a large number of data owners

# SECRECY AS A SERVICE



Secrecy

Data analysts

Data owners

# SECRECY AS A SERVICE



Data analysts

① **Submit query**

Secrecy

Data owners

# SECRECY AS A SERVICE

# SECRECY AS A SERVICE



Arithmetic sharing:  $s = s_1 + s_2 + s_3 \pmod{2^k}$

(for k-bit integers)

Boolean sharing:  $s = s_1 \oplus s_2 \oplus s_3$

(for k-bit strings)

# SECRECY AS A SERVICE

# SECRECY AS A SERVICE



Secrecy

Secrecy computing party

Secrecy computing party

Secrecy computing party

Data analysts

⑤ **Send result shares**

Data owners

13

# SECRECY AS A SERVICE



**Semi-honest model**

(parties are "honest but curious")

**Honest majority**

(can tolerate one compromised party)

Data analysts

Data owners

T. Araki, J. Furukawa, Y. Lindell, A. Nof, and K. Ohara. *High-Throughput Semi-Honest Secure Three-Party Computation with an Honest Majority*. CCS, 2016.

# FROM SECURE MPC PRIMITIVES TO RELATIONAL ANALYTICS



| SEMI-JOIN | ORDER-BY | COMPOSE | AGGREGATION |
| SELECT | DISTINCT | GROUP-BY | JOIN |

**Secure SQL operators**

| EQUALITY | INEQUALITY | CMP-SWAP | CONVERSION | ... |

**Complex operations**

| ADD (+) | MUL (×) | XOR (⊕) | AND (∧) |

**MPC primitives**

Arithmetic sharing: $s = s_1 + s_2 + s_3 \pmod{2^k}$          Boolean sharing: $s = s_1 \oplus s_2 \oplus s_3$

# SECRECY's CORE CONTRIBUTIONS

1. **Relational MPC primitives**

   - Amortize network I/O

   - Make secret-sharing competitive in WAN

# EXAMPLE: MESSAGE BATCHING IN SECRECY

$$k \; bits$$

| ID | Name | Timestamp |
|----|------|-----------|
| … | … | $t_1$ |
| … | … | $t_2$ |
| … | … | $t_3$ |
| … | … | $t_4$ |
| … | … | … |

$n$

*"Select all records with timestamp $t > \mathbb{T}$"*

# EXAMPLE: MESSAGE BATCHING IN SECRECY

$k$ bits

| ID | Name | Timestamp |
|----|------|-----------|
| ... | ... | $t_1$ |
| ... | ... | $t_2$ |
| ... | ... | $t_3$ |
| ... | ... | $t_4$ |
| ... | ... | ... |

$n$

$O(\log k)$

✉✉✉

+

$O(\log k)$

✉✉✉

+

$O(\log k)$

✉✉✉

+

...

*"Select all records with timestamp $t > \mathbb{T}$"*



Secret-sharing protocols exchange
many small messages between parties

Each inequality requires $O(\log k)$ communication rounds under MPC

# EXAMPLE: MESSAGE BATCHING IN SECRECY

$$k \; bits$$

| ID | Name | Timestamp |
|----|------|-----------|
| … | … | $t_1$ |
| … | … | $t_2$ |
| … | … | $t_3$ |
| … | … | $t_4$ |
| … | … | … |

$n$

*"Select all records with timestamp $t > \mathbb{T}$"*



Secrecy requires $\log k + 1$ communication rounds for the entire data table

(independent of the number of records)

# EFFECT OF MESSAGE BATCHING (LAN)



- Eager: Message batching disabled (one network I/O per row)

- Batched: Message batching enabled

Lower is better

# SECRECY's CORE CONTRIBUTIONS

1. **Relational MPC primitives**

   - Amortize network I/O

   - Make secret-sharing competitive in WAN

2. **Analytical cost model for MPC**

# SECRECY's CORE CONTRIBUTIONS

1. **Relational MPC primitives**

   - Amortize network I/O

   - Make secret-sharing competitive in WAN

2. **Analytical cost model for MPC**

   - Operation cost

   - Synchronization cost

   - Composition cost

# EXAMPLE: OPERATOR DECOMPOSITION IN SECRECY

$\gamma$ : aggregation

$\bowtie$ : join

🔒 MPC

$\gamma_{M.med}\text{COUNT(*)}$

$\bowtie_{M.id=P.id}$

M

P

```
SELECT M.med, COUNT(*)
FROM Medication as M, Patients as P
WHERE M.id = P.id
GROUP-BY M.med
```

*"Count the number of patients per prescribed medication"*

$\gamma$ : aggregation

$\bowtie$ : join

🔒 MPC

$\gamma_{M.med}\text{COUNT}(*)$

$\bowtie_{M.id=P.id}$

M

P

```
SELECT M.med, COUNT(*)
FROM Medication as M, Patients as P
WHERE M.id = P.id
GROUP-BY M.med
```

*"Count the number of patients per prescribed medication"*

*Applying* `GROUP-BY` *after the join will require materializing the cartesian product* $M \times P$

* Assuming the group-by operator is based on an odd-even circuit

$\gamma$ : aggregation

$\bowtie$ : join

$\gamma_{M.med}$COUNT(*)

MPC

$\bowtie_{M.id=P.id}$

M

P

MPC

$\gamma_{M.med}$SUM(cnt)

$\bowtie_{M.id=P.id}$ + partial aggregation

M

P

$O(n^2 \log^2 n)$ *operations / messages*

$O(\log^2 n)$ *rounds*    $O(n^2)$ *space*

$O(n^2)$ *operations / messages*

$\sim 4 \times$ fewer rounds     $O(n)$ *space*

# EFFECT OF JOIN–AGGREGATION DECOMPOSITION (LAN)



* Secrecy servers deployed on AWS EC2 r5.xlarge instances (us-east-2)

# SECRECY's CORE CONTRIBUTIONS
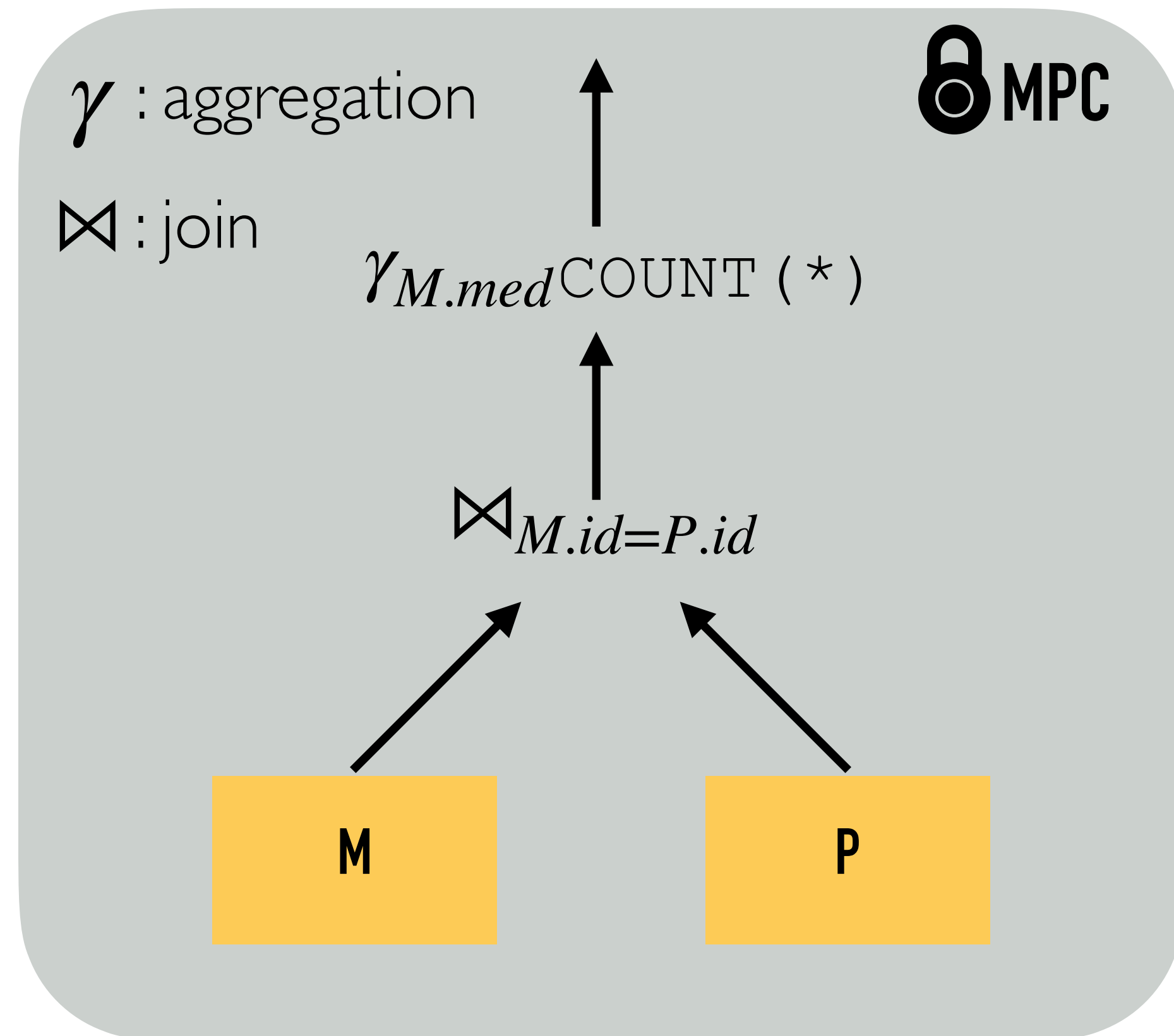
1. **Relational MPC primitives**

   - Amortize network I/O

   - Make secret-sharing competitive in WAN

2. **Analytical cost model for MPC**

   - Operation cost

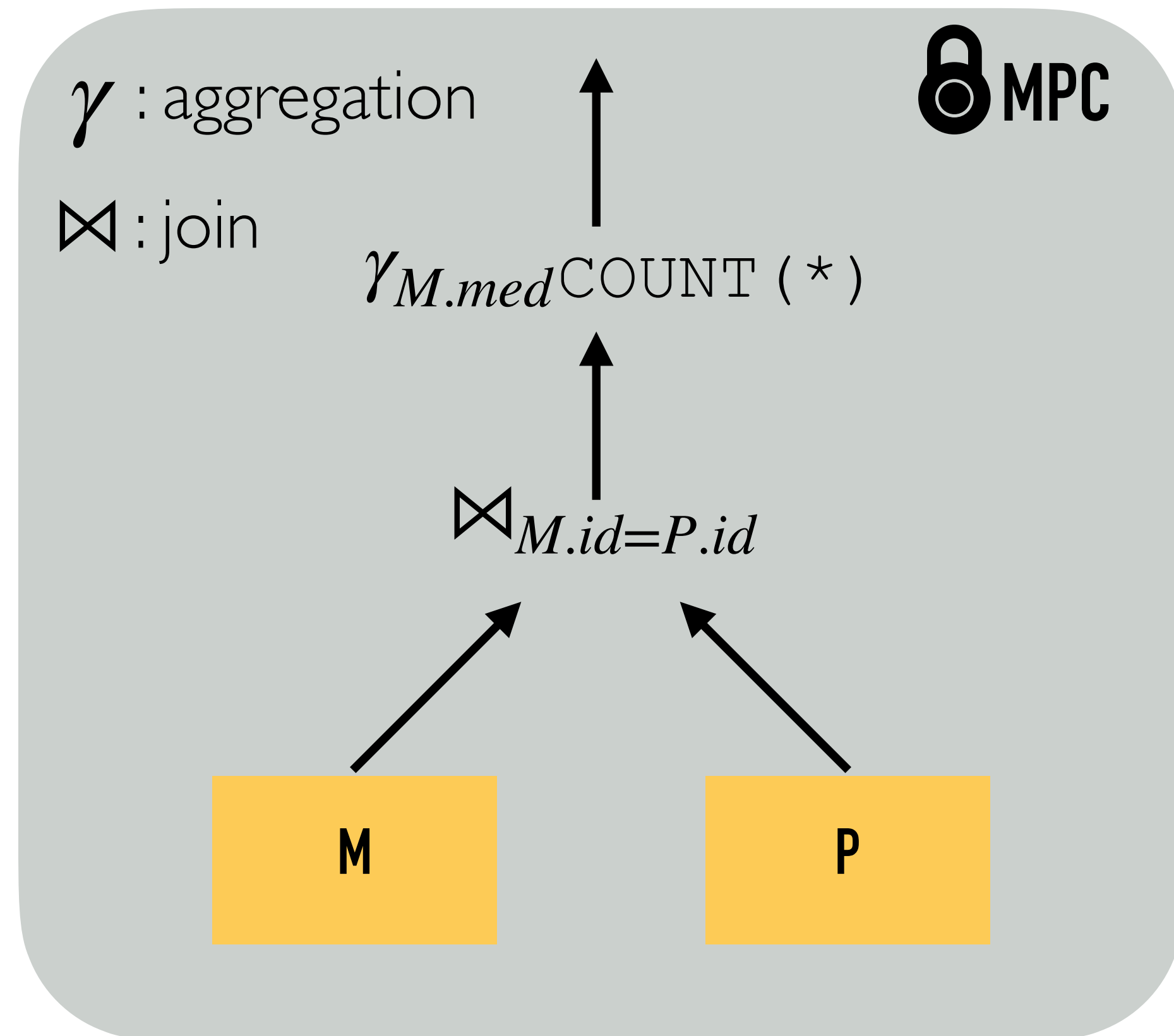   - Synchronization cost

   - Composition cost

3. **Vectorized MPC query processor**

   - Logical optimizations

   - System optimizations

   - Protocol-specific optimizations

# OPENING THE MPC BLACK BOXES



**Query Planner**

Secure

Plaintext

Input₁ Input₂

**Black-box MPC library e.g. EMP¹**

□ : data operator

**End-to-end secure MPC query engine**

Cross-layer optimizations

Input₁ Input₂

Supported optimizations:

- Logical (e.g. operator decomposition)

- Physical (e.g. message batching, operator fusion)

- Protocol-specific (e.g. dual sharing)

The Secrecy Framework

¹ X. Wang, A. J. Malozemoff, and J. Katz. *EMP-toolkit: Efficient MultiParty computation toolkit*, 2016. https://github.com/emp-toolkit

# EFFECT OF SECRECY OPTIMIZATIONS ON REAL QUERIES (MULTI-CLOUD)



Logical + System optimizations result in up to $2000 \times$ speedups

Legend: N... / Optimized

Bar values: Comorbidity 1.24x, RC. Diff 302.43x, Aspirin 91.11x, Pwd 1.46x, Credit 1.07x, TPC-H Q4 1.39x, TPC-H Q6 38.3x, TPC-H Q13 2088.17x

Y-axis: Time (s), $10^{-1}$ to $10^4$

\* Secrecy servers deployed in three clouds: AWS (Ohio), GCP (South Carolina), and Azure (Virginia)

\* Reported times are for 1000 rows per input table

\* Not optimized plans use message batching too (otherwise the cost of MPC is prohibitive)

# EFFECT OF SECRECY OPTIMIZATIONS ON REAL QUERIES (MULTI-CLOUD)



Protocol-specific optimizations result in up to 38 × speedups

Legend: Not optimized / Optimized

Y-axis: Time (s), $10^{-1}$ to $10^4$

Bars with labels:
- Comorbidity: 1.24x
- RC. Diff: 302.43x
- Aspirin: 91.11x
- Pwd: 1.46x
- Credit: 1.07x
- TPC-H Q4: 1.39x
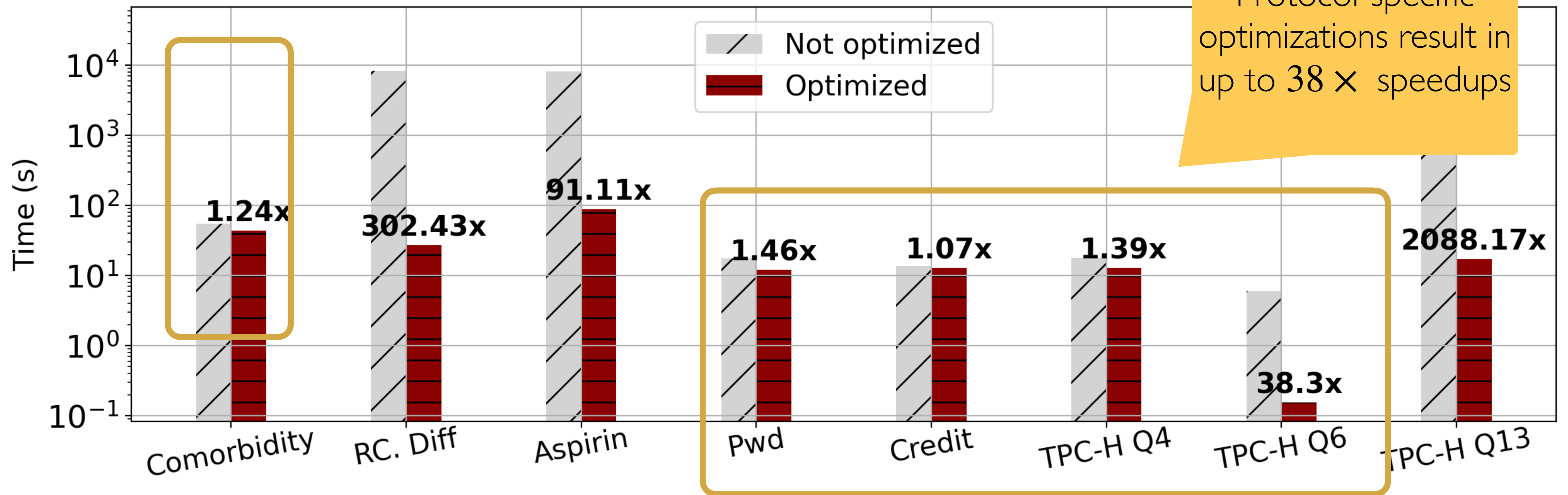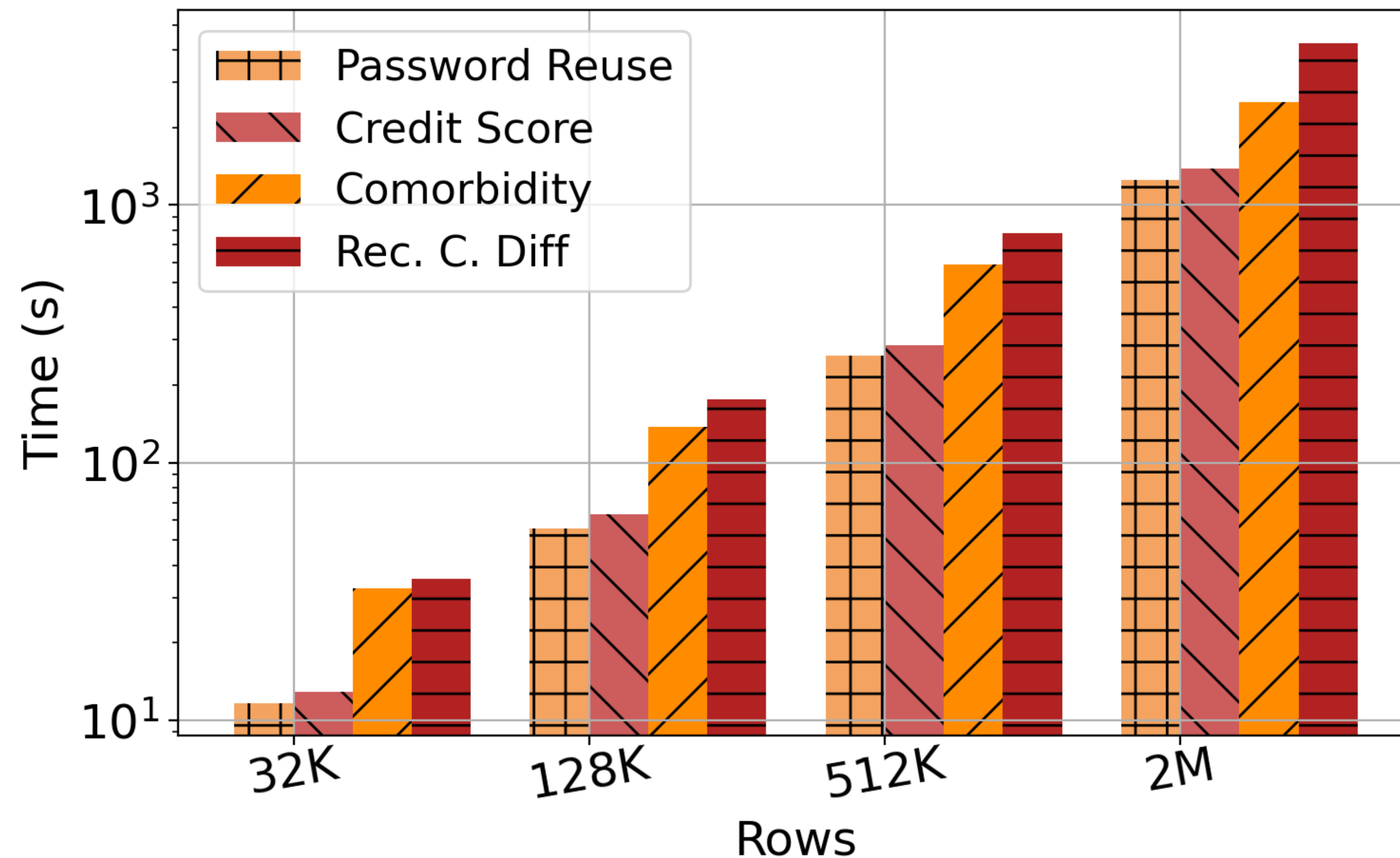- TPC-H Q6: 38.3x
- TPC-H Q13: 2088.17x

\* Secrecy servers deployed in three clouds: AWS (Ohio), GCP (South Carolina), and Azure (Virginia)

\* Reported times are for 1000 rows per input table

\* Not optimized plans use message batching too (otherwise the cost of MPC is prohibitive)

# SECRECY's SCALING BEHAVIOR (LAN)



Rec. C. Diff scales to 2 million rows in ~1.2h

```
WITH rcd AS (
    SELECT pid, time, row_no() OVER
    (PARTITION BY pid ORDER BY time)
    FROM diagnosis
    WHERE diag=cdiff)
 SELECT DISTINCT pid
 FROM rcd r1 JOIN rcd r2 ON r1.pid = r2.pid
 WHERE r2.time - r1.time >= 15 DAYS
 AND r2.time - r1.time <= 56 DAYS
 AND r2.row_no = r1.row_no + 1
```

*"Find the distinct ids of patients who have been diagnosed with cdiff and have two consecutive infections between 15 and 56 days apart"*
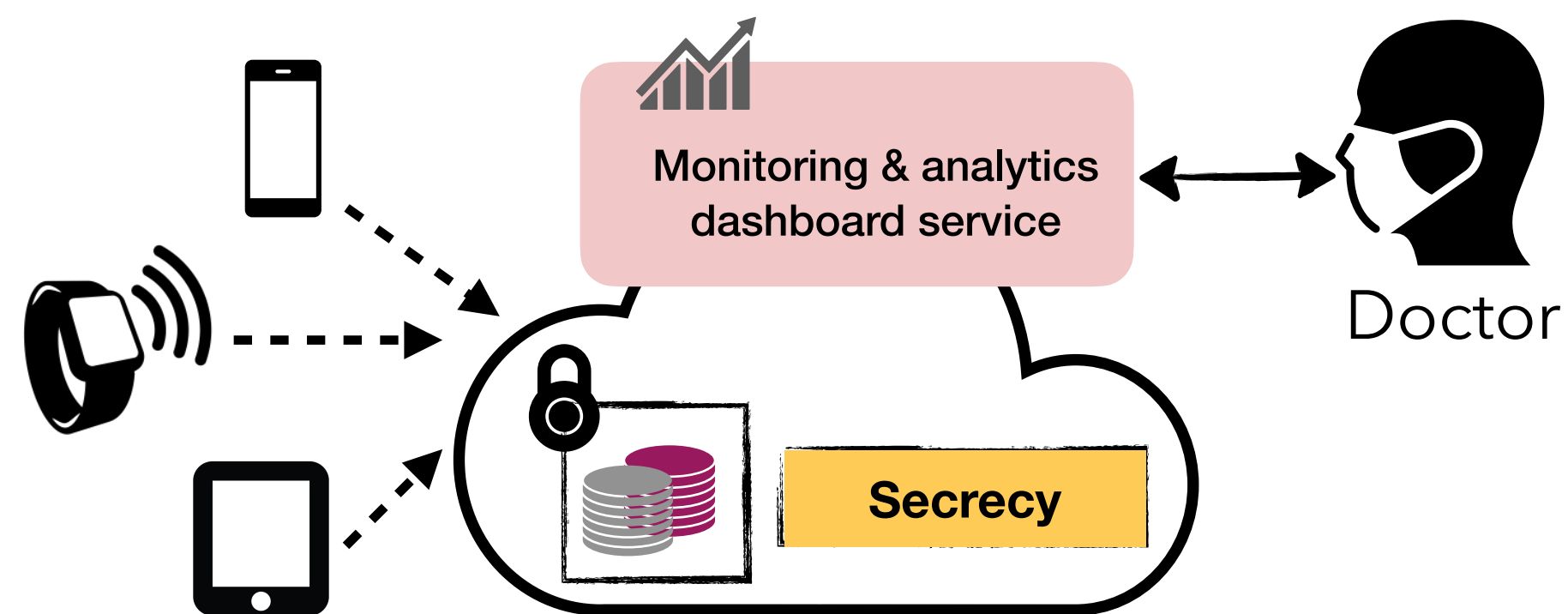
\* Secrecy servers deployed on AWS EC2 r5.xlarge instances (us-east-2)

\* Each server uses a single vCPU

# REAL–WORLD SECRECY USE CASES

## Secure digital health analytics[1]
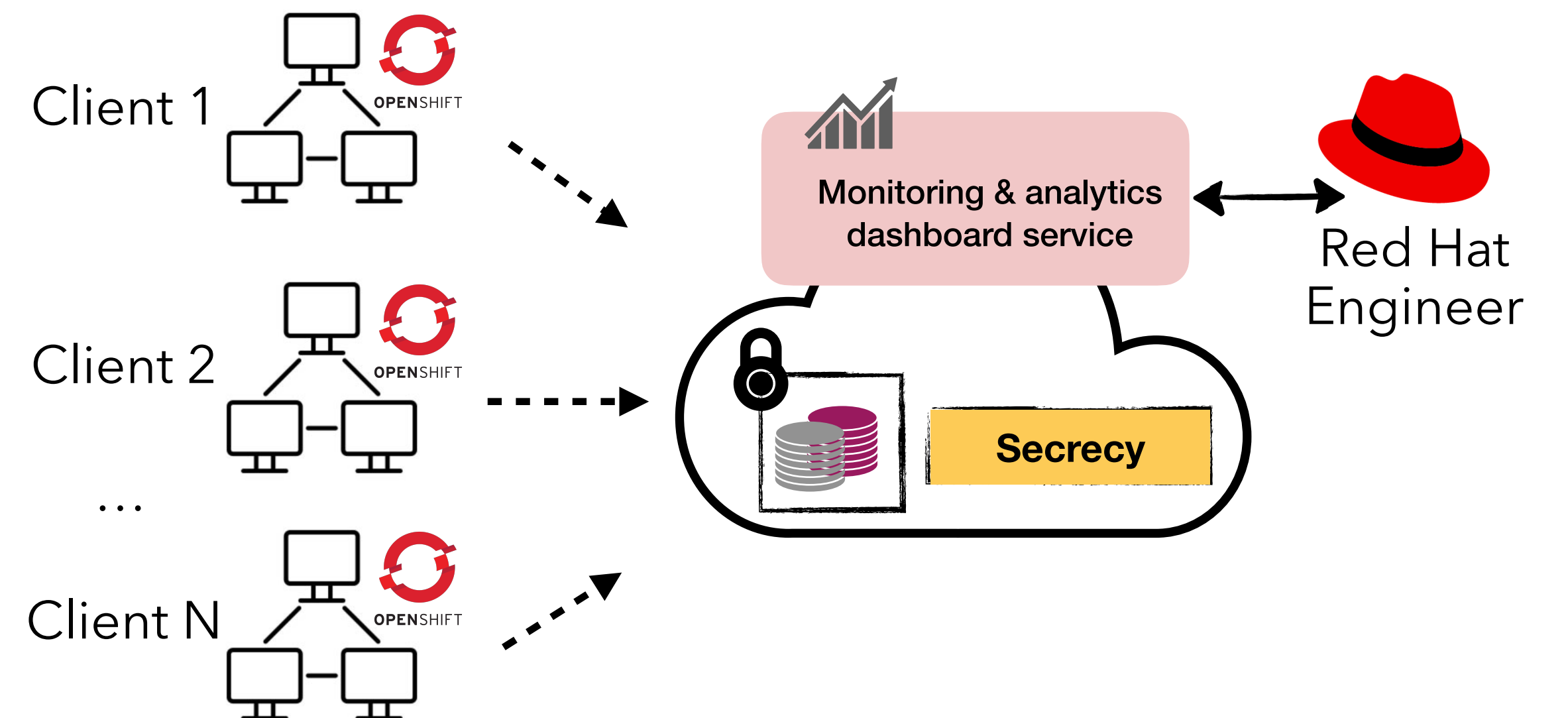
(BU Medical & Hariri Institute for Computing)
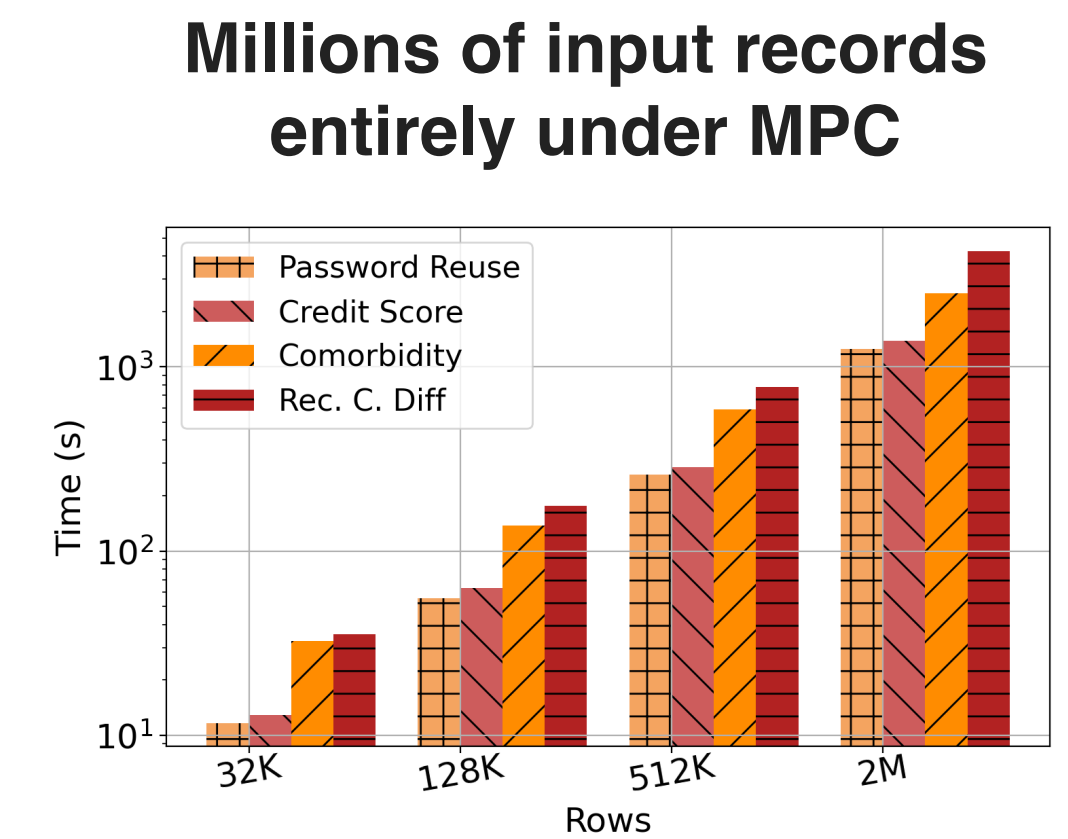


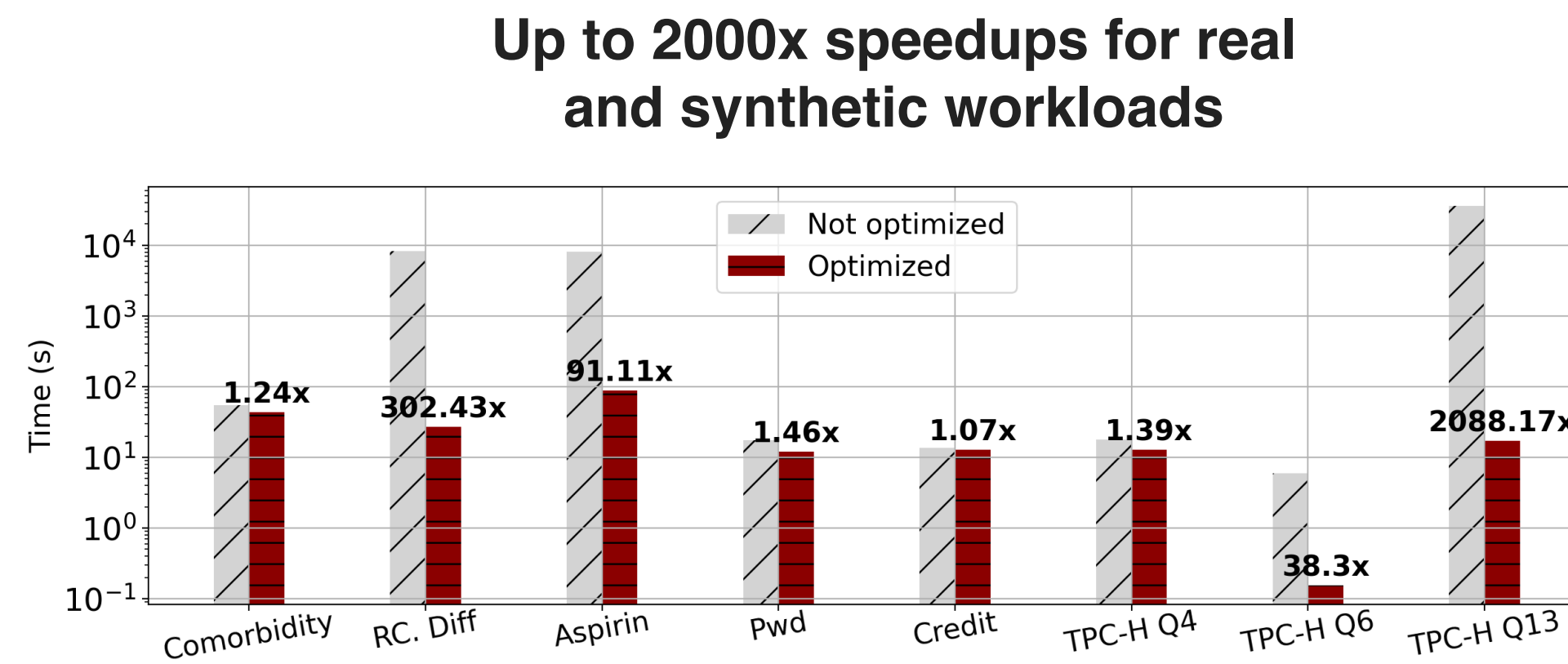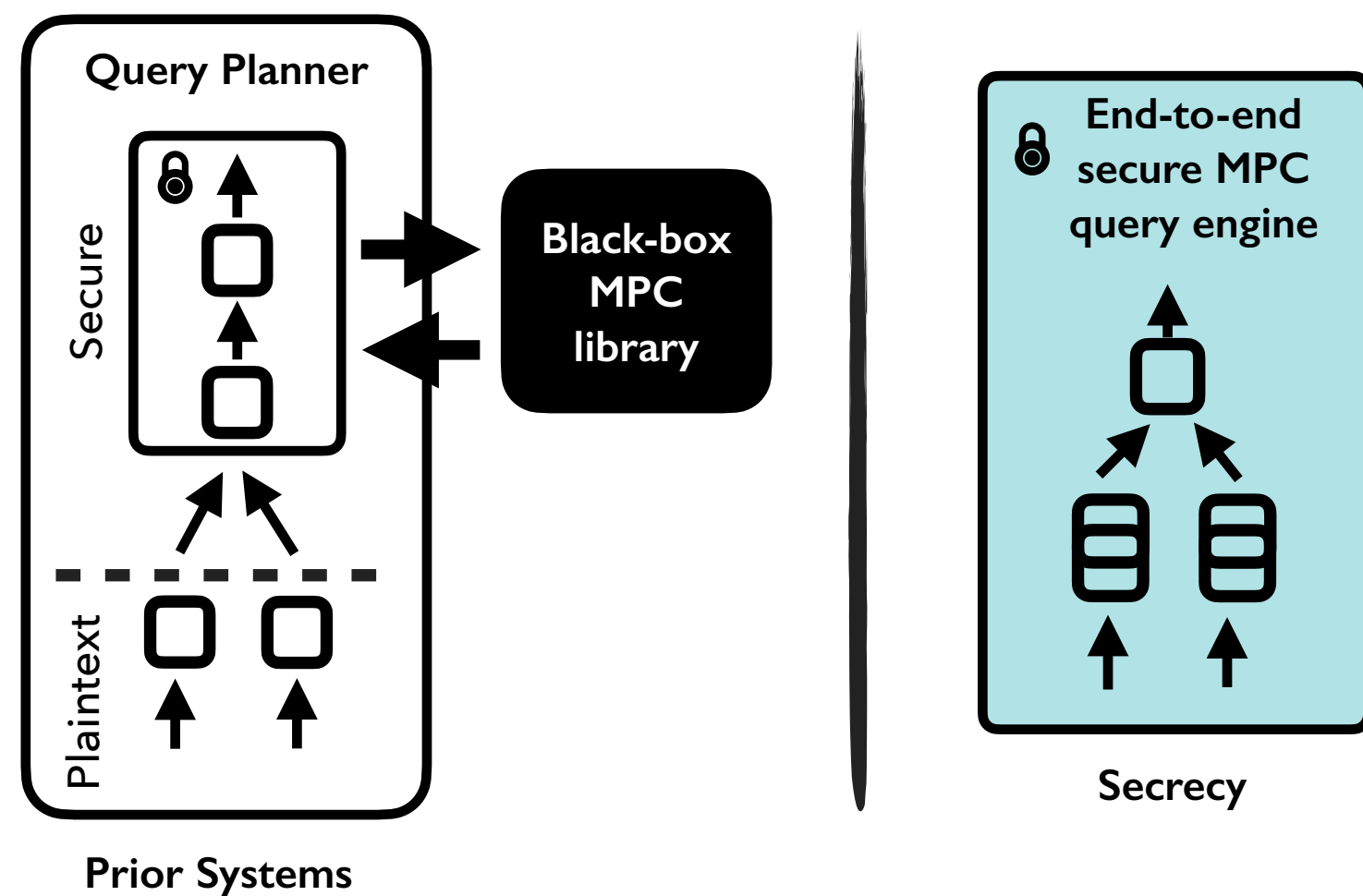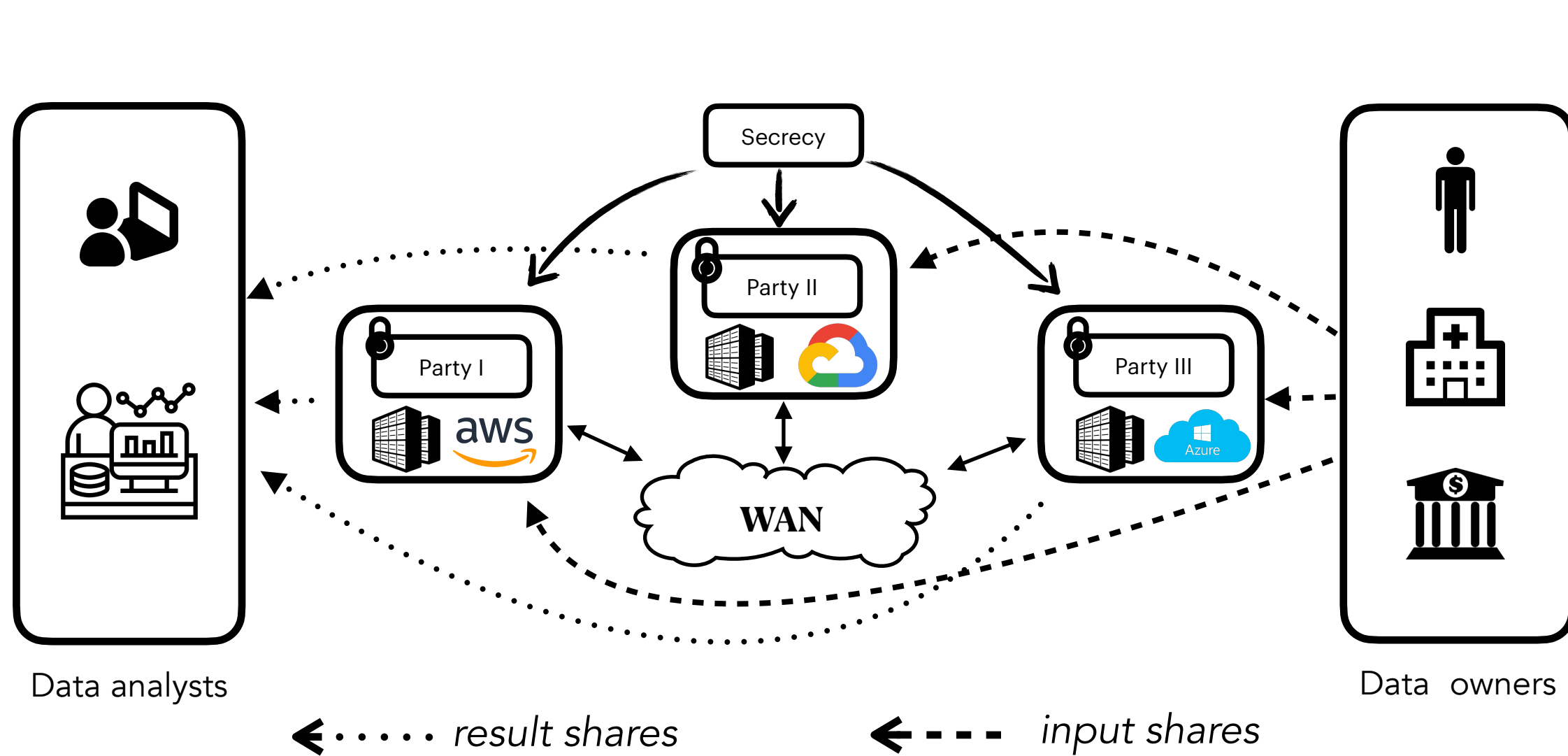## Secure cross-site analytics on OpenShift logs[2]

(BU Red Hat Collaboratory)

[1] https://www.bu.edu/hic/research/focused-research-programs/continuous-analysis-of-mobile-health-data-among-medically-vulnerable-populations/

[2] https://www.bu.edu/rhcollab/projects/security-privacy/secure-cross-site-analytics-on-openshift-logs/

# SECRECY SUMMARY



Data analysts

Secrecy
Party I
Party II
Party III
WAN
aws
Azure

Data owners

← · · · · result shares     ← - - - input shares

**Decoupling data owners from computing parties**

**No information leakage**

Outsourced MPC

End-to-end data protection

High expressivity

**No reliance on trusted execution environments**

General-purpose hardware

Efficient query execution

**Cost-based optimization**

**General and composable operators**

Query Planner
Secure
Plaintext
Black-box MPC library
**Prior Systems**

End-to-end secure MPC query engine
**Secrecy**

## Up to 2000x speedups for real and synthetic workloads



Legend: Not optimized · Optimized

1.24x  302.43x  91.11x  1.46x  1.07x  1.39x  38.3x  2088.17x

Comorbidity  RC. Diff  Aspirin  Pwd  Credit  TPC-H Q4  TPC-H Q6  TPC-H Q13

Time (s)

## Millions of input records entirely under MPC



Legend: Password Reuse · Credit Score · Comorbidity · Rec. C. Diff

32K  128K  512K  2M

Rows

Time (s)

Source code: https://github.com/CASP-Systems-BU/Secrecy

CASP Systems