

## Query Processing in Private Data Outsourcing Using Anonymization

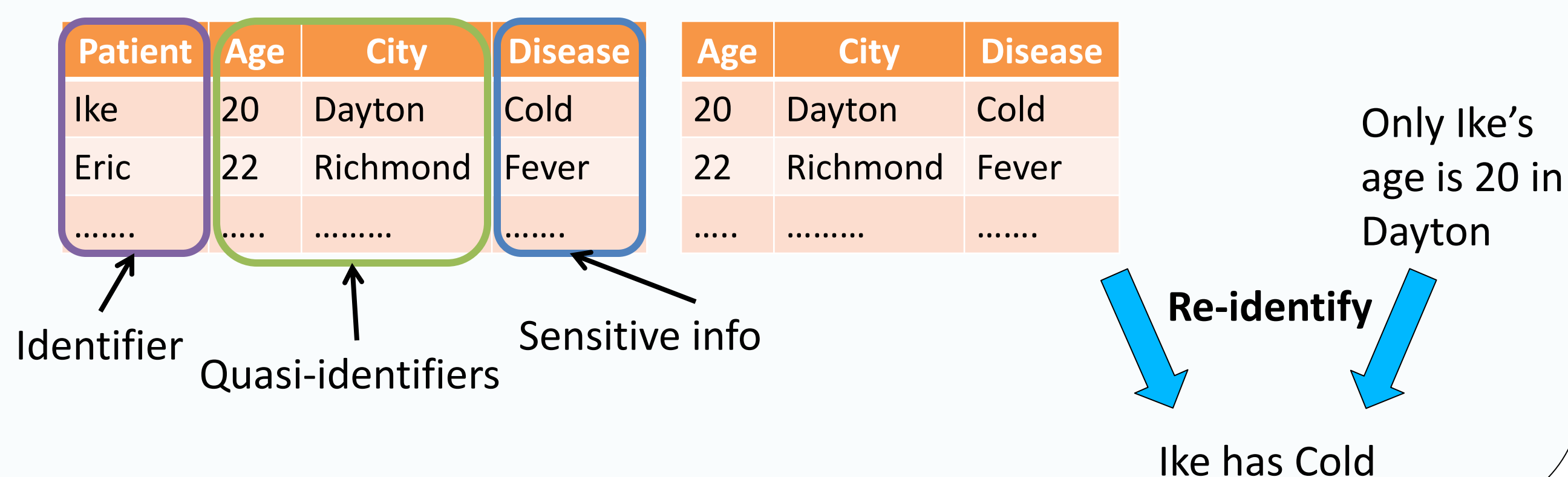
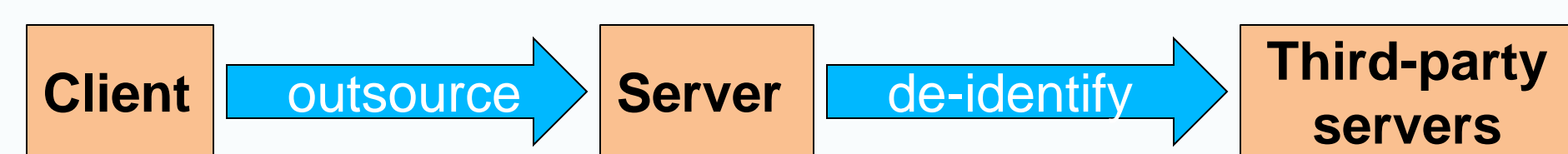
Ahmet Erhan Nergiz  
anergiz@cs.purdue.edu

Chris Clifton  
clifton@cs.purdue.edu

### Motivation

Data outsourcing benefits:

- Cheaper (affordable by small companies)
- High availability (24/7 access guarantee)
- Service quality (having experts on the field)



### Our Contribution

Encrypting whole database is not the answer

Value-added services: Data analysis, Data cleaning

"Sell" the data (or analysis of the data): Make the service free to the client!

Why not use existent work on data publishing (anatomy based on *l*-diversity)?

Patient	Age	City	Disease	Age	City	GID	GID	City
Ike	41	Dayton	Cold	41	Dayton	1	1	Cold
Eric	22	Richmond	Fever	22	Richmond	1	1	Fever
Olga	30	Lafayette	Flu	30	Lafayette	2	2	Flu
Kelly	35	Lafayette	Cough	35	Lafayette	2	2	Cough

Anonymize

### Related Work

Bucketization (Hacigumus et al.):

Patient	Age	City	Disease	Patient	Age	City	Disease	etuple
Ike	41	Dayton	Cold	23	6	3	2	E(tuple1)
Eric	22	Richmond	Fever	46	4	9	9	E(tuple2)
.....	.....	.....	.....	.....	.....	.....	.....	E(tuple3)

### How We Do It

Patient	Age	City	GID	SEQ	H(SEQ)	GID	Disease
Ike	41	Dayton	1	1	H <sub>k2</sub> (1)	1	Cold
Eric	22	Richmond	1	2	H <sub>k2</sub> (2)	1	Fever
Olga	30	Lafayette	2	3	H <sub>k2</sub> (3)	2	Flu
Kelly	35	Lafayette	2	4	H <sub>k2</sub> (4)	2	Cough
Faye	24	Richmond	3	5	H <sub>k2</sub> (5)	3	Flu
Mike	47	Richmond	3	6	H <sub>k2</sub> (6)	3	Fever
Jason	45	Lafayette	4	7	H <sub>k2</sub> (7)	4	Cough
Max	31	Lafayette	4	8	H <sub>k2</sub> (8)	4	Flu

Patient<sub>IT</sub> = Identifying table      Patient<sub>ST</sub> = Sensitive table

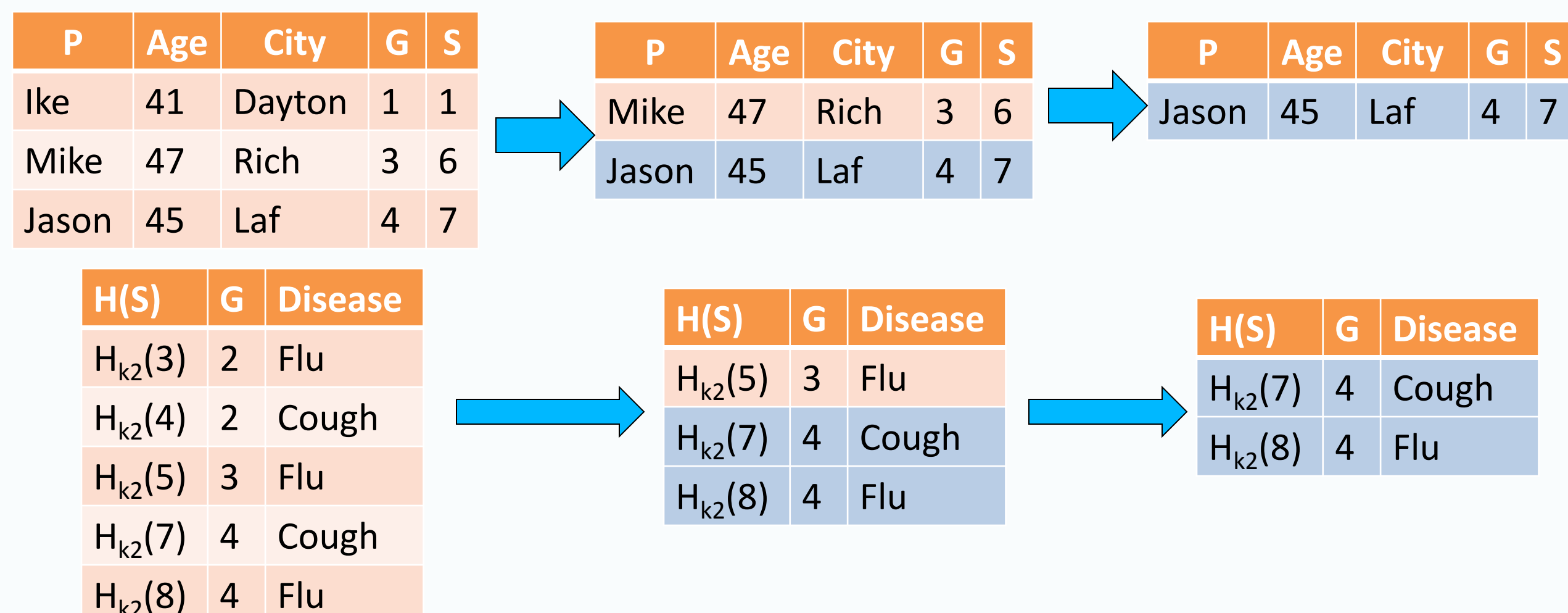
**Proposal:** Server does partial query processing on the anatomized database, client opens the encrypted link and joins IT and ST to complete the query

### Selection

Idea: Convert selection into subqueries that can be processed on the IT and ST

SELECT \* FROM Patient

WHERE Age > 40 AND (Disease = Flu OR Disease = Cough) AND (Disease = Cough OR Address = Lafayette)



Client joins the two tables and get:

Patient	Age	Address	G	S	H(S)	G	Disease
Jason	45	Lafayette	4	H <sub>k2</sub> (7)	H <sub>k2</sub> (7)	4	Cough

### Group By

SELECT AVG(Age) FROM Physician, Patient GROUP BY Gender, City

Doctor	Gender	G	S	H(S)	G	Patient	Patient	Age	City	G	S
Alice	Female	1	1	H <sub>k1</sub> (1)	1	Ike	Ike	41	Dayton	1	1
Carol	Female	1	2	H <sub>k1</sub> (2)	1	Eric	Eric	22	Richmond	1	2
Bob	Male	2	3	H <sub>k1</sub> (3)	2	Olga	Olga	30	Lafayette	2	3
Dave	Male	2	4	H <sub>k1</sub> (4)	2	Kelly	Kelly	35	Lafayette	2	4
Carol	Female	3	5	H <sub>k1</sub> (5)	3	Faye	Faye	24	Richmond	3	5
Alice	Female	3	6	H <sub>k1</sub> (6)	3	Mike	Mike	47	Richmond	3	6
Dave	Male	4	7	H <sub>k1</sub> (7)	4	Jason	Jason	45	Lafayette	4	7
Carol	Female	4	8	H <sub>k1</sub> (8)	4	Max	Max	31	Lafayette	4	8

Physician<sub>IT</sub>      Physician<sub>ST</sub>      Patient<sub>IT</sub>

Server processes the group-by partially and send the partial result to the client with the rest of the tuples:

Gender	City	AVG(Age)	COUNT(*)
Female	Dayton	41	1
Female	Richmond	31	3
Male	Lafayette	32.5	2

Client merges the two results to get the actual result

Gender	City	AVG(Age)
Female	Dayton	41
Female	Lafayette	31
Female	Richmond	31
Male	Lafayette	36.67

Client processes the group by query on the rest of the tuples by decrypting the link and joining the tables:

Gender	City	AVG(Age)	COUNT(*)
Male	Lafayette	45	1
Female	Lafayette	31	1

### Work In Progress

Implement it on top of a DBMS

Support data manipulation (insert, delete, update)

Challenges:

- Need incremental anatomization
- Previous works support only insert and delete w/o incremental anatomization
- Incremental anatomization needs to consider all inference channels existing between snapshots of the database

### Future Work

Integrate it to safe grouping for transactional datasets