

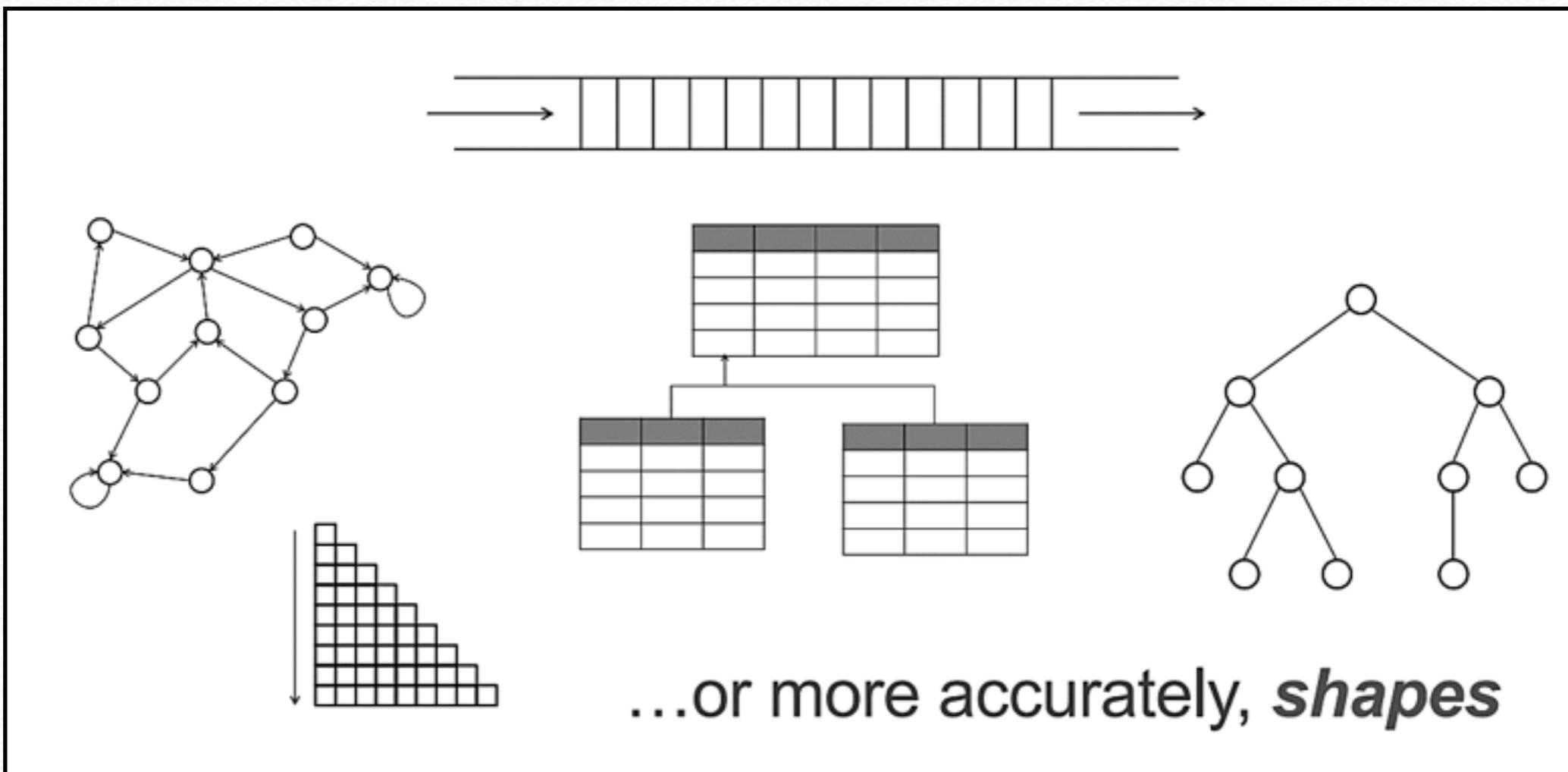
Polyglot Data Management

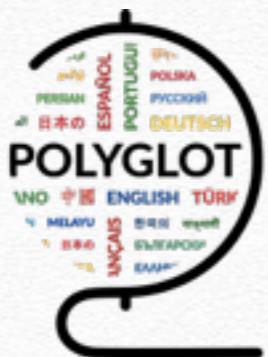
17/05/2017 - Big Data 2017



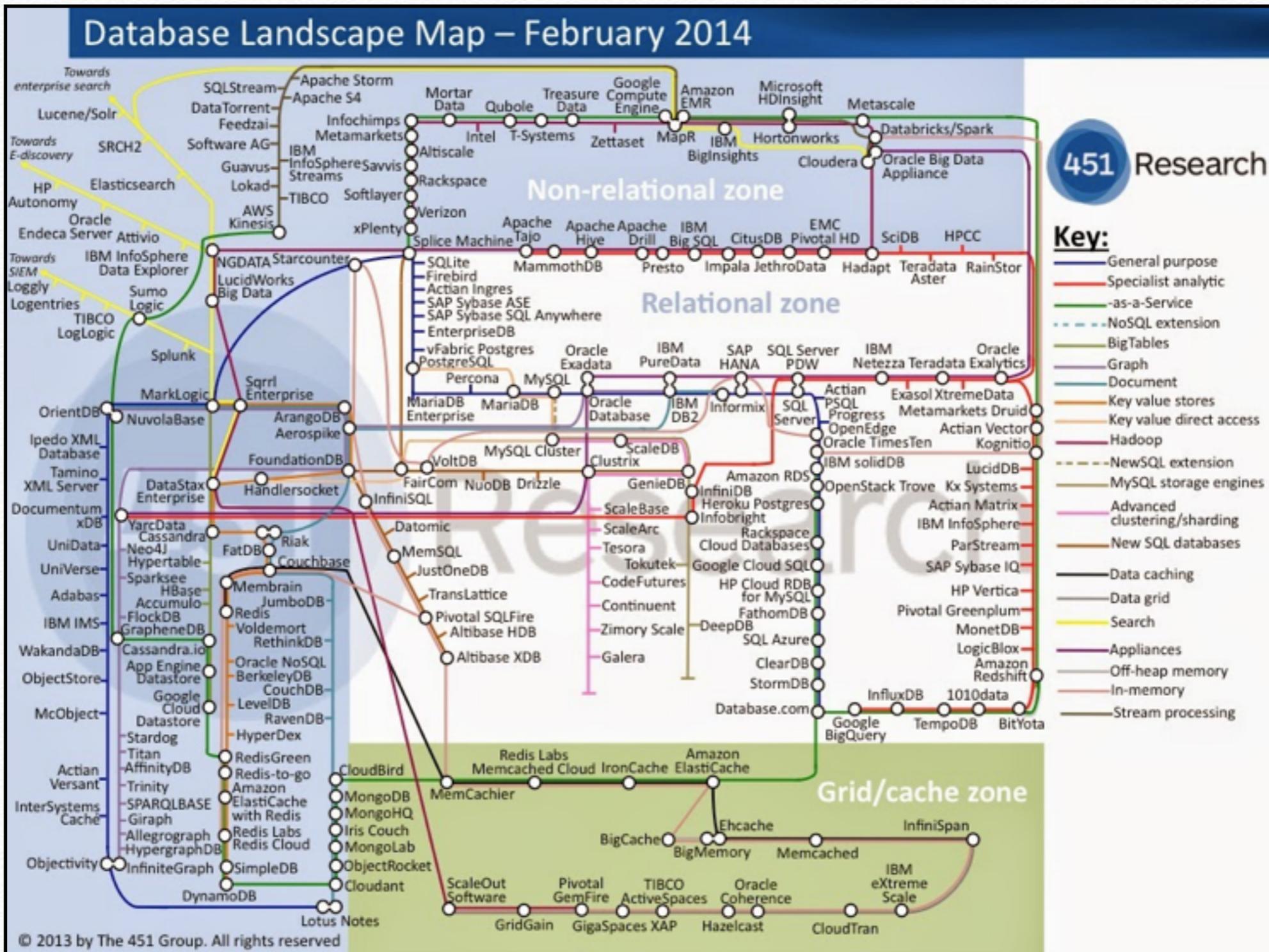
Different types of data shapes

- ❖ Data can have a variety of **shapes**. String cubes, graphs, relational tables, and trees are all examples of different data shapes

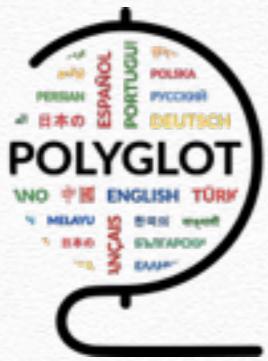




Different types of databases



451 Research

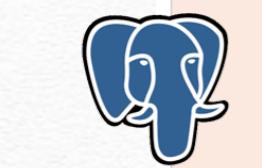


Polyglot Data Management

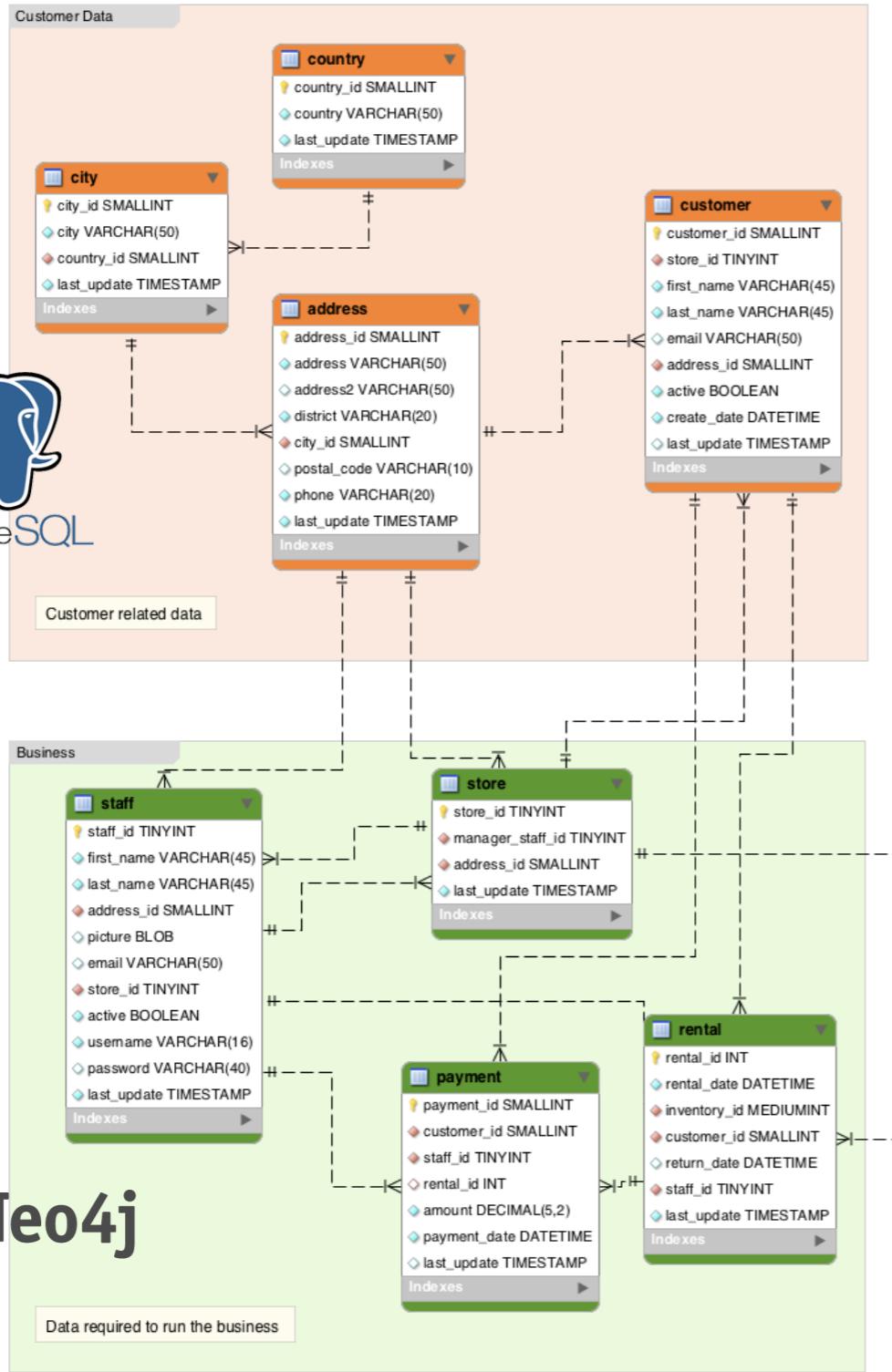
- ❖ A *polyglot data management platform* gives users the **flexibility** to choose the right tool for the job, instead of being forced into a solution that might not be optimal.
- ❖ The polyglot approach to data management must consider the **capabilities** associated with **data** as it flows through the information architecture from Acquisition to the General Core to Access.



Sample Scenario: a set of data containers



PostgreSQL

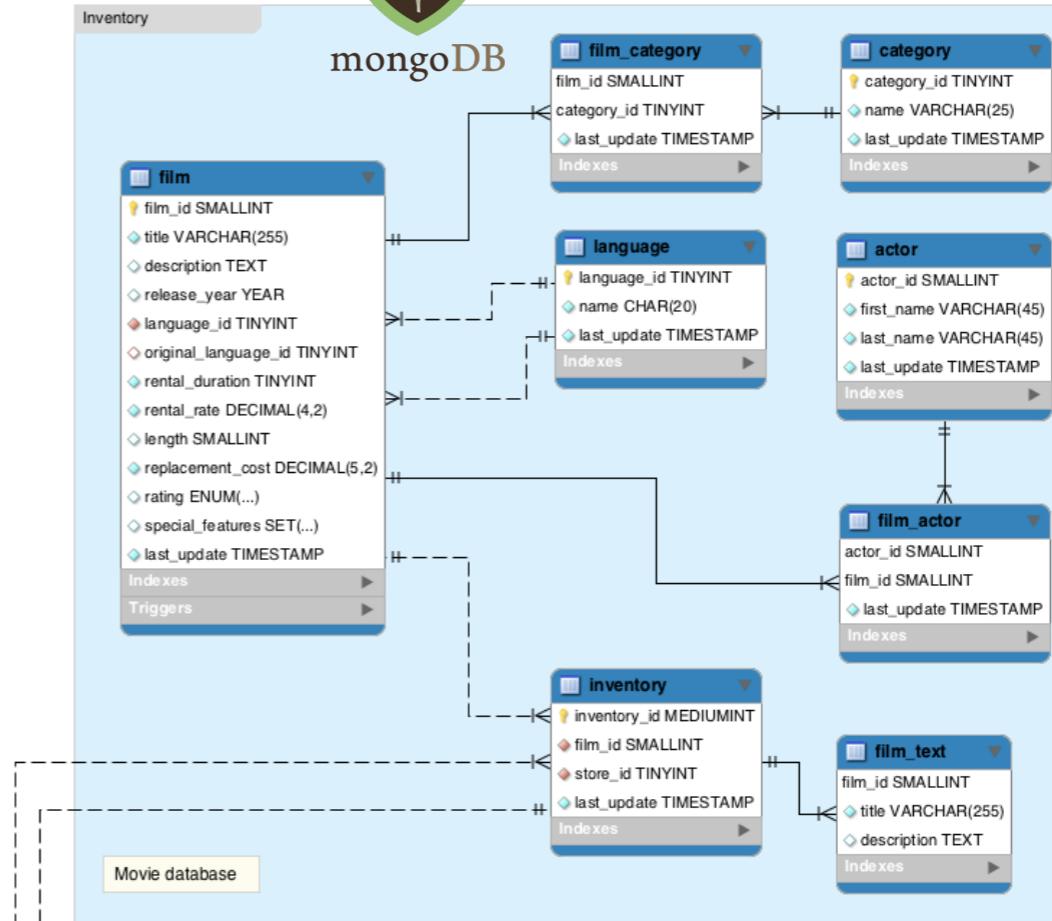


Neo4j

Data required to run the business



mongoDB



Views



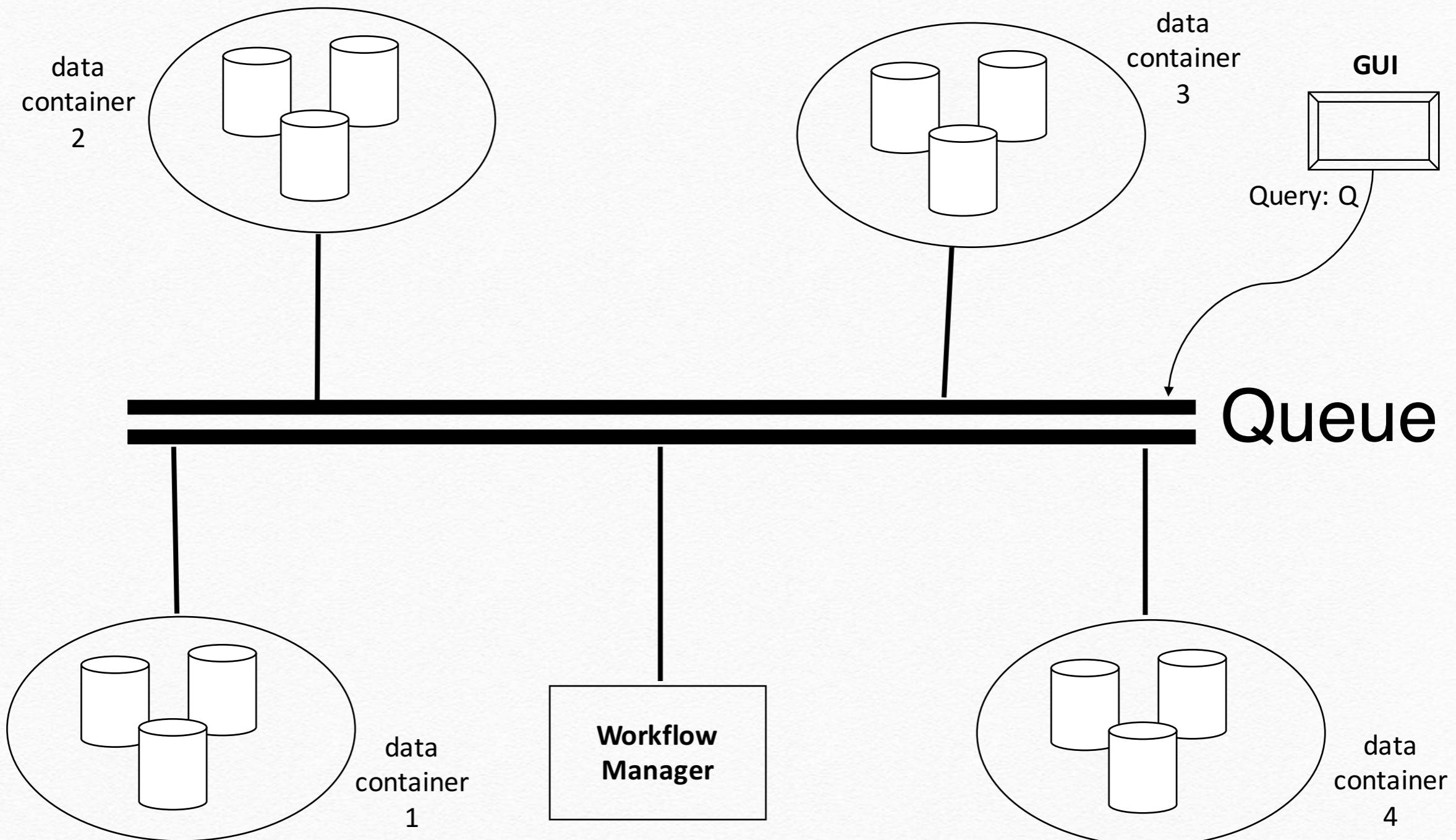
Special view on certain data used for appraisals

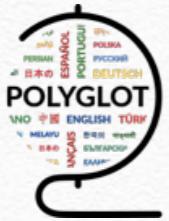


redis



Communication between data containers





Sample Scenario: technology plethora

- ❖ RDBMS (PostgreSql)



- ❖ Document Store (MongoDB)

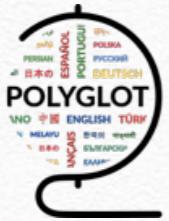


- ❖ Key-Value Store (Redis)



- ❖ GraphDB (Neo4J)





Sample Scenario: query language plethora

- ❖ RDBMS (SQL)



- ❖ Document Store (SQL binding, Java Object)



- ❖ Key-Value Store (Java Object)

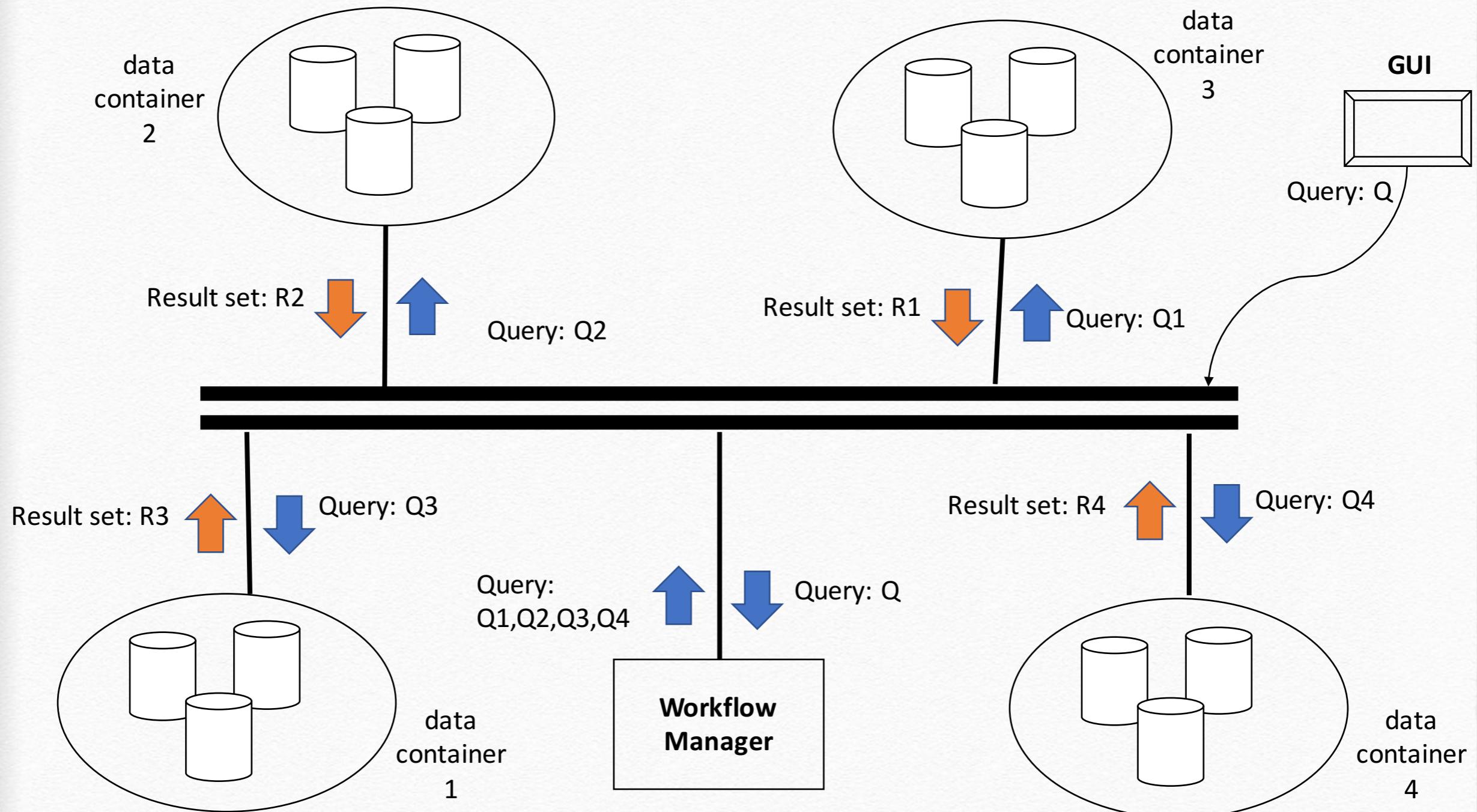


- ❖ GraphDB (Cypher)





USE CASE





Sample USE CASE

- ❖ Query Q: select last name of customers having rented the film with title “Titanic”
- ❖ Language: SQL

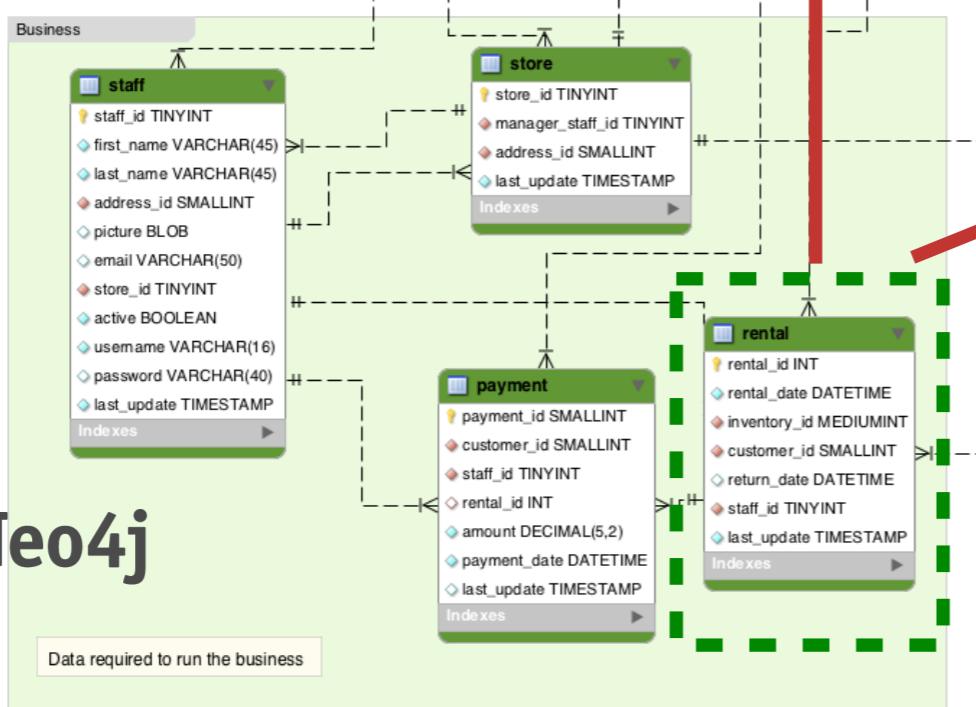
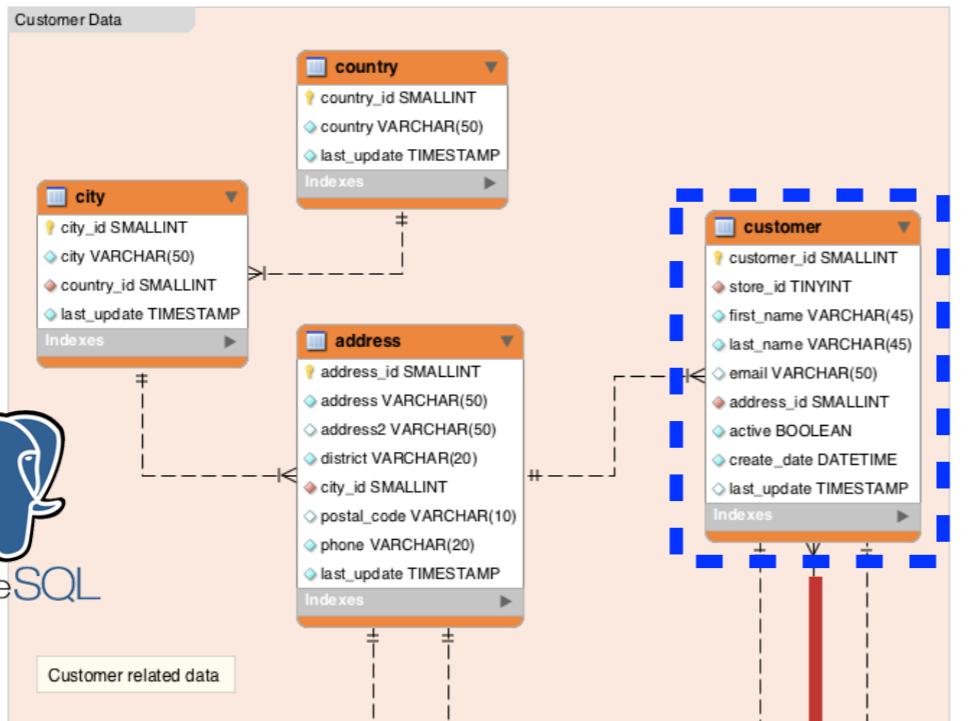
```
SELECT C.Last_Name  
FROM Customer C, Rental R, Inventory I, Film F  
WHERE C.customer_ID = R.customer_ID AND  
R.inventory_ID = I.inventory_id AND  
I.film_id = F.film_id AND F.title = "Titanic"
```



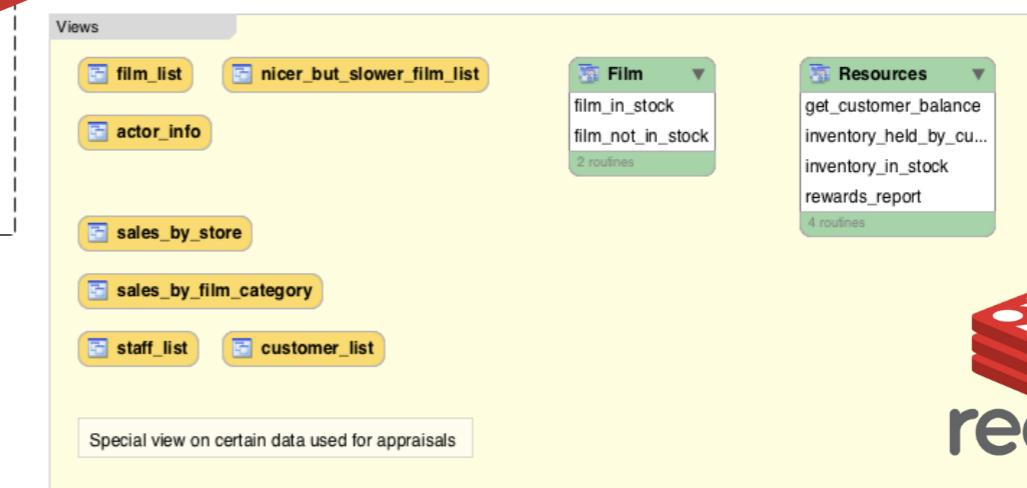
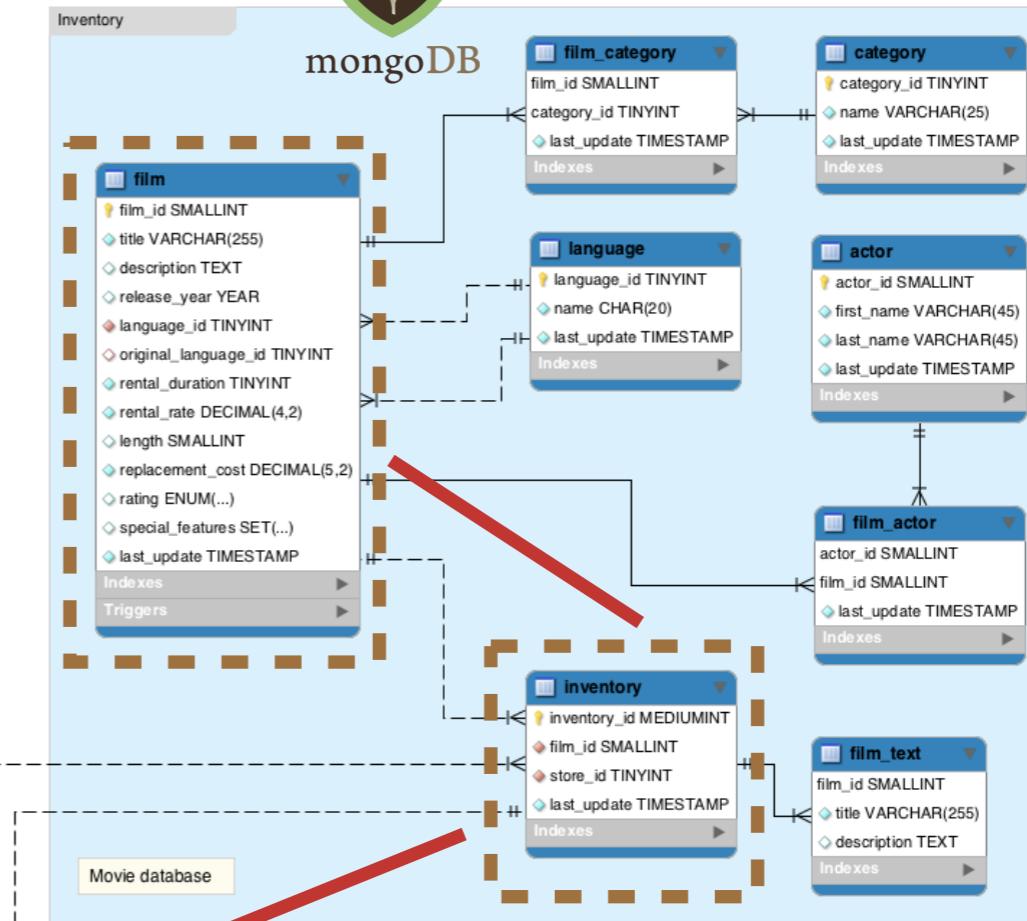
Sample USE CASE



PostgreSQL



Neo4j





Sample USE CASE



Q1

```
SELECT Last_Name FROM Customer WHERE customer_id = ?
```

Q2

```
MATCH (node1)  
WHERE node1.inventory_id = ?  
RETURN node1.customer_id
```



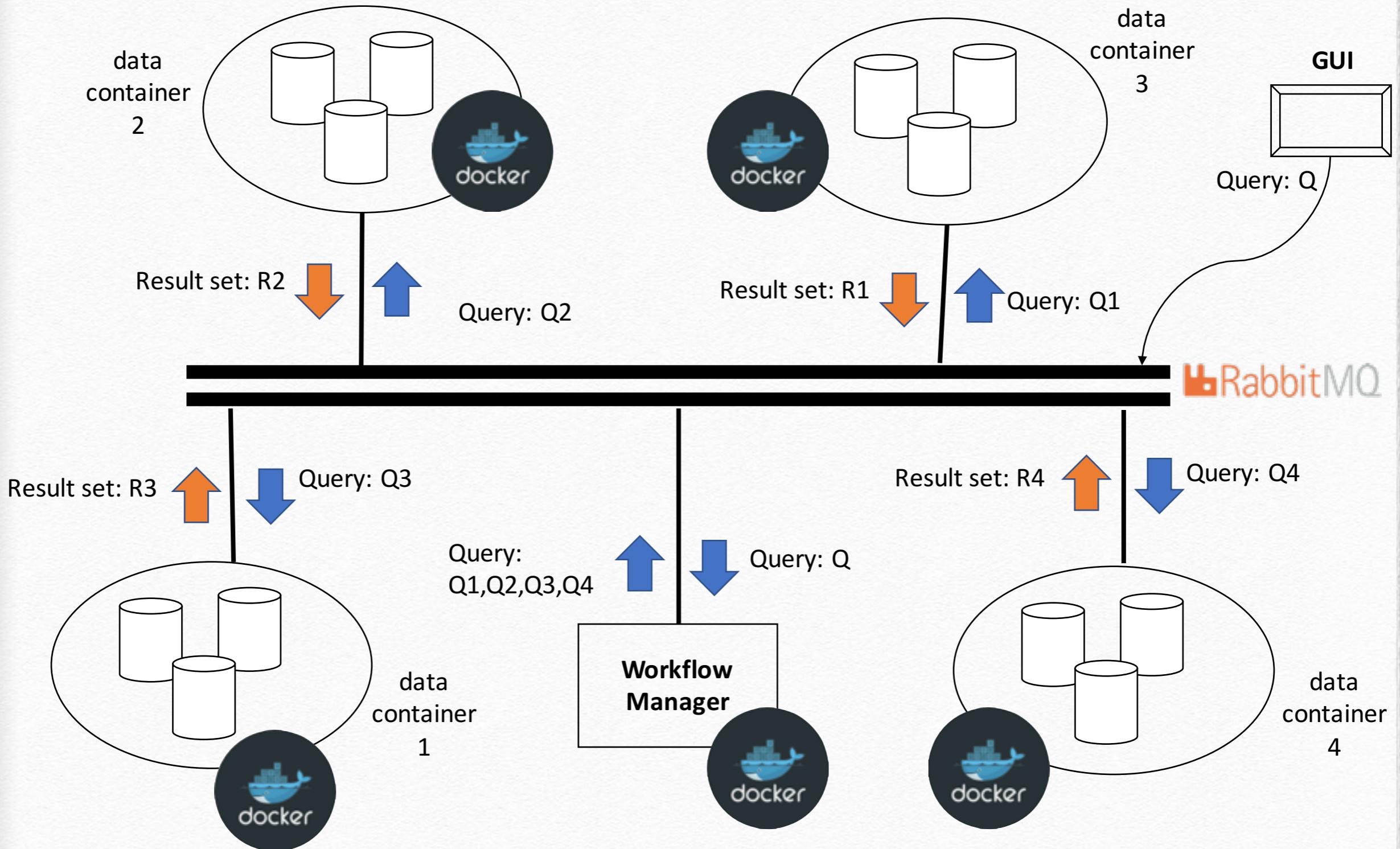
Q3

```
db.inventory.find(  
  { film: { title: "Titanic" } },  
  {inventory_id: 1})
```



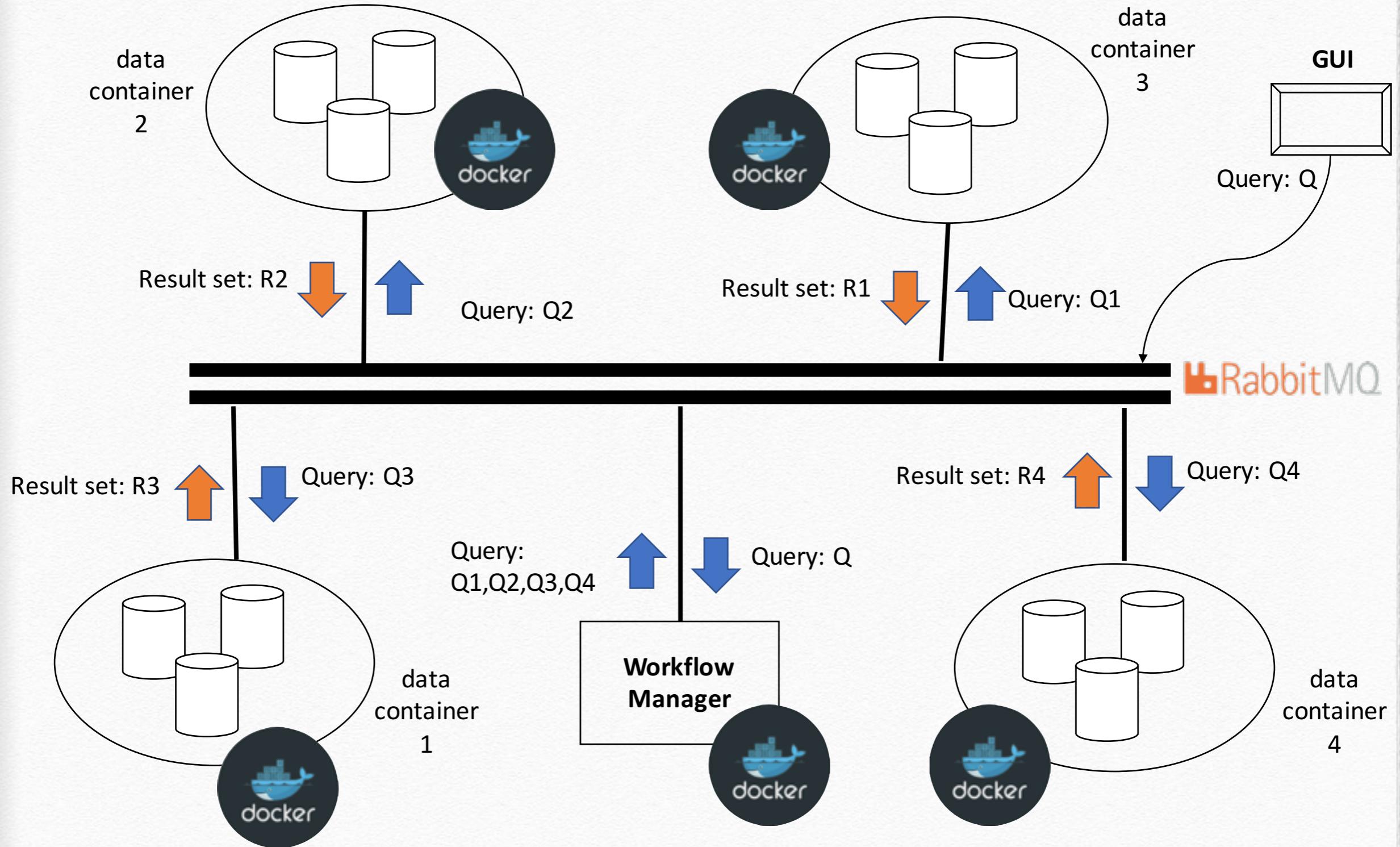


Micro-Services via Docker Containers



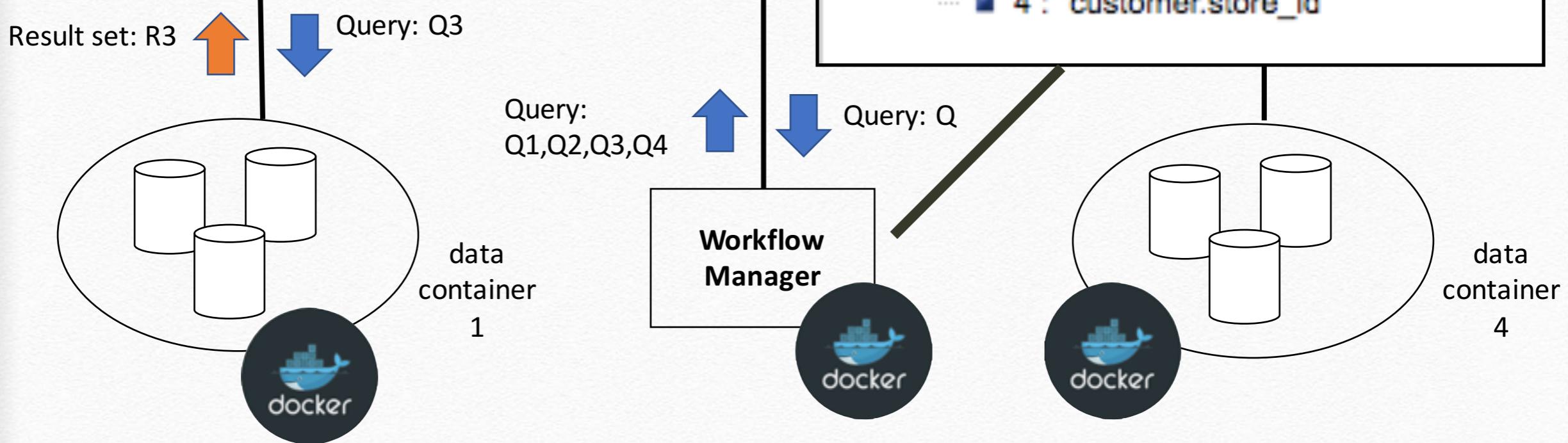
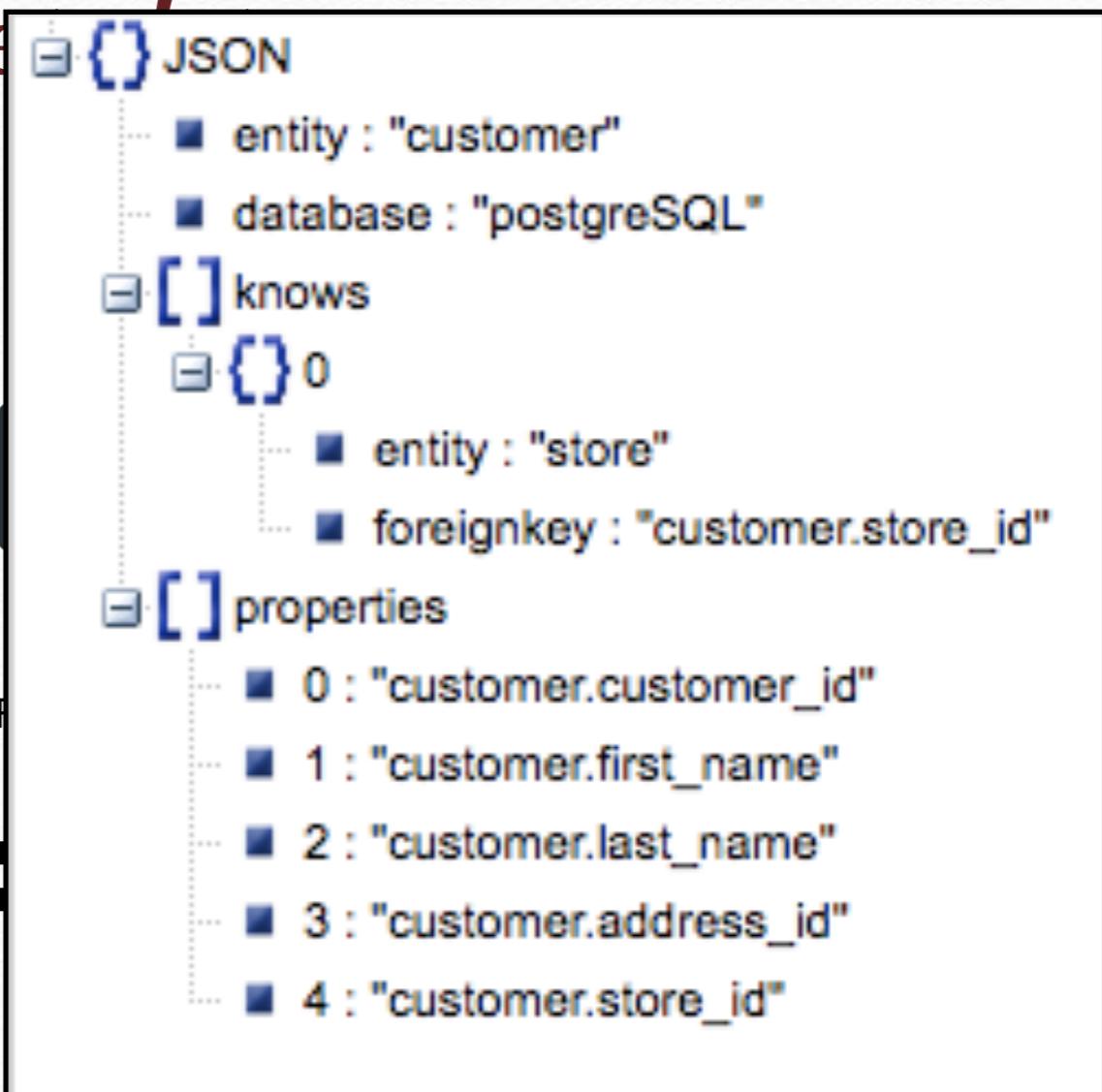
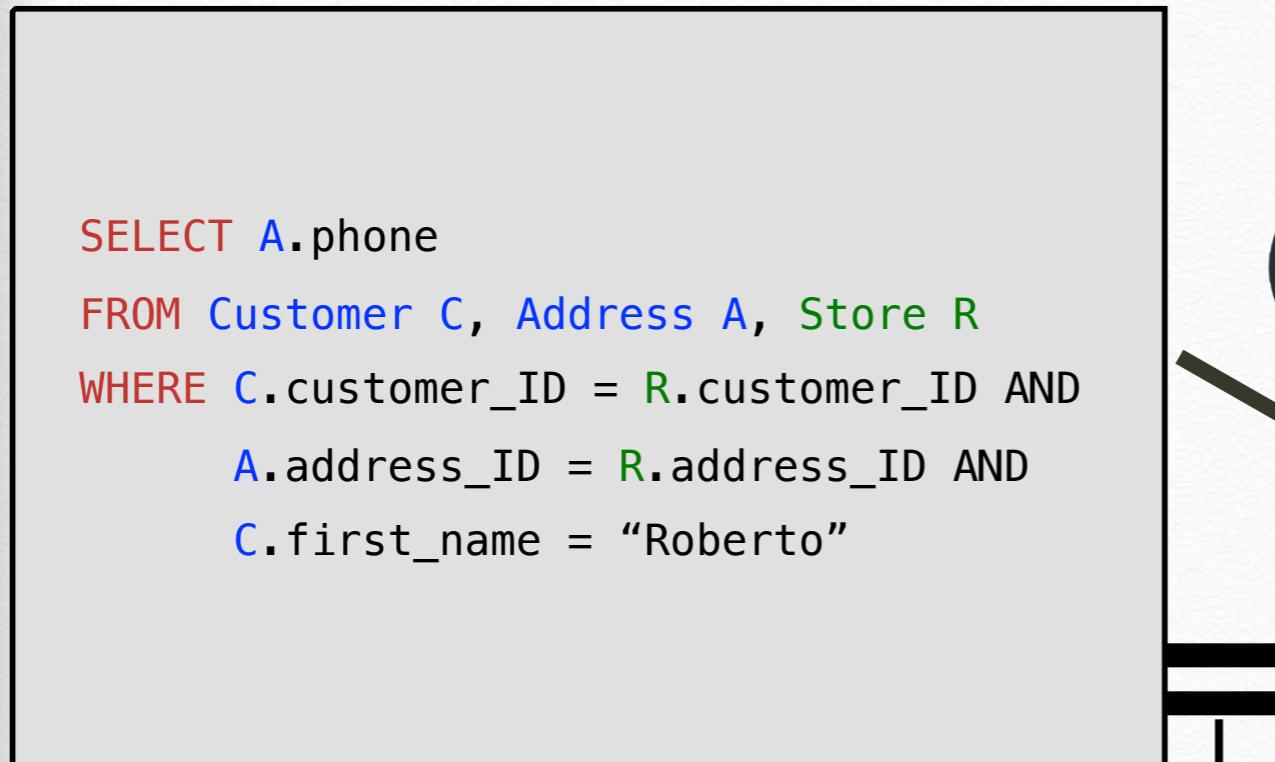


JSON Metadata



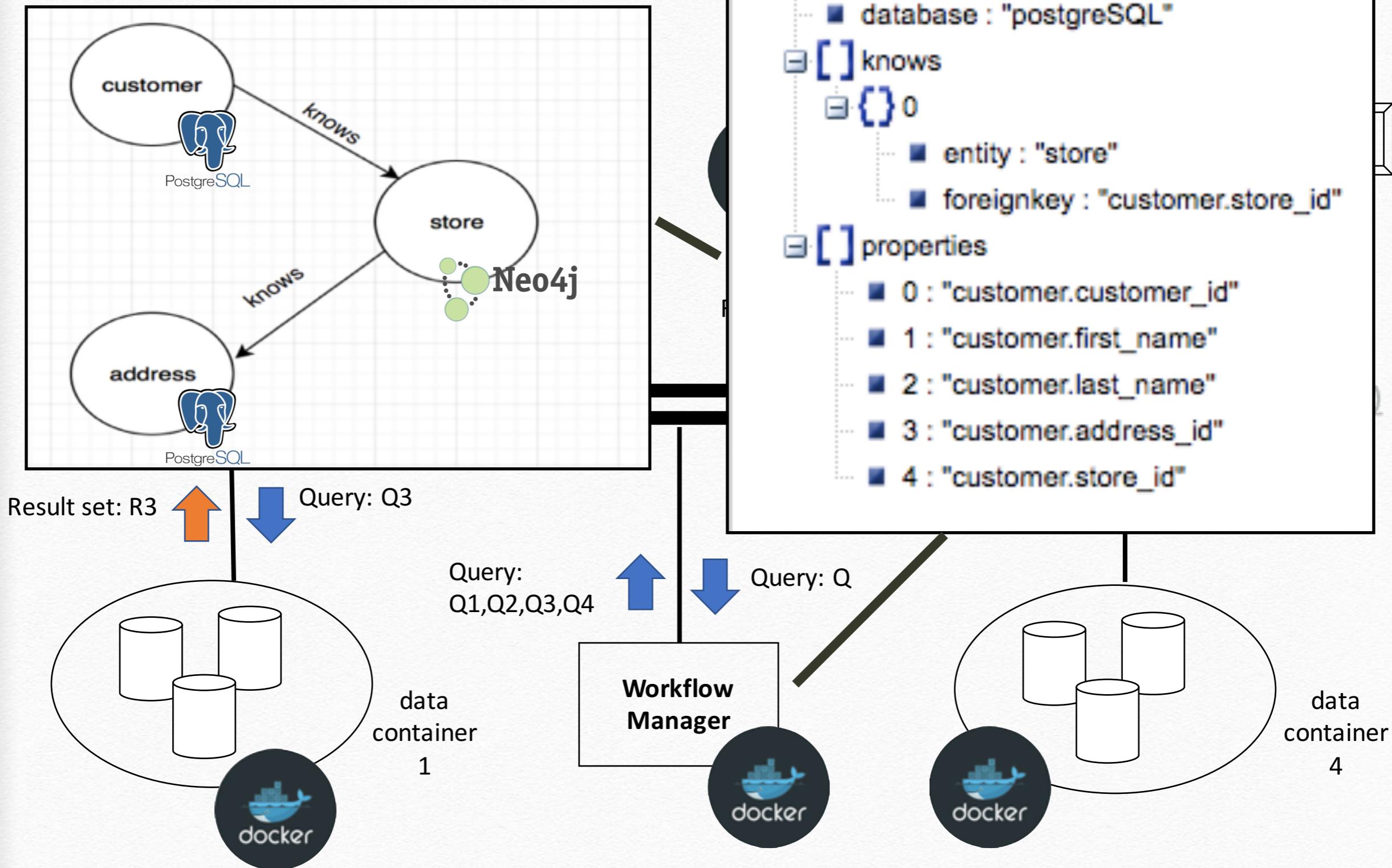


JSON Merge



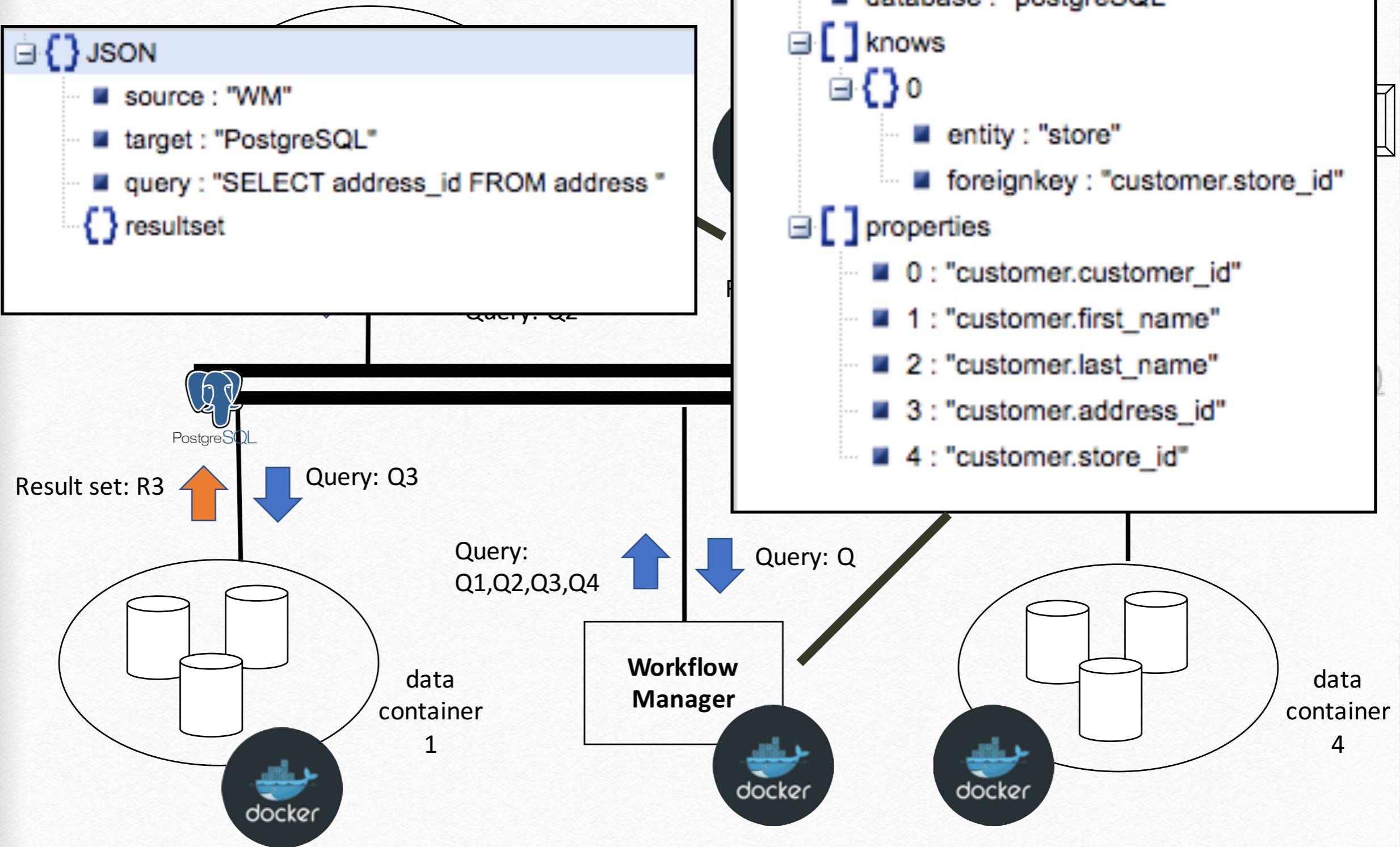


JSON Merge



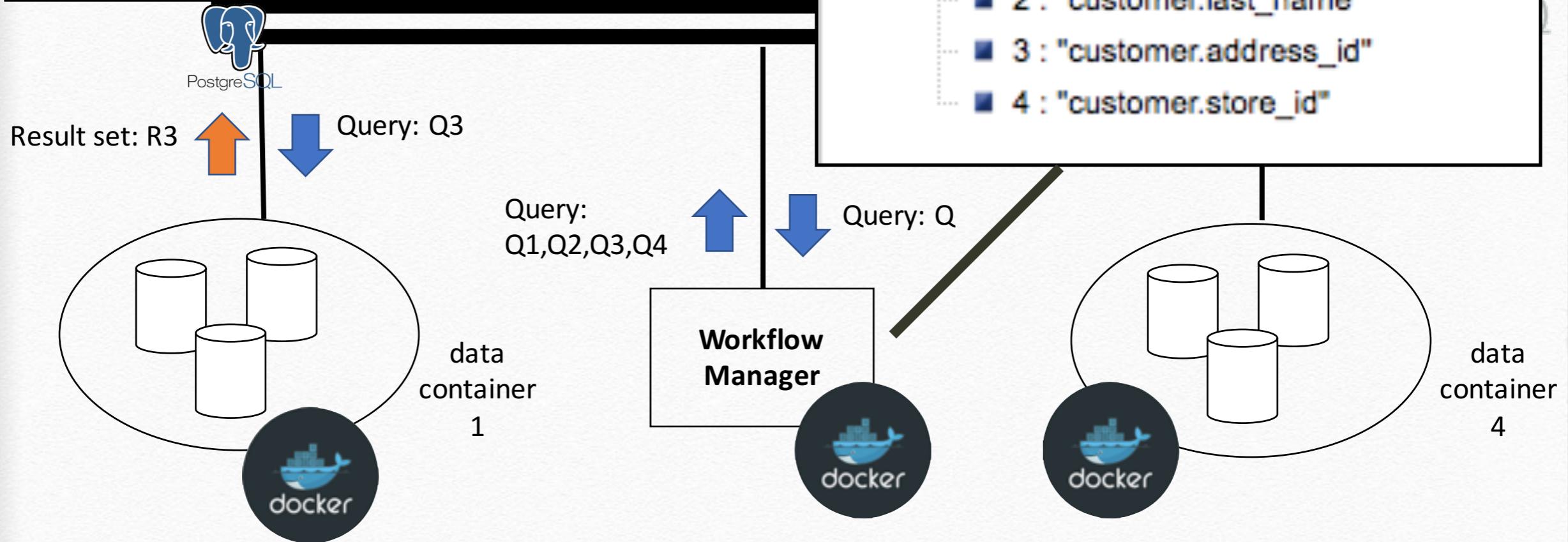
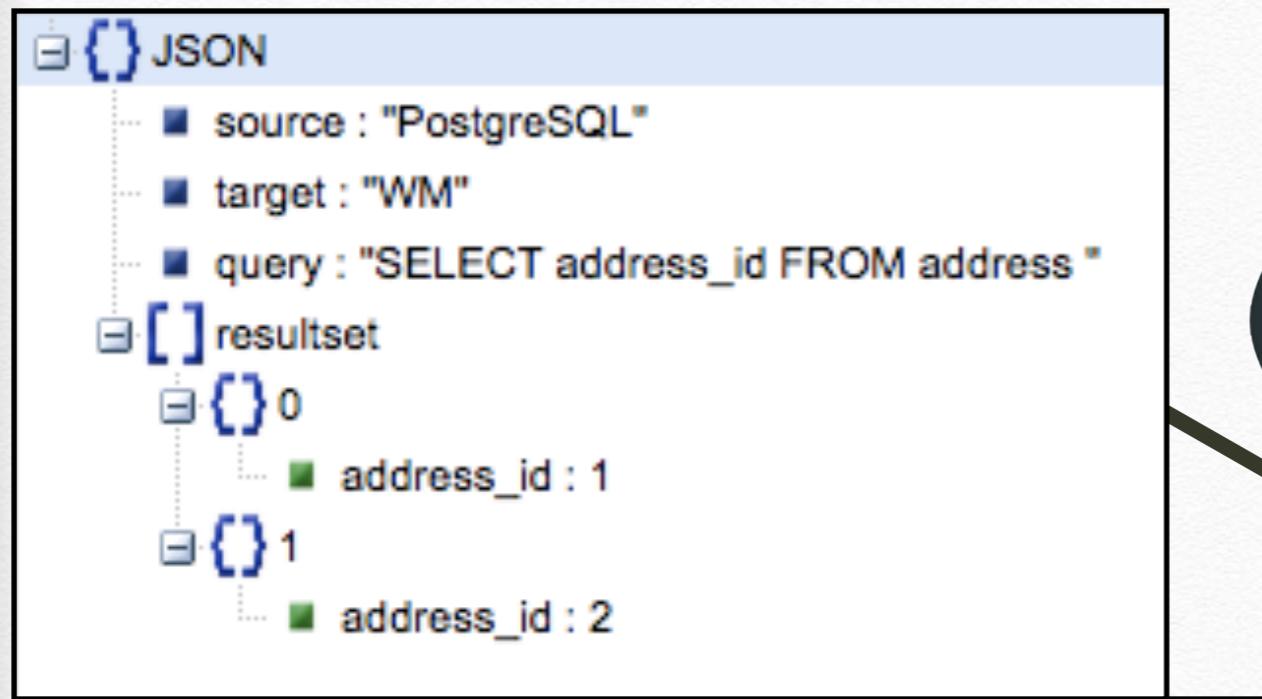


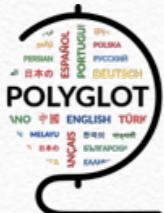
JSON Merge



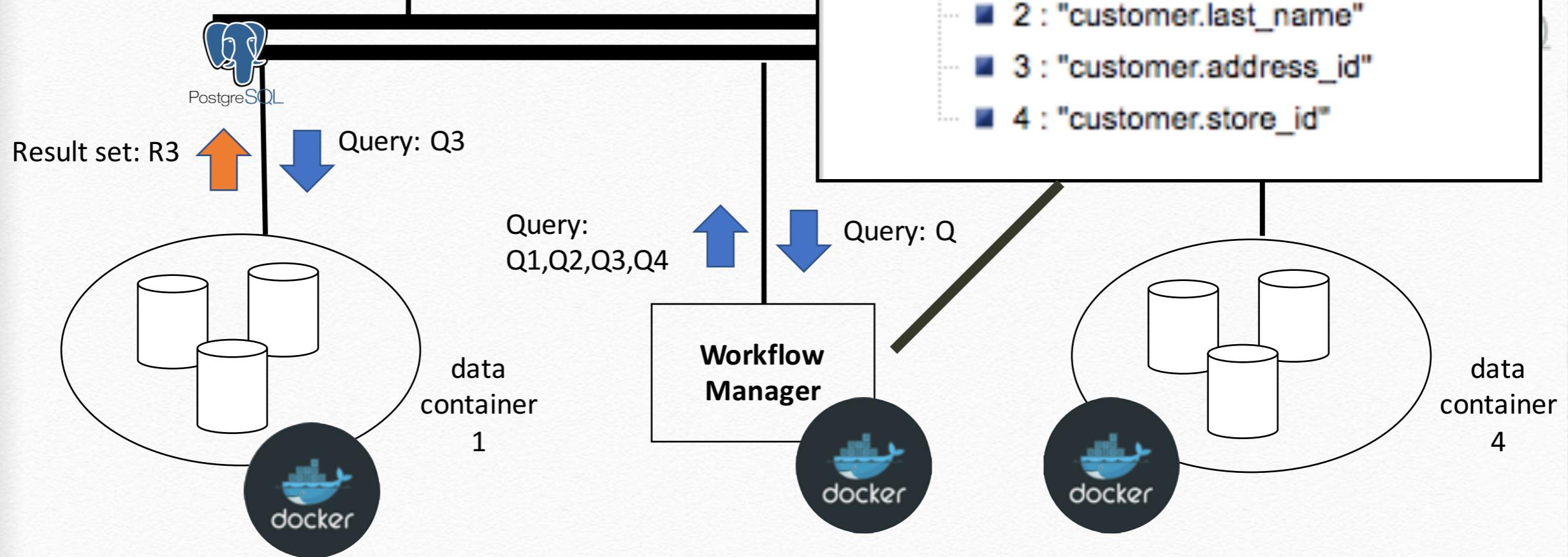
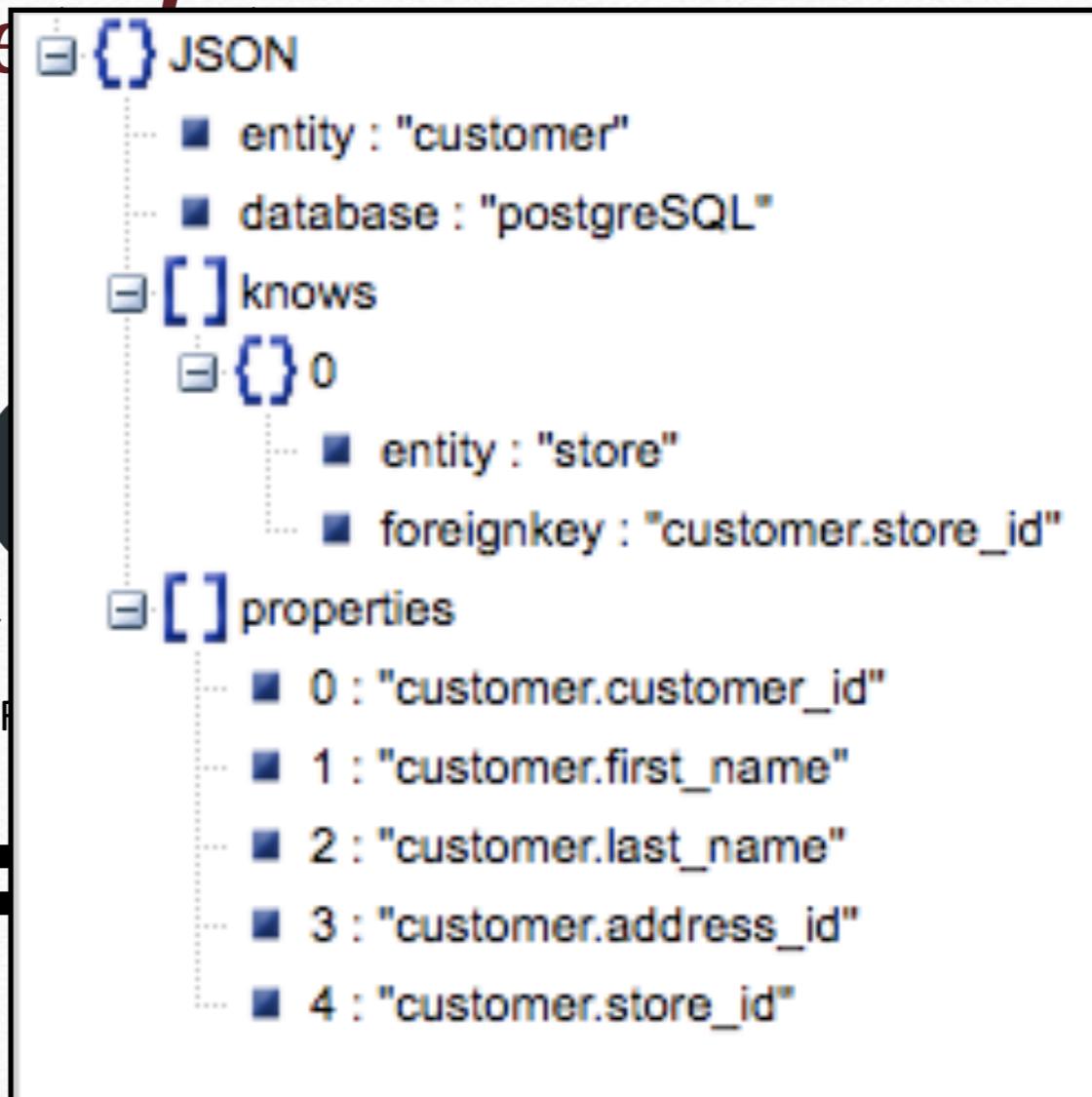
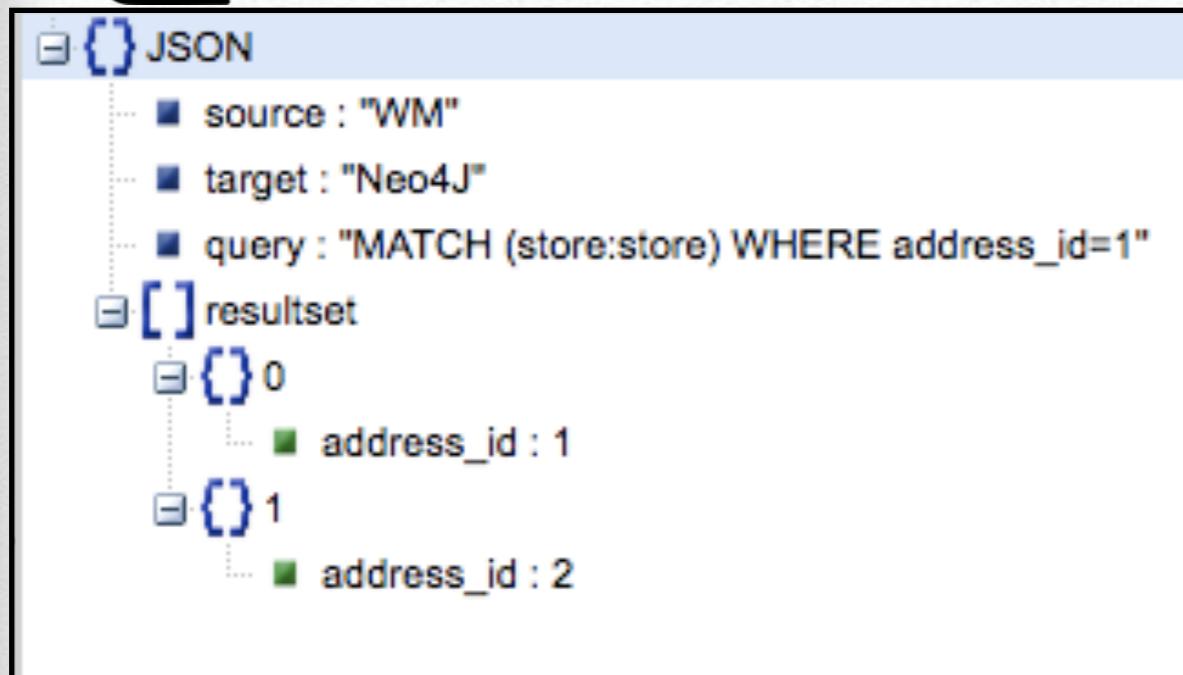


JSON Merge





JSON Merge





Prototype:

http://pamir.dia.uniroma3.it:8080/0Polystore1/

DATABASE

PostgreSQL ▾

customer ▾

address ▾

city ▾

country ▾

Neo4J ▾

store ▾

staff ▾

payment ▾

rental ▾

MongoDB ▾

Inventory ▾

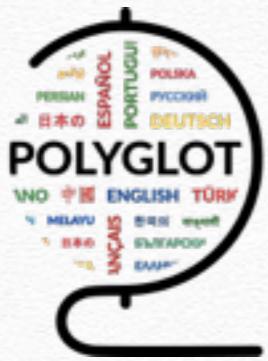
film ▾

Inserisci query nel linguaggio che conosci

Esempio: `SELECT * FROM customer` oppure `db.customer.find()` oppure `MATCH(customer:customer) RETURN customer.*`

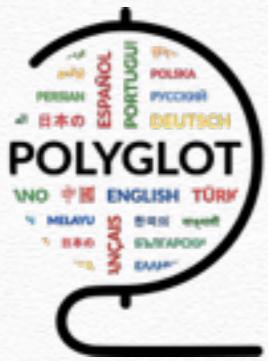
Query

Invia



Proposals: #1

- ❖ **Persist:** given input sources s_1, s_2, \dots, s_n , link and persist them in the system
- ❖ *case study:* Films and reviews
 - iMDB -> film details (MongoDB)
 - Wikipedia -> exploration of film details, e.g. plot, cast and so on, (PostgreSQL)
 - Youtube -> film trailer (Redis and/or Oracle NoSql)
 - Movielens -> film ratings (Neo4J)
 - Amazon -> film commerce (Cassandra)



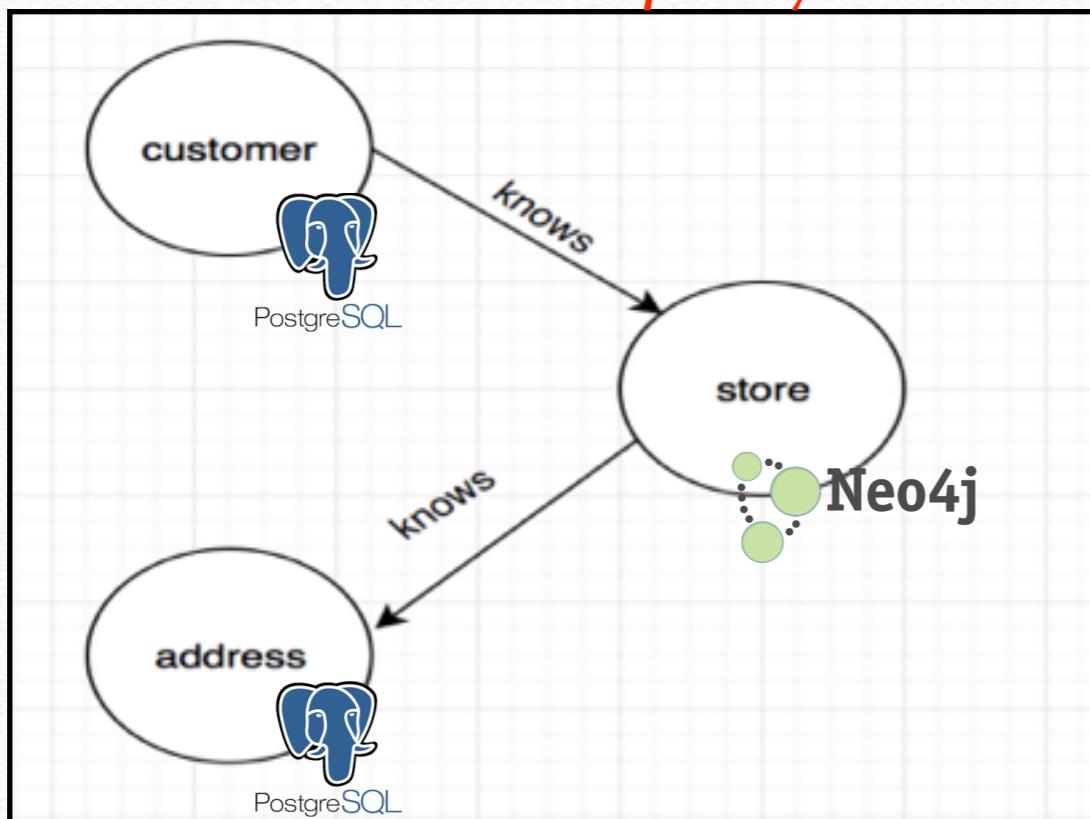
Proposals: #2



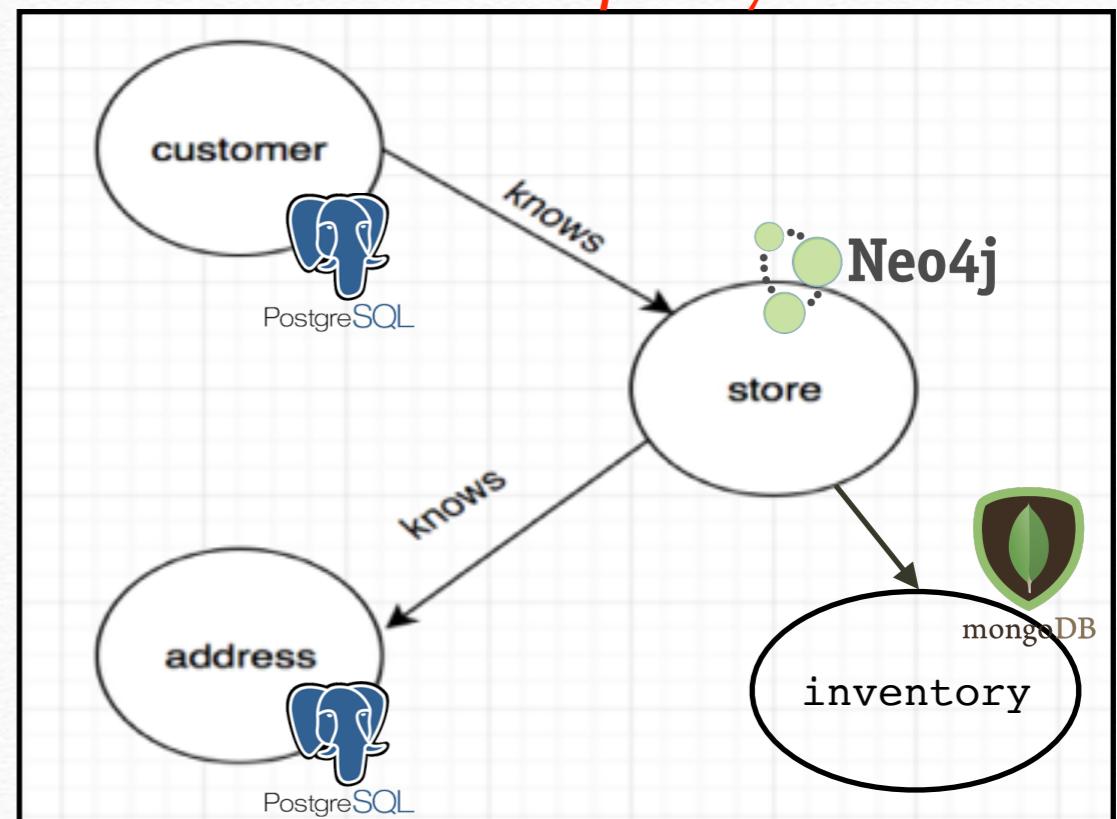
❖ Flow: execution plan improvement algorithm

❖ case study

chain query



star query





Proposals: #3

- ❖ **Query Analysis:** extended analysis of query language constructs used in the system and implementation in the query plan
- ❖ *case study:*
 - SQL -> OUTER JOINS, UNION, DIFFERENCE, NESTED QUERY
 - Cypher -> SUB-QUERY
 - Document Store -> NESTED COLLECTIONS

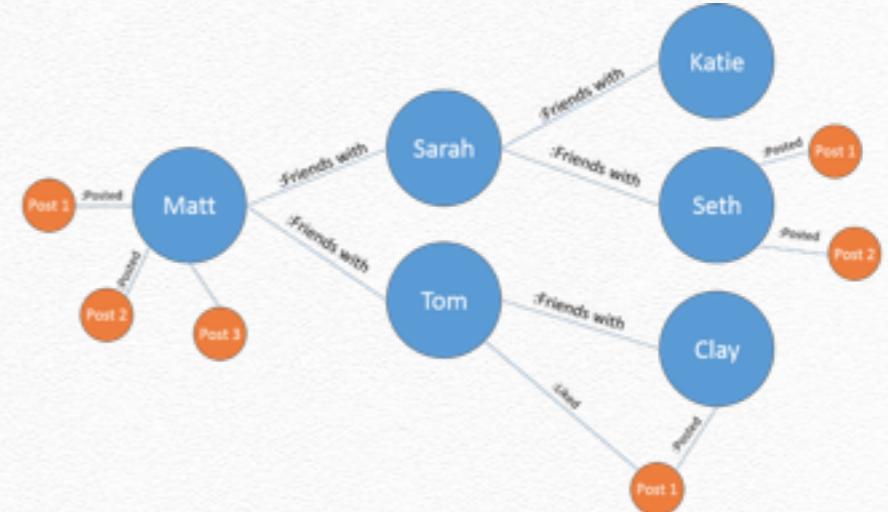


Proposals: #4

- ❖ Distributed Visual Exploration (data and schema):
- ❖ exporting of the system in a distributed environment
- ❖ front-end for exploring and querying data containers (i.e the exploration is depending on the selected query language)

	Results	Messages			
	LocationID	Name	CostRate	Availability	ModifiedDate
1	1	Tool Crib	0.00	0.00	2002-06-01 00:00:00.000
2	2	Sheet Metal Racks	0.00	0.00	2002-06-01 00:00:00.000
3	3	Paint Shop	0.00	0.00	2002-06-01 00:00:00.000
4	4	Paint Storage	0.00	0.00	2002-06-01 00:00:00.000
5	5	Metal Storage	0.00	0.00	2002-06-01 00:00:00.000
6	6	Miscellaneous Storage	0.00	0.00	2002-06-01 00:00:00.000
7	7	Finished Goods Storage	0.00	0.00	2002-06-01 00:00:00.000
8	10	Frame Forming	22.50	96.00	2002-06-01 00:00:00.000
9	20	Frame Welding	25.00	108.00	2002-06-01 00:00:00.000
10	30	Debur and Polish	14.50	120.00	2002-06-01 00:00:00.000
11	40	Paint	15.75	120.00	2002-06-01 00:00:00.000
12	45	Specialized Paint	18.00	80.00	2002-06-01 00:00:00.000
13	50	Subassembly	12.25	120.00	2002-06-01 00:00:00.000
14	60	Final Assembly	12.25	120.00	2002-06-01 00:00:00.000

```
{id: 1,  
name: Joe,  
url: '...',  
stream:  
[  
  {  
    user: 2,  
    title: 'today',  
    body: 'go fly a kite',  
    likes: [  
      {user: 3},  
      {user: 1},  
    ],  
  }  
]}  
}
```





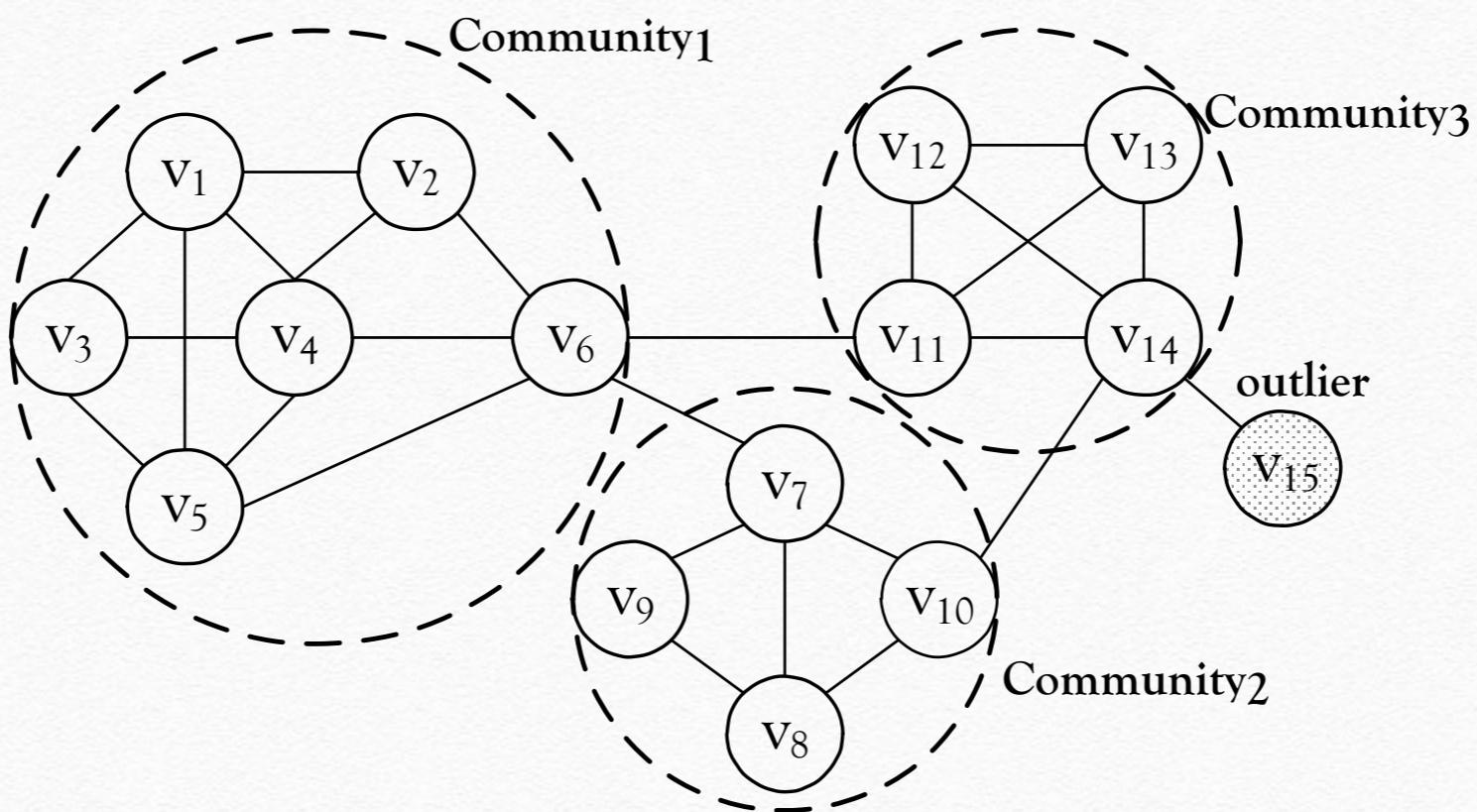
Community Profiling in Social Networks

17/05/2017 - Big Data 2017



Community in social media

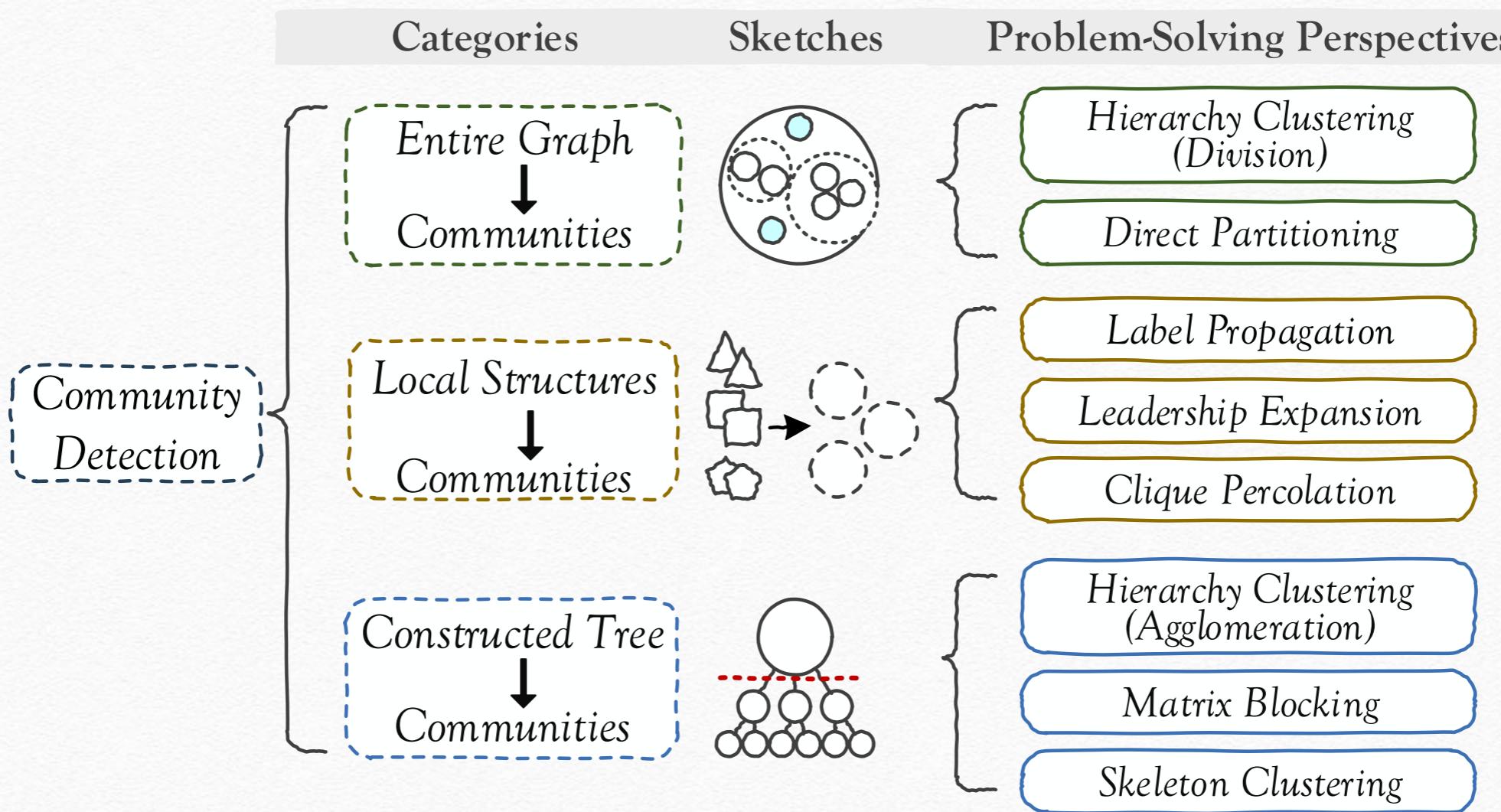
- ❖ **Community:** It is formed by individuals such that those within a group interact with each other more frequently than with those outside the group
- ❖ **Community detection:** discovering groups in a network where individuals' group memberships are not explicitly given





Community in social media

- ❖ **Community:** It is formed by individuals such that those within a group interact with each other more frequently than with those outside the group
- ❖ **Community detection:** discovering groups in a network where individuals' group memberships are not explicitly given





Community Profiling

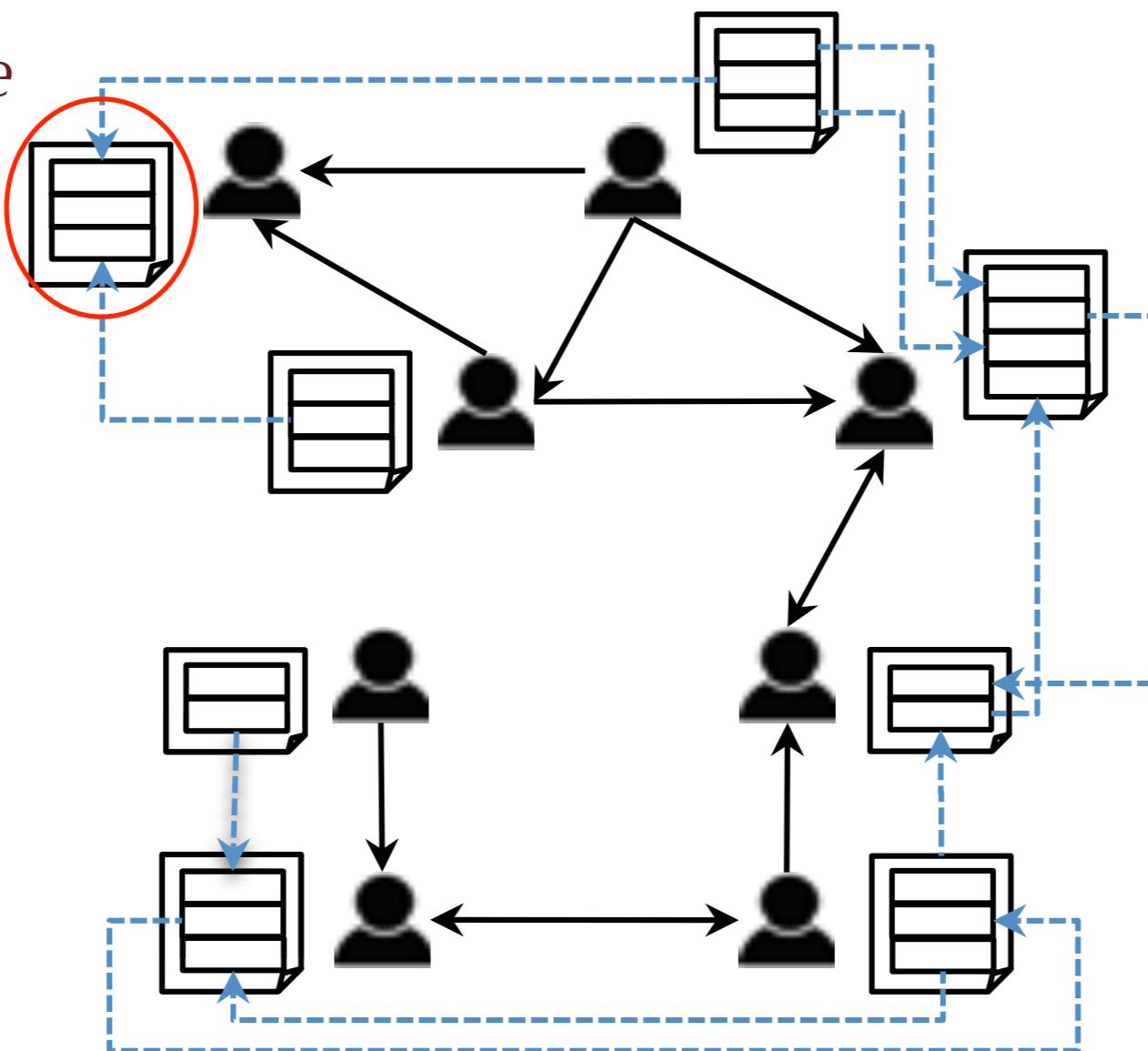
- ❖ **Positive Aspect:** community membership assists us to better understand the network structure.
- ❖ **Negative Aspect:** membership alone, without knowing **what a community is** and **how it interacts with others**, has only limited applications

- ❖ **Community profiling:** to characterize the intrinsic nature and extrinsic behavior of a community – thereby enabling useful community-level applications.



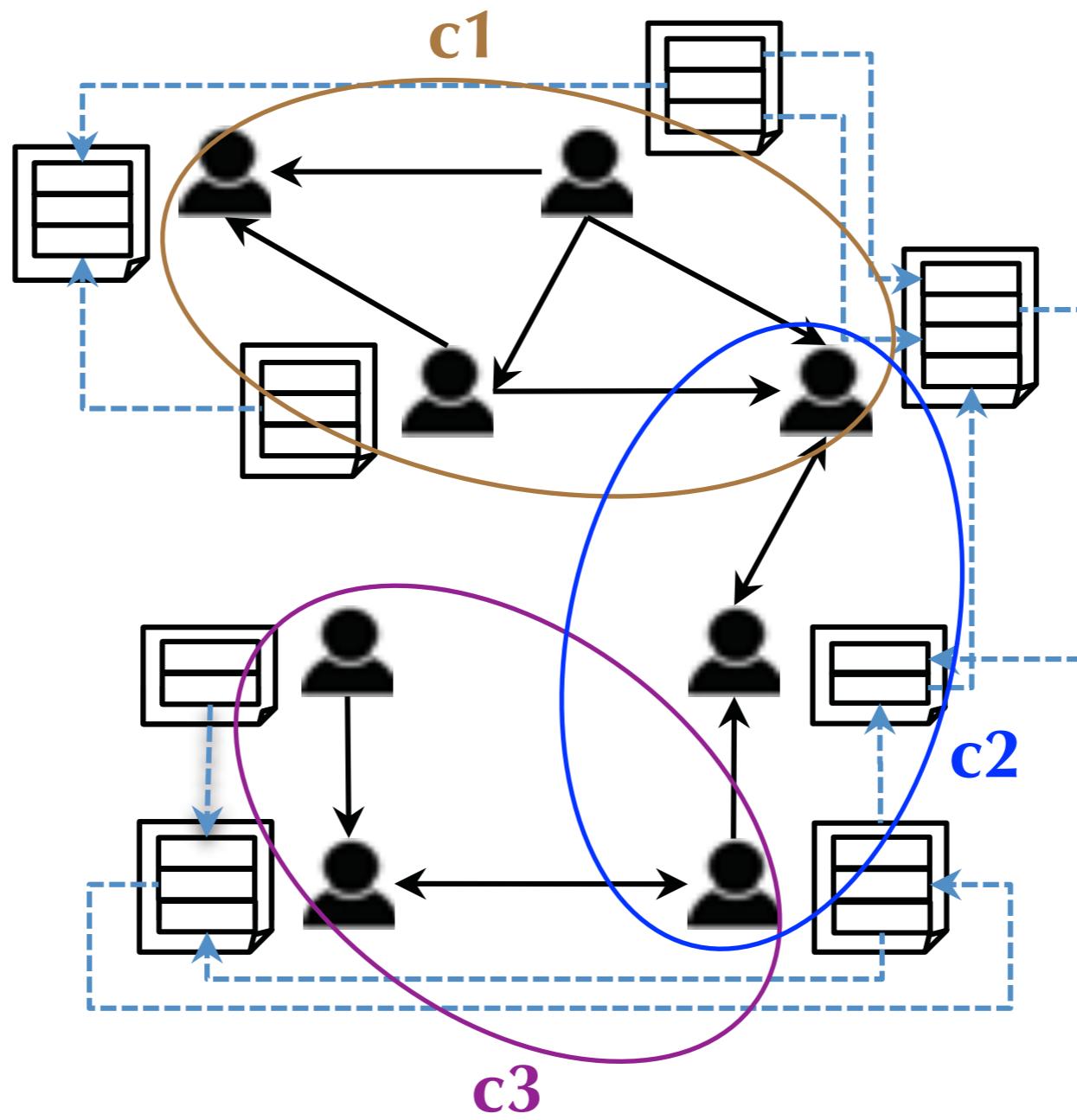
Community Profiling

user profile



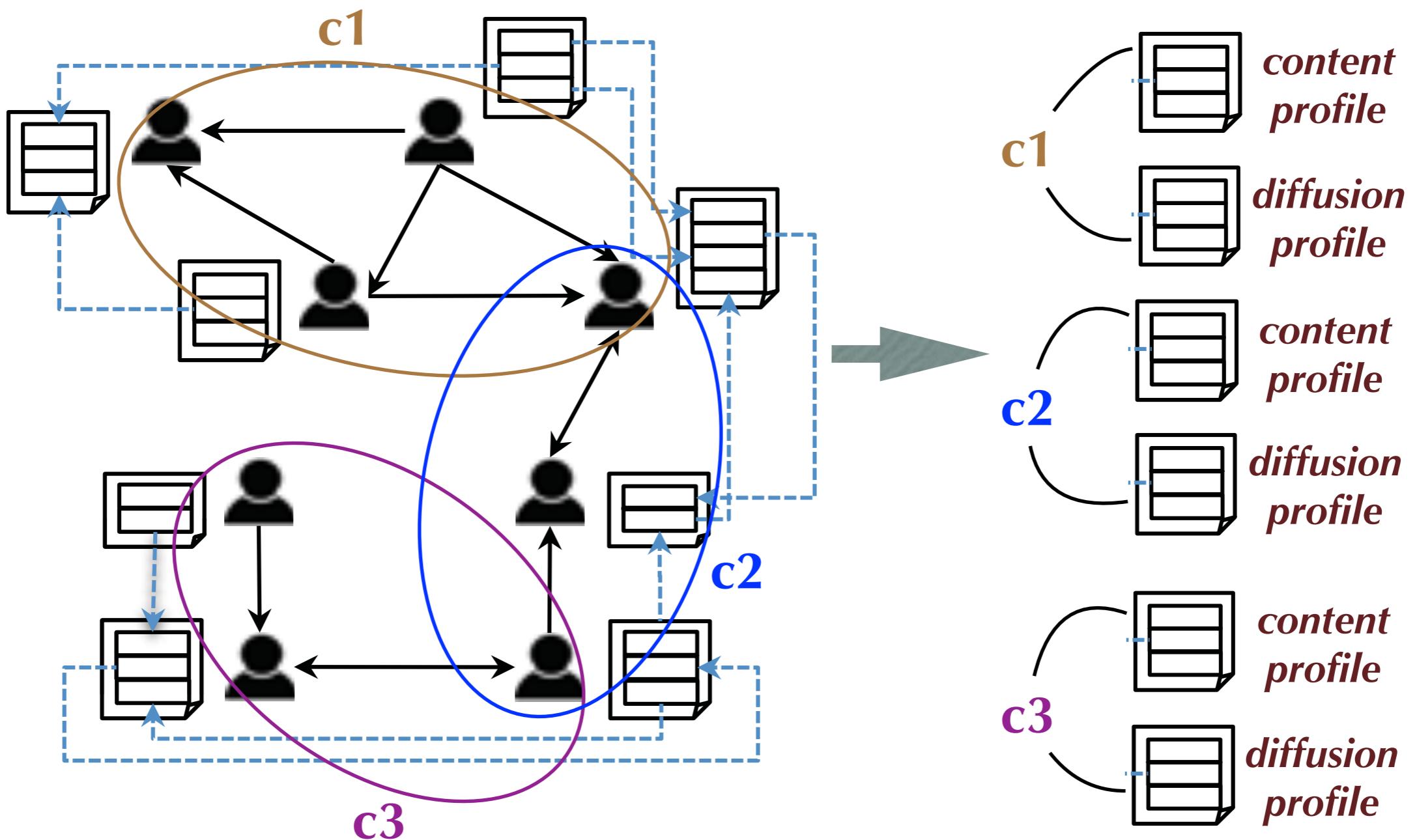


Community Profiling



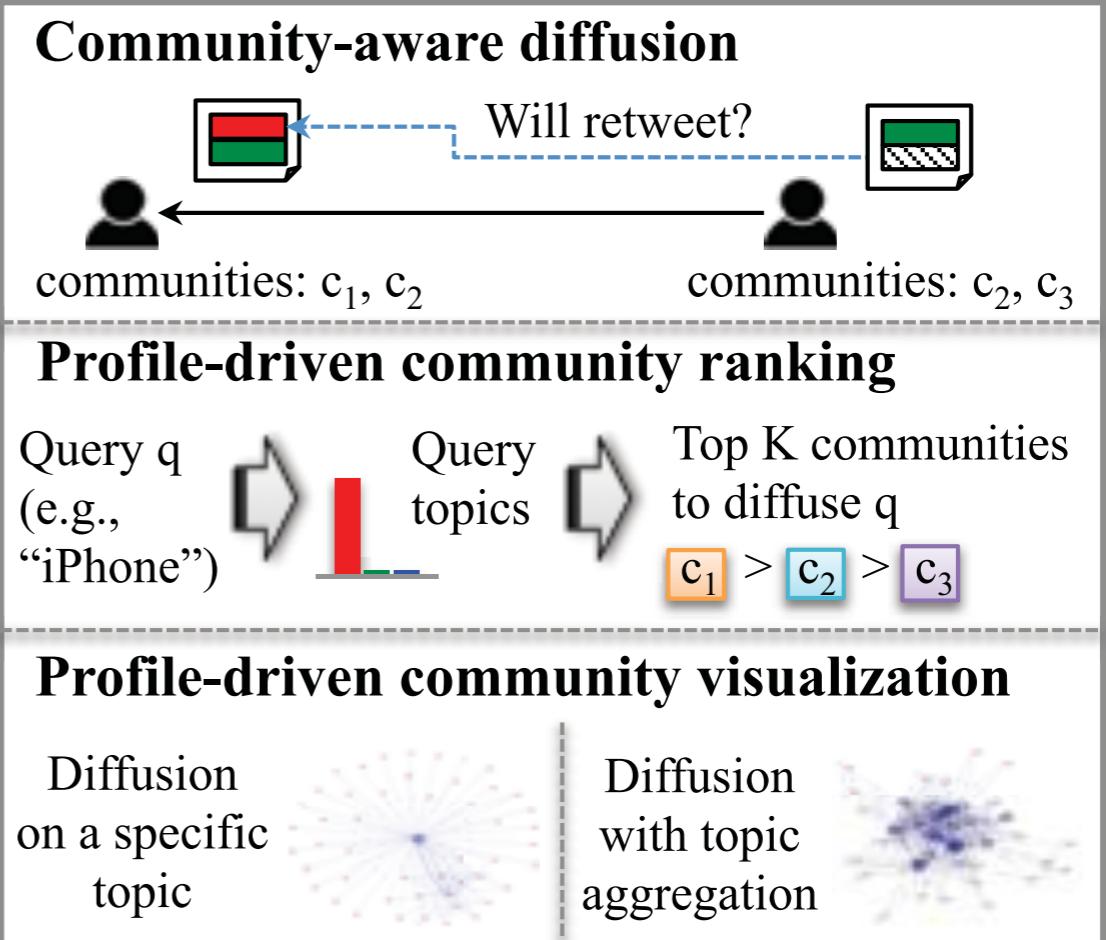
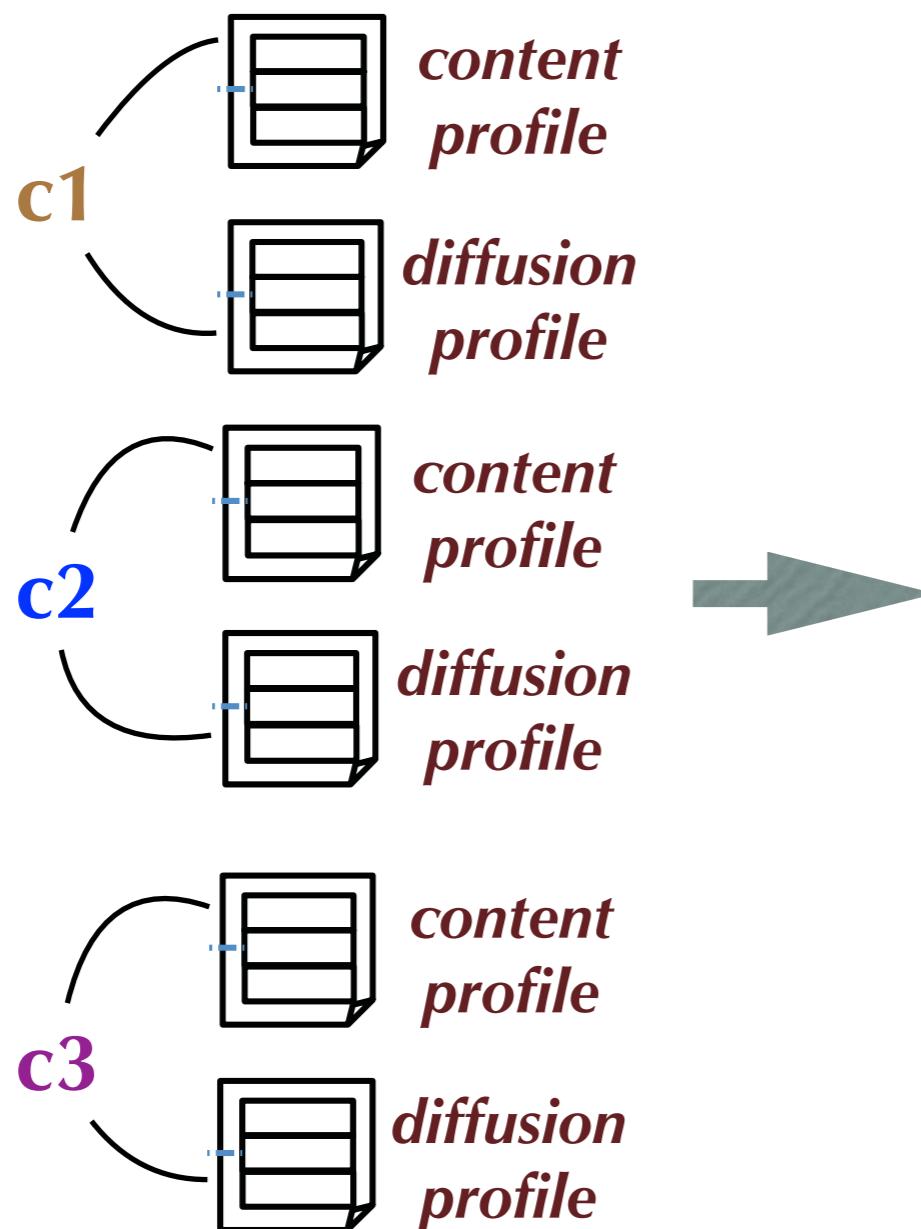


Community Profiling



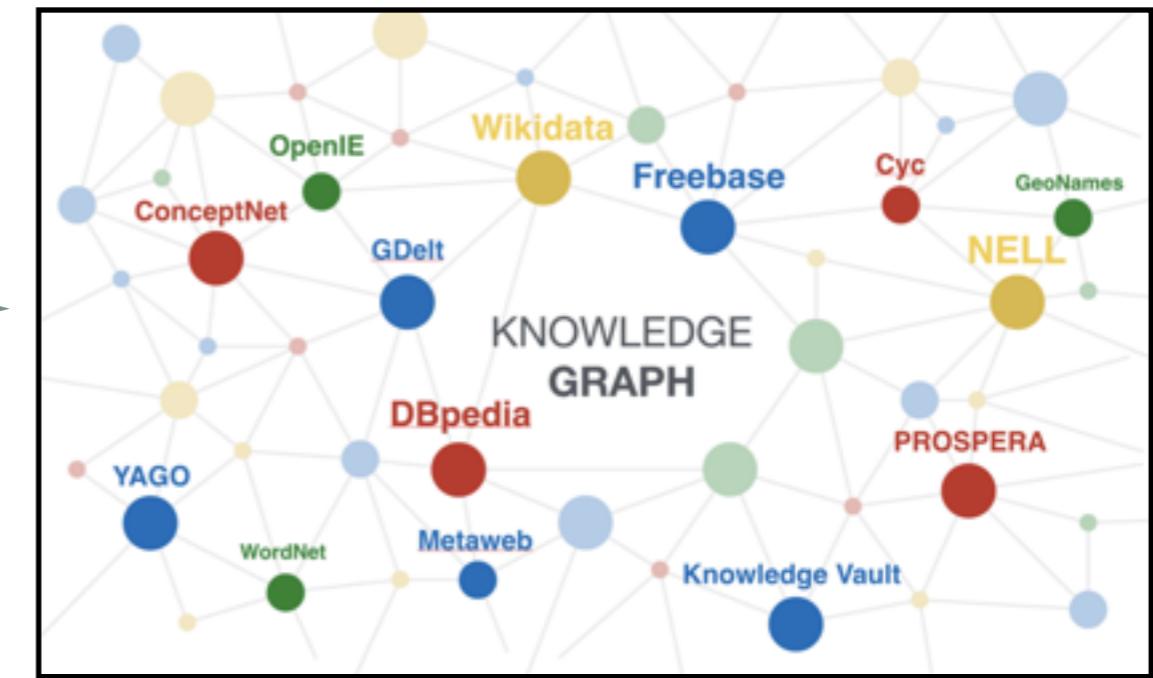
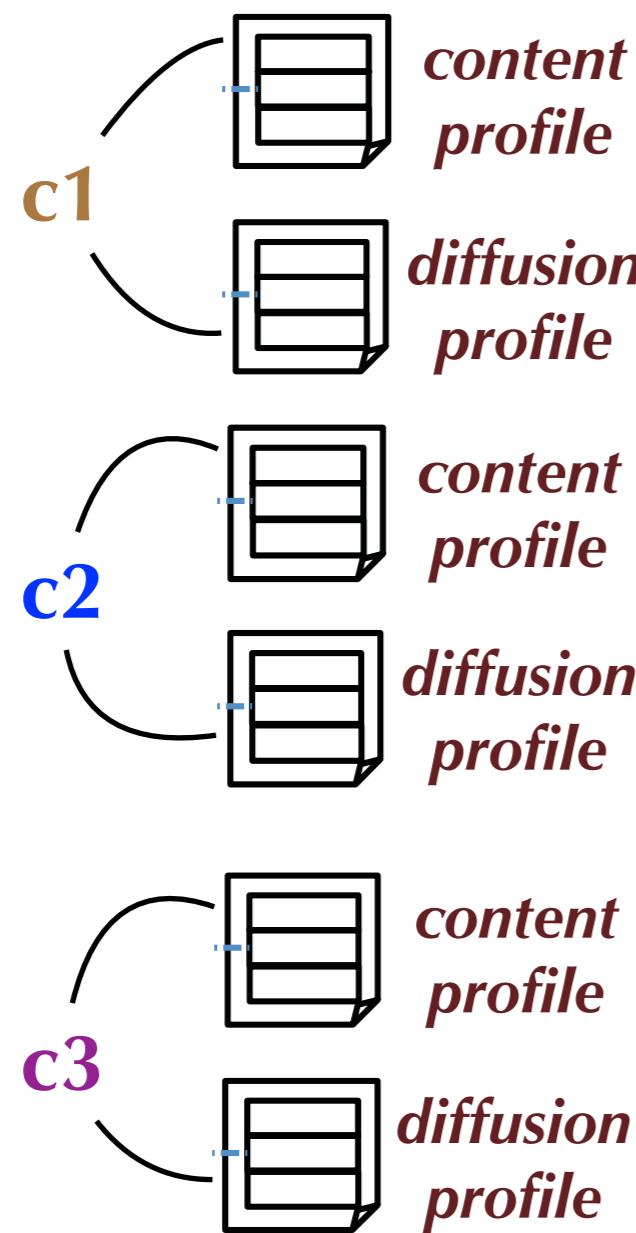


Community Profiling



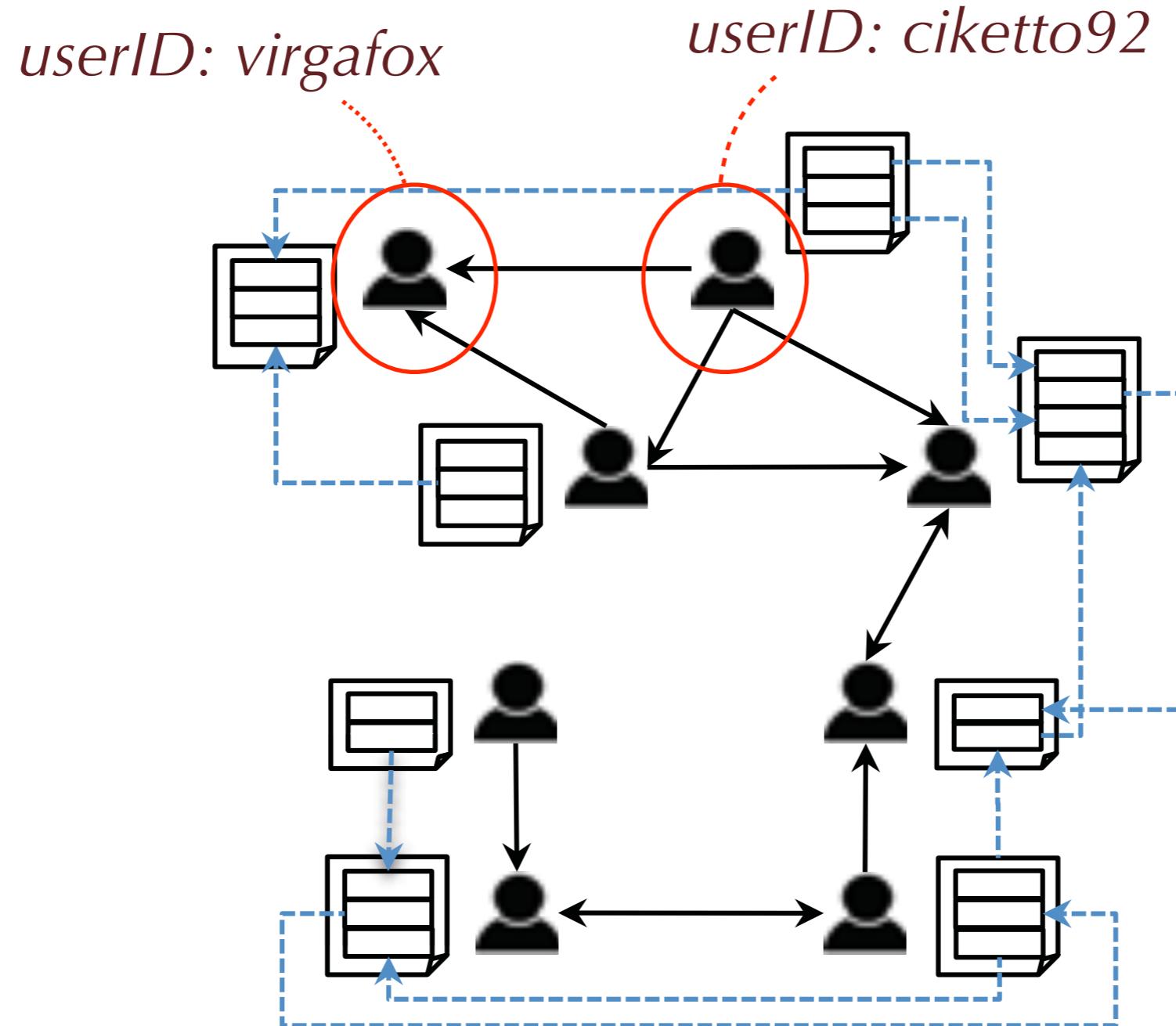


Community Profile Modeling





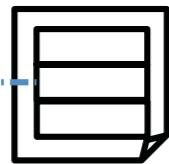
Case study: TWITTER (15.000 users)



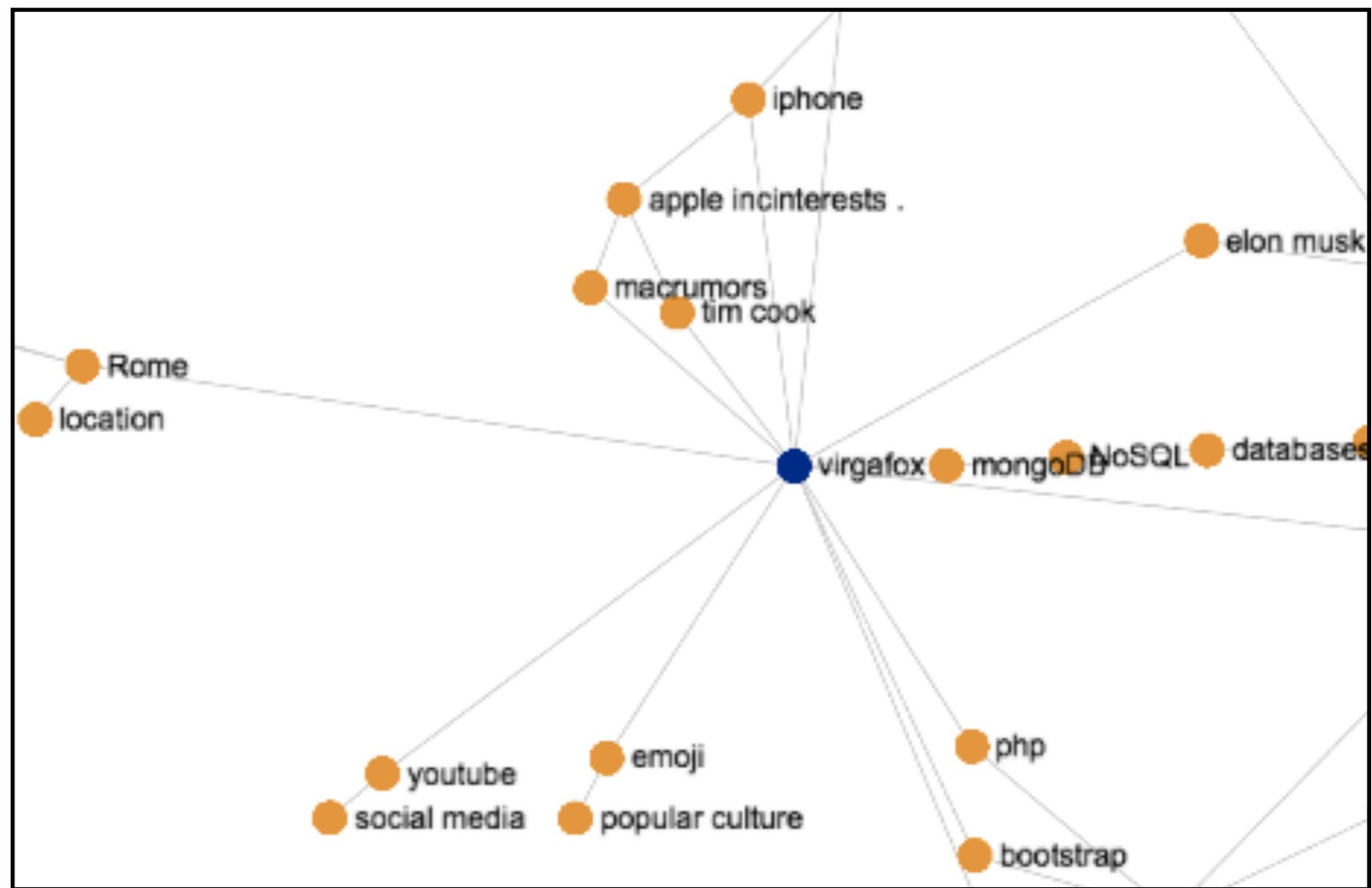
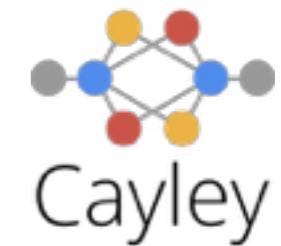


Case study: TWITTER (15.000 users)

userID: virgafox



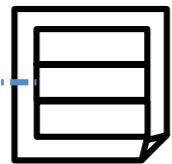
macromeasures



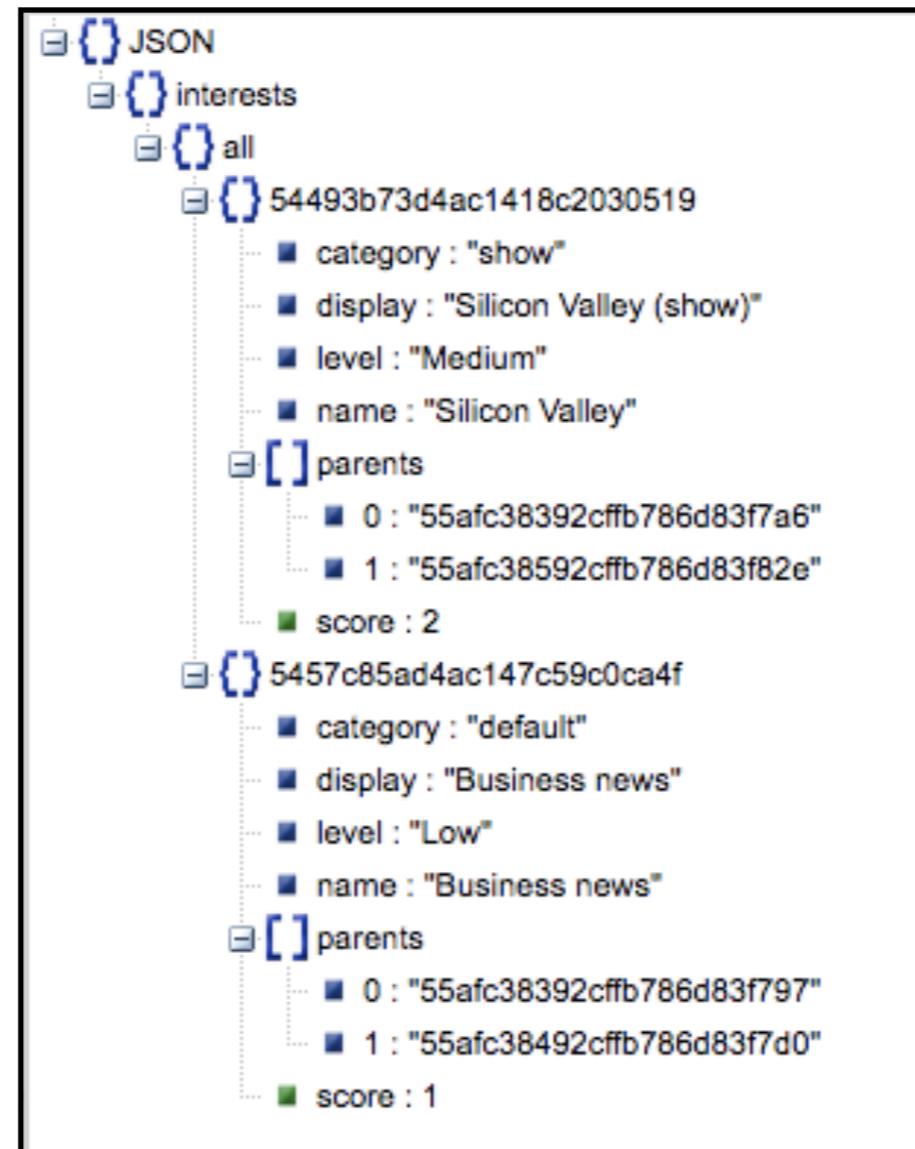


Case study: TWITTER (15.000 users)

userID: virgafox



macromeasures





Proposals: #1

- ❖ **Community Profile Algebra:** define and implement operators for profile management
- ❖ *case study:*
 - profile intersection
 - profile union
 - profile alignment
 - profile difference



Proposals: #2

❖ **Profile Modeling:** represent a user profile in terms of a knowledge graph referring to a different ontology

❖ *case study:*

ontology: Schema.org

social network: Twitter (15000 userIDs)



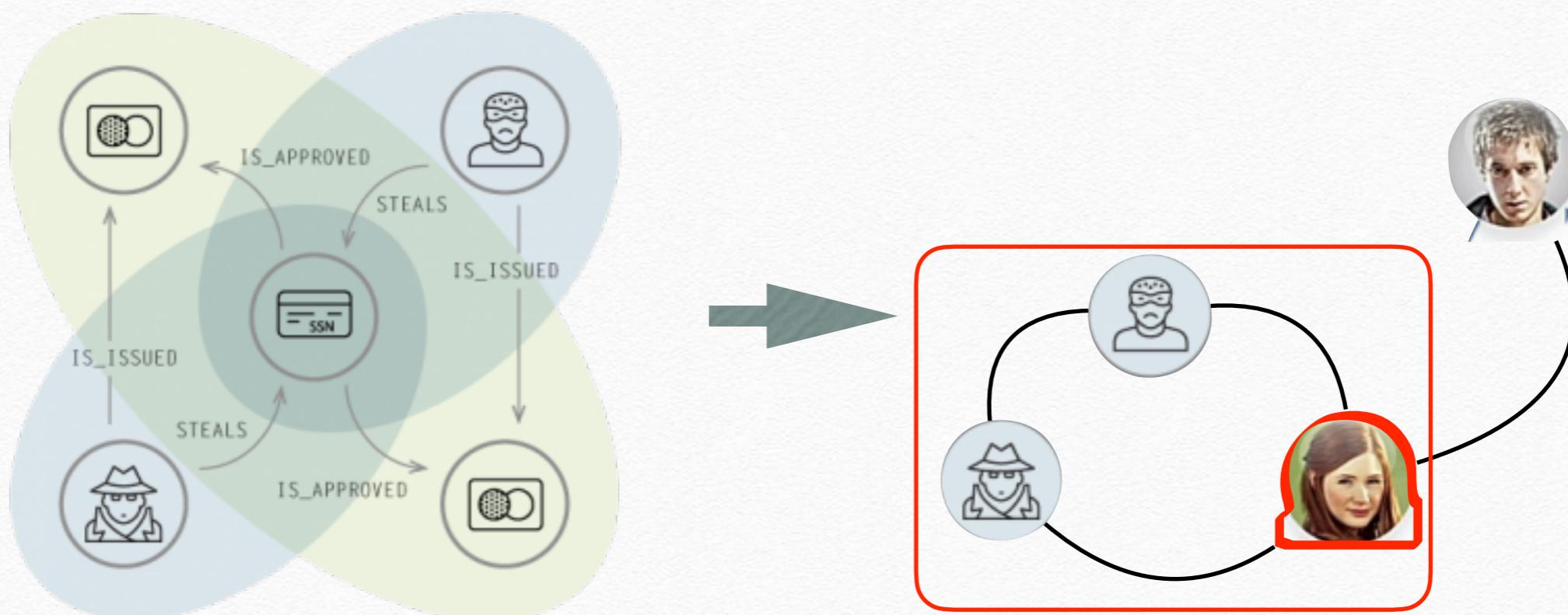
Fraud Detection

17/05/2017 - Big Data 2017



Fraud Detection via Community Detection

- ❖ **Detecting Fraud** : Fraud is an adaptive crime, so it needs special methods of intelligent data analysis to detect and prevent it. These methods exist in the areas of Knowledge Discovery in Databases (KDD), Data Mining, Machine Learning and Statistics
- ❖ **Community detection**: discovering groups in a network having relationships with defrauders

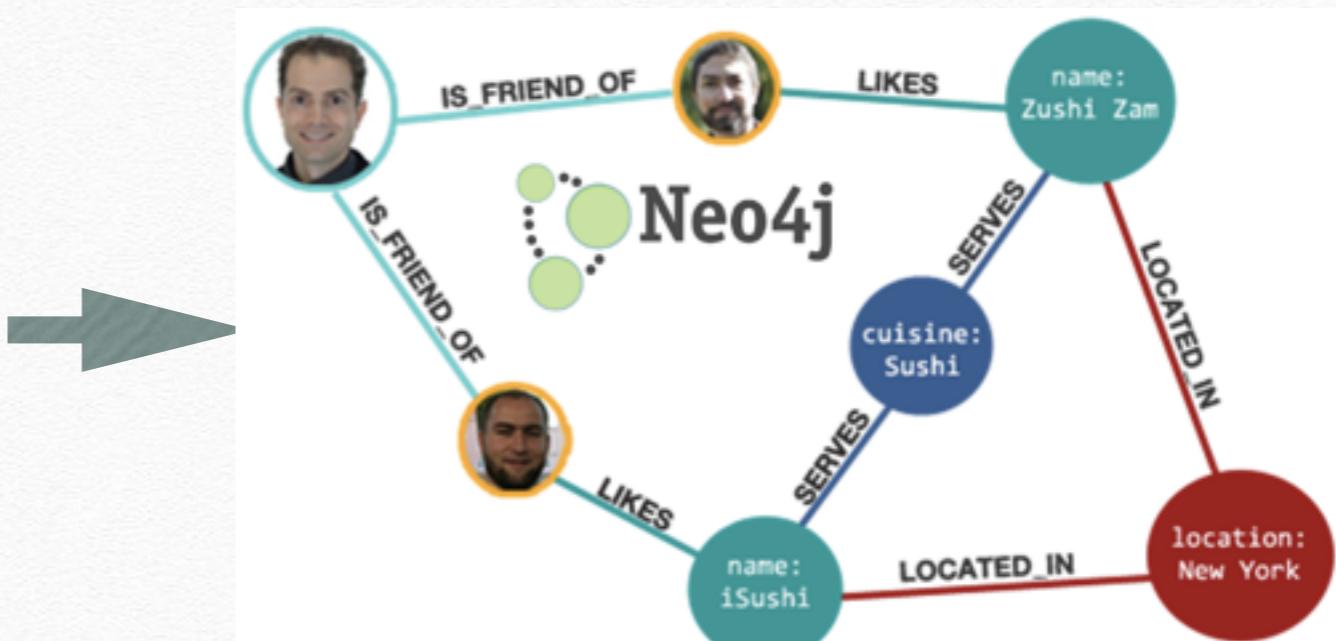




Fraud Detection via Community Detection

- ❖ **CASE STUDY:** Orders coming from an important phone company and requirements about fraud warnings
- ❖ **From Tables to Graph:** orders are structured in a relation DB;

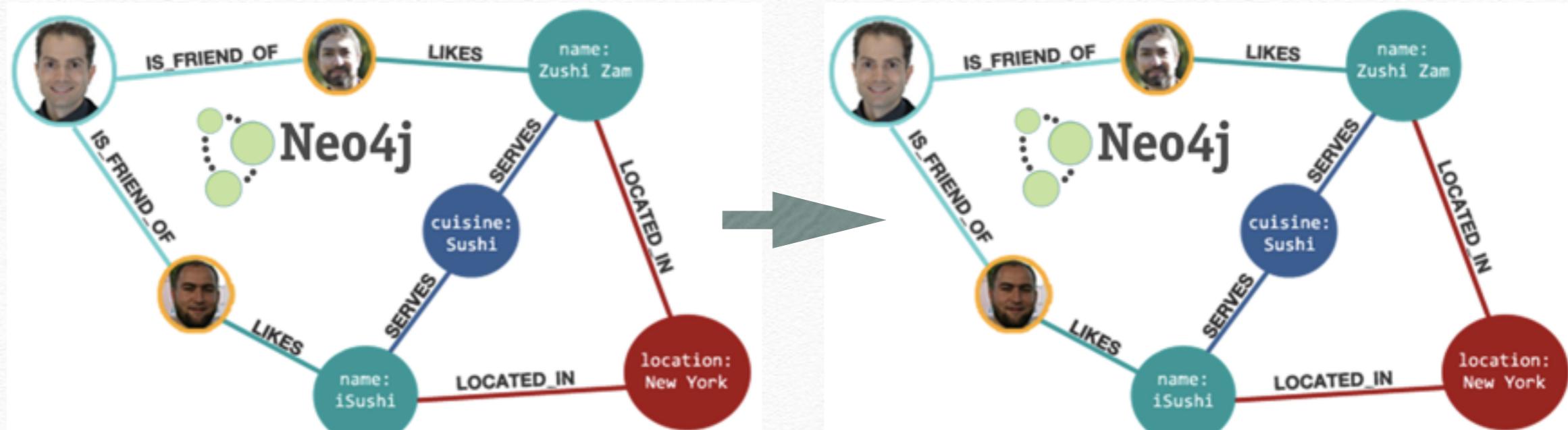
Name	FName	City	Age	Salary
Smith	John	3	35	\$280
Doe	Jane	1	28	\$325
Brown	Scott	3	41	\$265
Howard	Shemp	4	48	\$359
Taylor	Tom	2	22	\$250





Fraud Detection via Community Detection

- ❖ **CASE STUDY:** Orders coming from an important phone company and requirements about fraud warnings
- ❖ **From Graph to Graph:** requirements produce patterns (graph query) retrieving users with abnormal behavior





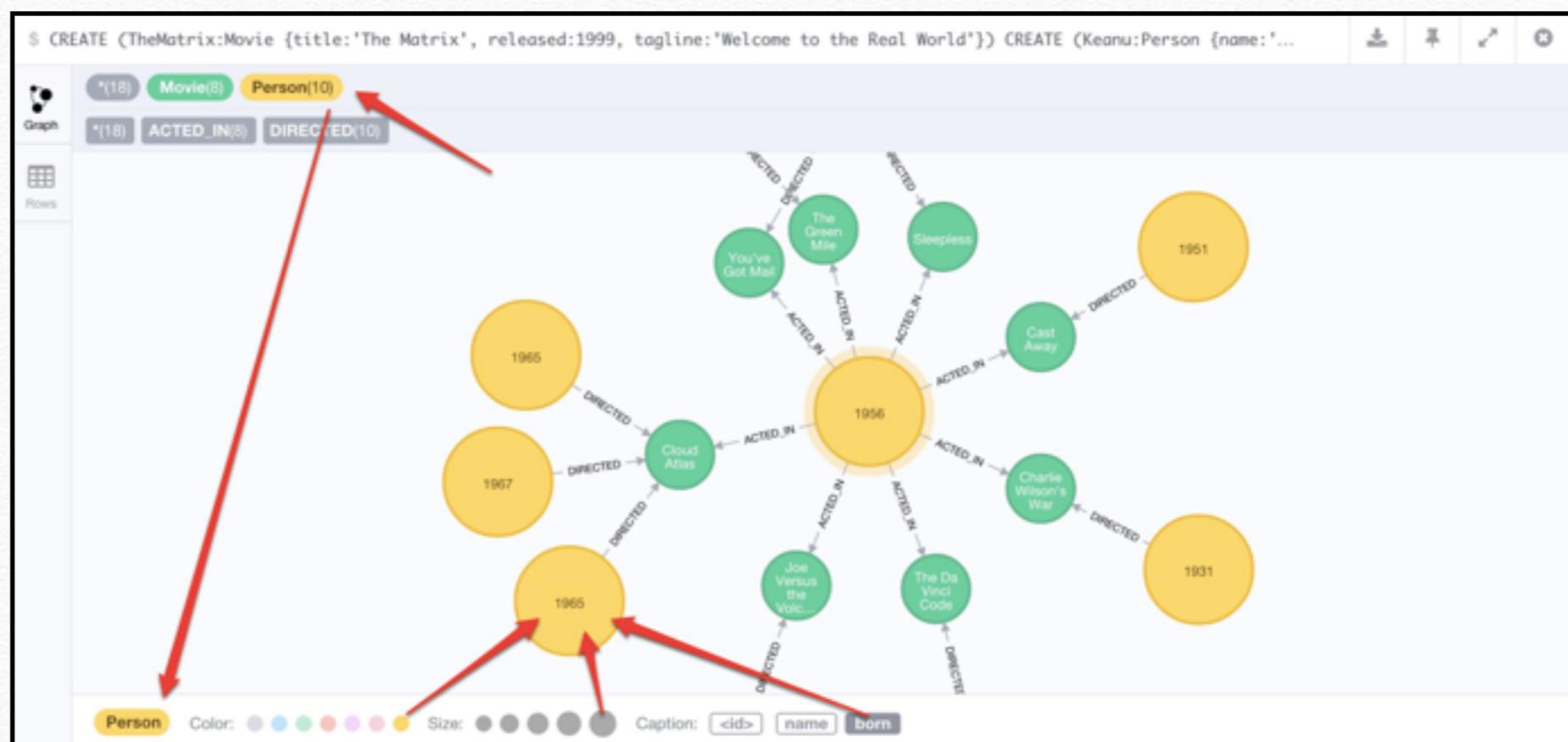
Proposals: #1

❖ **Visual Pattern computation:** define an on-click front-end to produce and to compute patterns on the graphed

❖ *case study:*

graph DBMS: Neo4J

multiple databases on Neo4J





Proposals: #2

❖ **Geo-referential Exploration:** usage of geo-referential front-end for pattern exploration

❖ *case study:*

geo-referential system: NanoCubes

language: Javascript and D3.js

