

# Exploring the Data Wilderness through Examples

*Davide Mottin, Matteo Lissandrini,  
Yannis Velegrakis, Themis Palpanas*

# Who we are



**Davide Mottin**

Graph Mining, Novel Query  
Paradigms, Interactive Methods

<https://mott.in>



**Matteo Lissandrini**

Knowledge Graphs , Novel Query  
Paradigms, Graph Mining

<http://people.cs.aau.dk/~matteo/>



**Yannis Velegrakis**

Big Data Management &  
Analytics, Information Integration

<https://velgias.github.io>



**Themis Palpanas**

Data Series Indexing & Mining,  
Data Management, Data Analytics

<http://www.mi.parisdescartes.fr/~themisp/>

**Slides.** <https://data-exploration.ml/>

# Our book is out!

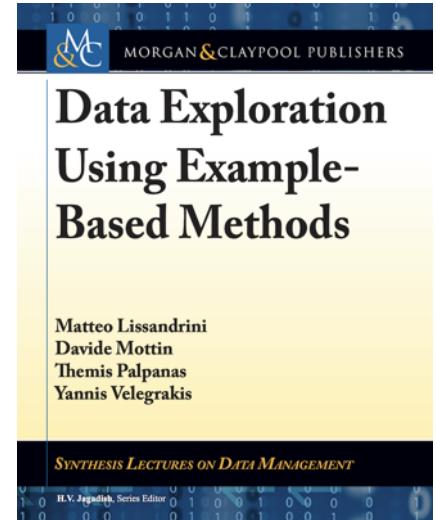
Check our book on Example-based methods!!!

M. Lissandrini, D. Mottin, T. Palpanas, Y. Velegrakis

## “Data Exploration Using Example-based Methods”

*Synthesis Lecture on Data Management*

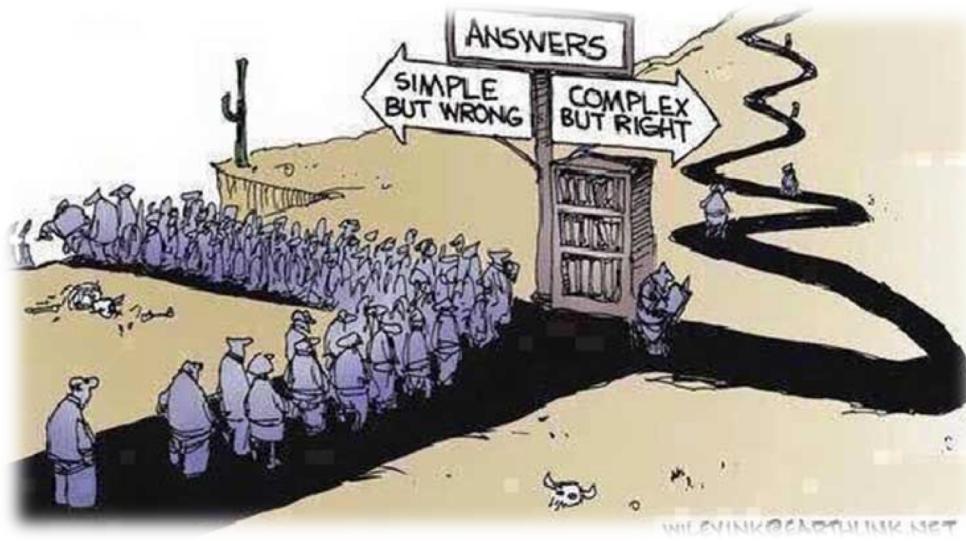
Morgan & Claypool eds. - 2018



# Big data – Easy value?



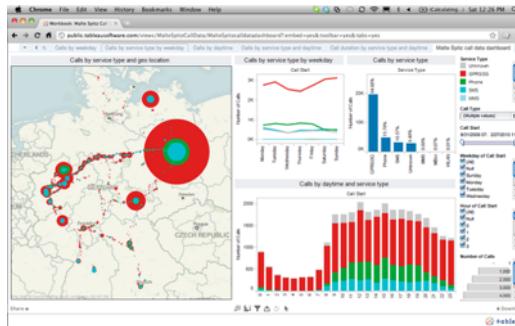
© marktoonist.com



# Exploring



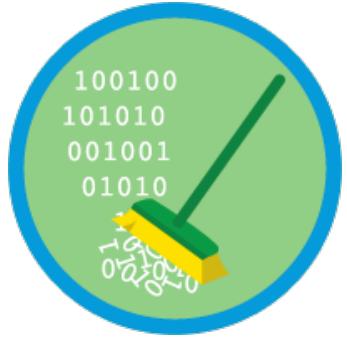
Traditional



On data



# Data exploration



Cleaning and profiling



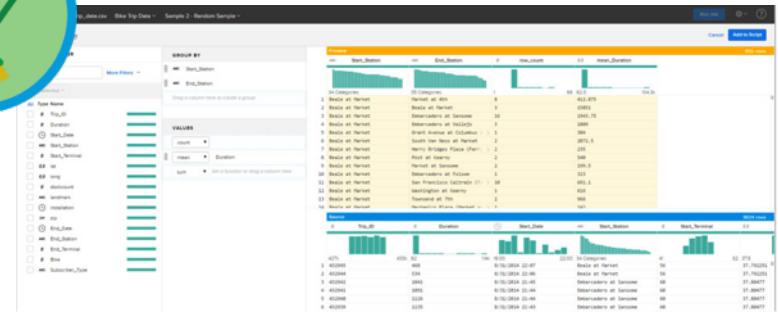
Visualization



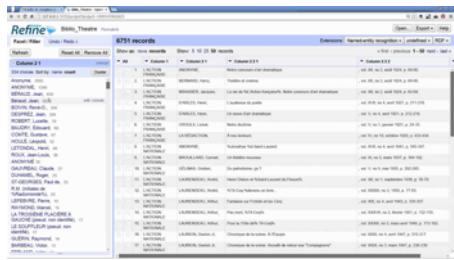
Analysis



# Data exploration software



Trifacta: data preparation



OpenRefine: data preparation and cleanup

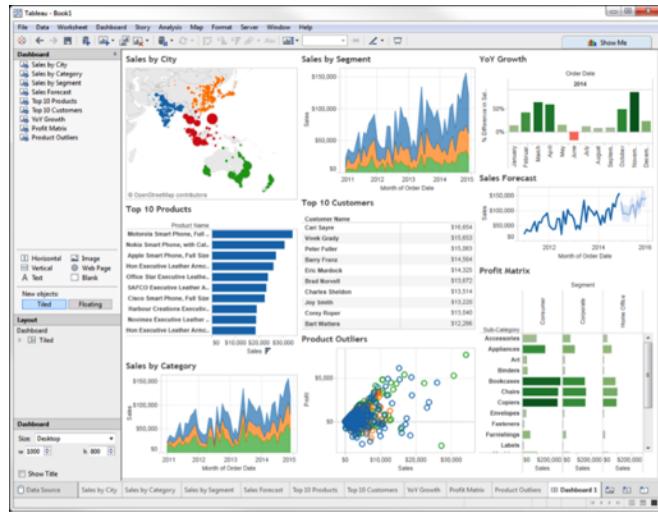


Tableau: analysis and statistics



# Traditional data exploration methods

[Idreos et al., 2015]

Efficiently extracting knowledge from data  
even if we do not know exactly what we are looking for

```
SELECT avg(system-stars)  
FROM Universe  
WHERE system-stars > 10  
GROUP BY galaxy
```



# Declarative Exploratory methods

```
SELECT galaxy_name  
FROM Universe.Galaxy
```

Simple query (exploratory)

Over generic  
100 billions results

```
SELECT g.galaxy_name, SUM(s.stars) as st_s  
FROM Universe.Galaxy AS g  
JOIN Universe.Systems AS s  
ON g.galaxy_name = s.galaxy_name  
WHERE  
    g.st_s > 100B  
    AND diameter > 100k AND diameter > 180k  
    AND has_black_hole = TRUE  
GROUP BY g.galaxy_name
```

Complex query  
(for data experts)

Specific  
Few results



# Examples as Exploratory Methods



Answers



# Historical perspective: Query-by-example

[Zloof et al. 1975]

Specify a query by example tables, or skeletons.

Incomplete values

Name	Stars	Diameter	Black_hole	Color	Life
P._	> 10B	>100k	TRUE	*	*
*	*	<180k	*	*	*

Unspecified values

Value conditions

Intuitive interface for simple queries

SQL not required

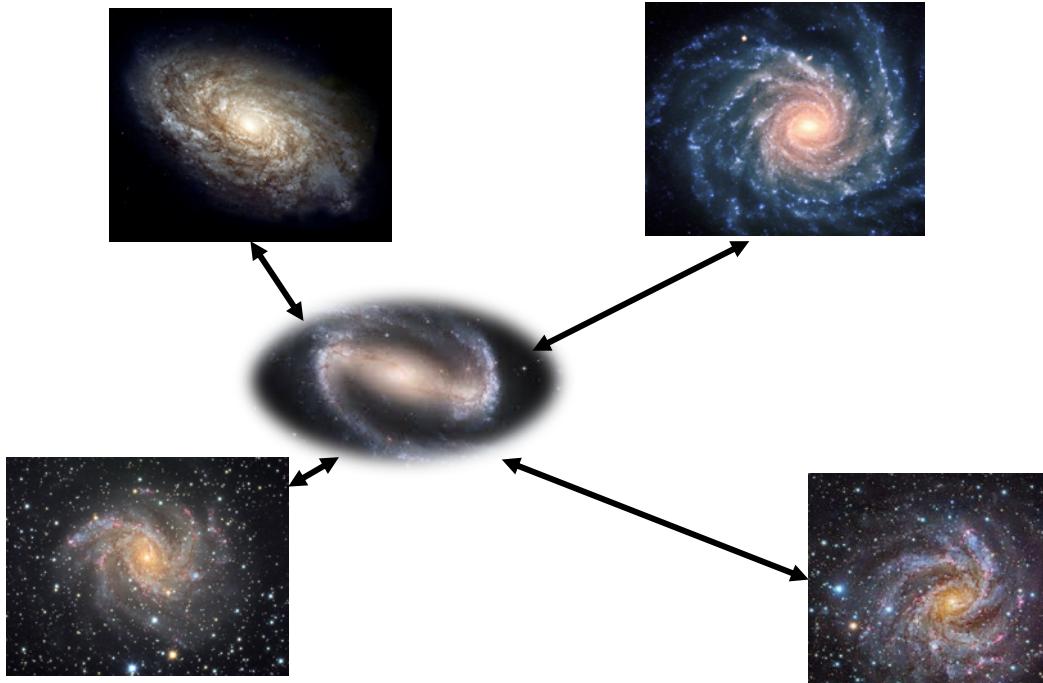
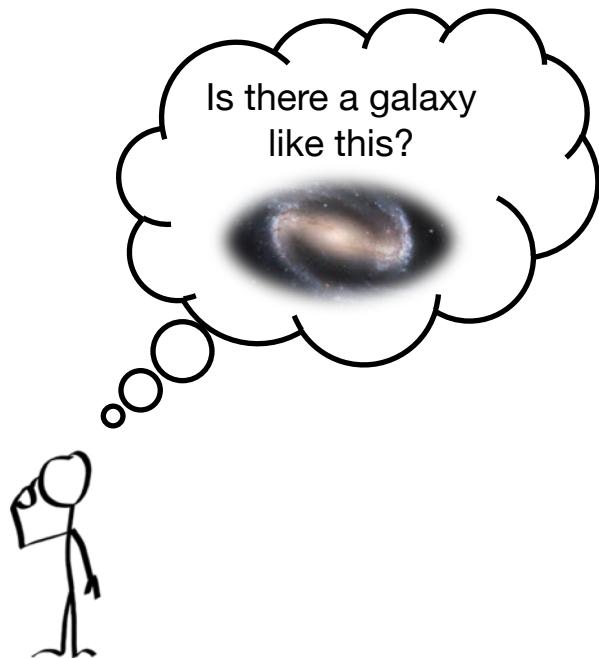
Restricted to SQL syntax but not explicitly

Not example-based



# Similarities are the key ...

If we knew how similar each item is with respect to any other for **each user**, we would know the answer to



# Similarities are the key ...

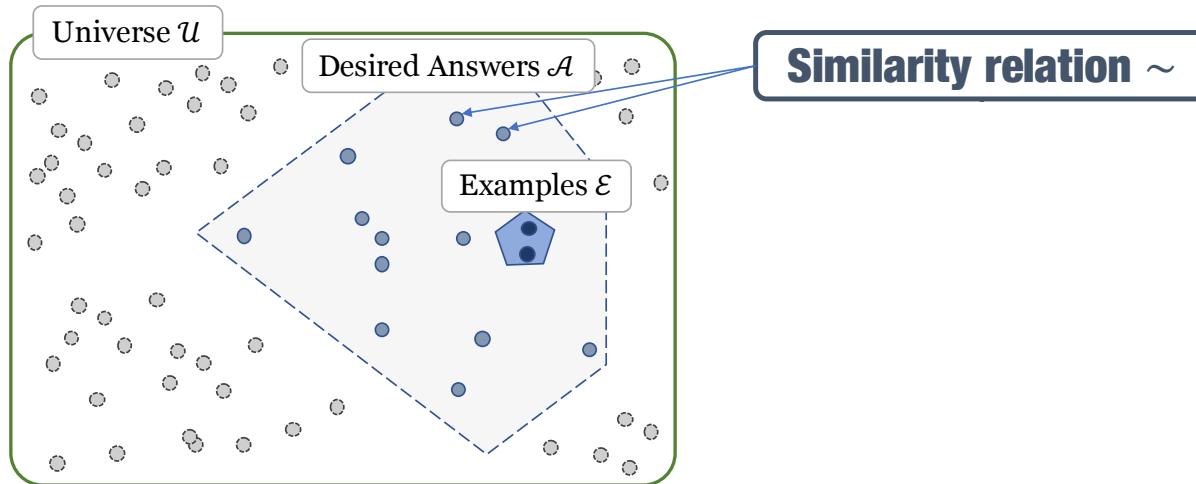
We define:

A universe  $\mathcal{U}$  of items

A similarity among items  $\sim$

A set of **input** examples  $\mathcal{E}$

A set of **output** user desired answers  $\mathcal{A}$



# The example-based problem

Given

a set of examples  $\mathcal{E}$  from a universe  $\mathcal{U}$

Find

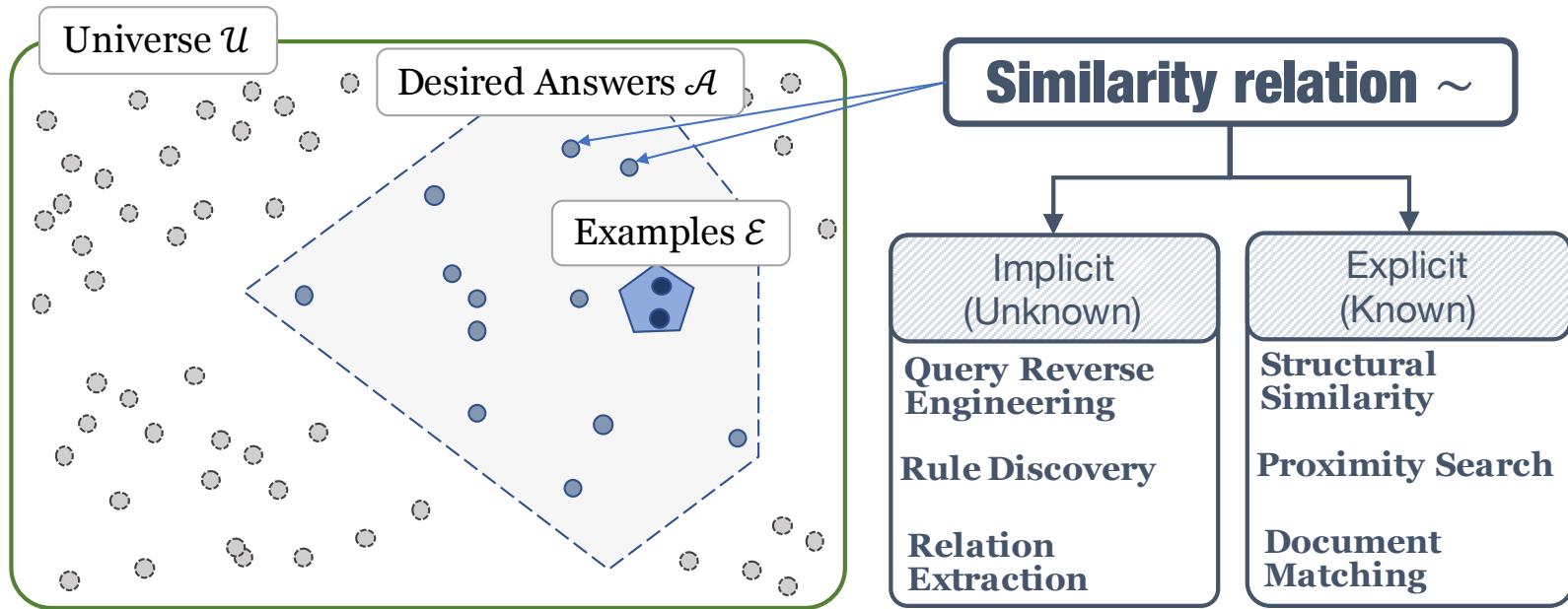
a similarity  $\sim$  such that

1.  $\mathcal{E}$  is part of the answers  $\mathcal{A}$  partially or totally
2. The answers in  $\mathcal{A}$  are the **most similar** to the examples in  $\mathcal{E}$  according to  $\sim$

How do we find  $\sim$  for each user?  
Do we need to know exactly  $\sim$ ?



# Example-based methods



# Tutorial's goals

- Exploratory methods using examples
- Algorithms for retrieving data without using query languages
- Interactive methods and user-in-the-loop feedback
- Machine learning for adaptive, online methods

But NOT

- Declarative query methods
- User interfaces and visualization
- Optimizations for fast data access
- Dynamic data



# Tutorial structure



Relational databases (50 min)



Textual data (30 min)



Graph and networks (50 min)

Machine learning  
(25 min)

Challenges and Remarks (15 min)

# Questions or comments? Please share!

<https://j.mp/ExploreSIGMOD>



# Example-based methods

- Reverse engineering queries
- Example-driven schema mapping
- Interactive data repairing



- Entity extraction by example text
- Web table completion using examples
- Search by example



- Community-based Node-retrieval
- Entity Search
- Path and SPARQL queries
- Graph structures as Examples



# Where we are



Relational databases



Textual data

Graphs and networks

Machine learning

Challenges and Remarks

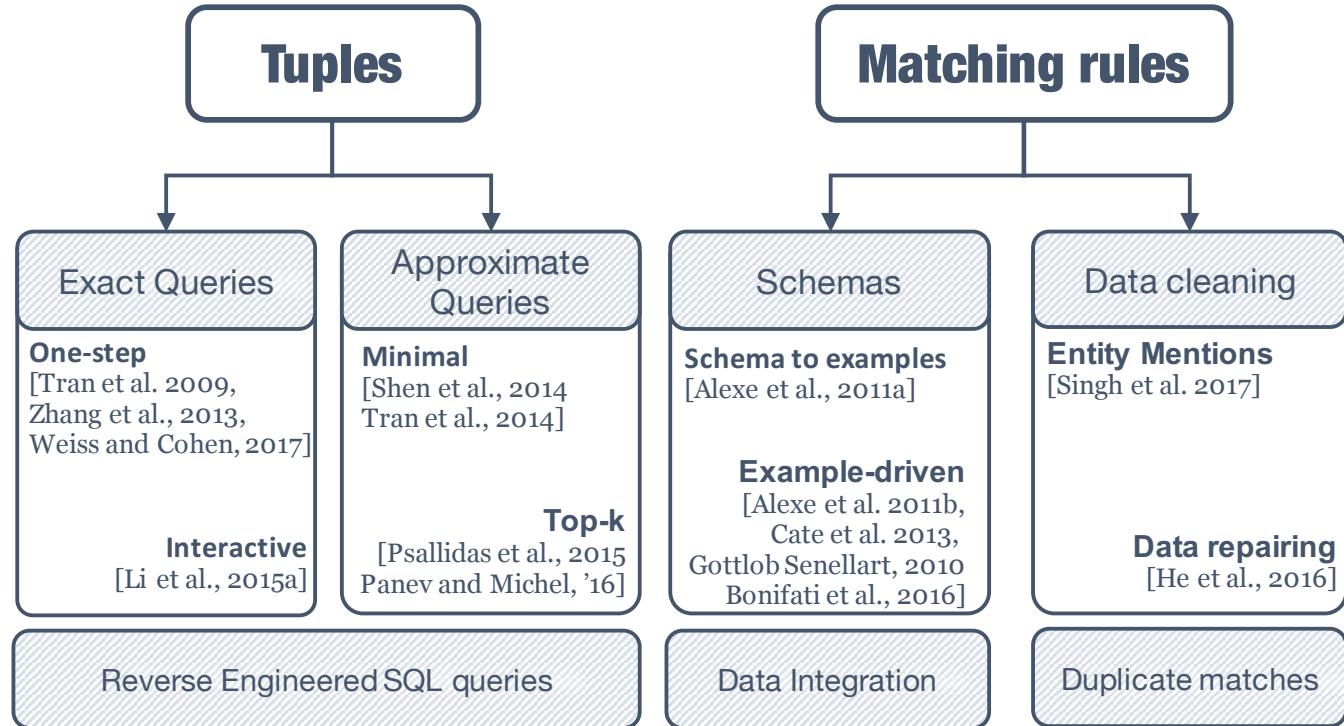
# What tasks?

SEARCHING FOR

BY FOCUSING ON

APPLYING

PRODUCES



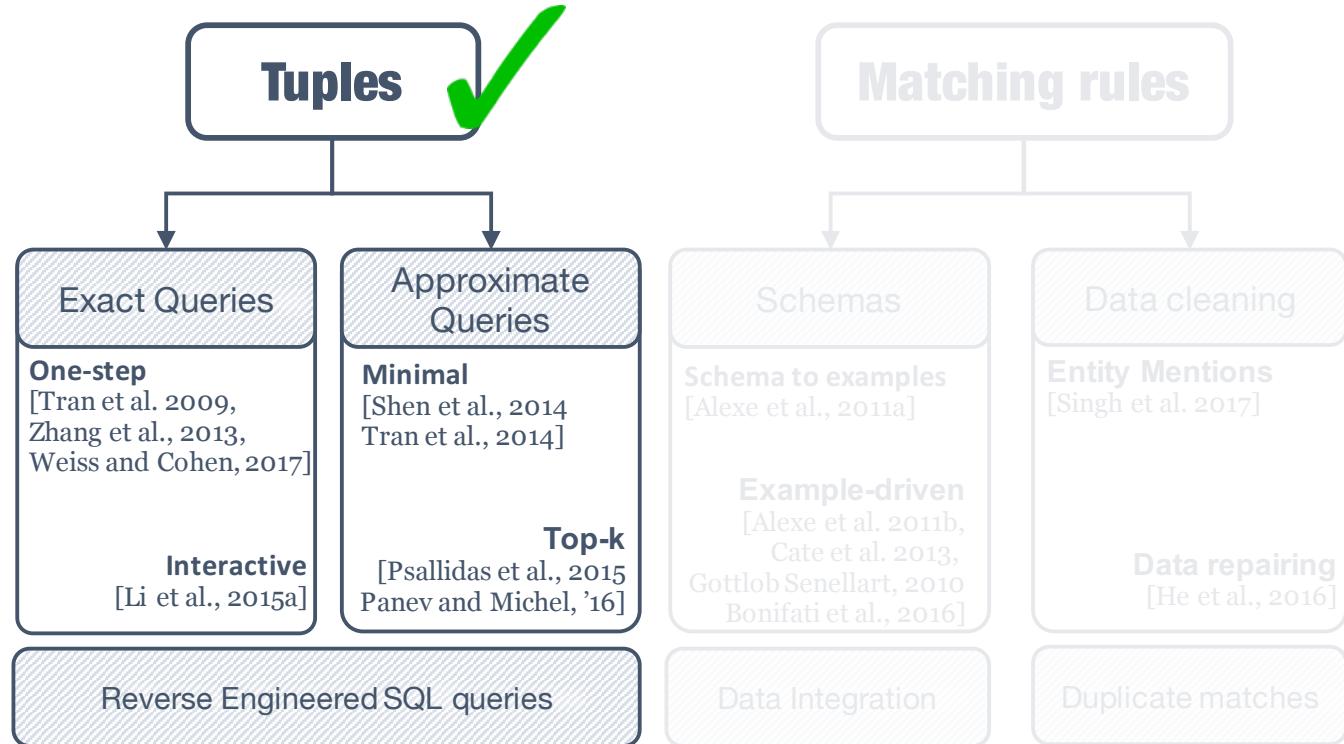
# Searching for ...

SEARCHING FOR

BY FOCUSING ON

APPLYING

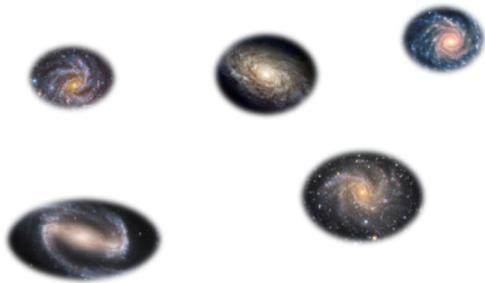
PRODUCES



# Reverse engineering queries (REQ)

Given a set of examples, find the query that generated that set of tuples

Example tuples

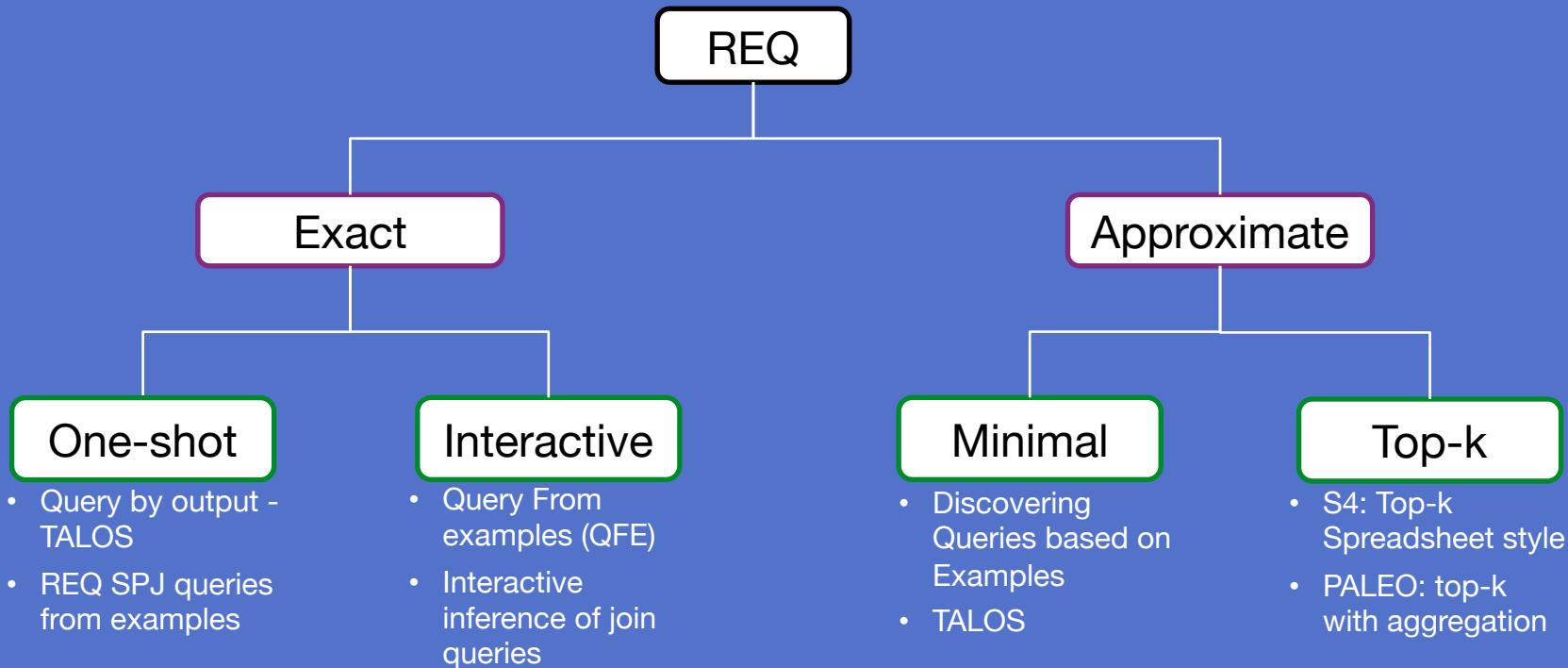


```
SELECT g.galaxy_name, SUM(s.stars) AS st_s  
FROM Universe.Galaxy AS g  
JOIN Universe.System AS s  
ON g.galaxy_name = s.galaxy_name  
WHERE  
    g.st_s > 100B  
    AND diameter > 100k AND diameter < 180k  
    AND has_black_hole = TRUE  
GROUP BY g.galaxy_name
```

How do you find such queries?

```
SELECT galaxy_name  
FROM Universe.Galaxy
```

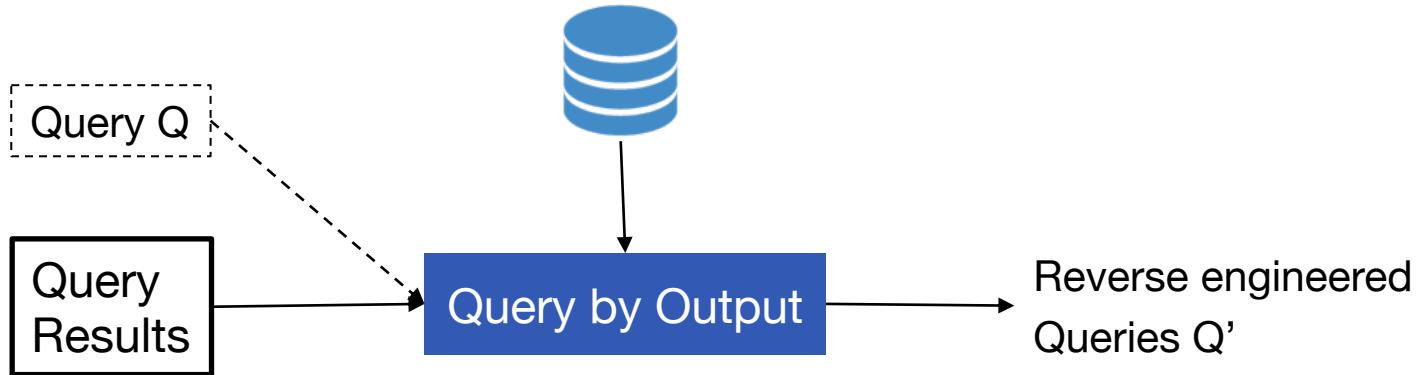
# Reverse engineering queries (REQ)



# Query by Output - TALOS

[Tran et al. 2013]

Main idea: Find the set of queries that exactly return a set of examples



Two queries  $Q$  and  $Q'$  are instance equivalent on a database  $D$ , if the results of  $Q$  are the same of the results of  $Q'$

# How many reverse queries?

Master						
	name	bat	throw	stint	HR	team
$t_1$	A	L	R	2	40	PIT
$t_2$	A	L	R	2	50	MT1
$t_3$	C	R	L	2	35	CHA
$t_4$	D	L	R	3	30	PIT
$t_5$	B	R	R	1	73	PIT
$t_6$	B	R	R	1	40	PIT
$t_7$	E	R	R	3	60	CHA

$r_1 \begin{array}{|c|c|} \hline B & PIT \\ \hline \end{array}$   
 $r_2 \begin{array}{|c|c|} \hline E & CHA \\ \hline \end{array}$

$Q(D)$

What queries generated  $Q(D)$ ?

Q1 = SELECT name, team FROM Master WHERE bat = 'R' AND throw = 'R'

Q2 = SELECT name, team FROM Master WHERE bat = 'R' AND weight > 35

Q3 = SELECT name, team FROM Master WHERE bat = 'R' AND stint <> 2

...

Equivalent  
Queries



# TALOS

[Tran et al. 2013]



pID	<b>name</b>	country	weight	bats	throws
P1	A	USA	85	L	R
P2	B	USA	72	R	R
P3	C	USA	80	R	L
P4	D	Germany	72	L	R
P5	E	Japan	72	R	R

(a) Master

pID	year	stint	<b>team</b>	HR
P1	2001	2	PIT	40
P1	2003	2	ML1	50
P2	2001	1	PIT	73
P2	2002	1	PIT	40
P3	2004	2	CHA	35
P4	2001	3	PIT	30
P5	2004	3	CHA	60

(b) Batting

<b>team</b>	year	rank
PIT	2001	7
PIT	2002	4
CHA	2004	3

(c) Team

Input

B	PIT
E	CHA

Master

Batting

**Join table**  
 $J = \text{Master} \bowtie \text{Batting} \bowtie \text{Team}$

$t_1$	name	bat	throw	stint	HR	team
$t_1$	A	L	R	2	40	PIT
$t_2$	A	L	R	2	50	MT1
$t_3$	C	R	L	2	35	CHA
$t_4$	D	L	R	3	30	PIT
$t_5$	B	R	R	1	73	PIT
$t_6$	B	R	R	1	40	PIT
$t_7$	E	R	R	3	60	CHA



Join graph computation



# TALOS

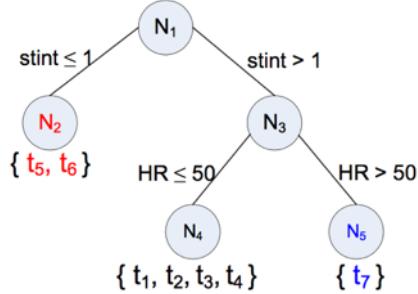
[Tran et al. 2013]

	B	PIT
	E	CHA
<i>t</i> <sub>1</sub>	A	L R 2 40 PIT
<i>t</i> <sub>2</sub>	A	L R 2 50 MT1
<i>t</i> <sub>3</sub>	C	R L 2 35 CHA
<i>t</i> <sub>4</sub>	D	L R 3 30 PIT
<i>t</i> <sub>5</sub>	B	R R 1 73 PIT
<i>t</i> <sub>6</sub>	B	R R 1 40 PIT
<i>t</i> <sub>7</sub>	E	R R 3 60 CHA

X X X ✓ ✓ ✓

Idea: treat the problem as a binary classification

1. **Strict**: all tuples must be captured
2. **At-Least-one**: one tuple for example must be captured



Decision tree

$$Gini(S) = 1 - (f_+^2 + f_-^2)$$

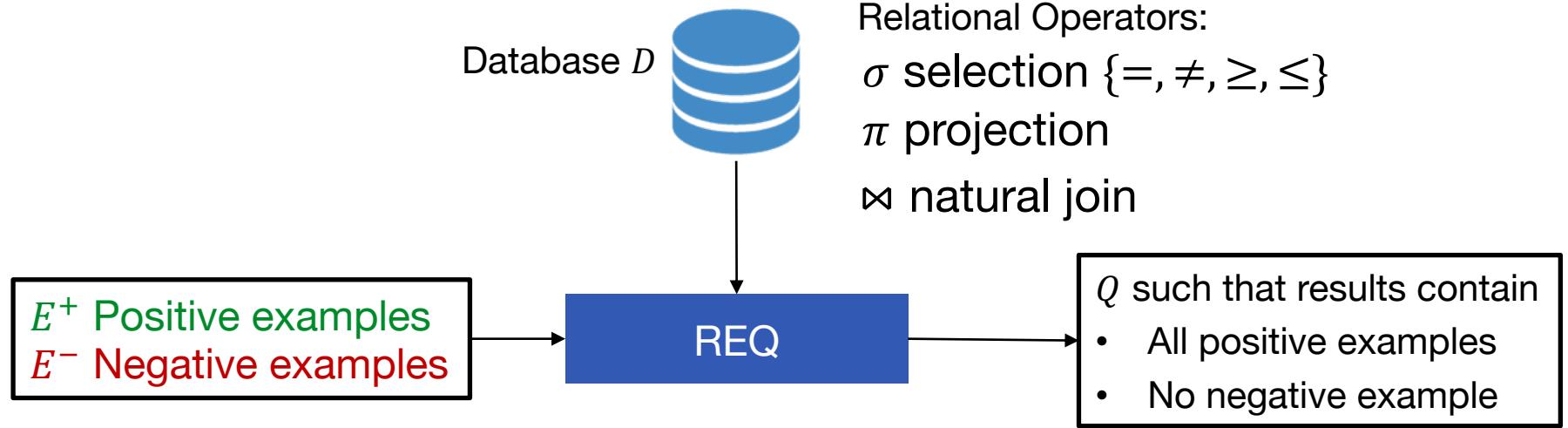
Positive and negative tuples in S

$$Gini(S_1, S_2) = \frac{(|S_1|Gini(S_1) + |S_2|Gini(S_2))}{|S_1| + |S_2|}$$



# How complex is exact REQ?

[Weiss et al., 2017]



How difficult is to find:  
A bounded size  $Q$ ? an unbounded  $Q$ ?

# Complexity - No parameters

[Weiss et al., 2017]

Operator	Unbounde d Queries	Bounded Queries
$\pi$	P	P
$\bowtie$	P	NPC
$\sigma$	P	NPC
$\sigma, \bowtie$	P	NPC
$\pi, \sigma$	NPC	NPC
$\sigma, \bowtie$	DP	DP
$\pi, \sigma, \bowtie$	DP	DP

Only projections: Easy

Unbounded selections: Easy

Unbounded selections: HARD

Combination of operators:  
**HARD!!!**



# Unbounded Select

[Weiss et al., 2017]

	A	B	C	D	E
<input checked="" type="checkbox"/>	1	2	3	4	5
<input checked="" type="checkbox"/>	1	3	2	3	4
	2	4	4	1	3
	5	3	2	4	2
<input checked="" type="checkbox"/>	4	2	3	1	2
	2	2	4	3	2
<input checked="" type="checkbox"/>	1	1	2	1	5
<input checked="" type="checkbox"/>	1	5	4	2	3

Possible queries?

- $A = 1 \text{ AND }$
- $B \geq 1 \text{ AND } B \leq 5 \text{ AND }$
- $C \geq 2 \text{ AND } C \leq 4 \text{ AND }$
- $D \geq 1 \text{ AND } D \leq 4 \text{ AND } D \neq 4$
- $E \geq 3 \text{ AND } E \leq 5 \text{ AND } E \neq 4$



# Bounded select

Reduction from  
Set Cover

NP-C

INPUT: Database D, Examples E, Query size k

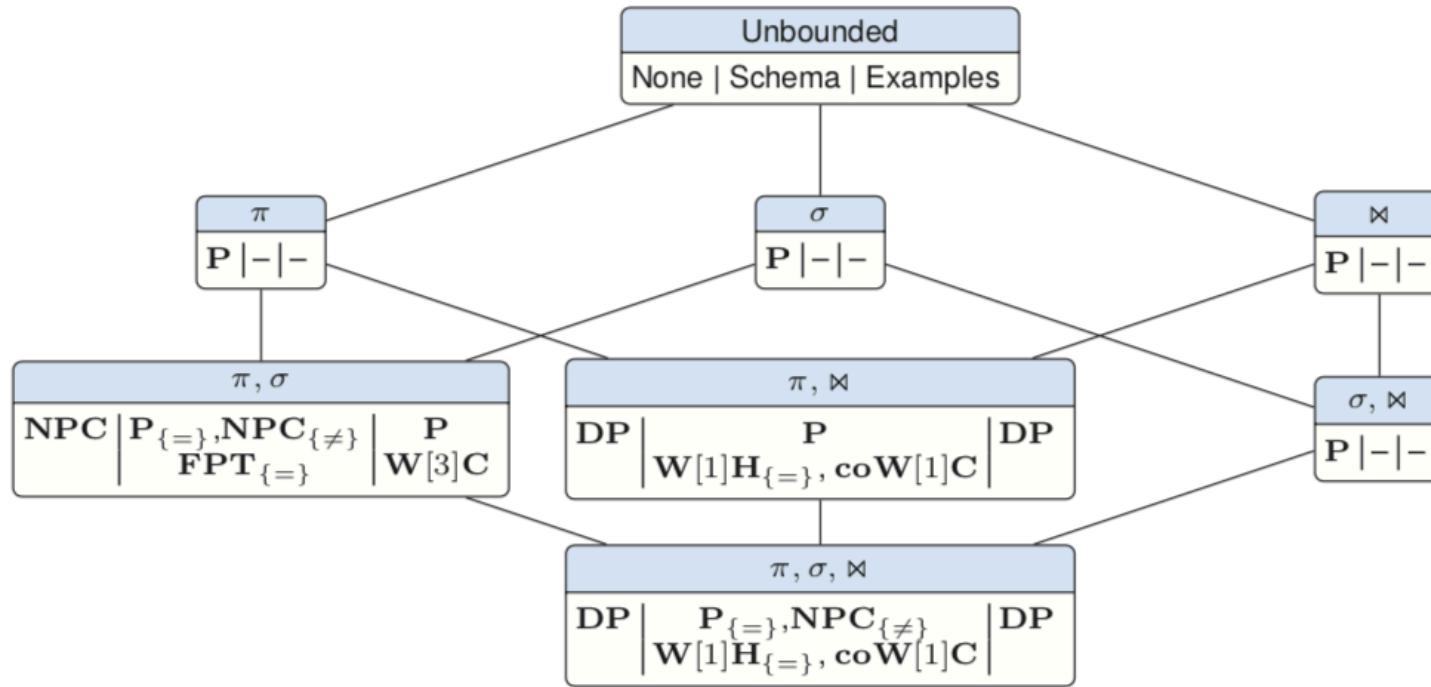
OUTPUT: Does there exist a query satisfying D and E, of size at most k?

$$U = \{1,2,3,4,5\} \quad S = \{ \{1,2,3\}, \{2,4\}, \{3,4\}, \{4,5\} \}$$

	$S_1$	$S_2$	$S_3$	$S_4$
✗	1	0	0	0
✗	1	1	0	0
✗	1	0	1	0
✗	0	1	1	1
✗	0	0	0	1
✓	1	1	1	1



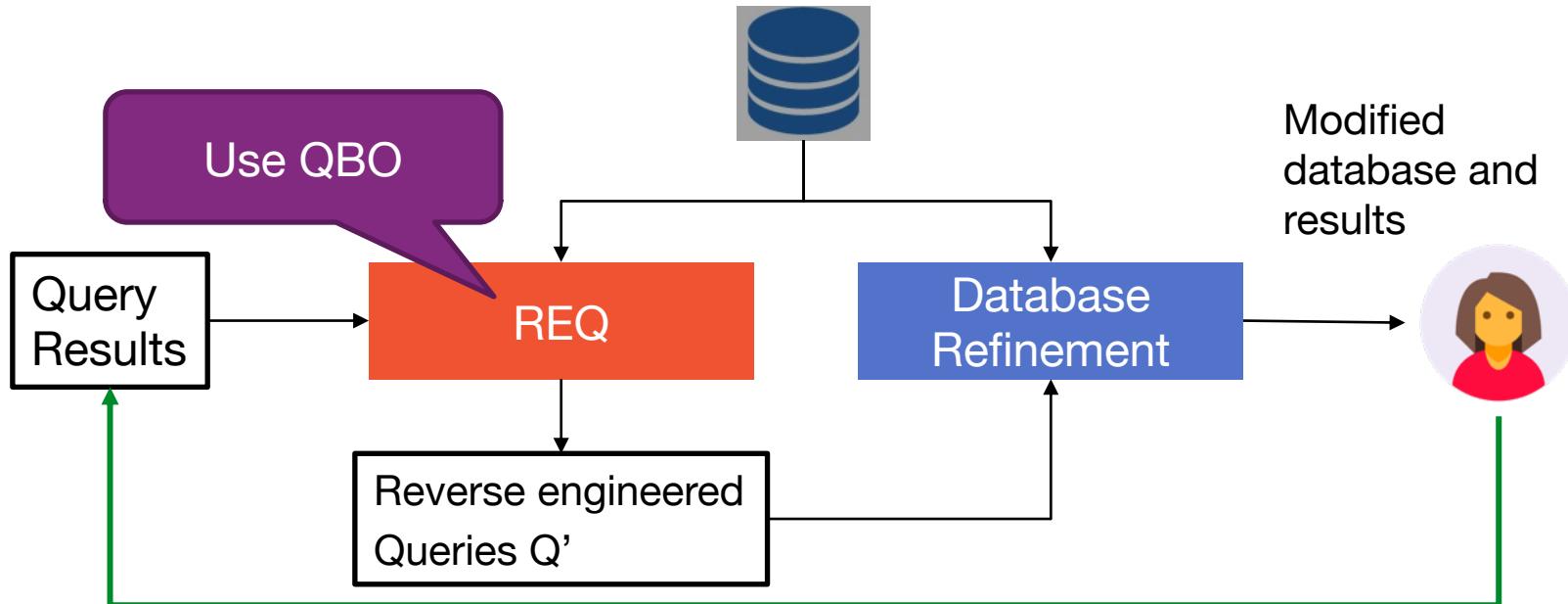
# Complexity - Parameters



# Interactive REQ – Query from Examples

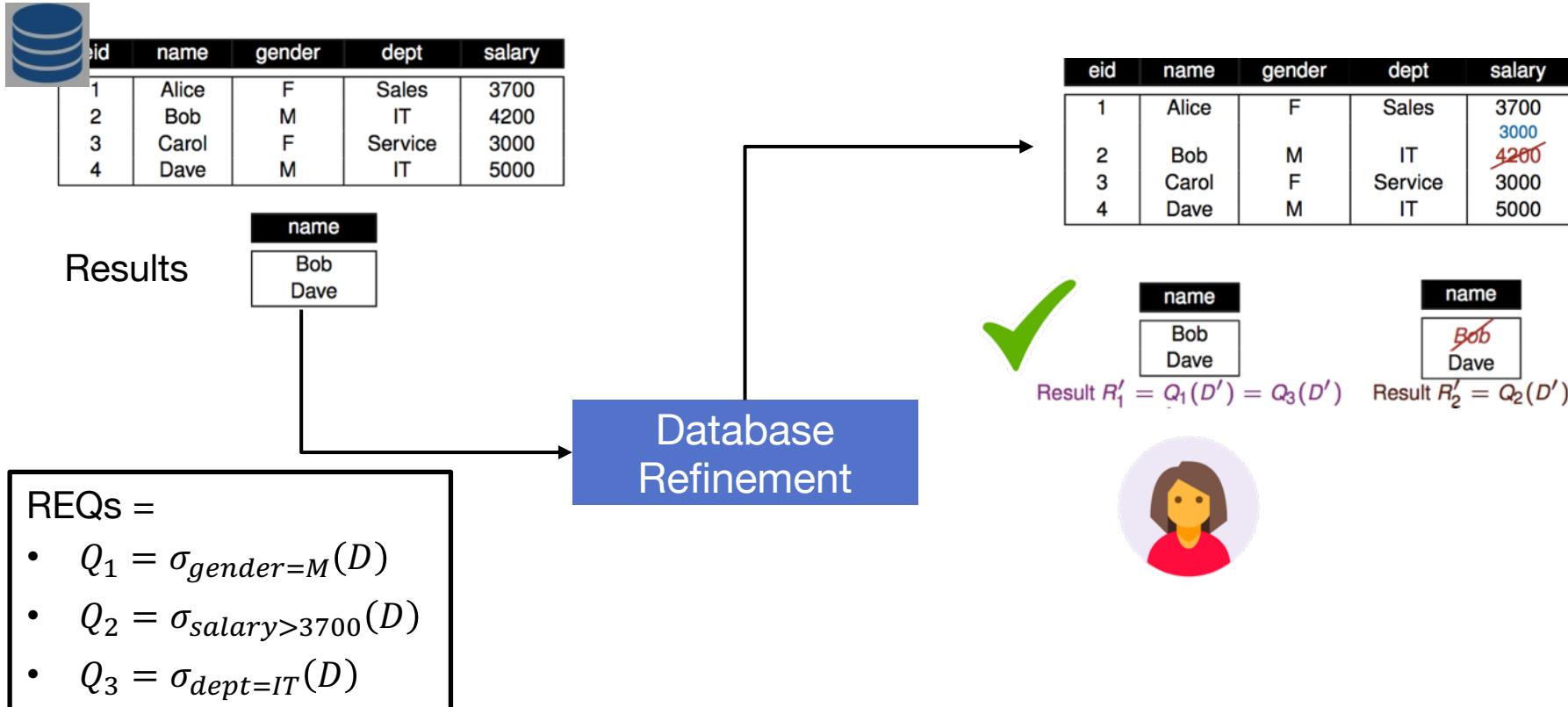
[Li et al., 2015]

Main idea: Interactively remove candidate queries proposing a new set of query results from a modified database



# Database Refinement

[Li et al., 2015]



# Cost model

[Li et al., 2015]

$$cost(D') = \boxed{edit(D, D') + \beta \cdot n} + \boxed{\sum_{i=1}^k edit(R, R_i)} + \boxed{N \cdot \frac{edit(D, D')}{\mu} + \beta} + \boxed{\frac{2}{k} \sum_{i=1}^k edit(R, R_i)}$$

Number of modified tables      Number of new result sets

DB cost      Results cost      Effort to examine  $D'$       Effort to examine new results

Current cost      Residual cost

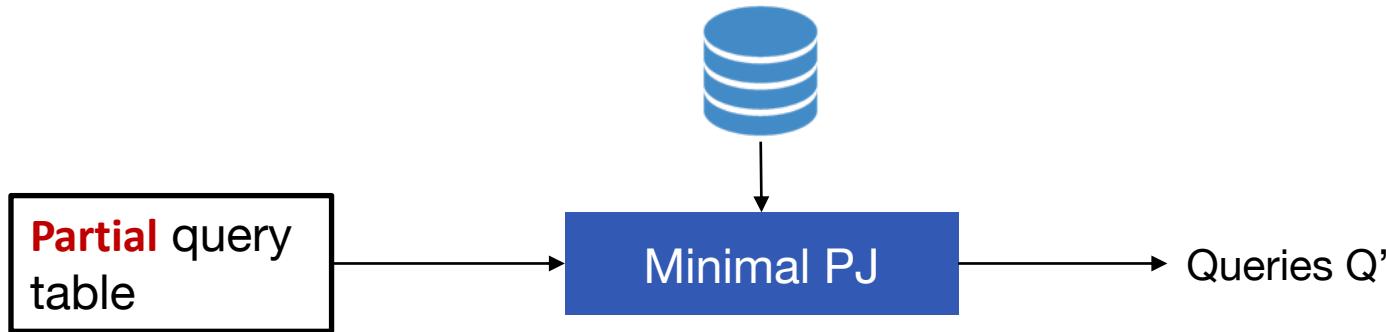
Main idea: Find a refined db  $D'$  and results  $R_1, \dots, R_k$  with:

1. Minimum number of results  $k$
2. Minimum differences in the database
3. The query are balanced (less interactions)

# Minimal Project Join REQ

[Shen et al., 2014]

Main idea: Find the set of queries that approximately return a set of examples



	A	B	C
1	Mike	ThinkPad	Office
2	Mary	iPad	
3	Bob		Dropbox

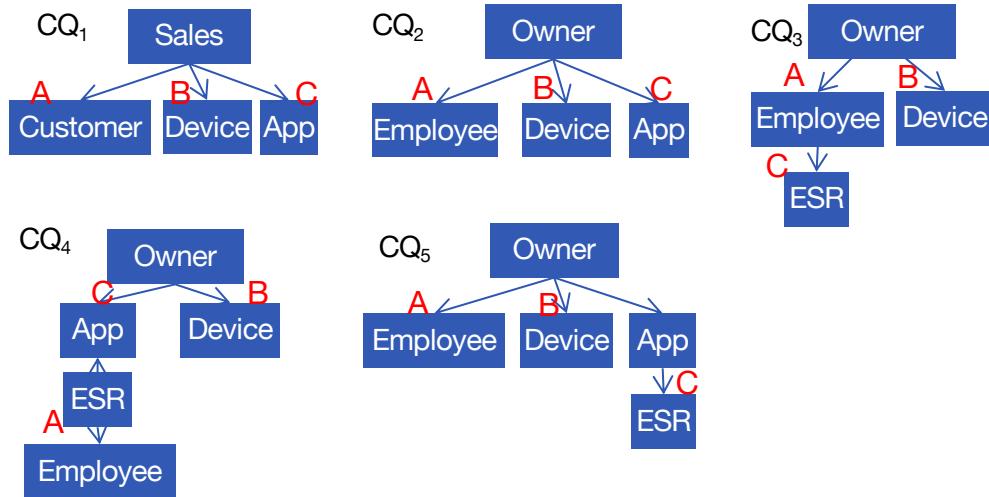
- **valid**: every tuple is present in query results
- **minimal**: any removal in query tree gets to an invalid query



# Candidate Query Generation

[Shen et al., 2014]

- Use candidate network generation algorithm  
(Hristidis 2002)



	A	B	C
1	Mike	ThinkPad	Office
2	Mary	iPad	
3	Bob		Dropbox

1. Generate join tree  $J$
2. Generate mapping  $\phi$
3. Check minimal:
  - Every leaf node contains a column that is mapped by an input column



# Validity verification

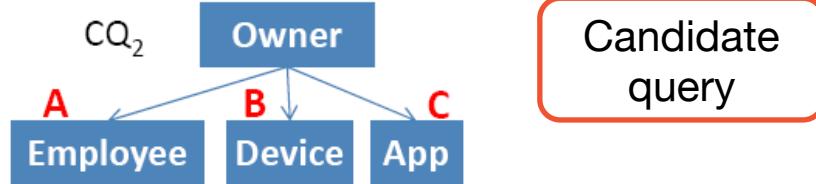
[Shen et al., 2014]

Naïve: check all candidate queries singularly if they return ALL examples

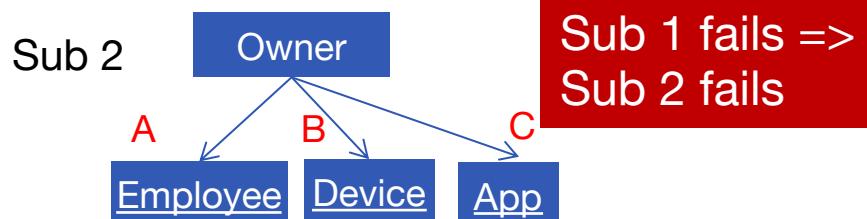
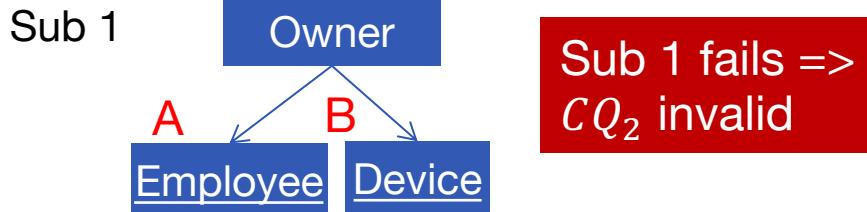
Better: exploit substructures in candidate queries for pruning

Best: adaptively select the substructures to have the min number of evaluations

NP-hard



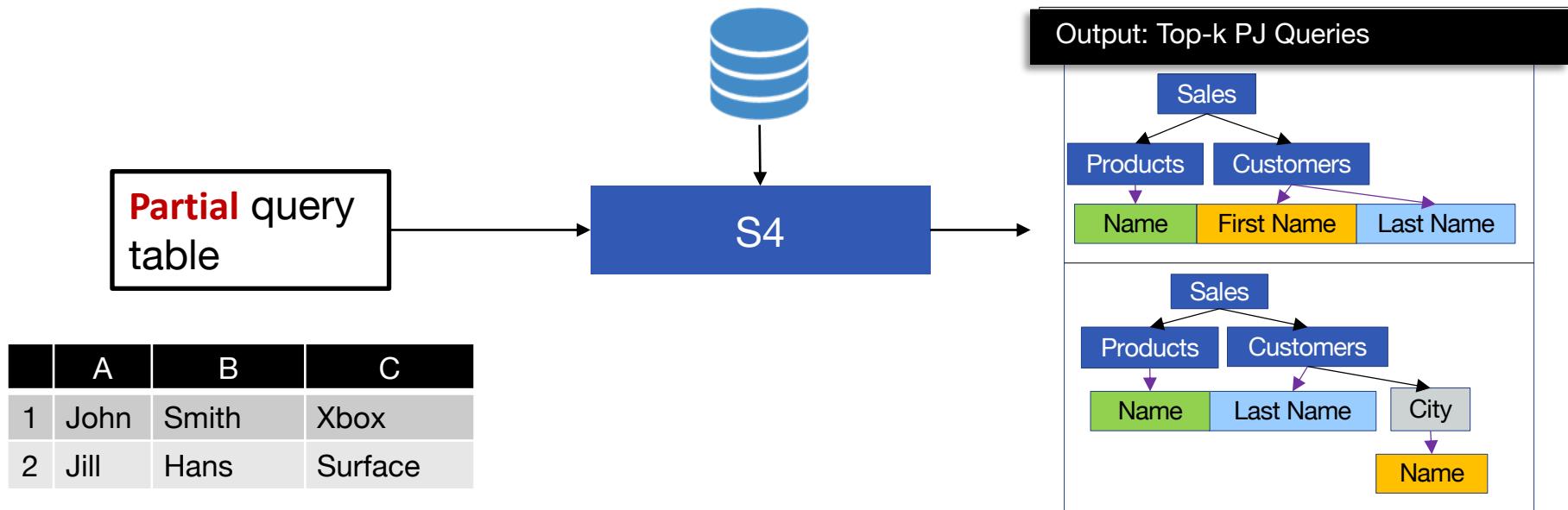
Substructures



# Minimal Project Join REQ

[Psallidas et al., 2015]

Main idea: Allow missing rows/columns and rank the k best queries



# Ranking score

[Psallidas et al., 2015]

Linear combination of row score and column score

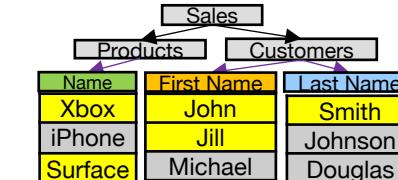
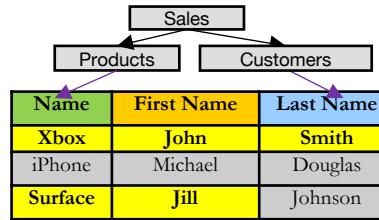
$$\frac{\alpha * \text{score}_{\text{row}}(Q) + (1 - \alpha) * \text{score}_{\text{col}}(Q)}{|Q|}$$

Row score

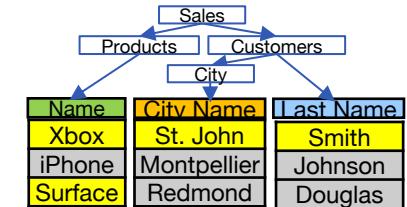
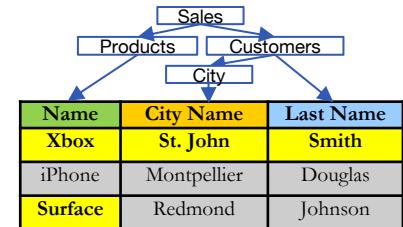
Row Score			
John	Smith	Xbox	3
Jill	Hans	Surface	2
			5
			4

Column score

Column Score	John	Smith	Xbox
	2	1	2
	2	1	1
			5



- $\alpha = 1$  penalizes missing rows
- $\alpha = 0$  penalizes missing columns



# S4 Optimizations

[Psallidas et al., 2015]

Upper bound

Row score is always bounded by the column score  
(row containment is more restrictive)  
Exploit inverted indexes on columns/rows

Early  
termination

Stop when current upper bound score is less than the k-th ranked  
evaluated query  
Scan queries on decreasing upper bound

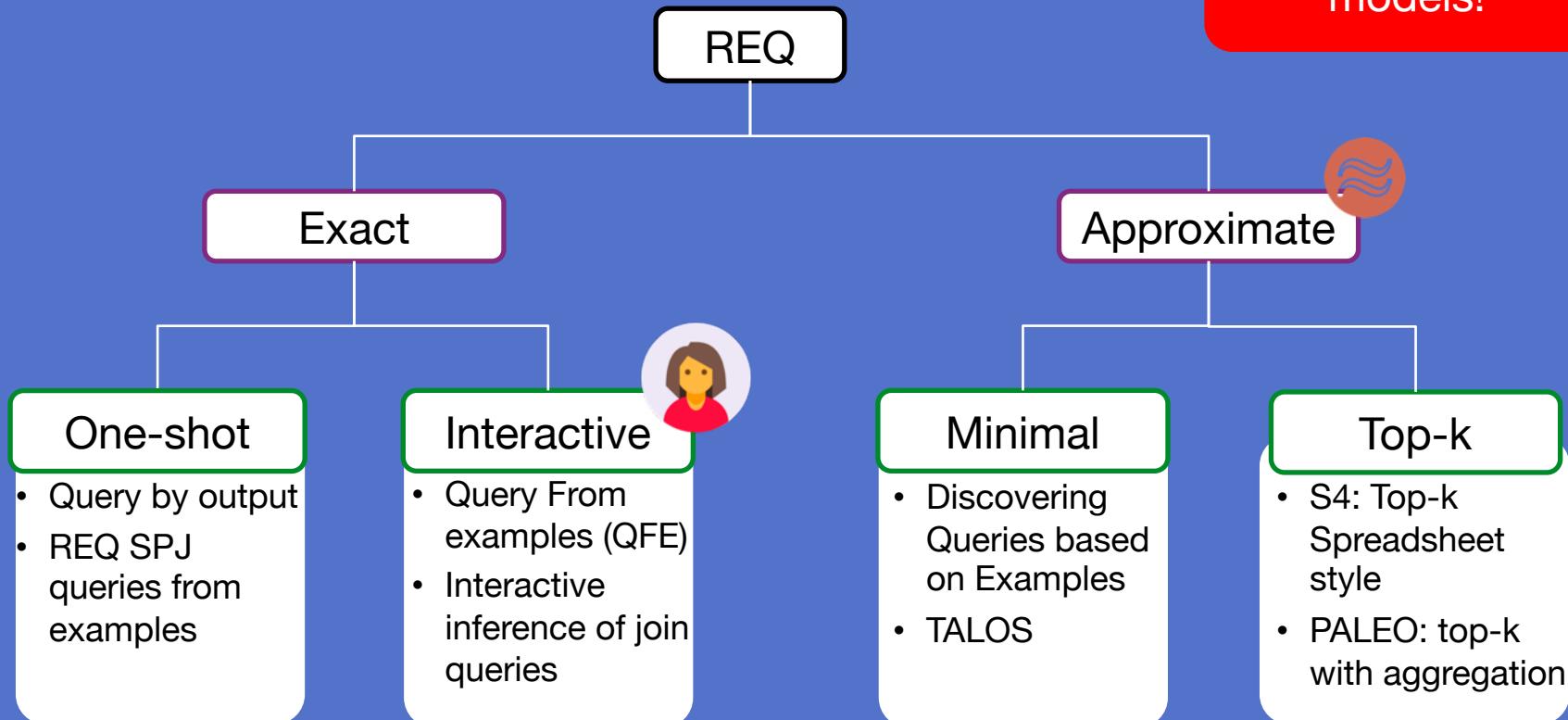
Caching

Reuse common subparts in the candidate queries



# Reverse engineering queries (REQ)

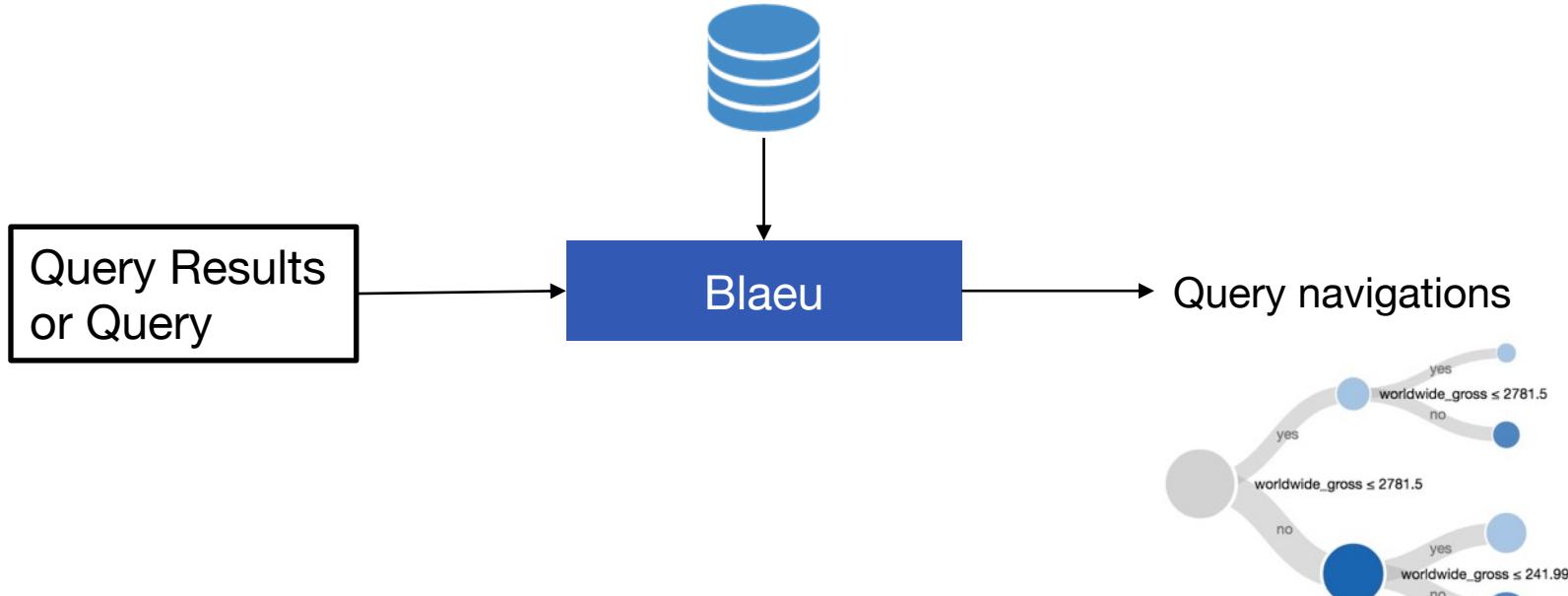
Lack of user models!



# Examples for query suggestion: Blaeu

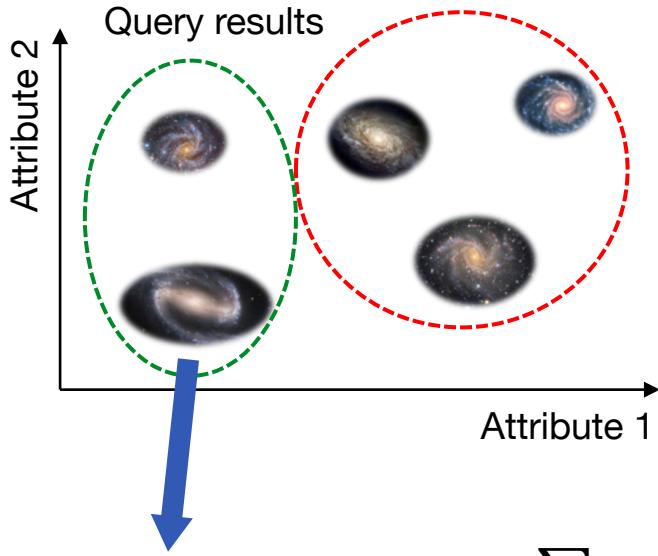
[Sellam et al., 2016]

Main idea: Allow interactive navigation of the query space in a hierarchy



# Examples for query suggestion: Blaeu

[Sellam et al., 2016]



$$u: DB \rightarrow \{-1,1\}, U(Q) = \sum_{t \in Q} u(t)$$

User utility

Given a result of an example query  $Q$ , explore the data through data maps = partitions

**Output:** Set of query refinements

**Problem:** User utility is unknown

- Cluster analysis for result exploration
- Zoom and projection operations
- User model



# Examples for query suggestion: Blaeu

[Sellam et al., 2016]

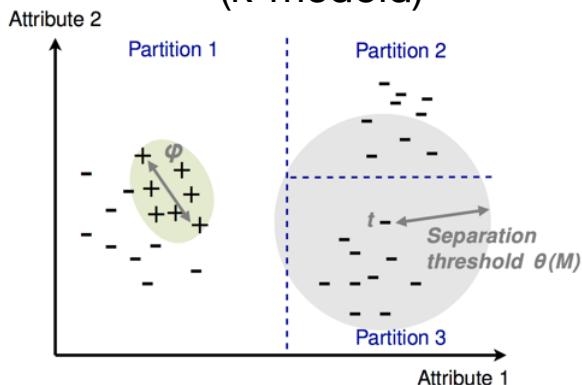
$$u: DB \rightarrow \{-1,1\}, U(C) = \sum_{t \in C} u(t)$$

Unknown User utility

Find the partition  $\mathcal{C} = \{C_1, \dots, C_n\}$  of the results of Q such that exists  $C_j \in \mathcal{C}: U(C_j) > U(Q)$

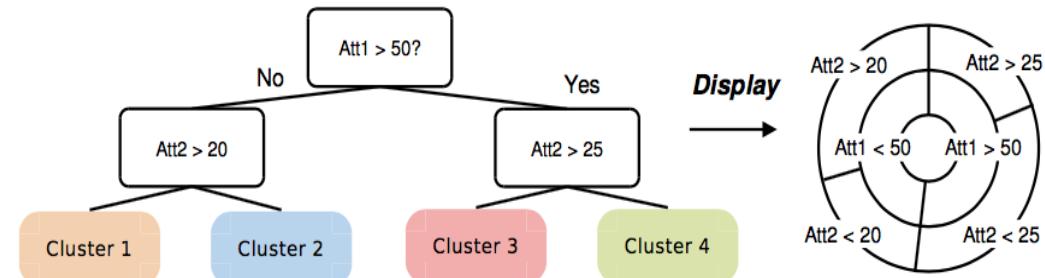
**Solution:** interesting tuples are close to each other within a maximum separation threshold  $\theta(\mathcal{C})$

Detect clusters  
(k-medoid)



Inference

Organize clusters



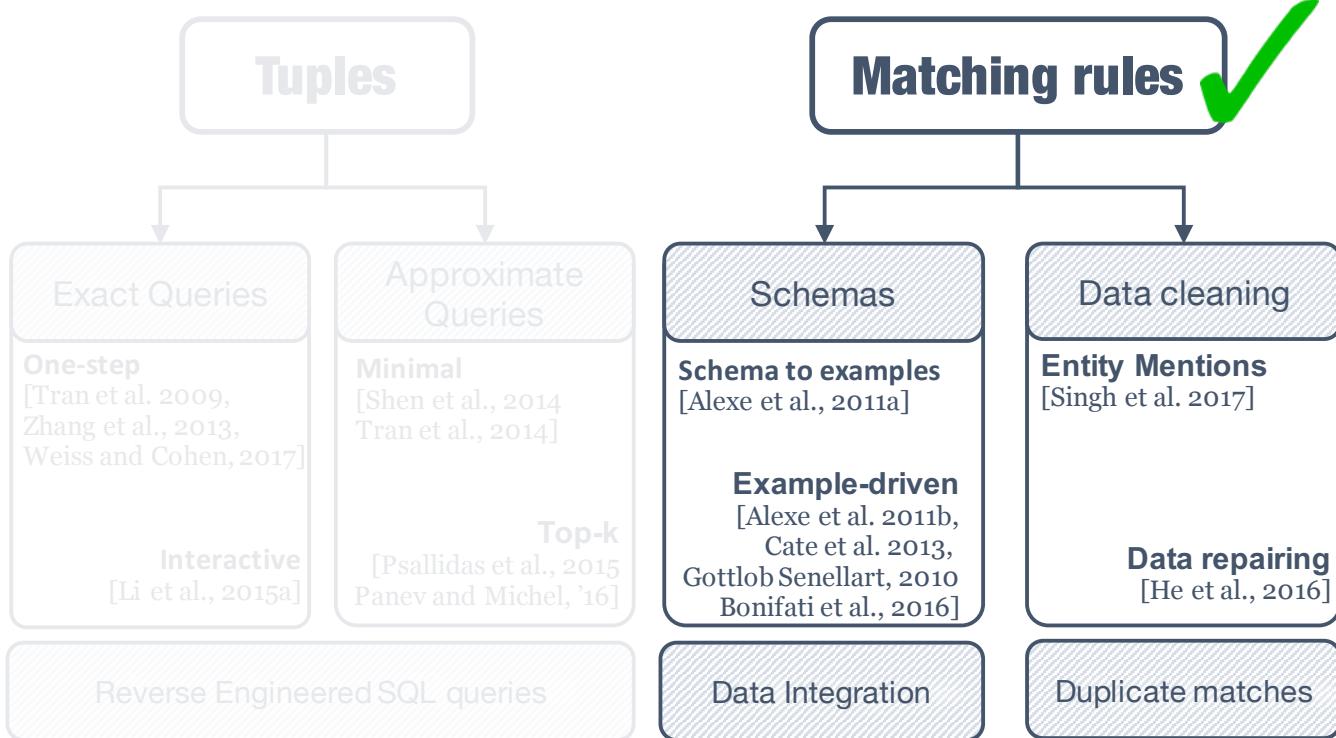
# Searching for ...

SEARCHING FOR

BY FOCUSING ON

APPLYING

PRODUCES



# Data Cleaning

- Often data have redundancy, wrong values, and missing values
- Different values can represent the same object (e.g., N.Y. and New York)
- Values can be simply wrong

**Data cleaning** refers to ways of making the data consistent and correct

<u><i>tid</i></u>	Date	Molecule	Laboratory	Quantity
<u>t1</u>	11 Nov	C <sub>16</sub> H <sub>16</sub> Cl	Austin	200
<u>t2</u>	12 Nov	statin	Austin	100
<u>t3</u>	12 Nov	C <sub>24</sub> H <sub>75</sub> S <sub>6</sub>	N.Y.	100
<u>t4</u>	12 Nov	statin	Boston	200
<u>t5</u>	13 Nov	statin	Austin	200
<u>t6</u>	15 Nov	C <sub>17</sub> H <sub>20</sub> N	Dubai	1000



<u><i>tid</i></u>	Date	Molecule	Laboratory	Quantity
<u>t1</u>	11 Nov	C <sub>16</sub> H <sub>16</sub> Cl	Austin	200
<u>t2</u>	12 Nov	C <sub>22</sub> H <sub>28</sub> F	Austin	100
<u>t3</u>	12 Nov	C <sub>24</sub> H <sub>75</sub> S <sub>6</sub>	New York	100
<u>t4</u>	12 Nov	statin	Boston	200
<u>t5</u>	13 Nov	C <sub>22</sub> H <sub>28</sub> F	Austin	200
<u>t6</u>	15 Nov	C <sub>17</sub> H <sub>20</sub> N	Dubai	100



Not in this tutorial

## Entity matching

*Find mentions of the same objects in different databases*

## Data repairing

*Update the values in the data to fix problems*



# Data repairing: rules

[He, J. et al. 2016]

A **rule** is a logical formula which determines how to change the value in a cell or a group of cells.

IF  $[X_1 = C_1 \dots X_n = C_n]$  UPDATE  $X_i$  to some value

- The update  $t_3[\text{Laboratory}] \leftarrow \text{"New York"}$  can be obtained by the rule
- IF [Laboratory = "N.Y."] UPDATE Laboratory to "New York"
- UPDATE Table  
SET Laboratory='New York'  
WHERE tid=t3

BUT it needs to be done for each cell!!

<u>tid</u>	Date	Molecule	Laboratory	Quantity
<u>t1</u>	11 Nov	C <sub>16</sub> H <sub>16</sub> Cl	Austin	200
<u>t2</u>	12 Nov	statin	Austin	100
<u>t3</u>	12 Nov	C <sub>24</sub> H <sub>75</sub> S <sub>6</sub>	N.Y.	100
<u>t4</u>	12 Nov	statin	Boston	200
<u>t5</u>	13 Nov	statin	Austin	200
<u>t6</u>	15 Nov	C <sub>17</sub> H <sub>20</sub> N	Dubai	1000



<u>tid</u>	Date	Molecule	Laboratory	Quantity
<u>t1</u>	11 Nov	C <sub>16</sub> H <sub>16</sub> Cl	Austin	200
<u>t2</u>	12 Nov	C <sub>22</sub> H <sub>28</sub> F	Austin	100
<u>t3</u>	12 Nov	C <sub>24</sub> H <sub>75</sub> S <sub>6</sub>	New York	100
<u>t4</u>	12 Nov	statin	Boston	200
<u>t5</u>	13 Nov	C <sub>22</sub> H <sub>28</sub> F	Austin	200
<u>t6</u>	15 Nov	C <sub>17</sub> H <sub>20</sub> N	Dubai	100

# Discovering rules

[He, J. et al. 2016]

## UPDATES:

$\Delta_1$ : t3[Laboratory]  $\leftarrow$  "New York"

$\Delta_2$ : t6[Quantity]  $\leftarrow$  100

$\Delta_3$ : t2[Molecule]  $\leftarrow$  "C22H28F"

### Some rules for $\Delta_1$ :

1. Change all Laboratory values to "New York" (t1 – t6)
2. Reformatting all "N.Y" to "New York"(t3)

### Some rules for $\Delta_2$ :

1. Update the quantity to 100 if the molecule is C17H20N and the date is 15 Nov (t6)

### Some rules for $\Delta_3$ :

1. Update to "C22H28F" if molecule is statin (t2,t4,t5)
2. Update to "C22H28F" if molecule is statin and Laboratory Austin (t2,t5)
3. Update to "C22H28F" if molecule is statin and lab is Austin and date is 12 Nov and quantity is 100 (t2)

<u><i>tid</i></u>	Date	Molecule	Laboratory	Quan ty
t1	11 Nov	<chem>C16H16Cl</chem>	Austin	200
t2	12 Nov	<chem>C22H28F</chem>	Austin	100
t3	12 Nov	<chem>C24H75S6</chem>	New York	100
t4	12 Nov	statin	Boston	200
t5	13 Nov	<chem>C22H28F</chem>	Austin	200
t6	15 Nov	<chem>C17H20N</chem>	Dubai	100



# Interactive data cleaning: problem

User validates rules, but has no capacity to validate all rules for each update.

- **Budget Repair Problem:** Given a set  $\mathcal{Q}$  of rules, a table  $T$  and a budget  $B$ , find  $B$  rules from  $\mathcal{Q}$  to maximize the number of repairs over  $T$
- Budget repair problem is an *online problem*

Corresponding *offline problem* is: given as input  $\mathcal{Q}$  rules where validity of each rule is known, select  $B$  rules from  $\mathcal{Q}$  to maximize the number of repairs over  $T$ . (**NP-Hard**)



# Search space

- Given two rules Q and Q', we say that  $Q \leq Q'$  if for all instance T over schema R,  $Q(T) \subseteq Q'(T)$
- In other words if every tuple updated by Q is also updated by Q'
- The relation  $\leq$  forms a **lattice**

<u>tid</u>	Date	Molecule	Laboratory	Quantity
t1	11 Nov	C <sub>16</sub> H <sub>16</sub> Cl	Austin	200
t2	12 Nov	statin → C <sub>22</sub> H <sub>28</sub> F	Austin	100
t5	13 Nov	statin	Austin	200

## For instance

$Q_1 = \text{Date} = "12 \text{ nov}" \text{ AND Laboratory} = "Austin"$   
 $\rightarrow \text{Molecule} = "C_{22}H_{28}F"$

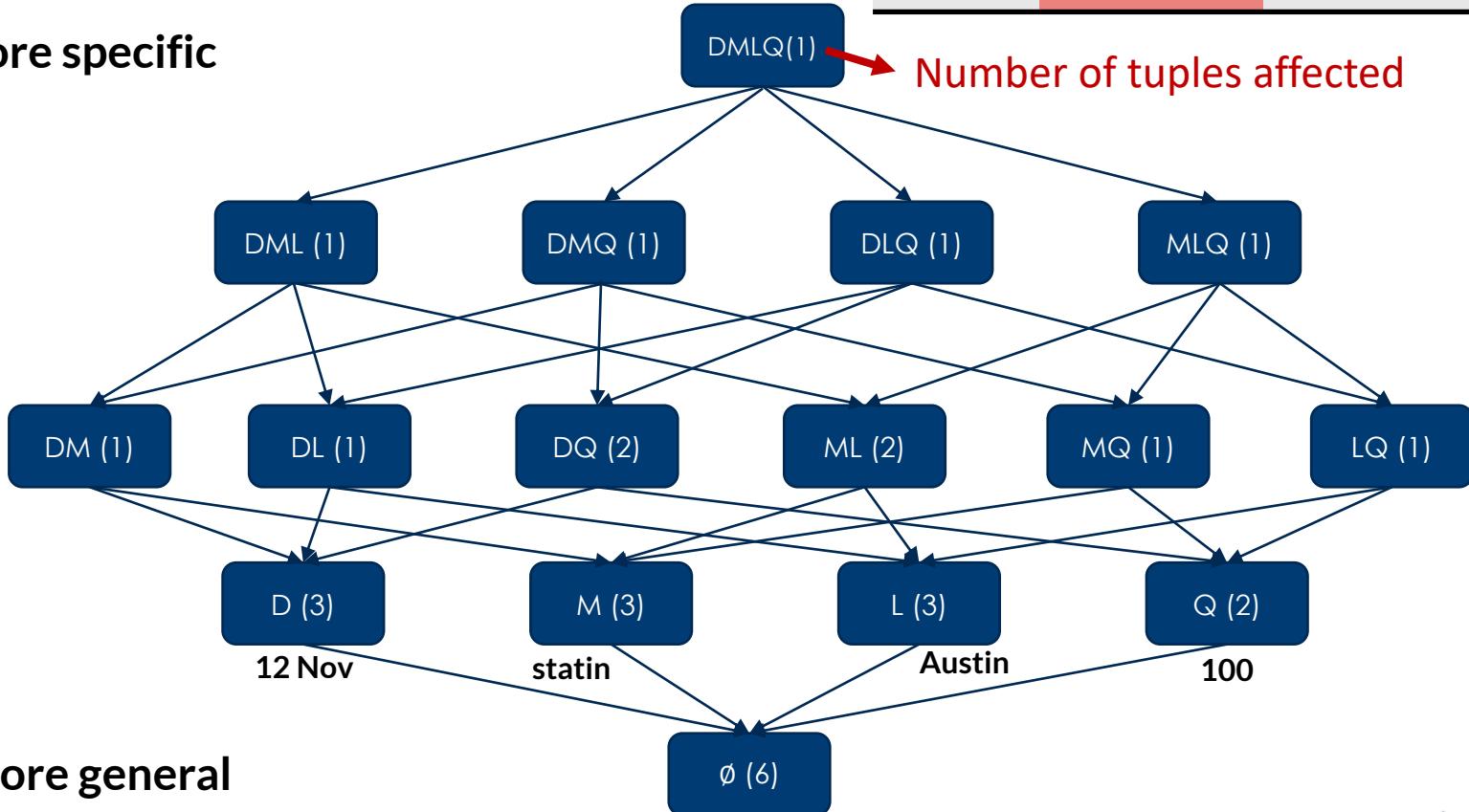
$Q_2 = \text{Molecule} = "statin" \text{ AND Laboratory} = "Austin"$   $\rightarrow \text{Molecule} = "C_{22}H_{28}F"$

$$Q_1 \leq Q_2$$



# Rule lattice

More specific

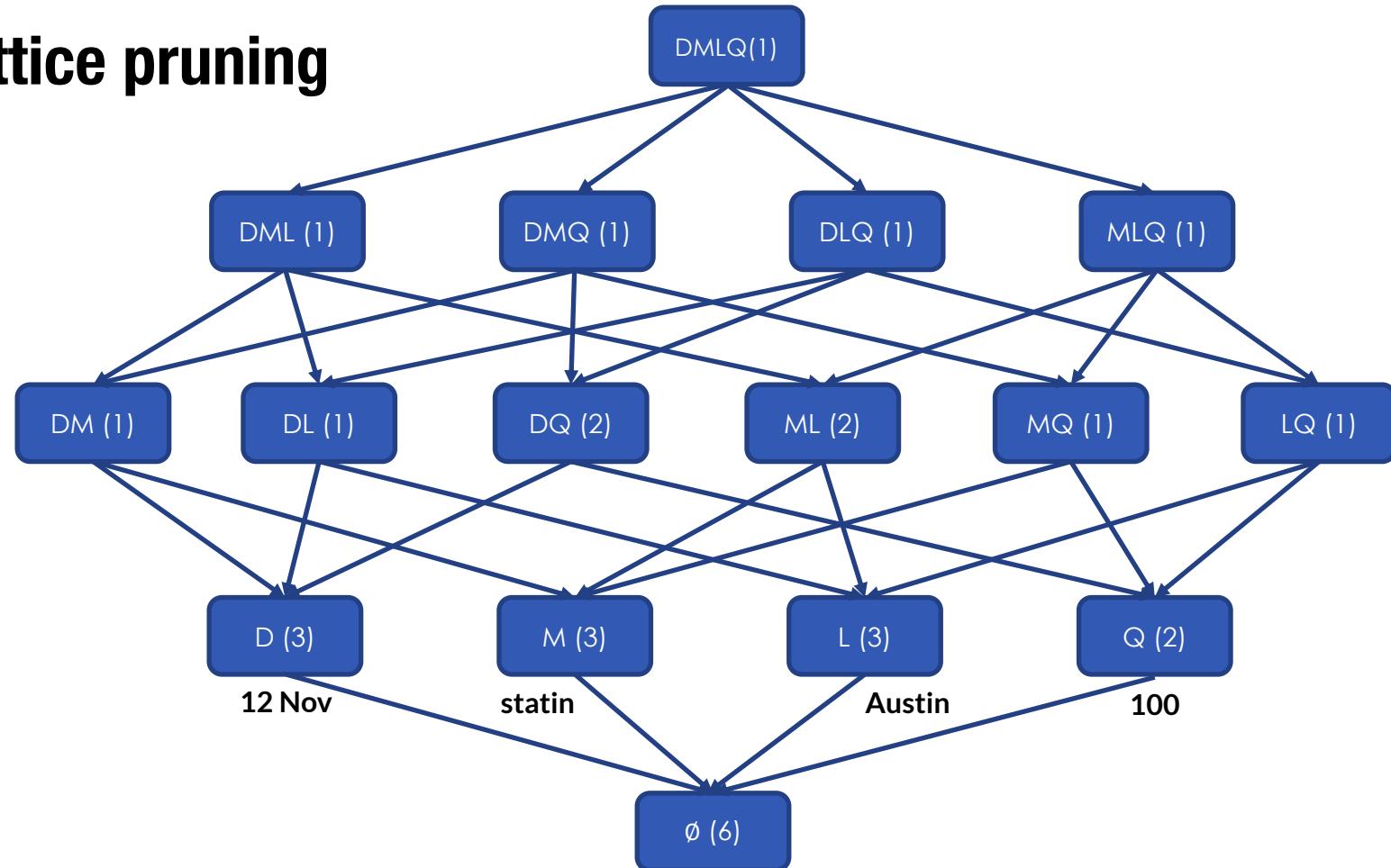


# Lattice pruning

1. If  $Q$  is valid,  $Q'$  is also valid if  $Q' \leq Q$
2. If  $Q$  is invalid,  $Q''$  is also invalid if  $Q \leq Q''$
3. If  $Q$  is valid, all  $Q'$  such that  $Q' \leq Q$  are valid.
4. If  $Q$  is invalid, all  $Q''$  such that  $Q \leq Q''$  are invalid.

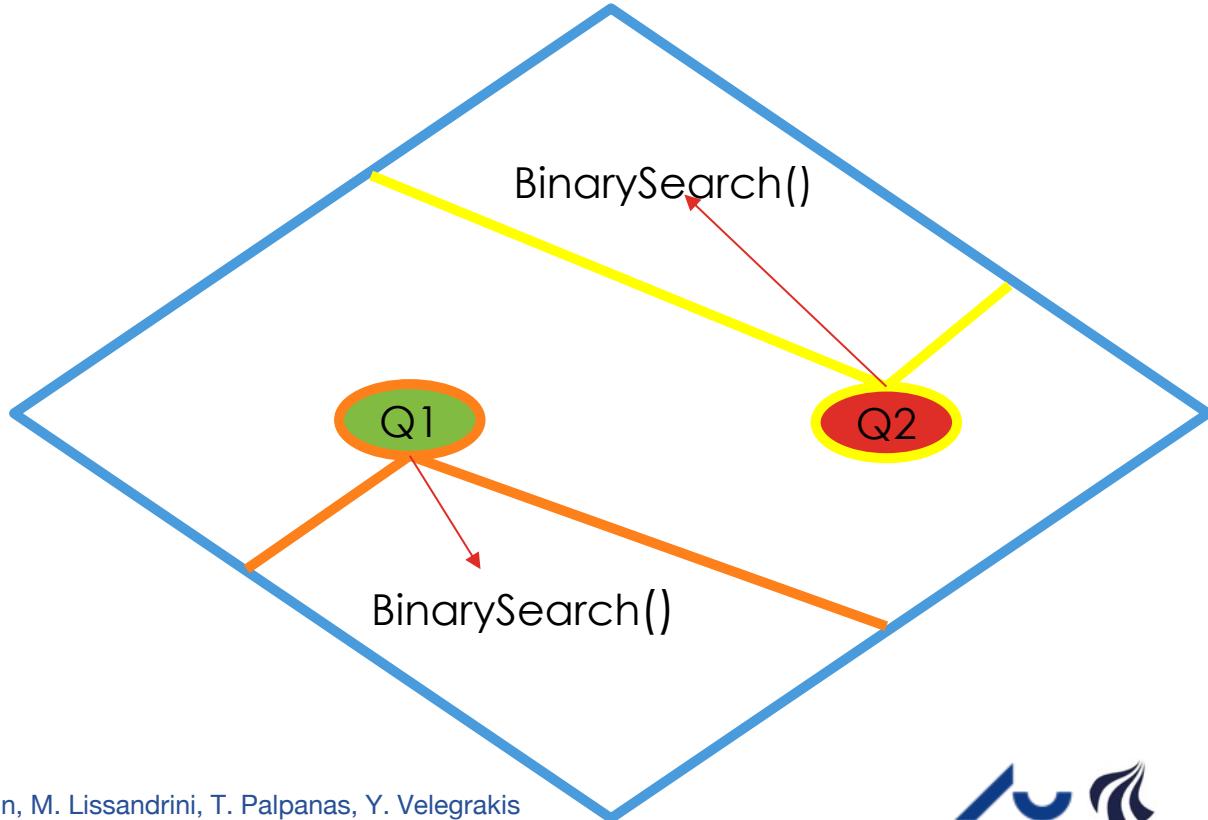


# Lattice pruning



# Dive search

- Binary Search over the lattice ,ordering with #affected tuples
- If  $T \rightarrow \text{BinarySearch}(Q_{\wedge})$
- If  $F \rightarrow \text{BinarySearch}(Q_{\vee})$

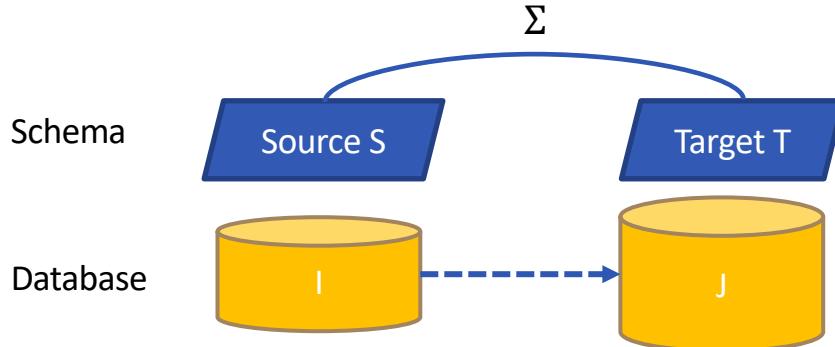


Algorithm complexity  
 $\mathcal{O}(B|Q| \log|Q|)$

# Schema mapping

- Schema mapping finds a way to represent items on one database to items on another database
- Finds a mapping  $\Sigma$  between two schemas such that a query on one database can be converted to a query on the other database
- Schema mappings in  $\Sigma$  are rules in first-order logic that specifies the relationships between schema S and T

$$\forall x \forall y S(x, y) \wedge U(x, z) \rightarrow \exists v T(v, y) \wedge T'(v, z)$$



# Schema mapping: issues

Mapping design is a hard task even for expert users

- Need for domain knowledge and schema understanding
- Need to be acquainted with the syntax and semantics of logic formula
- Need to be able to write queries or customized code



# Examples to the rescue

Not in this tutorial

## Example-driven

[Alexe et al. 2011, Gottlob and Senellar, 2010]

*Takes as input mapped tuples from source and target schema as examples*

"Schema Mapping and Data Examples" –  
[Tutorial at EDBT 2013 Kate, Kolaitis, Tan]

## Interactive Mapping Specification

[Bonifati et al., 2017]

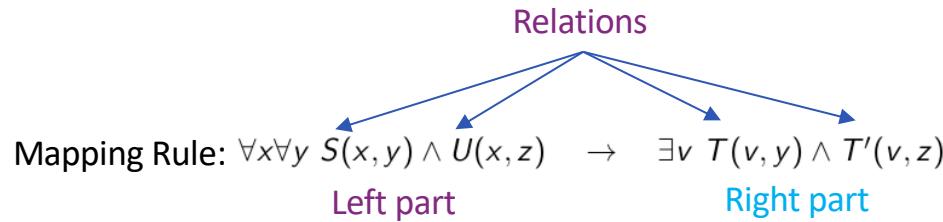
*Ask the user simple questions to refine the mapping*



# Mapping generation

Given a set of example target  $(E_S, E_T)$  and source tuples transform

- every **source** tuple into the **left part** of a mapping rule
- every **target** tuple into the **right part** of a mapping rule



# Mapping generation

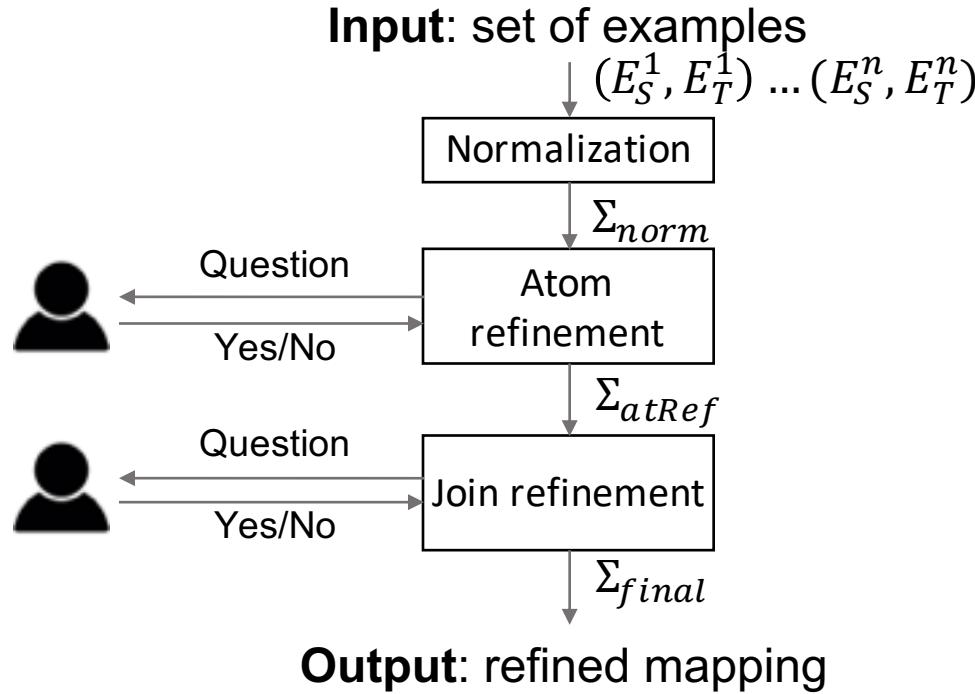
Company		Flight																										
E <sub>S</sub> :	<table border="1"><tr><th><i>IdCompany</i></th><th><i>Name</i></th><th><i>Town</i></th></tr><tr><td>'C1'</td><td>'AA'</td><td>'Paris'</td></tr><tr><td>'C2'</td><td>'Ev'</td><td>'Lyon'</td></tr></table>	<i>IdCompany</i>	<i>Name</i>	<i>Town</i>	'C1'	'AA'	'Paris'	'C2'	'Ev'	'Lyon'	<table border="1"><tr><th><i>Departure</i></th><th><i>Arrival</i></th><th><i>IdCompany</i></th></tr><tr><td>'Lyon'</td><td>'Paris'</td><td>'C1'</td></tr><tr><td>'Paris'</td><td>'Lyon'</td><td>'C2'</td></tr></table>	<i>Departure</i>	<i>Arrival</i>	<i>IdCompany</i>	'Lyon'	'Paris'	'C1'	'Paris'	'Lyon'	'C2'	<table border="1"><tr><th><i>Travel Agency</i></th></tr><tr><th><i>IdAgency</i></th><th><i>Name</i></th><th><i>Town</i></th></tr><tr><td>'A1'</td><td>'TC'</td><td>'L.A.'</td></tr></table>	<i>Travel Agency</i>	<i>IdAgency</i>	<i>Name</i>	<i>Town</i>	'A1'	'TC'	'L.A.'
<i>IdCompany</i>	<i>Name</i>	<i>Town</i>																										
'C1'	'AA'	'Paris'																										
'C2'	'Ev'	'Lyon'																										
<i>Departure</i>	<i>Arrival</i>	<i>IdCompany</i>																										
'Lyon'	'Paris'	'C1'																										
'Paris'	'Lyon'	'C2'																										
<i>Travel Agency</i>																												
<i>IdAgency</i>	<i>Name</i>	<i>Town</i>																										
'A1'	'TC'	'L.A.'																										
Firm		Departure																										
E <sub>T</sub> :	<table border="1"><tr><th><i>Id</i></th><th><i>Name</i></th><th><i>Town</i></th></tr><tr><td>'Id1'</td><td>'AA'</td><td>'Paris'</td></tr><tr><td>'Id2'</td><td>'Ev'</td><td>'Lyon'</td></tr><tr><td>'Id3'</td><td>'TC'</td><td>'L.A.'</td></tr></table>	<i>Id</i>	<i>Name</i>	<i>Town</i>	'Id1'	'AA'	'Paris'	'Id2'	'Ev'	'Lyon'	'Id3'	'TC'	'L.A.'	<table border="1"><tr><th><i>Town</i></th><th><i>IdFirm</i></th></tr><tr><td>'Lyon'</td><td>'Id1'</td></tr><tr><td>'Paris'</td><td>'Id2'</td></tr></table>	<i>Town</i>	<i>IdFirm</i>	'Lyon'	'Id1'	'Paris'	'Id2'	<table border="1"><tr><th><i>Town</i></th><th><i>IdFirm</i></th></tr><tr><td>'Paris'</td><td>'Id1'</td></tr><tr><td>'Lyon'</td><td>'Id2'</td></tr></table>	<i>Town</i>	<i>IdFirm</i>	'Paris'	'Id1'	'Lyon'	'Id2'	
<i>Id</i>	<i>Name</i>	<i>Town</i>																										
'Id1'	'AA'	'Paris'																										
'Id2'	'Ev'	'Lyon'																										
'Id3'	'TC'	'L.A.'																										
<i>Town</i>	<i>IdFirm</i>																											
'Lyon'	'Id1'																											
'Paris'	'Id2'																											
<i>Town</i>	<i>IdFirm</i>																											
'Paris'	'Id1'																											
'Lyon'	'Id2'																											

$m : \text{Company}(c1, aa, paris) \wedge \text{Company}(c2, ev, lyon) \wedge \text{TravelAgency}(a1, tc, la)$   
 $\wedge \text{Flight}(lyon, paris, c1) \wedge \text{Flight}(paris, lyon, c2)$   
 $\rightarrow \exists id1, id2, id3, Firm(id1, aa, paris) \wedge \text{Departure}(lyon, id1) \wedge \text{Arrival}(paris, id1)$   
 $\wedge \text{Firm}(id2, ev, lyon) \wedge \text{Departure}(paris, id2) \wedge \text{Arrival}(lyon, id2) \wedge \text{Firm}(id3, tc, la)$

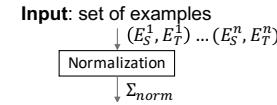


# Interactive Mapping

[Bonifati et al. 2017]



# Mapping normalization


$$\begin{aligned} m : & Company(c1, aa, paris) \wedge Company(c2, ev, lyon) \wedge TravelAgency(a1, tc, la) \\ & \wedge Flight(lyon, paris, c1) \wedge Flight(paris, lyon, c2) \\ \rightarrow & \exists id1, id2, id3, Firm(id1, aa, paris) \wedge Departure(lyon, id1) \wedge Arrival(paris, id1) \\ & \wedge Firm(id2, ev, lyon) \wedge Departure(paris, id2) \wedge Arrival(lyon, id2) \wedge Firm(id3, tc, la) \end{aligned}$$

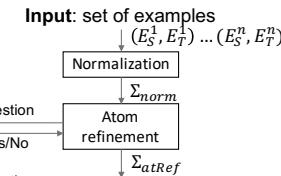
1. Split the mapping into an equivalent set of smaller rules

$$\phi = Company(c1, aa, paris) \wedge Company(c2, ev, lyon) \wedge Flight(lyon, paris, c1) \wedge Flight(paris, lyon, c2) \wedge TravelAgency(a1, tc, la)$$
$$\{\phi \rightarrow \exists id1, Firm(id1, aa, paris) \wedge Departure(lyon, id1) \wedge Arrival(paris, id1);$$
  
 ~~$\phi \rightarrow \exists id2, Firm(id2, ev, lyon) \wedge Departure(paris, id2) \wedge Arrival(lyon, id2),$~~   
 $\phi \rightarrow \exists id3, Firm(id3, tc, la)\}$ 

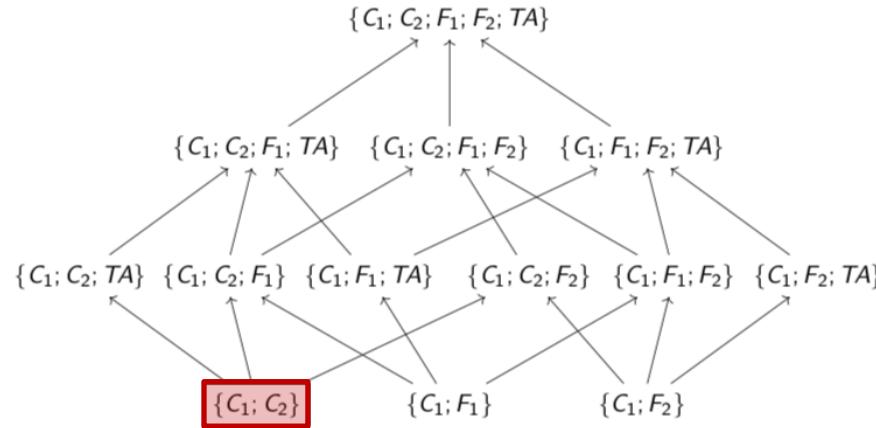
2. Redundancy suppression: either the first or the second rule is removed



# Atom Refinement



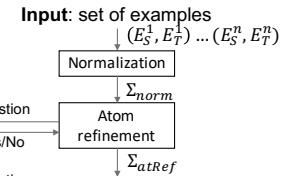
Ask the user and refine the left part of the rule



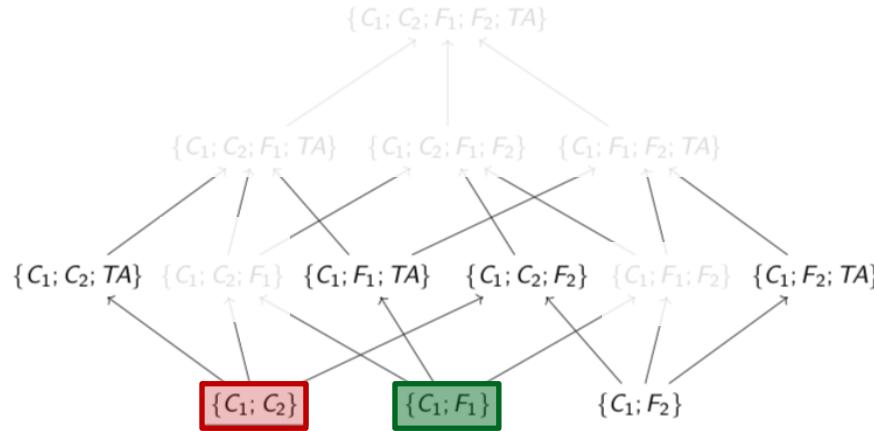
Are the tuples Company(c1,aa,paris); Company (c2, ev, lyon) enough to produce Firm(id, aa, Paris); Departure (Lyon, id); Arrival(Paris, id)?



# Atom Refinement

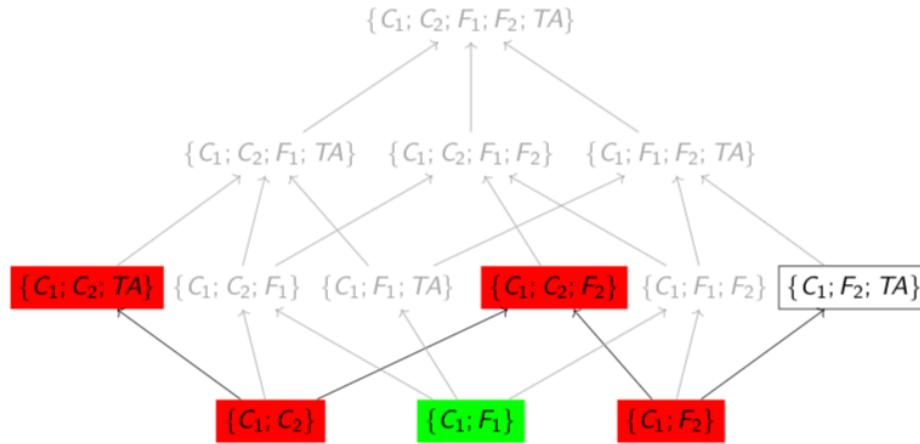
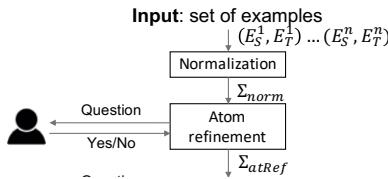


Ask the user and refine the left part of the rule

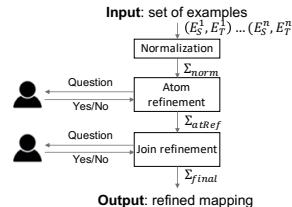


Are the tuples Company(c1,aa,paris); Flight (lyon, paris, c1) enough to produce Firm(id, aa, Paris); Departure (Lyon, id); Arrival(Paris, id)?

# Atom Refinement ...



# Join Refinement



For each rule generated by atom refinement, identify redundant joins entailed by multiple occurrences of a given variable

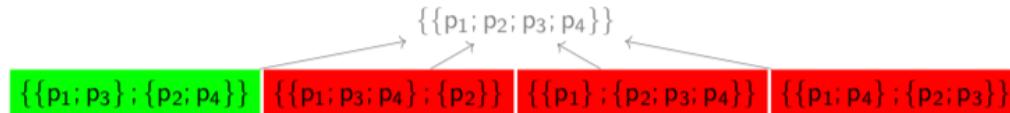
Exploration space:

- As with atom refinement, we can build a semi-lattice of partitions representing possible joins between variable occurrences.

User feedback:

- Similarly to atom refinement, the user is asked about the validity of small sets of tuples.

$$\begin{aligned} & \text{Company}(c1, aa, \textcolor{red}{paris}_1) \wedge \text{Flight}(\text{lyon}, \textcolor{blue}{paris}_2, c1) \\ \Rightarrow & \exists \text{id1}, \text{Carrier}(\text{id1}, aa, \textcolor{pink}{paris}_3) \wedge \text{Departure}(\text{lyon}, \text{id1}) \wedge \text{Arrival}(\textcolor{green}{paris}_4, \text{id1}) \end{aligned}$$



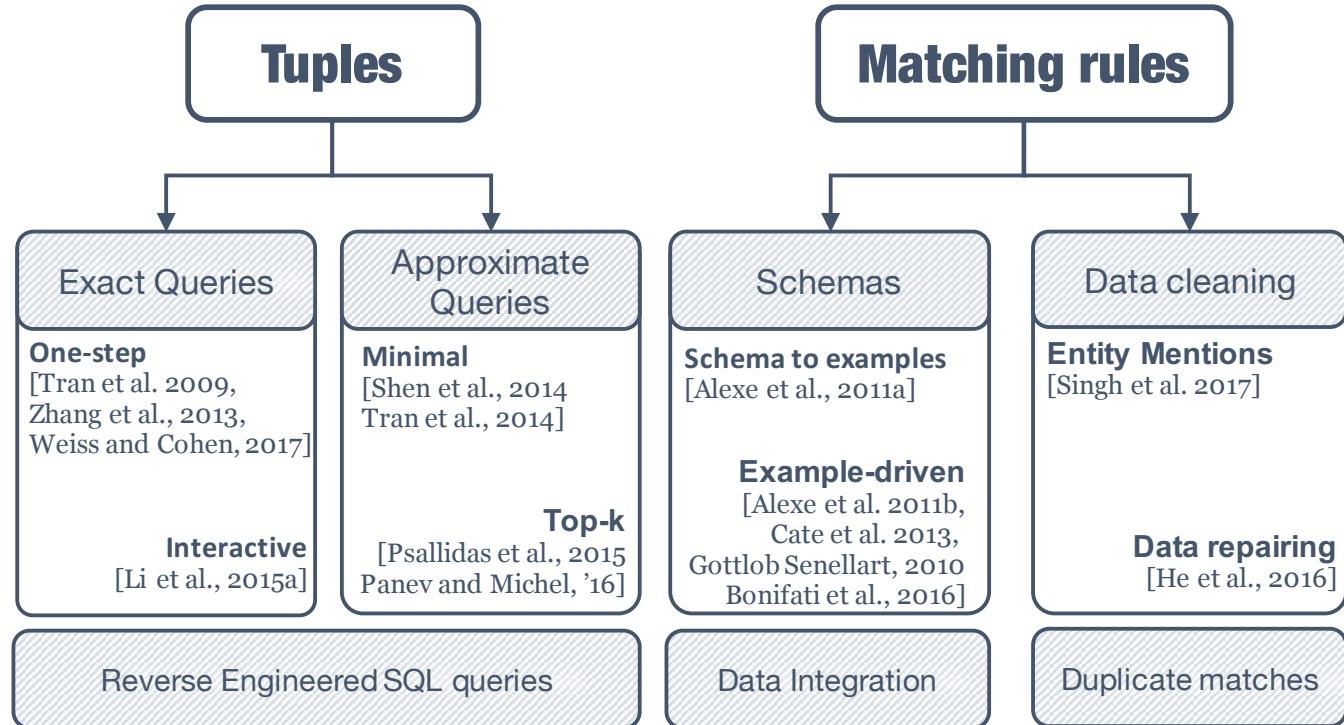
# Summary of relational search paradigms

SEARCHING FOR

BY FOCUSING ON

APPLYING

PRODUCES



<https://j.mp/ExploreSIGMOD>

# Where we are

Relational databases



Textual data



Machine learning

Graphs and networks

Challenges and Remarks

# SIMILARITY for DOCUMENTS

## Unstructured

<p><b>★★★★★ Super Mario Bros The Movie</b></p> <p>By Kay E. Platt on February 23, 2009</p> <p>Hello People, I am going to be reviewing a Movie that ruined my school reputation.... The Movie</p>
<p><b>★★★★★</b></p> <p>September 21, 2018</p> <p>Format: Prime Video</p> <p>Maybe don't name your musical "Rent" if you don't even have a single song about leasing law, property management procedures, or net lease calculations. As a real estate professional I am very disappointed and feel I was misled.</p>
<p><b>★★★★★ Don't be gullible</b></p> <p>January 9, 2019</p> <p>Format: Prime Video</p> <p>This movie is dumb. Neil Armstrong was not very smart at all and Ryan playing him is just wrong. This guy (Armstrong) was not a successor at all. I believe that there are some critical information add up to whether there was a d especially since what more then can't accomplish again. Why is quite advanced today. I feel possible to reach something that is moon and Mars when technology overall I do not give Armstrong many millions of Americans do.</p>
<p>complete as: recommend of 10, if you so you don't</p> <p><b>★★★★★ There are no magicians in this movie</b></p> <p>May 26, 2018</p> <p>Format: DVD</p> <p>I don't mean to give any spoilers away, but there are no magicians in this movie. Don't let the title fool you.</p>
<p>helpful</p>

## Semi-Structured

HR Information		Contact									
Position	Salary	Office	Extn.	Category	Structure	Country	City	Height (metres)	Height (feet)	Year built	Coordinates
Accountant	\$162,700	Tokyo	5407	Mixed use	Burj Khalifa	United Arab Emirates	Dubai	828.1	2,717	2010	25°11'50.0"N 139°16'26.6"E
Chief Executive Officer (CEO)	\$1,200,000	London	5797	Self-supporting tower	Tokyo Skytree	Japan	Tokyo	634	2,080	2011	35°42'38.5"N 139°48'39"E
Junior Technical Author	\$86,000	San Francisco	1560	Guyed steel lattice mast	KVLY-TV mast	United States	Blanchard, North Dakota	628.8	2,063	1963	47°20'32" N 97°17'25" W
Software Engineer	\$1	Abraj Al Bait Towers	2,000	Clock building	Abraj Al Bait Towers	Saudi Arabia	Mecca	601	1,972	2011	21°25'08"N 39°48'35"E
Software Engineer	\$2	Lotte World Tower	1,000	Pre-Sales Support	Lotte World Tower	South Korea	Seoul	555.7	1,823	2017	37°30'45"N 127°0'10"E
Integration Specialist	\$3	One World Trade Center	1,000	Sales Assistant	One World Trade Center	United States	New York, NY	541	1,776	2013	40°42'46.8"N 74°0'48.6"W
Software Engineer	\$1	Military structure	1,000	Senior Javascript Developer	Large masts of INS Kattaborman	India	Tirunelveli	471	1,545	2014	8°22'42.5"N 77°44'38.45"E ; 8°22'30.15"N 77°45'21.07"E
Pre-Sales Support	\$1	United States	1,000	Company	Contact	Country					21°25'11.87"N ; 158°08'53.67"E ; 21°25'13.38"N 158°09'14.35"E ; 3°09'27.45"N 101°42'40.7"E ; 3°09'29.45"N 101°42'43.4"E
Sales Assistant	\$1	Lahaina, Hawaii	458	Launch Pad	Maria Anders	Germany					1972
Senior Javascript Developer	\$4	Malaysia	1,000	Pay Talk	Francisco Chang	Mexico	Kuala Lumpur	452	1,482	1998	
		United States	1,000	Earn More	Roland Mendel	Austria	New York	425.5	1,396	2015	
		United States	1,000	Island Trading	Helen Bennett	UK					

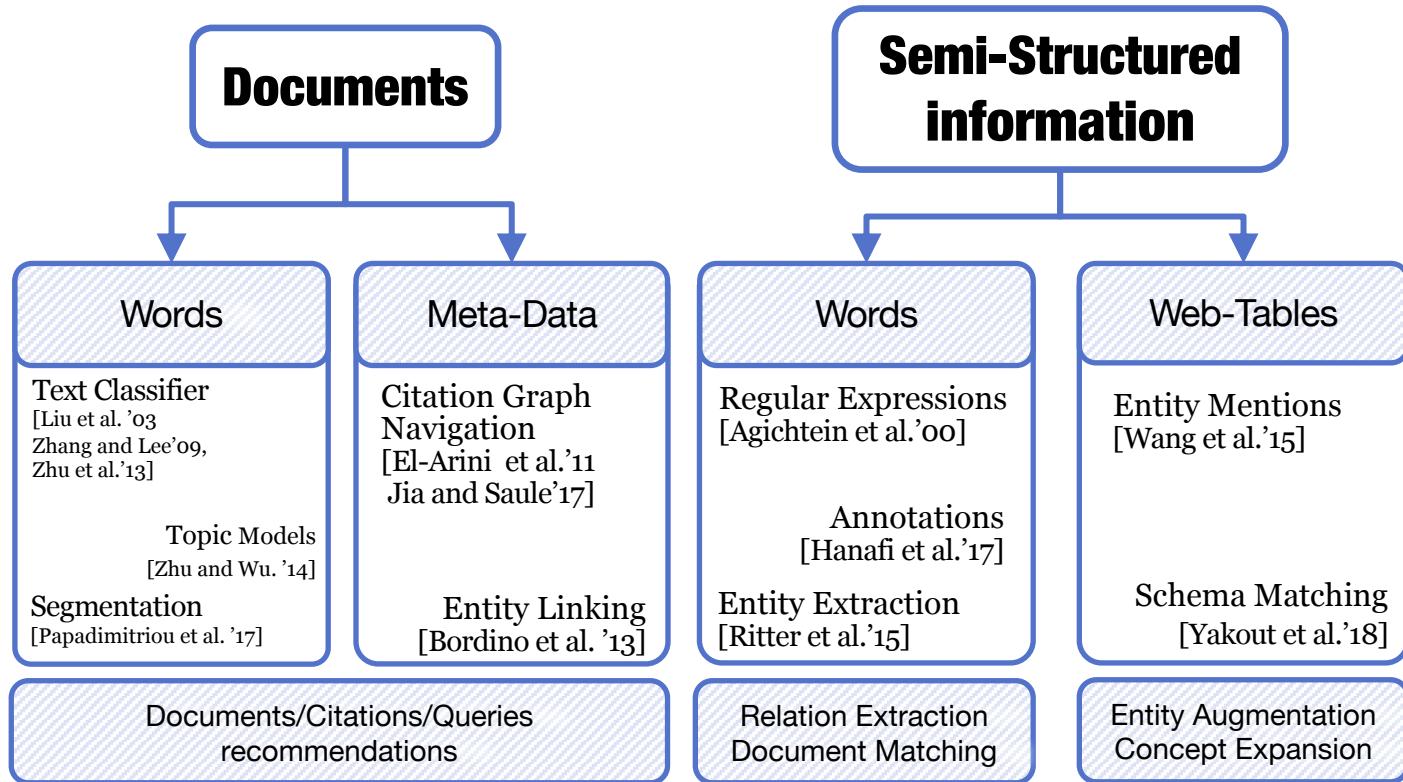


**SEARCHING FOR**

**BY LOOKING AT**

**APPLYING**

**PRODUCES**



# Document Search

Keyword Queries  
& Relevance

**Keyword query: search text with text**

“Action movie with magic”

Search documents containing those exact words

... a live action movie...

.... there is plenty of action...

... packed with action...

... Magic Mike is comedy movie ...

... in Harry Potter magic is everywhere..

Is this enough?

***Identify “relevant words”  
and “relevant documents”***

<https://data-exploration.mli>



★★★★★ Super Mario Bros The Movie

By Kay E. Platt on February 23, 2009

Hello People, I am going to be reviewing a Movie that ruined my school reputation.... The Movie itself is OK....

These famous actors who are chosen to play mario and luigi are acting in this movie, OK.. So I was in first grade when I watched this on VHS, and then my best friend Louis who was sitting next to me at story time was talking to me and then th

★★★★★

September 21, 2018

Format: Prime Video

Maybe don't name your musical "Rent" if you don't even have a single song about leasing law, property management procedures, or net lease calculations. As a real estate professional I am very disappointed and feel I was misled.

★★★★★ There are no magicians in this movie

May 26, 2018

Format: DVD

I don't mean to give any spoilers away, but there are no magicians in this movie. Don't let the title fool you.

★★★★★ Don't be gullible

January 9, 2019

Format: Prime Video

This movie is dumb. Neil Armstrong was not very smart at all and Ryan playing him is just wrong. This guy (Armstrong) was not a successor at all. I believe that there are some critical information that doesn't quite add up whether there was a

especially since what more then accomplish again. Why is quite advanced today. I feel like to reach something that is Mars and Mars when technology I do not give Armstrong millions of Americans do.

★★★★★ pokémon

January 17, 2013

Verified Purchase

Format: VHS Tape

I will watch this while wearing pokémon clothes, sitting with my pokedoll, listening to the theme song, while playing pokémon on my ds.



# Document Search

## Relevant Keywords

Relevance: which keywords are more helpful in describing the content of the document?

Relevance  $\neq$  Frequency

What keywords are more likely to be used to describe the document we want and not other documents

1. Term-frequency: how many times the term appears in the document
2. Document-frequency: In how many documents the term appears

TF-IDF: Term Frequency  
Inverse Document Frequency



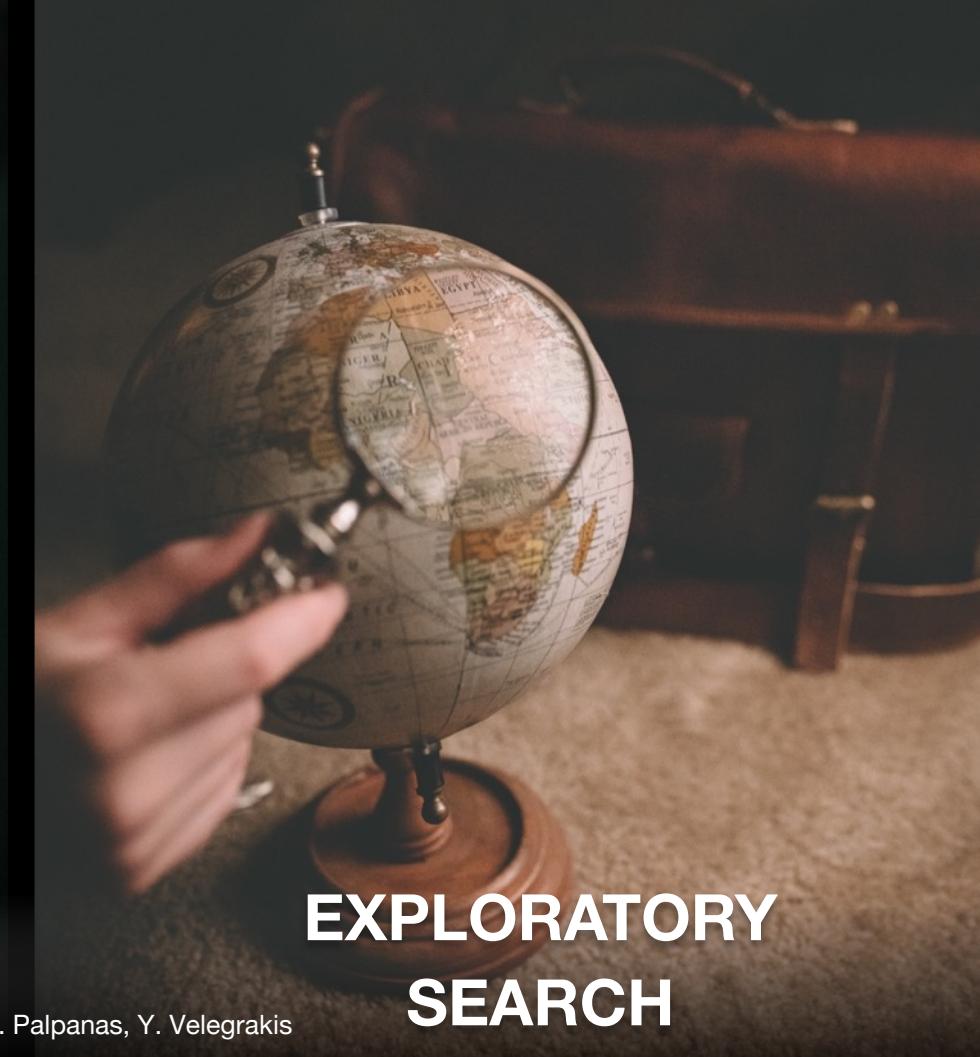
Frequent	TF - IDF	Frequent	TF - IDF	Frequent	TF - IDF
film	maui	film	parrs	film	sulley
moana	te	the	syndrome	sulley	waternoose
the	moana	incredibles	violet	monsters	boo
million	fiti	bird	omnidroid	the	cda
disney	cravalho	pixar	parr	mike	randall
maui	goddess	release	mirage	monster	scarer
				anisan	fizt
				eelen	story
					celia

FEW SELECTED KEYWORDS  
IN THE USER QUERY  
*What keywords to choose?*



## TRADITIONAL SEARCH

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis



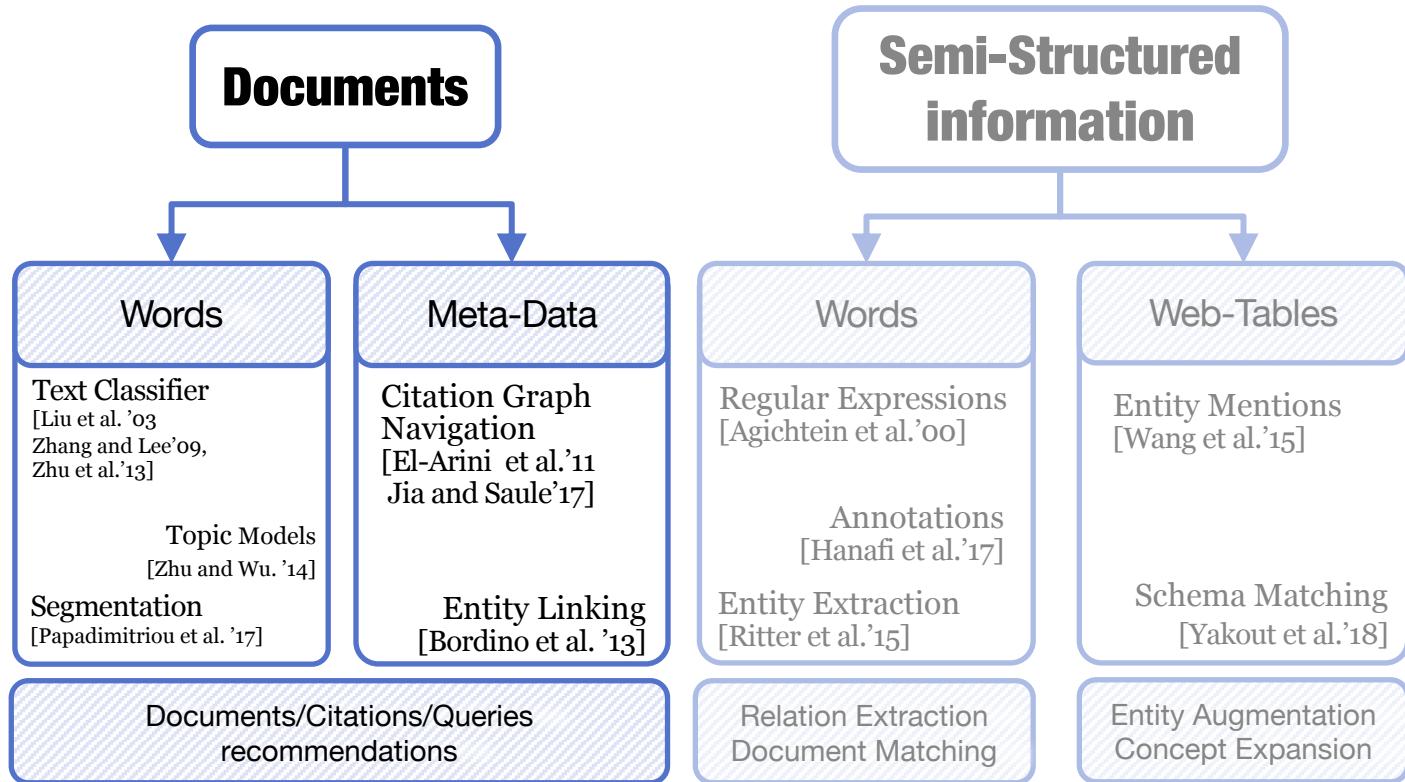
## EXPLORATORY SEARCH

**SEARCHING FOR**

**BY LOOKING AT**

**APPLYING**

**PRODUCES**



# Documents as Examples

Liu et al. [2003]

Exemplar documents

Set of exemplar documents

rather than a set of keywords.

An entire document may contain more information!

It also contains more noise

Identify what makes them special, i.e., relevant

## Example-based Document Search

Given a corpus of documents  $D$ ,  
and a small set of relevant documents ( $D_{\text{rel}}$ ),  
identify a set of answer documents  $D_A$   
such that  $D_{\text{rel}} \subseteq D_A \subseteq D$ .

Model as a classification problem!

Find me movies like these:



Monsters, Inc.

2001 · Fantasy/Adventure · 1h 32m

Monsters Incorporated is the largest scare factory in the monster world, and James P. Sullivan (John Goodman) is one of its top scarers. Sullivan is a huge, intimidating monster with blue fur, large purple spots and horns. His scare assistant, best friend and roommate is Mike Wazowski (Billy Crystal), a green, opinionated, feisty little one-eyed monster. Visiting from the human world is Boo (Mary Gibbs), a tiny girl who goes where no human has ever gone before.

The Incredibles

2004 · Action/Adventure · 1h 56m

In this lauded Pixar animated film, married superheroes Mr. Incredible (Craig T. Nelson) and Elastigirl (Holly Hunter) are forced to assume mundane lives as Bob and Helen Parr after all super-powered activities have been banned by the government. While Mr. Incredible loves his wife and kids, he longs to return to a life of adventure, and he gets a chance when summoned to an island to battle an out-of-control robot. Soon, Mr. Incredible is in trouble, and it's up to his family to save him.

## PROBLEM: MISSING NEGATIVE CLASS

**Few positive examples** and  
**a large set of unknowns**.

What features can discriminate relevant and irrelevant?

Would be better to have some negative examples



# Text Classifiers

Liu et al. [2003]

## Using Positive and Unlabeled Examples

### Positive Unlabeled learning

- a corpus of documents  $D$ ,
- 2 Classes: relevant  $T$  & irrelevant  $\perp$
- relevant documents ( $D_{rel}$ )  
 $\forall d \in D_{rel}. \text{class}(d) = T$
- Unlabeled documents  $U = D - D_{rel}$

### Goal:

train a classifier  $C : D \rightarrow \{T, \perp\}$ ,  
to predict class( $u$ )  $\forall u \in U$ .

### Missing:

To train  $C$  we need examples  
for the negative class  $\perp$

---

### Algorithm 4.9 Document Classification with Positive and Unlabeled Data

---

**Input:** Relevant Documents  $D_{rel} \subseteq \mathcal{D}$ , Unlabeled Documents  $U \subseteq \mathcal{D}$

**Output:** Classifier  $C$

- 1:  $D_{neg} \leftarrow \text{getNegativeSample}(U)$   $\triangleright$  See Li and Liu [2003], Liu et al. [2002], Yu et al. [2002]
  - 2:  $C \leftarrow \text{trainClassifier}(D_{rel}, D_{neg}, U \setminus D_{neg})$   $\triangleright$  E.g., Expectation Maximization, SVM, or Rocchio
  - 3: **return**  $C$
- 



# Inferring Negative Examples (I)

Liu et al. [2003]

Assign a label to Unlabeled data:

how to determine a negative sample set without asking the user

## 4 Alternative approaches

- **Naïve Bayes** (McCallum et al. [1998])
  - All unlabeled data are assumed negatives
  - NB-Classifier estimates  $P(c|d)$  based on  $P(w|c)$  with  $c \in \{T, \perp\}$ ,  $d \in D$ , and words  $w \in W$
- The **Rocchio** technique (Raskutti et al. [2002])
  - $\forall d \in D$   $\vec{d}$  is the TF-IDF vector representation
  - Build prototype vectors  $\vec{c}_T$  for documents in  $D_{rel}$
  - and  $\vec{c}_\perp$  for documents in  $U$
  - Compare each  $\forall d \in U$  with  $\vec{c}_T$  and  $\vec{c}_\perp$
  - assign the class of the most similar vector

### Goal:

Determine set of elements to be regarded as reliable negatives (RN)

$$\vec{c}_T = \alpha \frac{1}{|\mathbf{D}_{rel}|} \sum_{d \in \mathbf{D}_{rel}} \frac{\vec{d}}{\|\vec{d}\|} - \beta \frac{1}{|\mathbf{U}|} \sum_{d \in \mathbf{U}} \frac{\vec{d}}{\|\vec{d}\|}$$

Train a “simplistic” classifier



# Inferring Negative Examples (II)

Liu et al. [2003]

Assign a label to Unlabeled data:

how to determine a negative sample set without asking the user

## 4 Alternative approaches

- The **Spy** technique (Liu et al. [2002])
  - Extract a sample S from the positive example
  - Merge S in U (deploy the spies!)
  - Build NB classifier with EM
  - Determine threshold t such that all spies are correctly classified
  - Document above the threshold are considered negative
- **1-DNF\*** technique (Yu et al. [2002]).

*\*Disjunctive Normal Form*

- *Positive Example Based Learning*
- Get words  $W_f \subset W$ .  $\text{freq}(w, D_{\text{rel}})/|D_{\text{rel}}| > \text{freq}(w, U)/|U|$
- Remove from U all documents containing any word in  $W_f$

### Goal:

Determine set of elements to be regarded as reliable negatives (RN)

Train a “simplistic” classifier



# Training the Expert Classifier

Exploit the partial-supervision

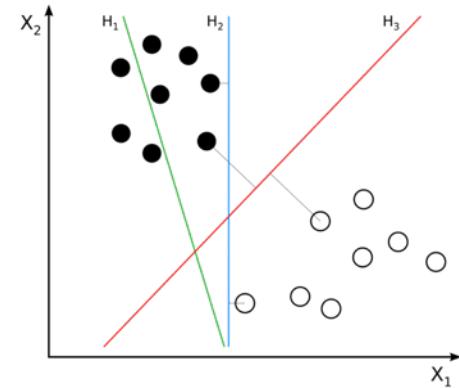
Liu et al. [2003]

## Expert Classifier

Builds on the result of the first step to train a much more sophisticated and precise classifier.

- **1-shot approach**
  - Use  $D_{rel}$  and RN and train a classifier (SVM or EM)
- **Iterative approach**
  - Use  $D_{rel}$  and RN and train a classifier  $C_i$
  - Use  $C_i$  and extract new negative documents Q
  - Add Q to RN, train a new classifier  $C_{i+1}$
  - Continue until no more negative documents are retrieved

[*Optionally*] evaluate the last trained classifier over  $D_{rel}$  and discard it if it performs poorly



Methods perform **poorly** when the initial set of documents is very small

The Rocchio approach + EM is best for this case

Advanced models with TF-IDF or Topic models  
Zhu et al. [2013] - Zhu and Wu [2014]

Beware of Class Imbalance!  
SMOTE: Synthetic Minority Over-sampling Technique



# Document Segmentation

Intention-based relatedness

Papadimitriou et al. [2017]

## Model documents as Composite Objects

Do not perform matching across the posts as a whole but across **fragments** of them that are **written for the same intention**

### Intuition:

Different parts of the document

Have different Purposes:

- Provide background information
- Describe Problem
- Ask question...

I have an HP system with a RAID 0 controller and 4 disks in form of a JBOD. I would like to install Hadoop with a replication 4 HDFS and only 320GB of disk space used from every disc. **Do you know whether it would perform ok or whether the partial use of the disk would degrade performance.** Friends have downloaded the Cloudera distribution but it didn't work. It stopped since the web site was suggesting to have 1TB disks. I am asking because I do not want to install Linux and then realize that my **hardware configuration is not the right one.**

Doc A

Extra RAID disk drives seem to be the solution to my problem but does adding RAID drives requires a **reformat and rebuild of the system to improve performance?**

Doc C

My boss gave me yesterday an HP Pavilion computer with Intel Matrix Storage System, a 320GB drive and Linux pre-installed. I am thinking to add an extra disk drive using a RAID 0 or 1. Can I do it without having to rebuild the entire system? I have already looked at the HP official web site for how to use a JBOD. But I have not found anything related to it.

Doc B

My HP Pavilion stops working after 15 min of activity. I called our technical department but no luck. Despite the many calls, I did not manage to find a **person with adequate knowledge to find out what is wrong.** All they said is bring it to up and we will see, which frustrated me. At the end I had the brilliant idea to move it to a cooler place and voila. No more p

Doc D

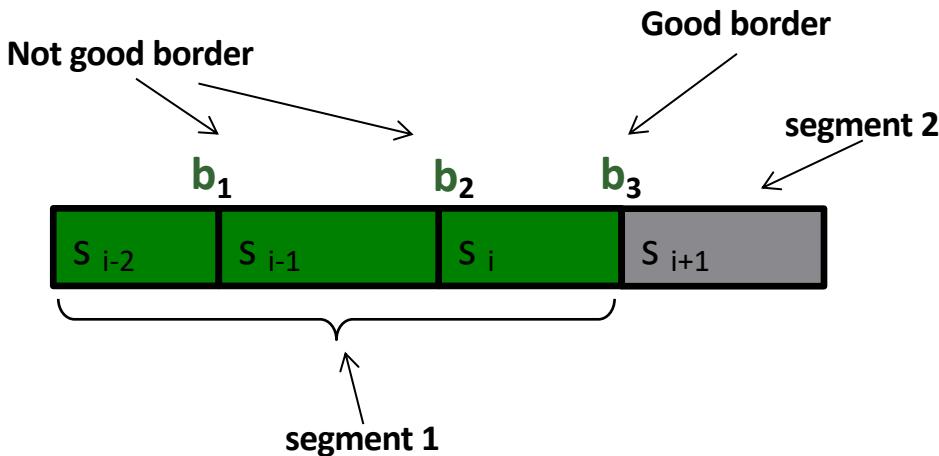


# Segmentation

Papadimitriou et al. [2017]

## Boundaries

Use text characteristics and identify points in which a significant variation of these characteristics occurs, and place a segmentation border there.



## Communication means & Text Features

Tense ( $CM_{tense}$ )	present	past	future
Subject ( $CM_{subj}$ )	I/we	you	it/they/(s)he
Style ( $CM_{qneg}$ )	interrog.	negative	affirmative
Status ( $CM_{pasact}$ )	passive	active	
Part of Speech ( $CM_{pos}$ )	verb	noun	adj./adverb

0 I have an HP system with a RAID 0 controller and 4 disks in form of a JBOD. 75 I would like to install Hadoop with a replication 4 HDFS and only 320GB of disk space used from every disc. 182 Do you know whether 201 it would perform ok or whether the partial use of the disk 259 would degrade performance. 285 Friends have downloaded the Cloudera distribution but 338 it didn't work. 355 It stopped since 371 the web site was suggesting to have 1TB disks. 418 I am asking because 436 I do not want to install Linux and then realize that 488 my hardware configuration is not the right one. 535

## Bottom-up approach

1. Start with single words as segments
2. Compute a **Diversity Index** in each segment
3. Merge segments with low diversity



# Intention Clustering & Matching

Matching among segments with the same intention

Clusters are based on intentions

Given a document  $d_q$ ,

1. the system will **segment**  $d_q$ ,
2. identify for each segment the **segments in the same cluster**
3. **aggregate the similarity** of those segments into a score for each document.

C1

I am asking because I do not want to install Linux and then

I have already looked at the HP official web site for how to use a JBOD. But I have not found anything related to it.

I have an HP system with a RAID 0 controller and 4 disks in Extra RAID disk drives seem to be the solution to my problem but does adding RAID drives requires a reformat and rebuild of the system to improve performance?

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

8

Papadimitriou et al. [2017]

Friends have downloaded the Cloudera distribution but it didn't work. It stopped since the web site was suggesting to have 1TB disks.

C3

I am thinking to add an extra disk drive using a RAID 0 or 1. Can I do it without having to rebuild the entire system?

Do you know whether it would perform ok or whether the

Despite the many calls, I did not manage to find a person with adequate knowledge to find out what is wrong.

Explore based on related topics related to common goals

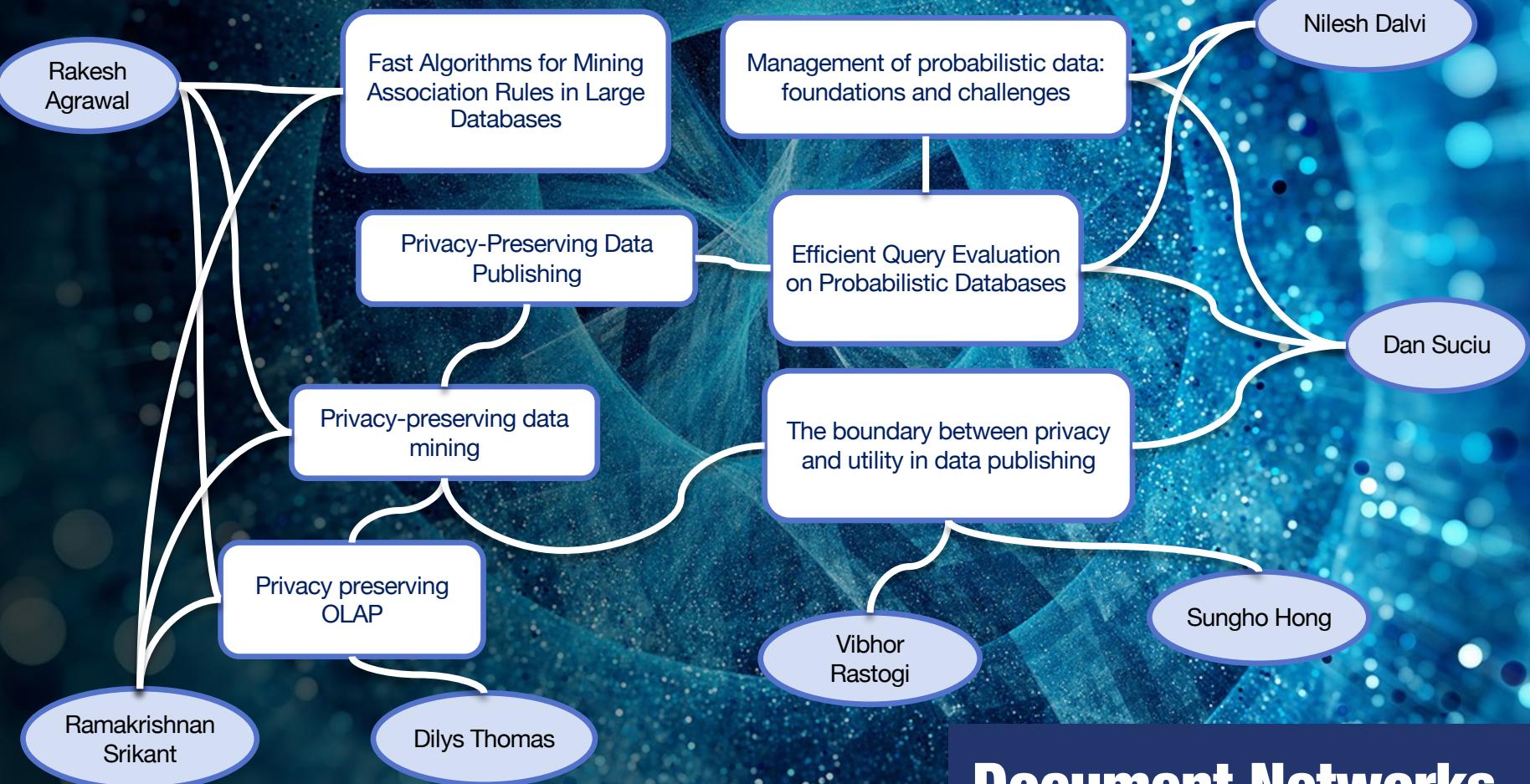
C2

All they said is bring it to up and we will see, which frustrated me. At the end I had the brilliant idea to move it to a cooler place and voila. No more problems.

My HP Pavilion stops working after 15 min of activity. I called our technical department but no luck.

Linux pre-installed.





# Document Networks

# Influence in Citation Networks

Document relevance based on influence

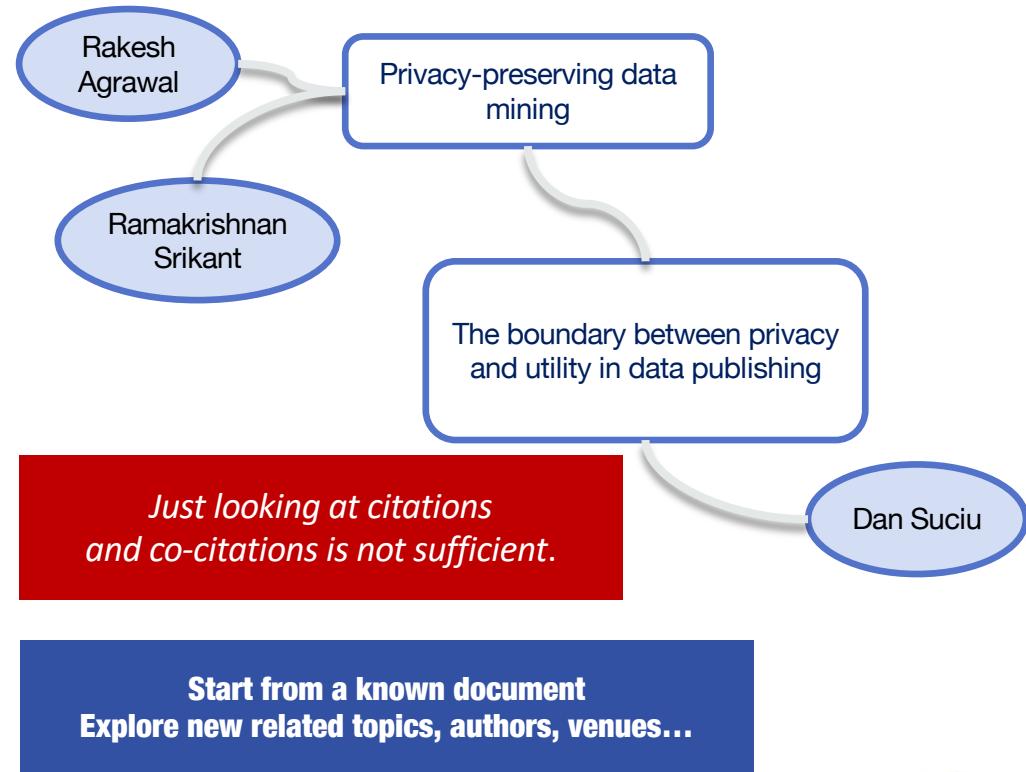
El-Arini and Guestrin [2011]  
Jia and Saule [2017]

## Citation Network

- Nodes are Authors and Papers
- Edges are Authorship and Citations
- Influence is based on connecting Paths

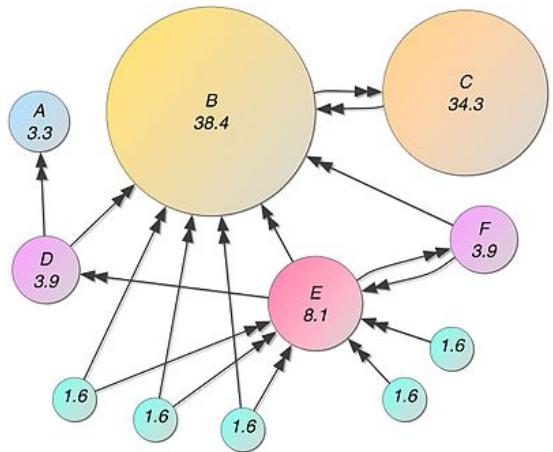
## Advance Models

- El-Arini and Guestrin [2011] :
  - Condition influence on topics  
*Iterate for each topic T: Select topic T, keep only papers relevant for T, compute connecting Paths.*
  - Weight edges with Influence-Probability
- Jia and Saule [2017]
  - Enrich graph with Keywords & Venues



# Traverse Document Networks

How to navigate links and connections



## Global Page Rank

Starting from a random node,  
traversing randomly, random  
restart point

## Personalized Page Rank

- Start from seed nodes, i.e. the documents  $D_{\text{rel}}$
- Navigate towards locally connected nodes

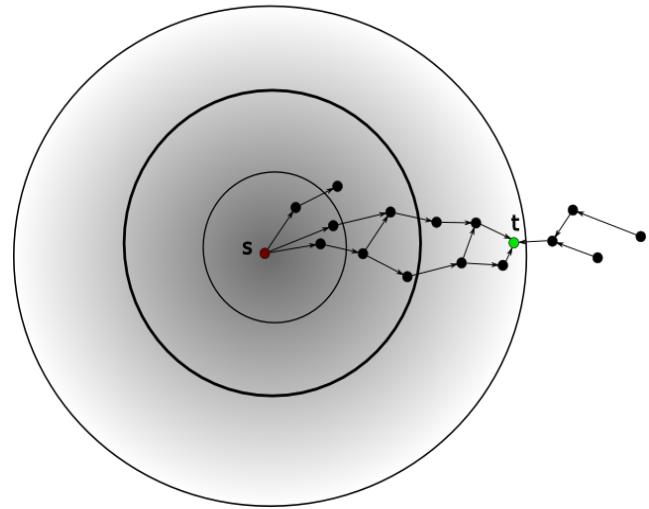
**Example based Exploration implies locality**

**CHALLENGE:**  
*Identify meaningful transition probabilities*

*E.g., El-Arini and Guestrin [2011]*

El-Arini and Guestrin [2011]

Jia and Saule [2017]



## Personalized Page Rank

Starting from a limited set of nodes,  
traversing randomly,  
restart point is one in the initial set.  
Bound not to travel too far



# Serendipitous Search

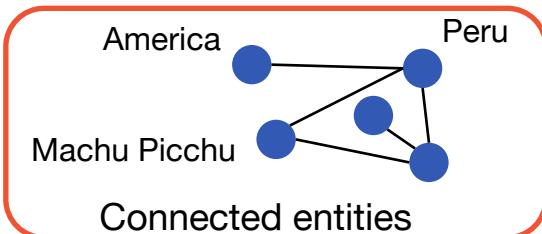
Bordino et al. [2013]

Enhance document links with Entities and Query-logs

**Input:** Query/Document  
**Output:** Queries



Document



**Serendipity**  
*Related topics potentially come to mind after consulting the page.*

rafting excursion down the urubamba river  
el dorado temple of sun  
indios quechuas  
map of peru  
sapa inca

Exploit “lateral connections”  
in User Search Behaviors

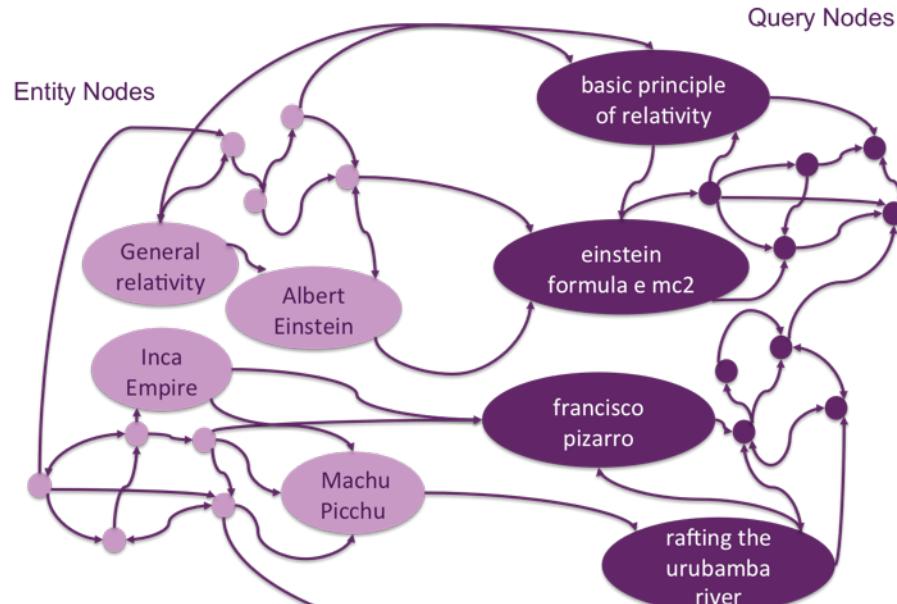
Searches related to  
Document content



# Entity Query Graph

Bordino et al. [2013]

Entity-Query graph from queries to entities and back



## EQGraph Weighted Edges

1. query to query:

$$w_Q(q_i \rightarrow q_j) = w_{QFG}(q_i \rightarrow q_j)$$

Queries in the same session

2. entity to query

$$w_{EQ}(e \rightarrow q) = \frac{f(q)}{\sum_{q_i | e \in X_E(q_i)} f(q_i)}$$

Frequency-based approach

3. entity to entity

$$w_E(e_u \rightarrow e_v) = 1 - \prod_{i=1, \dots, r} (1 - p_{q_{i_s} \rightarrow q_{i_t}}(e_u \rightarrow e_v))$$

The more queries entities share the higher the probability

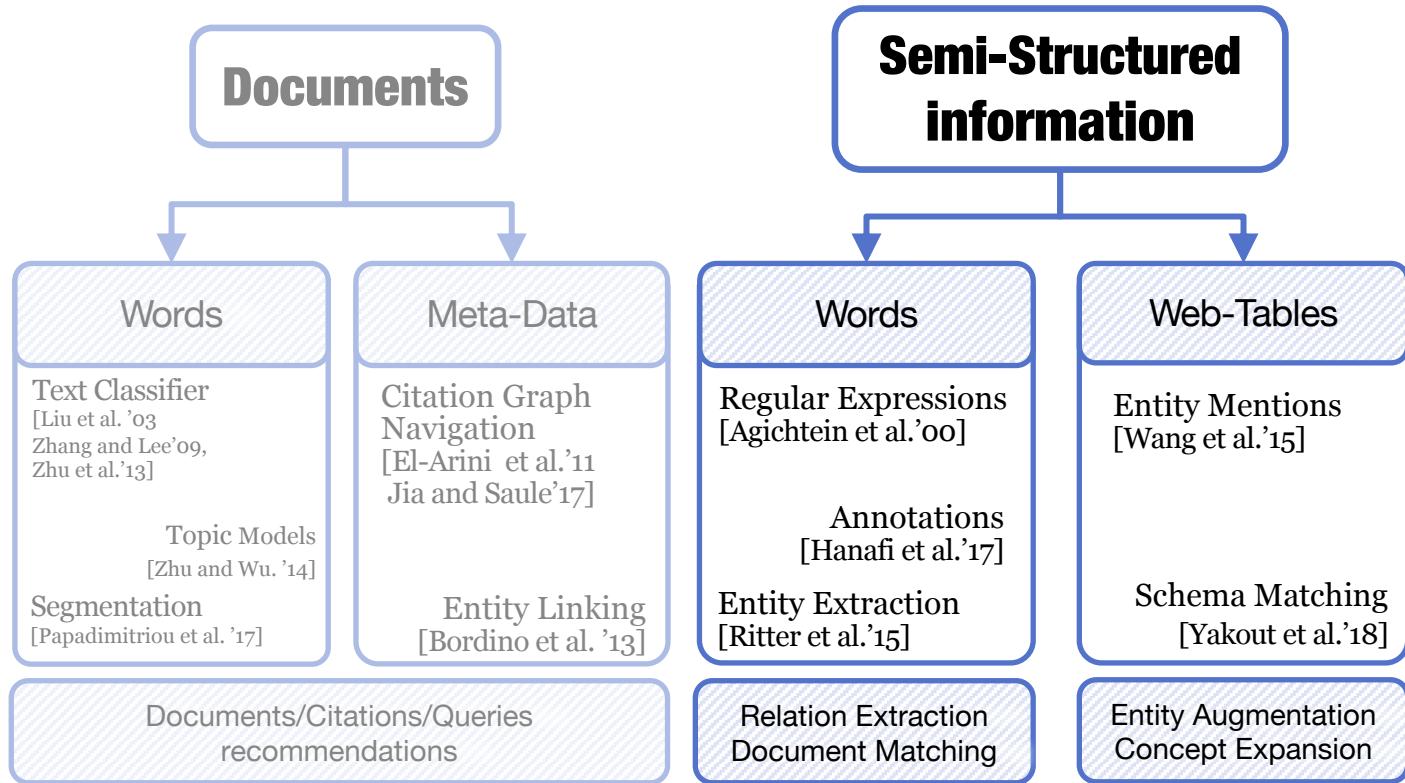
Based on query to query edges

**SEARCHING FOR**

**BY LOOKING AT**

**APPLYING**

**PRODUCES**



# Entity Mentions & Web-Tables

Documents & semi-structured information

In fact, the Chinese market has the three most influential names of the retail and tech space – Alibaba, Baidu, and Tencent (collectively touted as BAT), and is betting big in the global AI in retail industry space. The three giants which are claimed to have a cut-throat competition with the U.S. (in terms of resources and capital) are positioning themselves to become the 'future AI platforms'. The trio is also expanding in other Asian countries and investing heavily in the U.S. based AI startups to leverage the power of AI. Backed by such powerful initiatives and presence of these conglomerates, the market in APAC AI is forecast to be the fastest-growing one, with an anticipated CAGR of 45% over 2018 - 2024.

To further elaborate on the geographical trends, North America has procured more than 50% of the global share in 2017 and has been leading the regional landscape of AI in the retail market. The U.S. has a significant credit in the regional trends with over 65% of investments (including M&As, private equity, and venture capital) in artificial intelligence technology. Additionally, the region is a huge hub for startups in tandem with the presence of tech titans, such as Google, IBM, and Microsoft.

HR Information		Contact			
Position	Salary	Office	Extn.		
Accountant	\$162,700	Tokyo	5407		
Chief Executive Officer (CEO)	\$1,200,000				
Junior Technical Author	\$86,000				
Software Engineer	\$132,000				
Software Engineer	\$206,850				
Integrator	\$270,000				

Structure	Country	City	Height (metres)	Height (feet)	Year built	Coordinates
Burj Khalifa	United Arab Emirates	Dubai	828.1	2,717	2010	25°11'50.0"N 55°16'26.6"E
Tokyo Skytree	Japan	Tokyo	634	2,080	2011	35°42'36.5"N 139°48'39"E
KVLY-TV mast	United States	Blanchard, North Dakota	628.8	2,063	1963	47°20'32"N 97°17'25"W
Abraj Al Bait Towers	Saudi Arabia	Mecca	601	1,972	2011	21°25'08"N 39°49'35"E
Lotte World Tower	South Korea	Seoul	555.7	1,823	2017	37°30'45"N 127°6'10"E
One World Trade Center	United States	New York, NY	541	1,776	2013	40°42'46.8"N 74°0'48.6"W
Large masts of INS Kattabomman	India	Tirunelveli	471	1,545	2014	8°22'42.52"N 77°44'38.45"E ; 8°22'30.13"N 77°45'21.07"E
Lualualei VLF transmitter	United States	Lualualei, Hawaii	458	1,503	1972	21°25'11.87"N 158°08'53.67"W ; 21°25'13.38"N 158°09'14.35"W
Country		Malaysia	Kuala Lumpur	452	1,482	1998
Germany		Austria	Munich	425.5	1,396	2015
Pay Talk	Francisco Chang	Mexico				
Earn More	Roland Mendel	Austria				
Island Trading	Helen Bennett	UK				

# Entity-relation tuples

Example-based extraction of Entity mentions and Relations

Brin [1998]

Agichtein and Gravano [2000]

## Search for Information **WITHIN** Documents

Explore new Entities  
and new ways to express relations

### 1. Example

⟨ Google ; Menlo Park ⟩

### 2. Match

Google founded in Menlo Park...

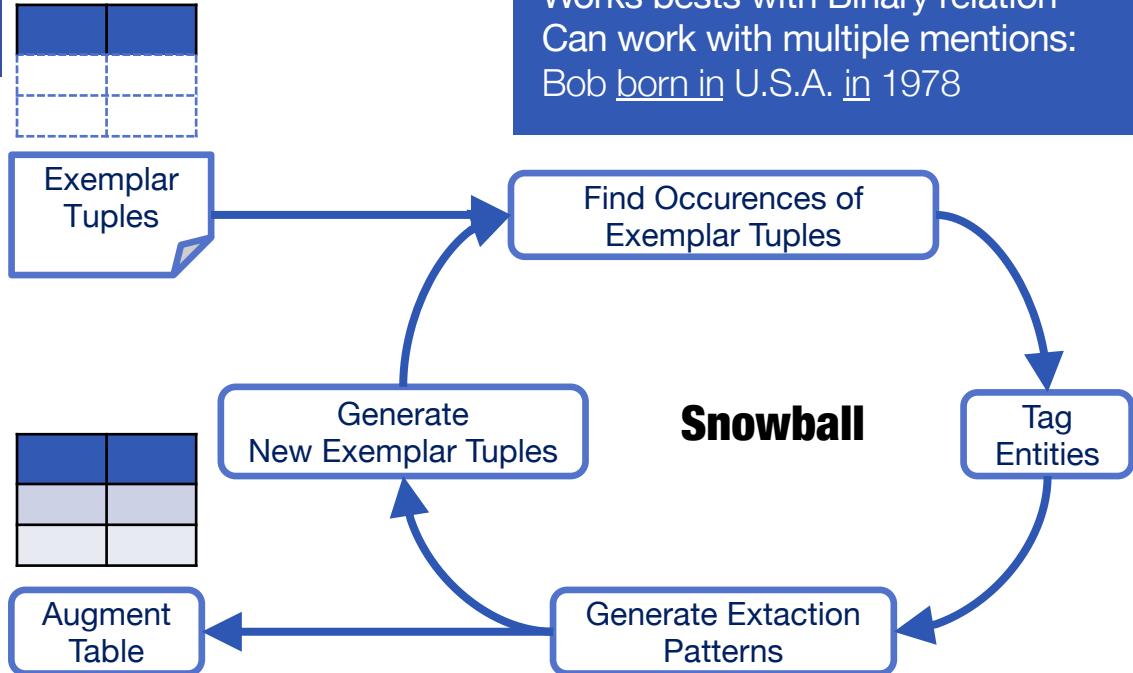
### 3. Extract Pattern

... [X] founded in [Y] ...

### 4. Extract New Mentions & Patterns

Apple founded in Cupertino ...

Apple headquarters in Cupertino



# Entity-relation tuples

Example-based extraction of Entity mentions and Relations

Brin [1998]

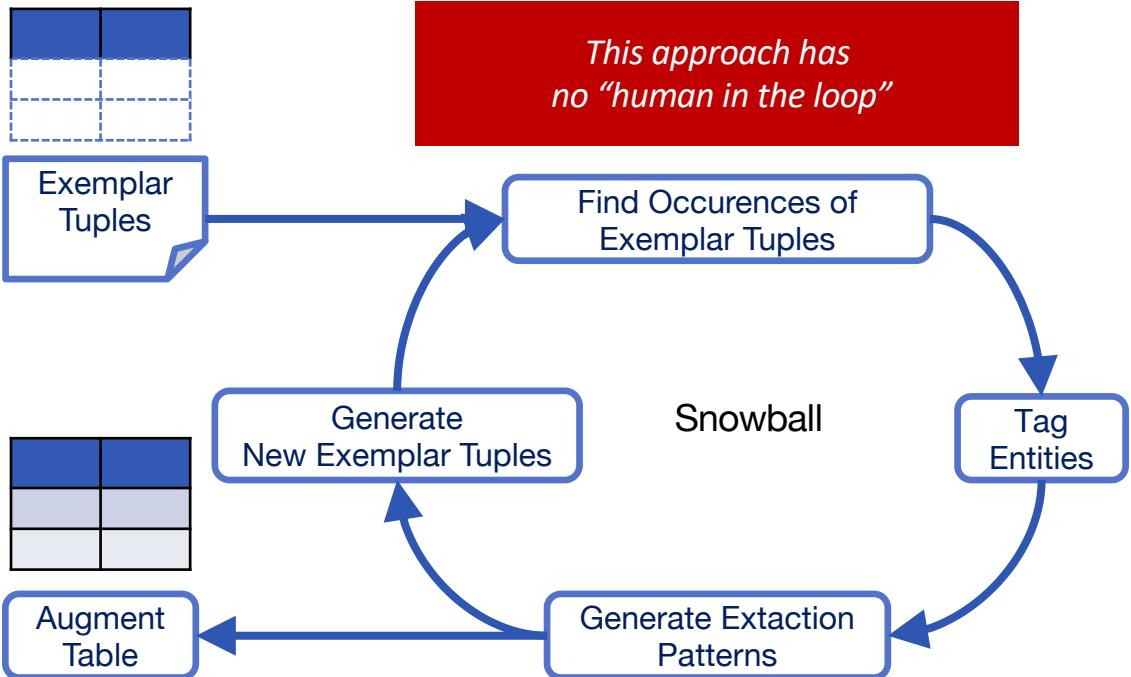
Agichtein and Gravano [2000]

*How to validate the new rules extracted automatically?*

1. **Compare extracted rules with known tuples:** confidence of R is based on how many known tuples extracts

2. **Compare extracted tuples with known rules:** confidence of T is based on how many known rules also extract T

New extracted Rules and Tuples should not create contradictions



# Entity-extraction by Example

Hanafi et al., [2017]

Learn extraction rules from example

Allow to match from text  
both Positive and Negative examples



**Goal:** Supervised Extraction

definition) increased 9.6 percent, the number of murders increased 6.2 percent, aggravated assaults increased 2.3 percent, the number of rapes (revised definition) rose 1.1 percent, and robbery violations were up 0.3 percent.  
Violent crime increased in all but two city groupings. In cities with populations from 50,000 to 99,999 inhabitants, violent crime was down 0.3 percent, and in cities with 500,000 to 999,999 in population, violent crime decreased 0.1 percent. The largest increase in violent crime, 5.3 percent, was noted in cities with 250,000



SEER

**Output:** Extraction rules

P: Percentage = 1.0 = 1.0  
D: {5, 6} = 0.4      D: {percent, %} = 0.4 = 0.4  
R: [0-9]+ = 0.2      D: {percent, %} = 0.4 = 0.3

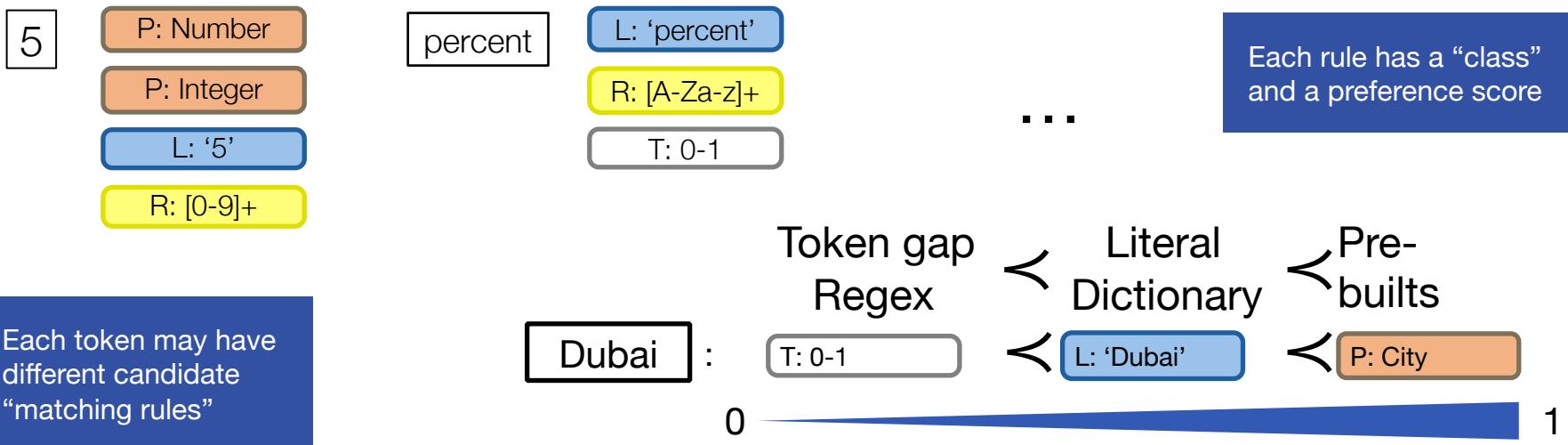
# Matching Rules

From string tokens to “semantics”

Hanafi et al., [2017]

Example: 5 percent up in Dubai

**Intuition:** Exploit a vocabulary of simple specialized patterns with known semantics

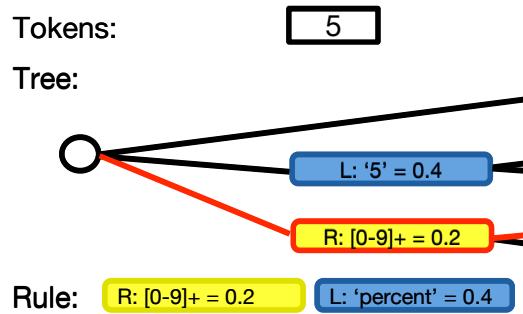


# Merging Rules

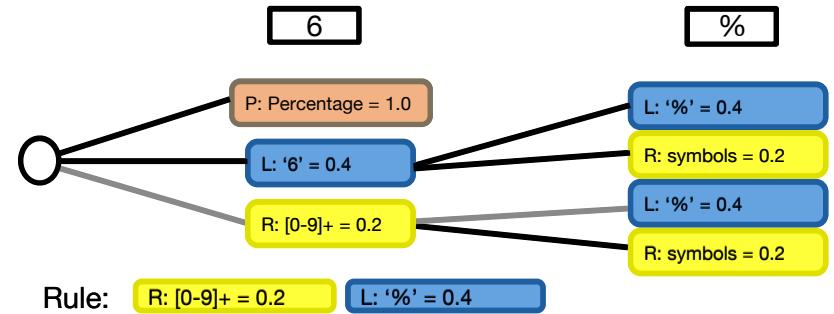
Hanafi et al., [2017]

Reconcile multiple interpretations

Example: 5 percent

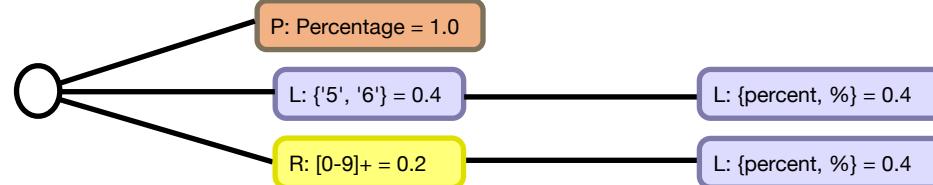


Example: 6 %



Intersection: [5 percent, 6%]

Consider also Negative Examples to prune candidates



# Web Tables

Semi-structured data on the web

<https://en.wikipedia.org/wiki/Denmark#Regions>

## Regions

The governing bodies of the regions are the [regional councils](#), each with forty-one councillors elected for four-year terms headed by regional district chairmen ([regionsrådsformanden](#)), who are elected by the council.<sup>[79]</sup> The areas of responsibility of the regional councils are the [national health service](#), [social services](#) and [regional development](#).<sup>[79][80]</sup> Unlike the counties they represent, the regions are allowed to levy taxes and the health service is partly financed by a national health care contribution until 2018 ([sundhedsafdelen](#)) from both government and municipalities.<sup>[18]</sup> From 1 January 2019 this contribution will be abolished, as it is being replaced by a general tax instead.

The area and populations of the regions vary widely; for example, the [Capital Region](#), which encompasses the Copenhagen metropolitan area with the exception of the subtracted province East Zealand but includes the [Baltic Sea](#) island of [Bornholm](#), has a population three times greater than that of [North Denmark Region](#), which covers the more sparsely populated area of northern Jutland. Under the county system created in 1970, the most densely populated municipalities, such as [Copenhagen Municipality](#) and [Frederiksberg](#), had been given a status equivalent to that of the regions, making them first-level administrative divisions. These [sui generis](#) municipalities were incorporated into the new regions under the 2007 reforms.

Danish name	English name	Admin. centre	Largest city (populous)	Population (January 2017)	Total area (km <sup>2</sup> )
Hovedstaden	Capital Region of Denmark	Hillerød	Copenhagen	1,807,404	2,568.29
Midtjylland	Central Denmark Region	Viborg	Aarhus	1,304,253	13,095.80
Nordjylland	North Denmark Region	Aalborg	Aalborg	587,335	7,907.09
Sjælland	Region Zealand	Sorø	Roskilde	832,553	7,268.75
Syddanmark	Region of Southern Denmark	Vejle	Odense	1,217,224	12,132.21

Source: [Regional and municipal key figures](#)

Google food calories

All Images Videos Maps News More Settings Tools

About 317.000.000 results (0,57 seconds)

Food Group	Carbohydrates (Grams)	Calories
Milk (higher % of simple carbohydrates; less nutrient dense)		
Chocolate milk (1 cup)	26	208
Low fat (2%) milk	12	121
Pudding (any flavor) (1/2 cup)	30	161
67 more rows		

[Carbohydrate and Calorie Content of Foods By Item | MomsTeam](#)  
<https://www.momsteam.com/nutrition/.../carbohydrate-and-calorie-content-of-foods>

Google country capitals

All Images Maps Videos News More Settings Tools

About 2.460.000.000 results (0,77 seconds)

According to [countries-of-the-world.com](#)  
View 40+ more

Country	Capital city
Albania	Tirana
Algeria	Algiers
Andorra	Andorra la Vella
Angola	Luanda
Angola	Luanda

List of world capitals

Country	Capital city
Albania	Tirana
Algeria	Algiers
Andorra	Andorra la Vella
Angola	Luanda

109 more rows

List of world capitals by countries  
<https://www.countries-of-the-world.com/capitals-of-the-world.html>

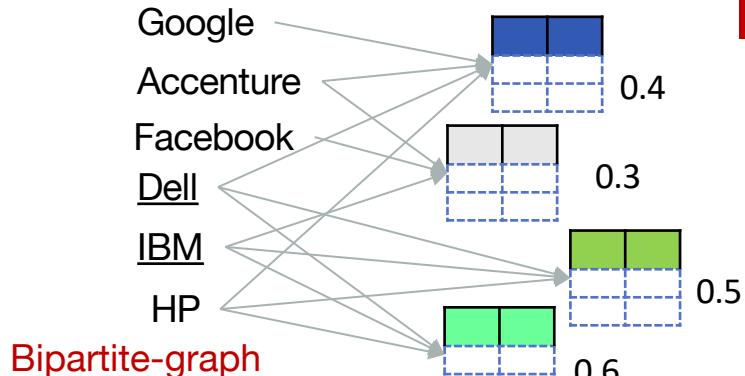


# Entity List Expansion

Wang et al. [2015]

Augmentation: identify entities to complete the list

1. Input: Incomplete list + Keyword query
2. Retrieve tables from paged based on the keyword query
3. Assign Score to tables based on relevance
4. Extract entity mentions from tables
5. Analyze Entity mention co-occurrence
6. Pick "co-occurring" Entities



*Problem: entities may appear together for different reasons*

Score Propagation

*Problem: Here PPR Causes concept drift*

Heuristic Propagation

**Goal:** Given some seed entity mentions, retrieve more entities of the same type

Incomplete table

IT Company
Dell
IBM
Lenovo
....?

Augmented table

IT Company
Dell
IBM
Lenovo
Apple
Samsung
HP
Acer



# Web-Table Completion

Identify relevant content, retrieve missing information

Yakout et al. [2012]

**Goal:** Retrieve missing attribute values

**Intuition:** If there is a structure, we can match it!

Model	Brand
S80	Benq
A10	
GX-1S	
T1460	

Incomplete table

Model	Brand	Part No	Mfg
S80	Nikon	DSC W570	Sony
Easyshare CD44	Kodak	DSC W570	Benq
DSC W570	Sony	T1460	Optio E60
Optio E60	Pentax	Optio E60	Pentax
		S8100	Nikon
		S8100	Nikon

## Web tables



Complete table

**Extra Input:** target attribute name or example of completing attribute



# Table Correlation Graph

Yakout et al. [2012]

Schema matching for web-page and web-tables  
Binary-relations only

## Determine Table Match

Direct Match between Q(K,A) and T(K,B)

K=entity name column

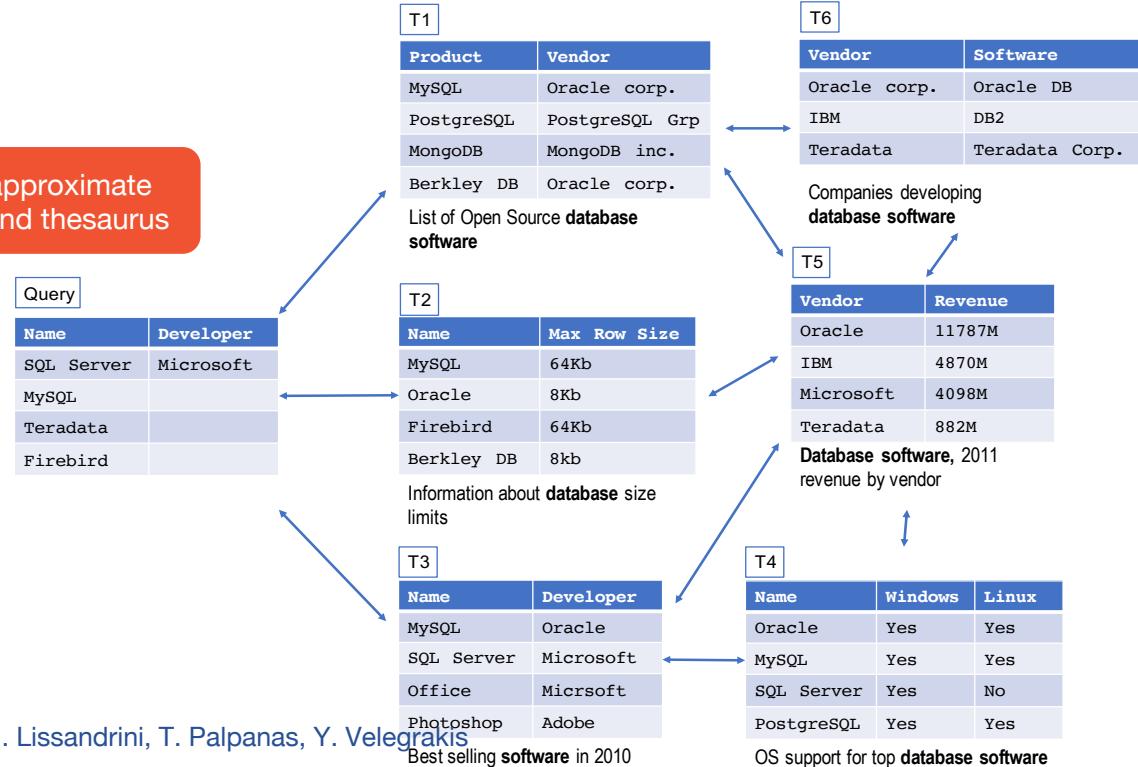
A,B = entity attribute column

$$S_{DMA}(T) = \begin{cases} \frac{|T \cap_K Q|}{\min(|Q|, |T|)} & \text{if } Q.A \approx T.B \\ 0 & \text{otherwise} \end{cases}$$

Can use approximate matching and thesaurus

**Problem:** considers only direct links between Q and T

**Goal:** Retrieve missing attribute values



# Table Correlation Graph

Yakout et al. [2012]

Schema matching for web-page and web-tables  
Binary-relations only

## Determine Table Match

Holistic Match

1. Assign Direct Match Score from Query to Tables
2. Scores >0 are starting nodes

## Build Classifier using

- Context similarity
- Table-to-content similarity
- URL similarity
- Tuples Similarity

the model predicts the match between two tables with a probability

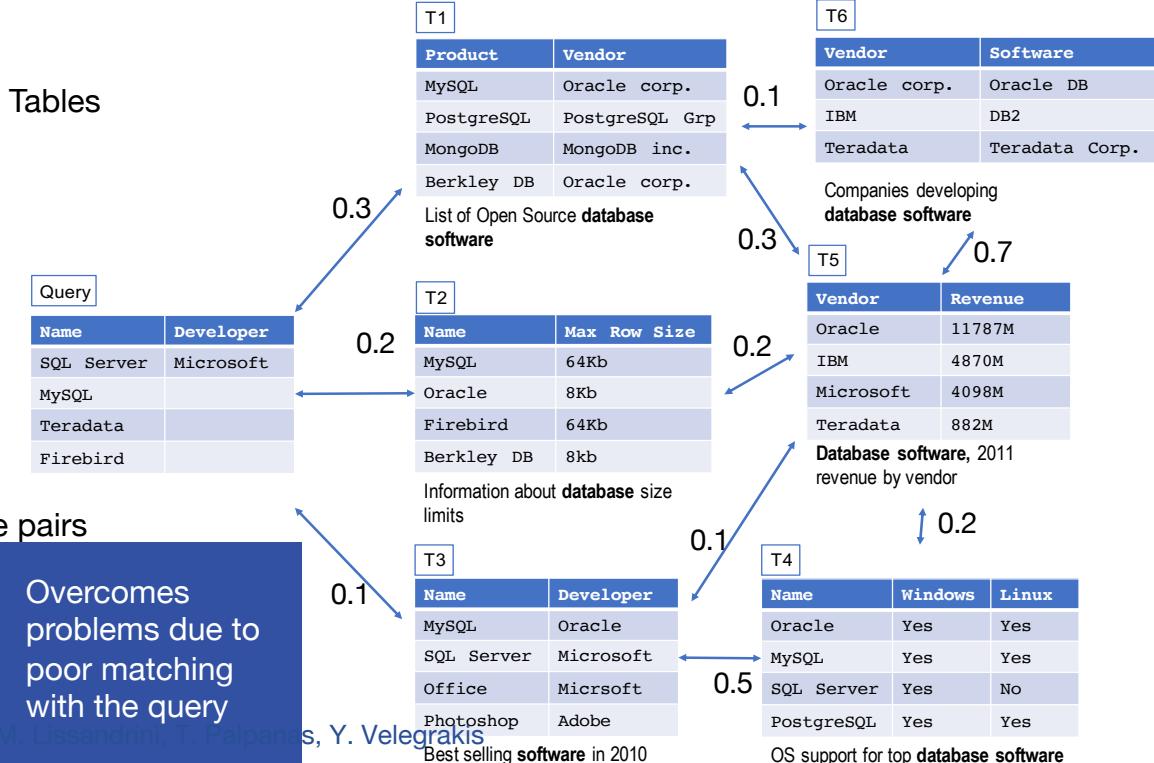
3. Use classifier to add weight to other table pairs
4. Use starting node and execute PPR
5. Use PPR scores to rank matching tables

Overcomes problems due to poor matching with the query

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

Best selling software in 2010

**Goal:** Retrieve missing attribute values

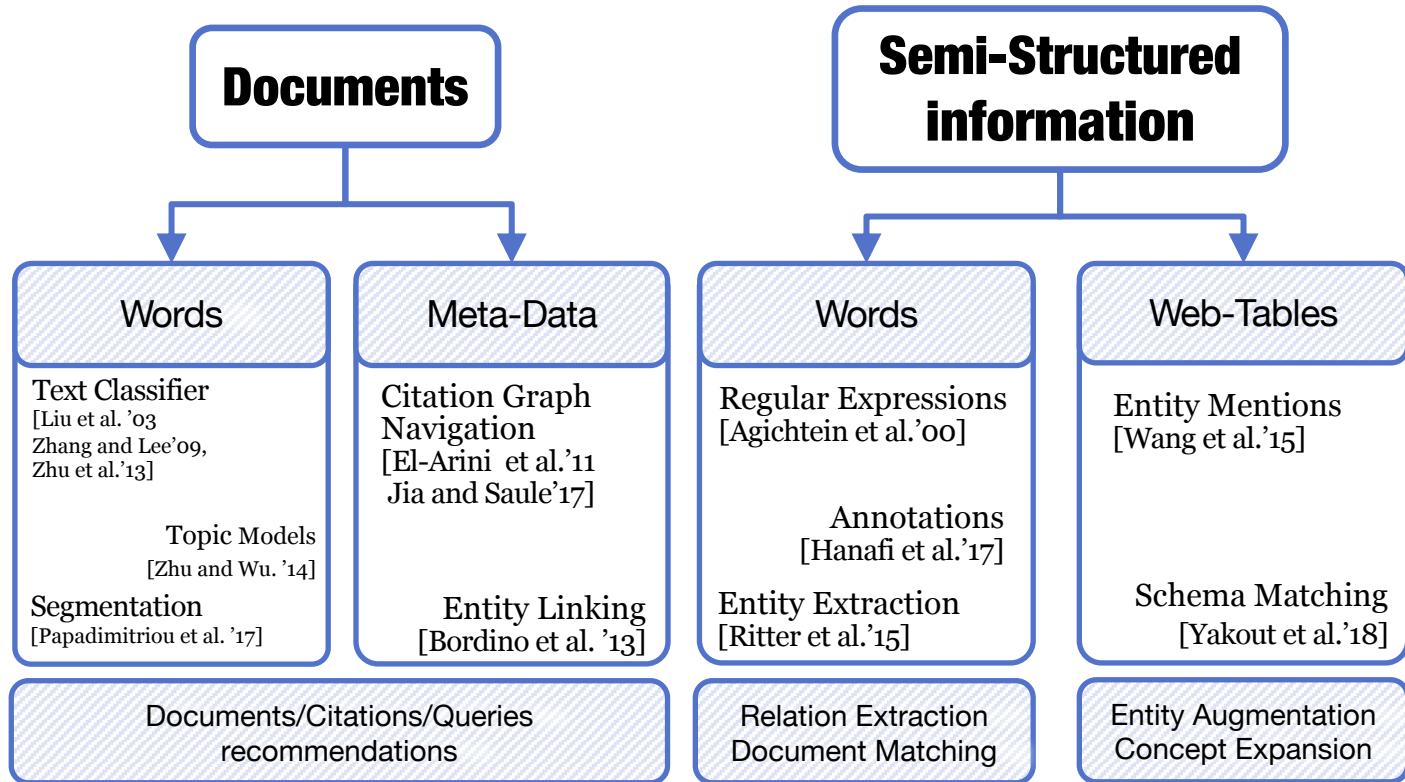


**SEARCHING FOR**

**BY LOOKING AT**

**APPLYING**

**PRODUCES**



<https://j.mp/ExploreSIGMOD>

# Where we are

Relational databases

Textual data

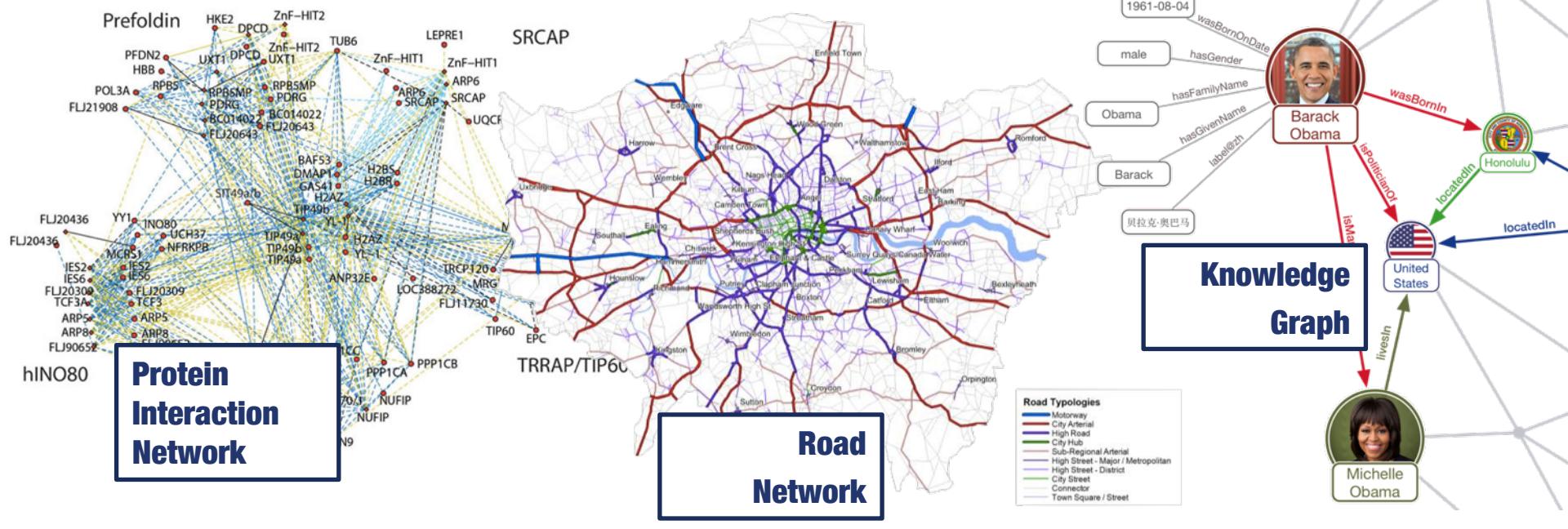


Graphs and networks



Machine learning

Challenges and Remarks



# Graphs are Everywhere

# Graphs

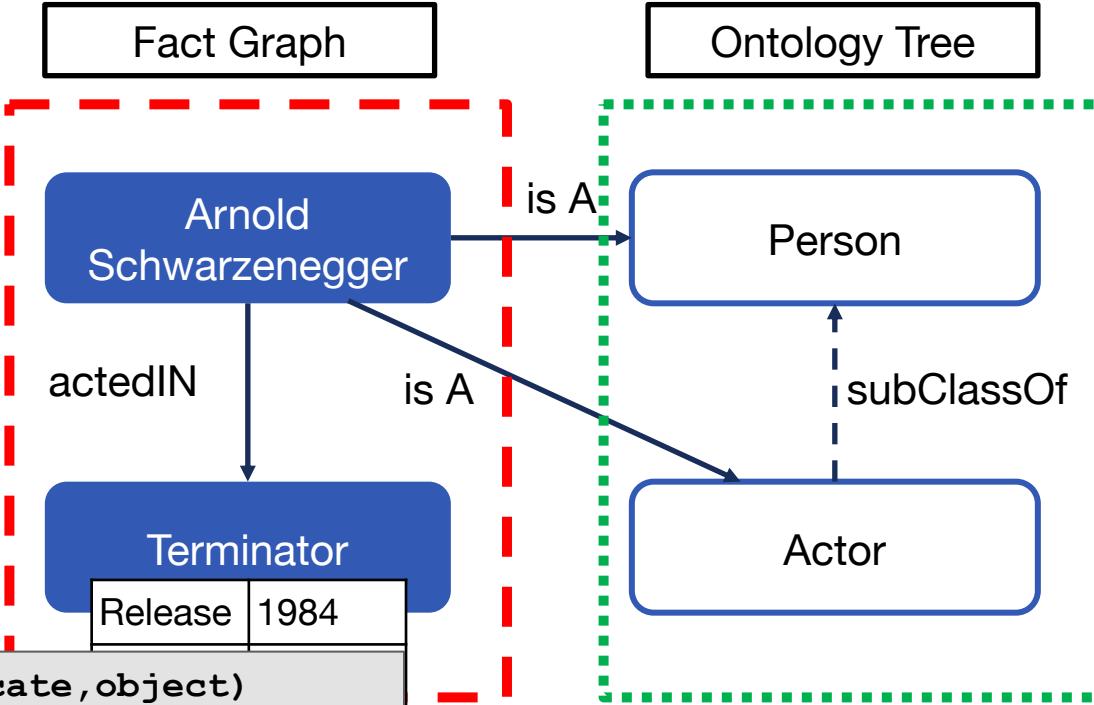
Connected Data

Edge-labelled  
Multigraphs  
 $G: \langle V, E, L, \ell \rangle$

Attributes:  
 $V/E: \langle \text{key}, \text{value} \rangle$

**RDF**  $(\text{subject}, \text{predicate}, \text{object})$

$(\text{Arnold\_Schwarzenegger}, \text{isA}, \text{Person})$   
 $(\text{Actor}, \text{subClassOf}, \text{Person})$   
 $(\text{Arnold\_Schwarzenegger}, \text{actedIn}, \text{Terminator})$



**The Structure of the Graph  
Is as important as the Data-values**

# Exemplar Queries

Mottin et al. [2014,2016]

Example-driven graph search

**Input:**  $Q_e$ , an example element of interest

**Output:** set of elements in the desired result set

Nodes/Entities  
Edges/Facts  
Structures

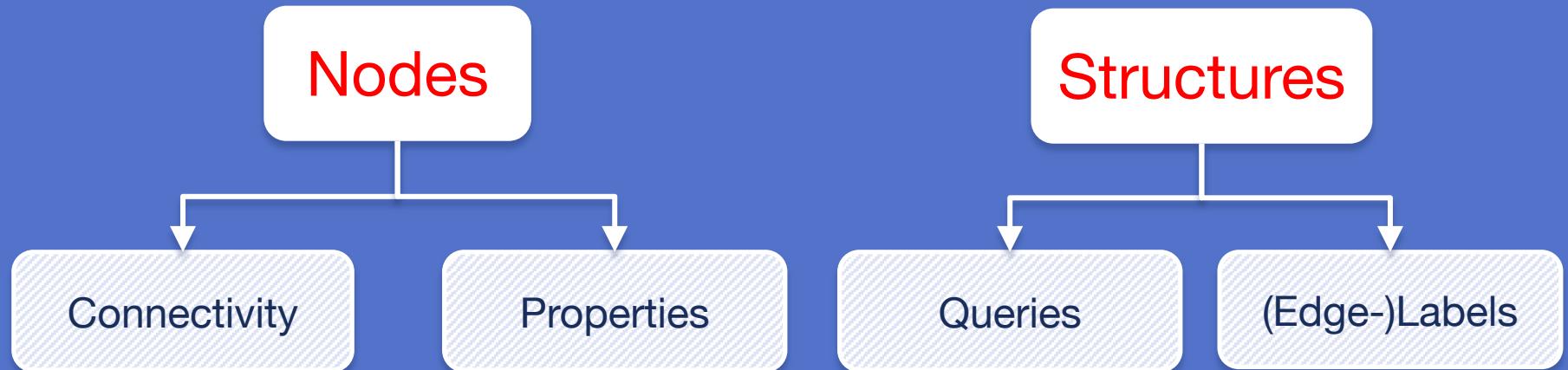
## Exemplar Query Evaluation

- evaluate  $Q_e$  in a database D, finding a sample S
- find the set of elements A **similar** to S given a **similarity relation**
- [OPTIONAL] return only the subset A<sup>R</sup> that are relevant

Usually requires an intermediate step:  
User input (keywords) → Element in the graph



# SIMILARITY for GRAPHS

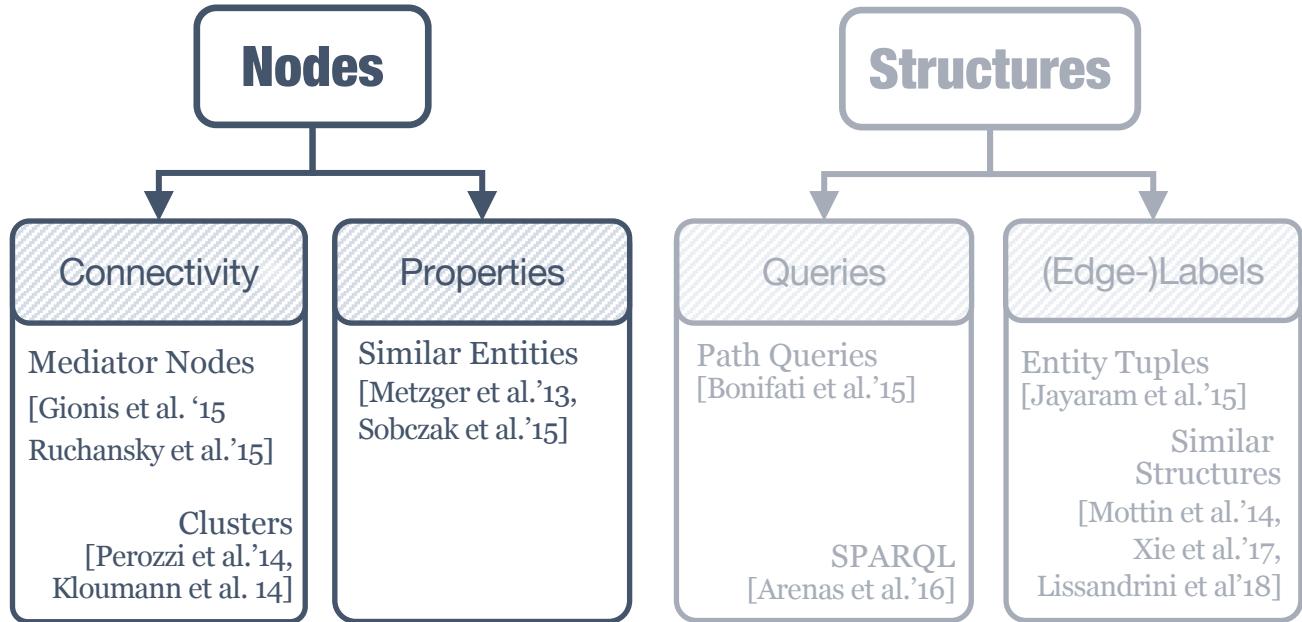


**CHALLENGE: DISCOVER USER PREFERENCE**

**CHALLENGE: EFFICIENT SEARCH**



# SEARCHING FOR BY LOOKING AT PRODUCES



# Seed Set Expansion

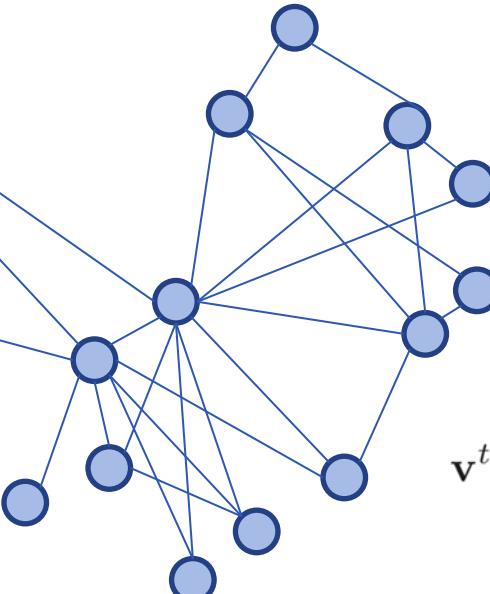
Kloumann and Kleinberg [2014]

Nodes connected  
by a community



Communities can be extremely large  
Identify “central nodes”  
or “the core subgraph”

Given a graph  $G$ , and a set of **query nodes**  $V_Q \subseteq V_G$ ,  
**retrieve all other nodes**  $V_C \subseteq V_G$ ,  
where  $C$  is a community in  $G$ , and  $V_Q \subseteq V_C$ .



Solution: PPR

$$\mathbf{v}^{t+1} = (1 - \alpha)\mathbf{M} \cdot \mathbf{v}^t + \alpha\mathbf{v}^0$$

# The Minimum Wiener Connector Problem

Ruchansky et al. [2015]

**Model:** Unlabeled Undirected Graph

**Query:** A set of Nodes  $Q$

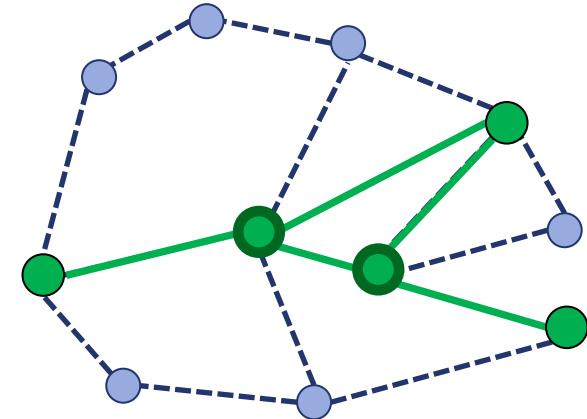
**Similarity:** Shortest-Path distance

**Output:** A Set of Connector Nodes  $H$

“explains” connections in  $Q$

Connectors:  
Nodes with HIGH closeness  
to ALL the inputs

Similar to a Steiner-Tree but  
overall pairwise distances are optimized



Case: Infected Patients  
→ Culprit/Other Infected

Case: Target Audience  
→ Influencers



# The Minimum Wiener Connector Problem

Ruchansky et al. [2015]

**Model:** Unlabeled Undirected Graph

**Query:** A set of Nodes Q

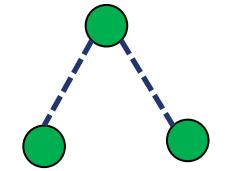
**Similarity:** Shortest-Path distance

**Output:** A Set of Connector Nodes H

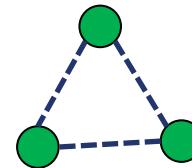
“explains” connections in Q

minimize the sum of pairwise  
shortest-path-distances  
between nodes in the connector H

Called: Wiener Index.  
tradeoff between size  
and average distance



$$W=1+2+1 = 4$$



$$W=1+1+1 = 3$$

Sometimes The Best  
Solution is  
NOT A Tree

NP-Hard

$$\min \sum_{(u,v) \in H} d(u, v)$$

d(u, v) is the shortest-path distance



# Approximate minimum Wiener Index Connector

Ruchansky et al. [2015]

CHOOSE  $r \in Q$  &  $\lambda \in [1, \log_{(1+\beta)} |V|]$

All Pairwise Distances

→ Distances from a root  $r$

Measure distance in  $H$

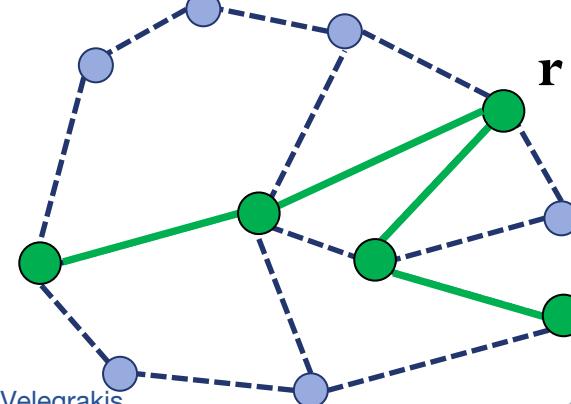
→ Precomputed distance in  $G$

Edge Weights

$$w(u, v) = \lambda + \frac{\max\{d_G(r, u), d_G(r, v)\}}{\lambda}$$

Approximated with  
Edge-Weighted SteinerTree

Enumerate Candidate Solutions  
for  $r \in Q$  &  $\lambda$   
and keep best tree



# Focused Clustering and Outlier Detection

Similarity based on attributes

**Model:** Unlabeled Undirected Graph with Node Attributes

**Query:** A set of Nodes Q

**Similarity:** To Be Inferred

*based on Attribute Values & Connectivity*

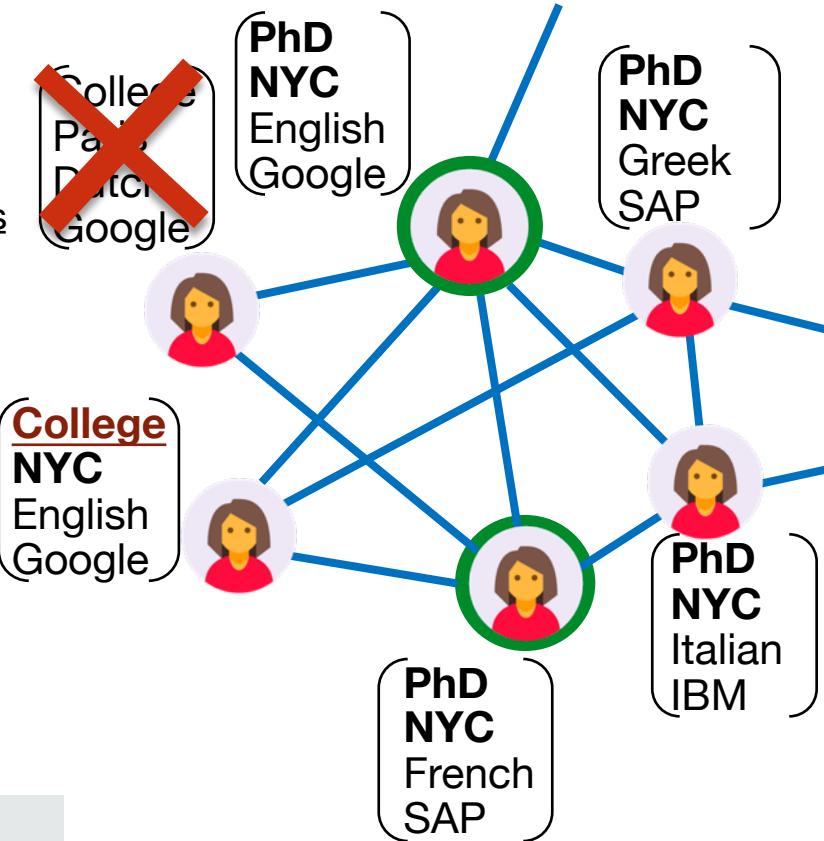
**Output:** Clusters of Nodes: Dense & Coherent  
+ Outliers

Case: Target Users → Community with same interests

Case: Products → Co-purchased products with similar features

D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

Perozzi et al. [2014]



# Focused Clustering

Infer User Focus

TASK: Infer “FOCUS”, important attributes

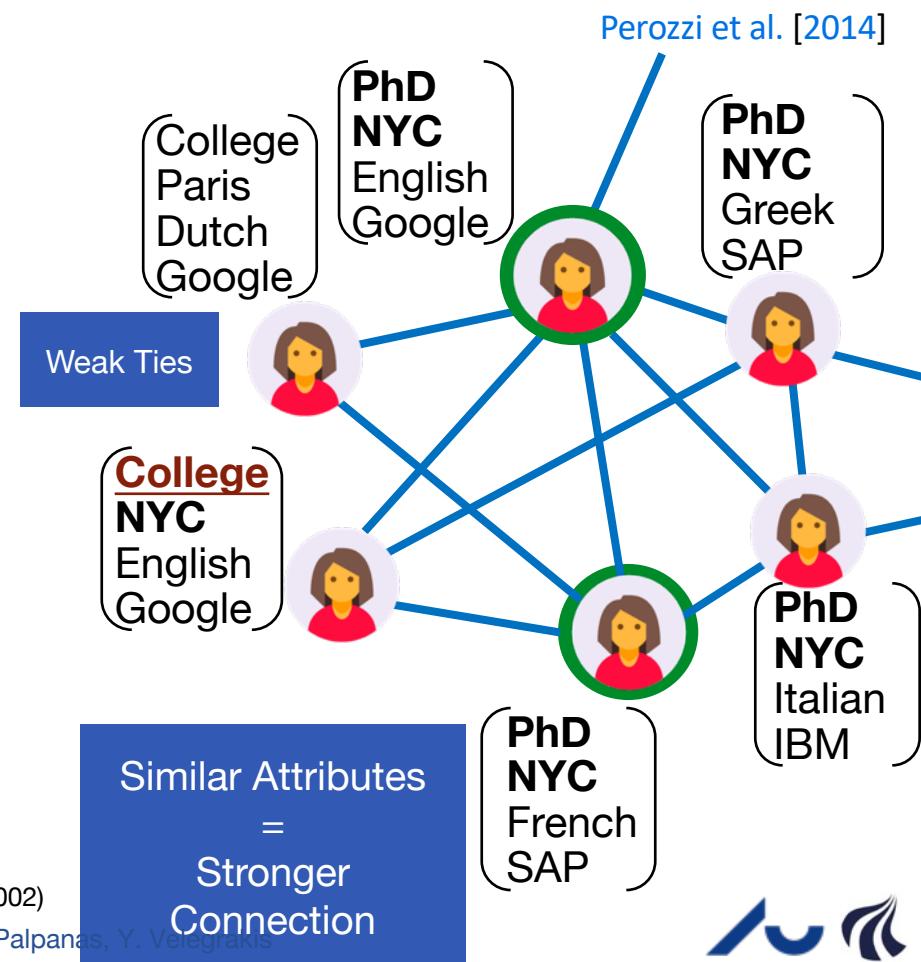
attribute weights  $\beta$

$$\begin{pmatrix} \text{PhD} \\ \text{NYC} \\ \text{English} \\ \text{Google} \end{pmatrix} \quad \begin{pmatrix} \text{PhD} \\ \text{NYC} \\ \text{French} \\ \text{SAP} \end{pmatrix} \xrightarrow{\hspace{1cm}} \begin{pmatrix} 0.5 \\ 0.5 \\ 0 \\ 0 \end{pmatrix}$$

1. Set of similar pairs, PS (from Q)
2. Set of dissimilar pairs, PD (random sample)
3. Learn a distance metric between PS and PD

$$\min_{\mathbf{A}} \sum_{(u,v) \in P_S} (f_i - f_j)^T \mathbf{A} (f_i - f_j) - \gamma \log \left( \sqrt{(f_i - f_j)^T \mathbf{A} (f_i - f_j)} \right)$$

( Distance Metric Learning, inverse Mahalanobis distance: Xing, et al 2002)



# Focused Clustering

Perozzi et al. [2014]

Prune the Graph and keep dense communities

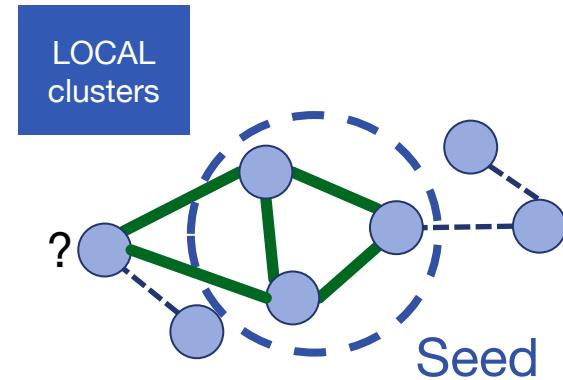
## TASK: Extract Clusters on Focused Graph

attribute weights  $\beta \rightarrow$  Edge Weight

### 1. Find Starting Set of Small Candidate Clusters

1.a Drop low-weight edges

1.b Extract Strongly Connected Component  $C_1, C_2, \dots$



### 2. Grow Clusters around Candidates

2.a Compute conductance of  $C$ :  $\phi^{(w)}(C, G)$

2.b Select node to add to  $C'$ : best improvement to  $\Delta\phi^{(w)}(C, C')$  (greedy)

2.c Prune Underperforming nodes

### 3. Detect Outliers: High unweighted conductance

w.r.t. low weighted conductance

**Weighted Conductance:**  
ratio between the weighted sum of edges crossing the boundaries of the cluster and the weighted sum of those residing within it.

**Performant Strategy:**  
Start with local solution and expand around them to avoid complete scans of the graph

# iQBEEs: Entity Search by Example

Knowledge Graph Search

**Model:** Knowledge Graph (Edge-labels)

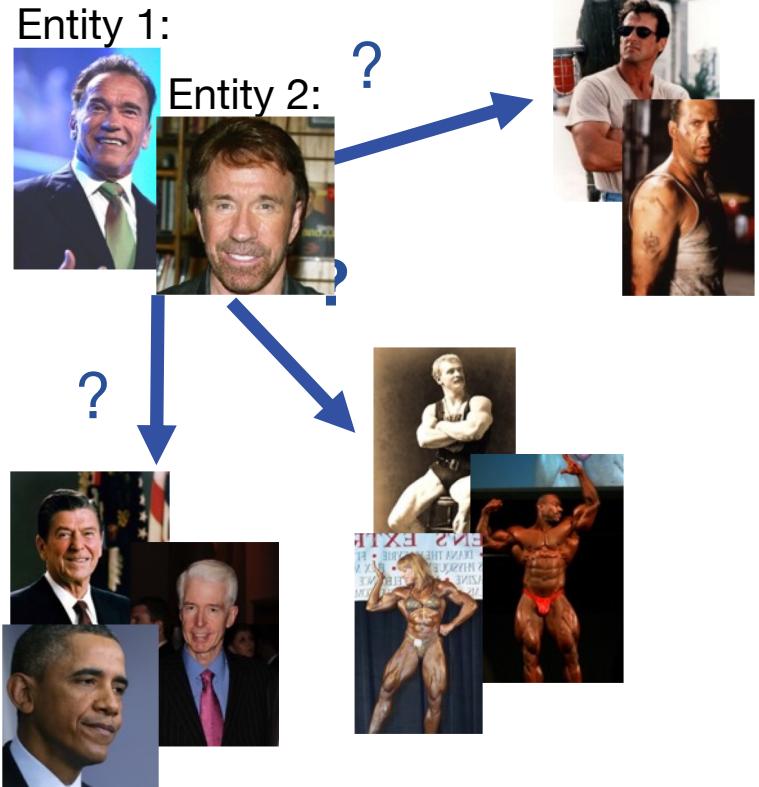
**Query:** A set of Entities Q

**Similarity:** Shared semantic properties

**Output:** A Set of Similar Entities (ranked)

Case: Products → Products with similar aspects

Case: Social Media → User recommendation



# Maximal Aspects

Selecting Features of Entity Similarity



?x sport BodyBuilding  
?x type AmericanActor



Is not maximal if  
Adding any aspect  
 $\rightarrow E(A)=\{\text{Arnold}\}$

**1. Prune  
generic  
aspects**

?x type AmericanActor  
?x governorOf California



Include  
Typical Types

**2. Rank  
Set of  
aspects**

?x hasHeight 1.88m  
?x type Entity



Use most  
Specific Type

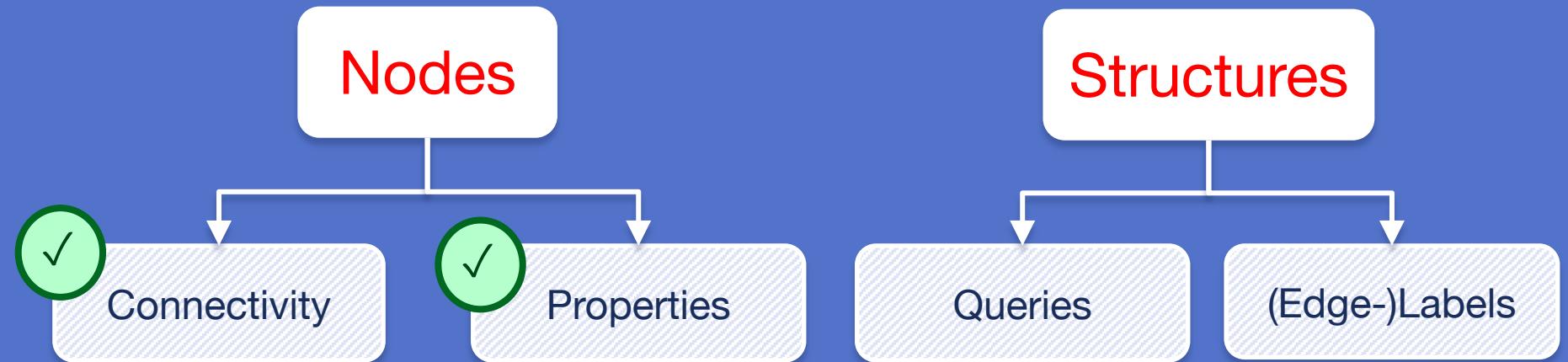
?x type AmericanActor  
?x actedIn TheExpendables  
?x type ActionActor



REPEATABLE  
Update Q

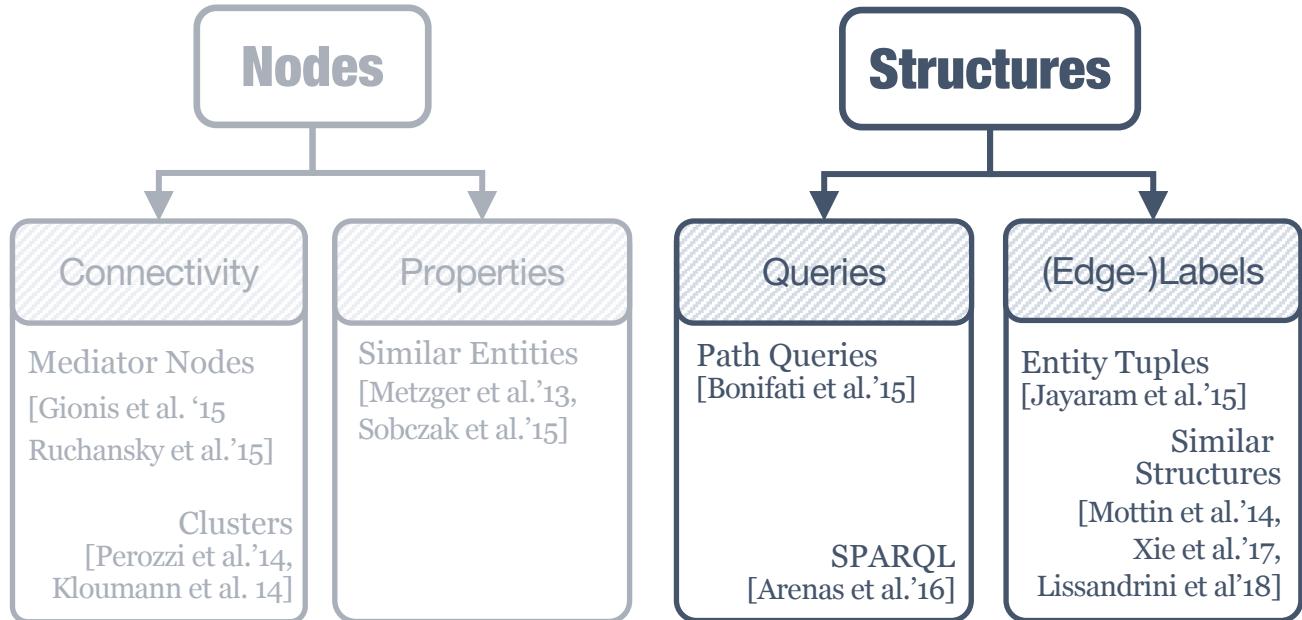


# SIMILARITY for GRAPHS



**Queries can retrieve  
both Nodes and Structures**

# SEARCHING FOR BY LOOKING AT PRODUCES



# Learning Path Queries on Graphs

Queries from Examples

**Model:** Edge Labeled Graph

**Query:** 2 sets of Entities  $Q^+$ ,  $Q^-$   
Positive, Negative

**Similarity:** Common Path Query (RegExp)

$$q := \epsilon \mid a(a \in \Sigma) \mid q_1 + q_2 \mid q_1 \cdot q_2 \mid q^* \\ (\text{bus}|\text{tram})^*+ \text{Cinema}$$

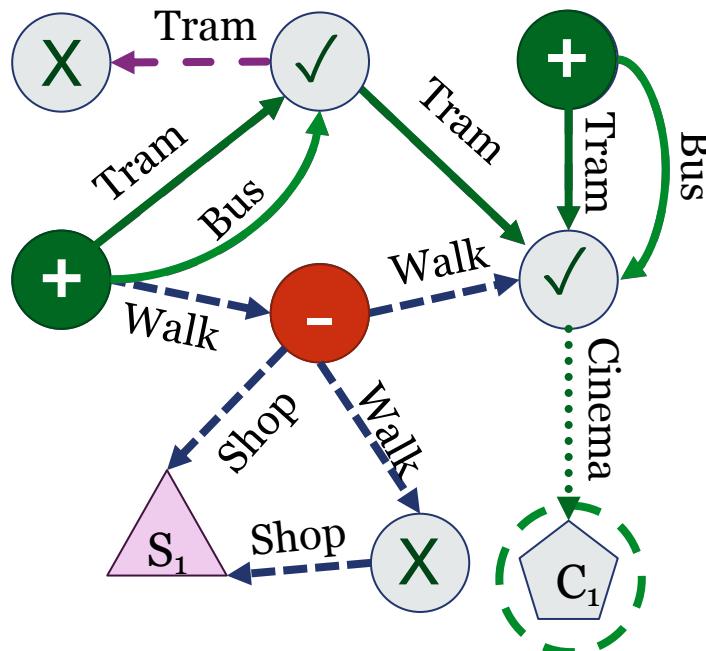
**Output:** Set of Nodes satisfying paths for  $Q^+$

but not paths for  $Q^-$

Case: Proteins → Similar interactions/co-expression

Case: Tasks Initiator → Similar Processes/Behaviours

Bonifati et al. [2015]



MONADIC: only starting nodes  
extensible to

BINARY/ N-ARY : path from X to Y

# Learnability of Path Queries

Bonifati et al. [2015]

When is possible and How

Query: 2 sets of Entities  $Q^+, Q^-$

Sometimes Positive & Negative Examples Cannot be reconciled!

## Consistency:

1. Select Smallest Consistent Path

$$\forall v \in Q^+. paths_G(v) \not\subseteq paths_G(Q^-)$$

2. Loops cause infinite paths? Fix Maximal Length K

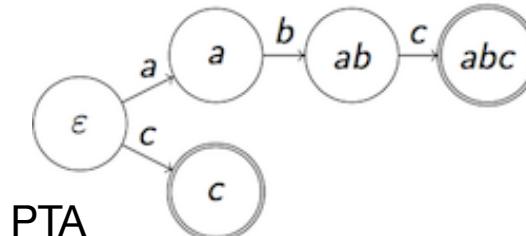
When to use Kleene star \* ?

$$C \mid (A.B.C) \rightarrow (A.B)^*.C$$

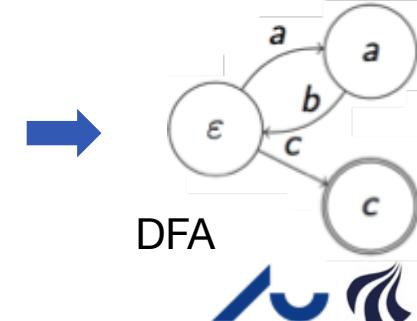
3. Generalize SCP

- a) Construct Prefix Tree Acceptor
- b) Generalize into DFA with Merge

Can be INTERACTIVE! The system presents to the user nodes to label as Positive/Negative



PTA



DFA



# Reverse engineering SPARQL queries

Arenas et al. [2016]

Knowledge Graph Search

**Model:** Knowledge Graph (Edge-labels)

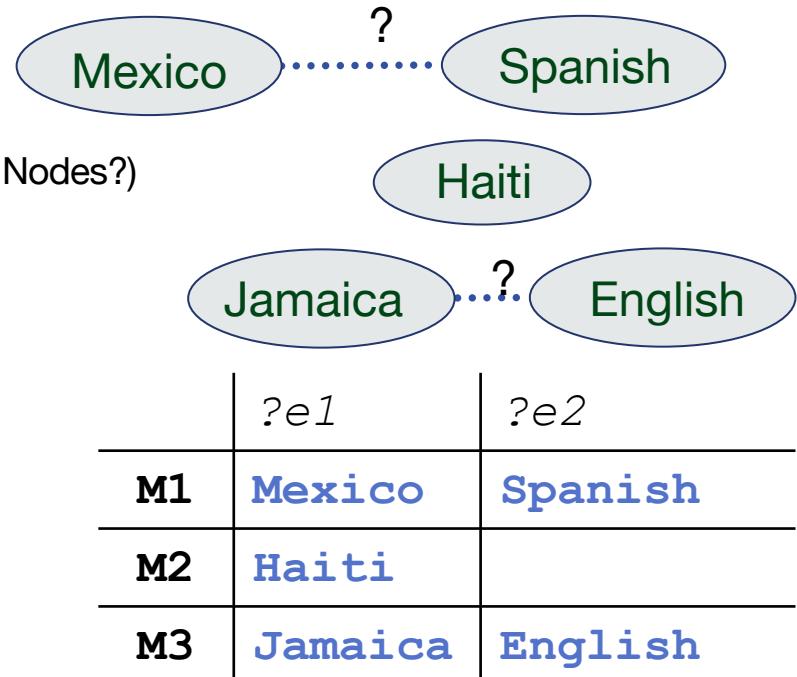
**Query:** Set of Answers → Not Graphs but Tuples (of Nodes?)

**Similarity:** common AND/OPT/FILTER query

**Output:** a SPARQL query / query results

Case: Open Data → Query Unknown Schema

Case: Novice User → Avoid SPARQL



**MATCH** (?X, is\_a, Country)  
**OPT** (?X, has\_language, ?Y)



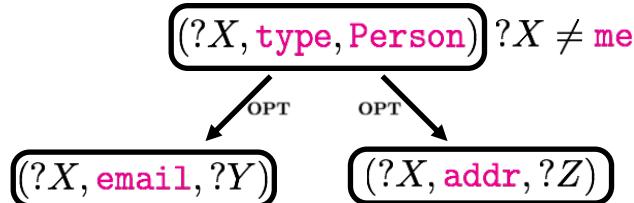
# Reverse engineering SPARQL queries

Arenas et al. [2016]

Challenges and Complexity

**Query:** Set of Variable Mappings

	?X	?Y	?Z
M1	John		
M2	Mary	mary@email.eu	
M3	Lucy		Roses Street



Incomplete Mappings are  
treated as OPTIONAL  
Typical of RDF queries

Enumerate all possible  
SPARQL queries satisfied  
by the mappings

INTRACTABLE  
 $\Sigma_2^p$ -complete

Build tree-shaped  
SPARQL queries IMPLIED  
by the mappings



# Reverse engineering SPARQL queries

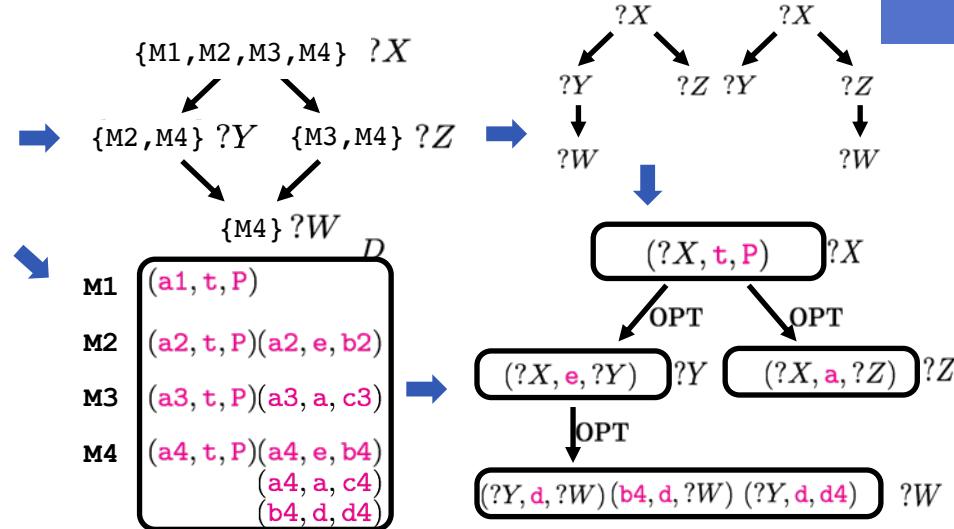
Arenas et al. [2016]

Challenges and Complexity

Query: Set of Variable Mappings  $\Omega$

	$\Omega$			
	?X	?Y	?Z	?W
M1	a1			
M2	a2	b2		
M3	a3		c3	
M4	a4	b4	c4	d4

Greedy: keep just enough to cover all variables



- 3 Instantiations:
1. Only Positive Examples
  2. Positive & Negative
  3. Exact Result only

# Graph Exemplar Queries

Mottin et al. [2016]

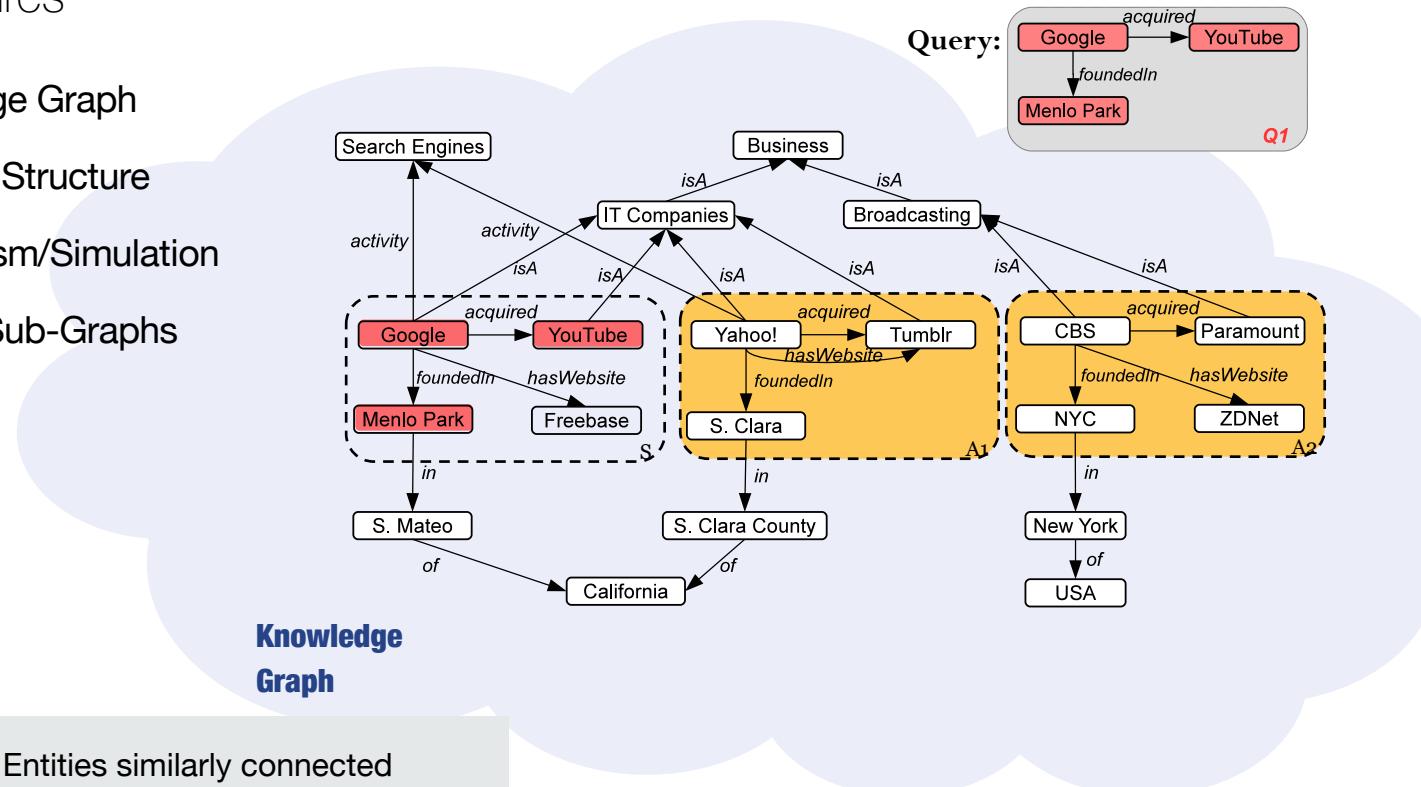
Search for Structures

**Model:** Knowledge Graph

**Query:** Example Structure

**Similarity:** Isomorphism/Simulation

**Output:** A set of Sub-Graphs

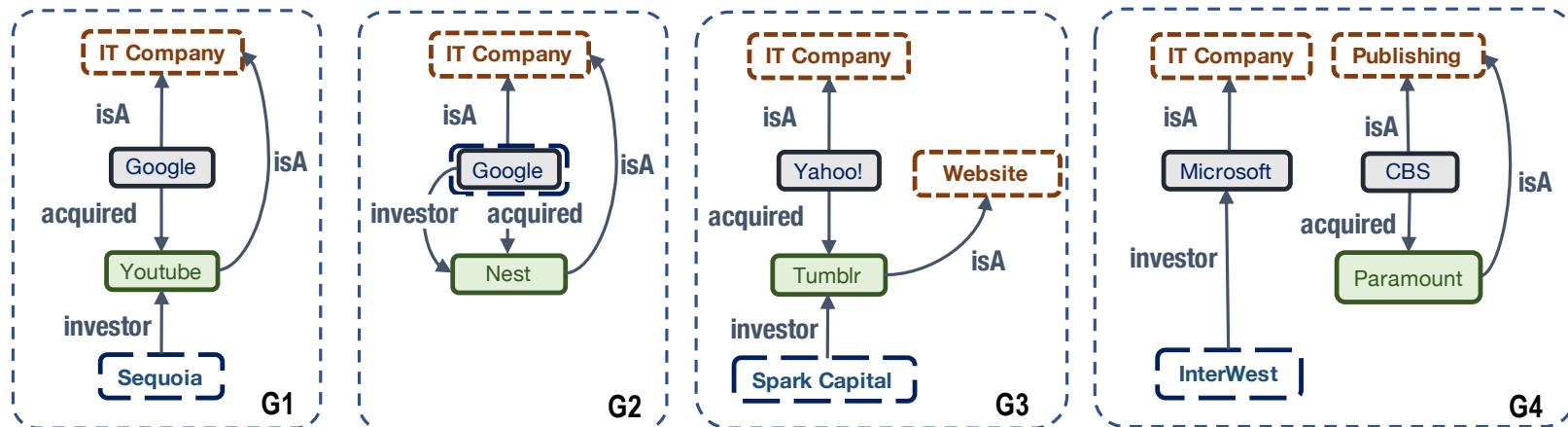


Case: Rich Schema → Find Entities similarly connected

# Graph Isomorphism vs. Simulation Variants

Structural Congruence/Similarity

Isomorphism requires an bijective function  
Simulation requires only a surjective relation  
Preserves only Parent → Child relationships



Example of Simulating ( $G_1 \sim \{G_2, G_3, G_4\}$ ) and Strong Simulating Graphs ( $G_1 \approx G_2$ )

Strong simulation: Capturing topology in graph pattern matching  
– Shuai Ma et al., 2014

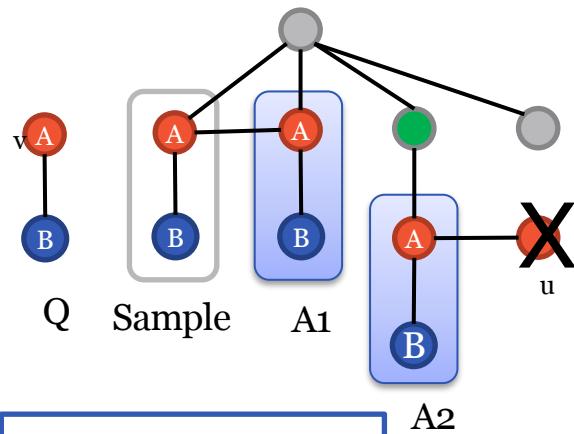


# Computing Exemplar Queries (i)

Mottin et al. [2016]

Fast Structure Matching

Reduce Search Space:  
Removes nodes that cannot be part of a solution



NP-complete  
(subgraph isomorphism)  
 $O(|V|^4)$  (simulation)

**Exact Pruning technique:**

- Compute the neighbor labels of each node
- Prune nodes not matching query nodes neighborhood labels
- Apply iteratively on the query nodes

$$W_{n,a,i} = \{n_1 | l(n_1, n_2) = a \vee n_2 \in N_{i-1}(n)\}$$

neighborhood ( $v$ ) =  $\{(B,1)\}$   
neighborhood ( $u$ ) =  $\{(A,1)\}$

No Match

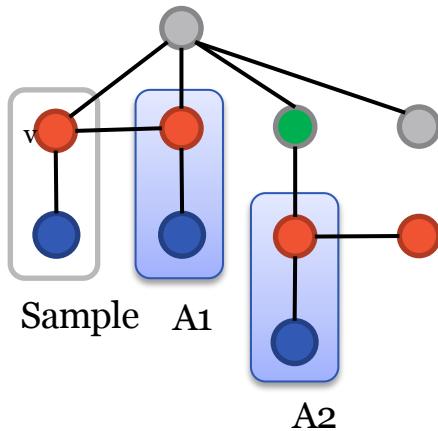


# Computing Exemplar Queries (ii)

Mottin et al. [2016]

Prune Irrelevant Answers

Reduce Search Space:  
Removes nodes that are likely to be less  
relevant



NP-complete  
(subgraph isomorphism)  
 $\Theta(|V|^4)$  (simulation)

## Approximation:

- Nodes closer to the sample are more important
- Use **Personalized PageRank** with a weighted matrix

$$\mathbf{v} = (1 - c)A\mathbf{v} + cp$$

- Weight edges: frequency of the edge-label

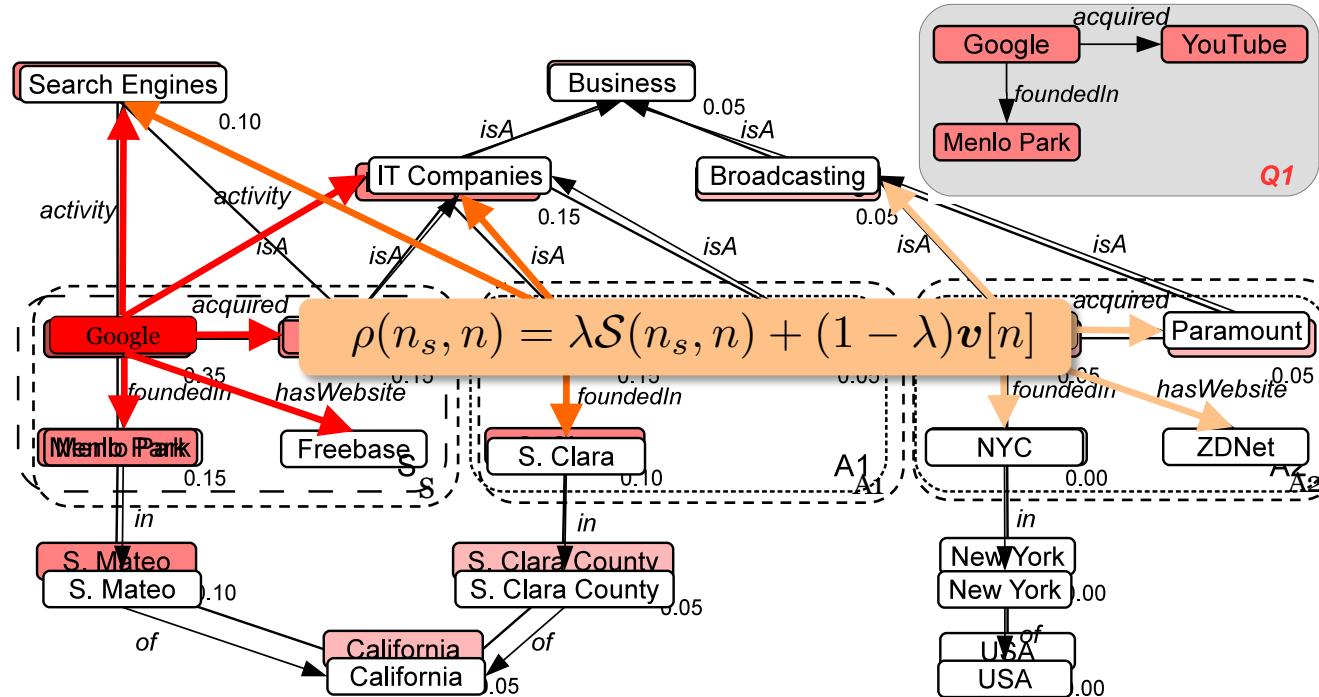
$$I(e_{ij}^\ell) = I(\ell) = \log \frac{1}{P(\ell)} = -\log P(\ell)$$

$$P(\ell) = \frac{|E^\ell|}{|E|}$$

# Ranking Results

Score Relevance of Answers

Mottin et al. [2016]



## Combination of two factors

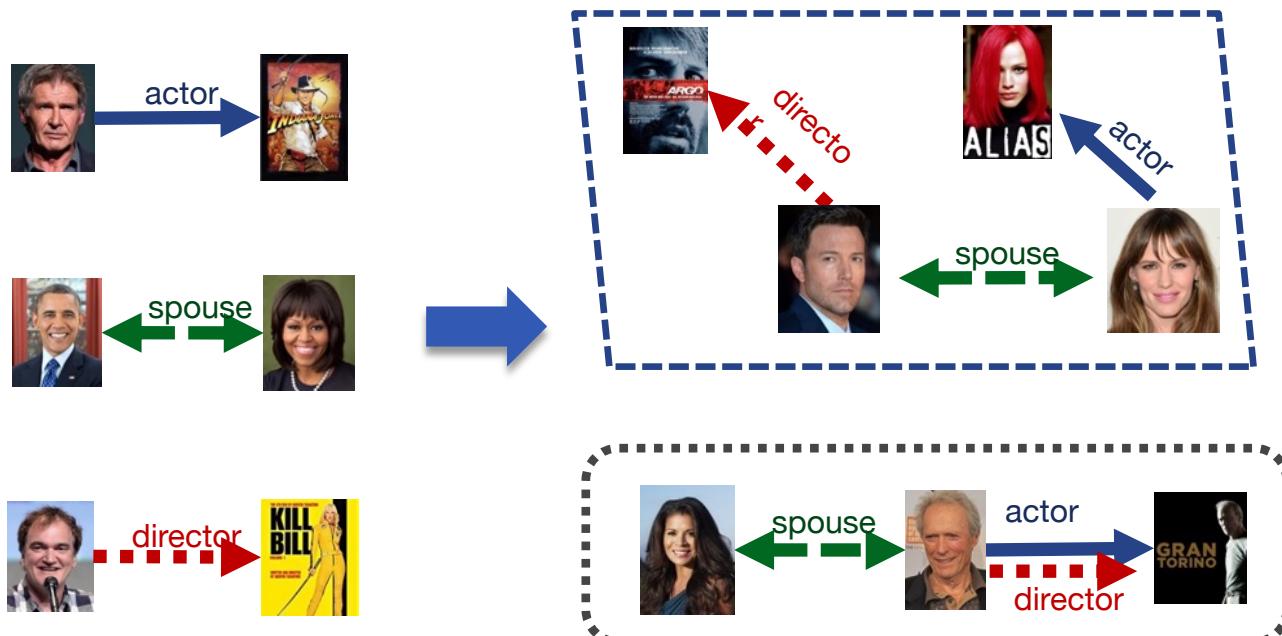
1. Structural: similarity of two nodes in terms of neighbor relationships
2. Distance-based: the PageRank already computed



# Search with Multiple Examples

Lissandrini et al. [2018]

Combining partial answers



- Multiple Simple Examples
- Each Example describes an Aspect
- Results are Combinations of aspects
- Results have possibly Multiple Structures

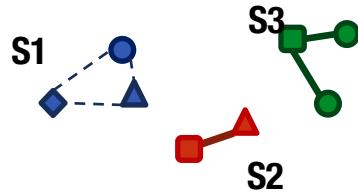
Case: Unknown Structures → Find Complex Connections with Simpler Components



# Search Framework

Lissandrini et al. [2018]

Pruning and Partial matching



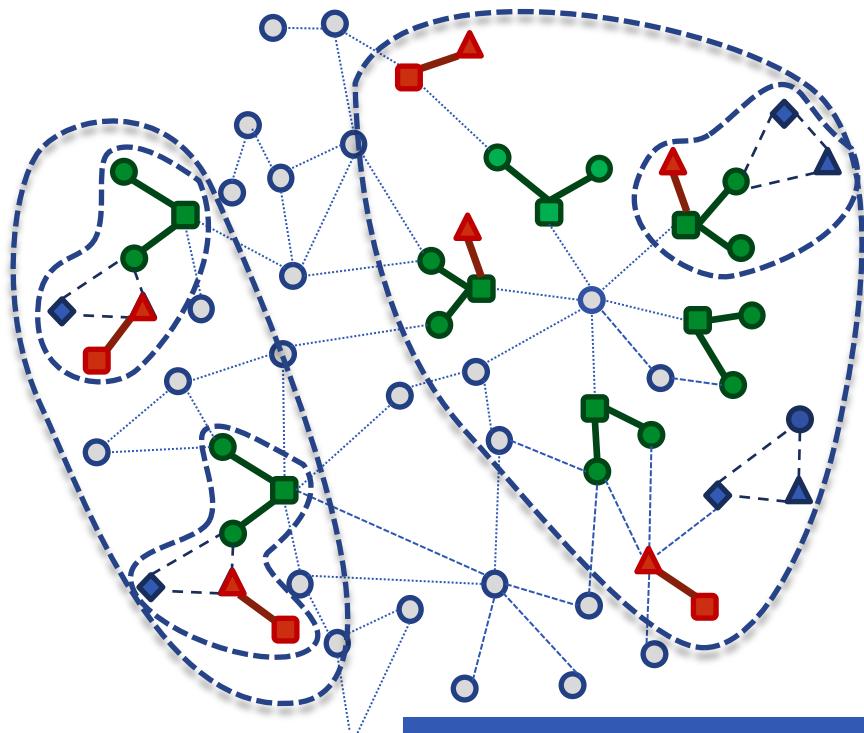
## Multi-exemplar Answering

**Input:** Database  $G : \langle V, E, \ell \rangle$

**Input:** Samples  $S : \langle s_1, \dots, s_m \rangle$

**Output:** Answers  $\mathcal{A}$

- 1:  $\mathcal{G} \leftarrow \text{PARTIAL}(G, S)$  ←
- 2:  $\mathcal{A} \leftarrow \text{SEARCH}(\mathcal{G}, S)$
- 3: **return**  $\mathcal{A}$



Exploit Localized Search

# Graph Query by Example(GQBE)

Jayaram et al. [2015]

Search for example Tuples

Model: Knowledge Graph

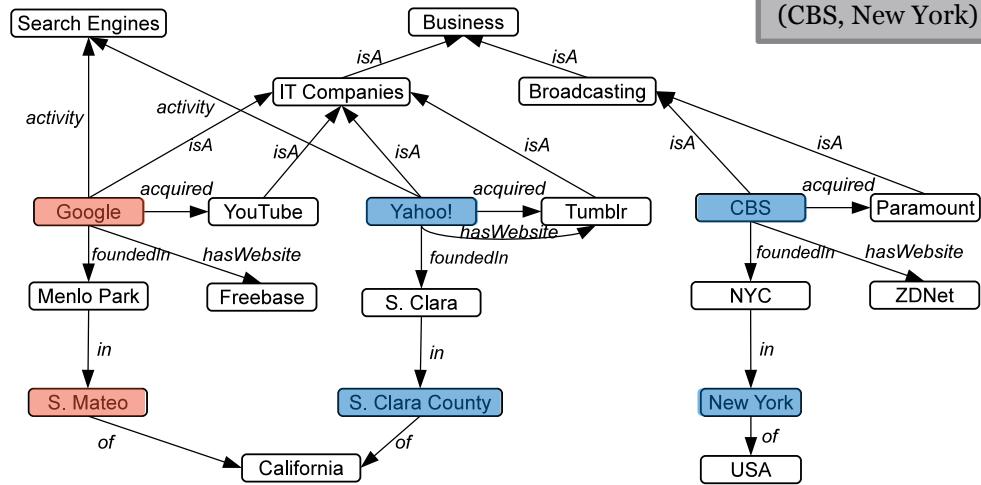
Query: Entity Tuples

Similarity: ~Isomorphism

Output: A set of Tuples

In GQBE Input is a set of (disconnected) entity mention tuples

$Q = (\text{Google}, \text{S. Mateo})$   
Results =  
 $(\text{Yahoo}, \text{S. Clara})$   
 $(\text{CBS}, \text{New York})$



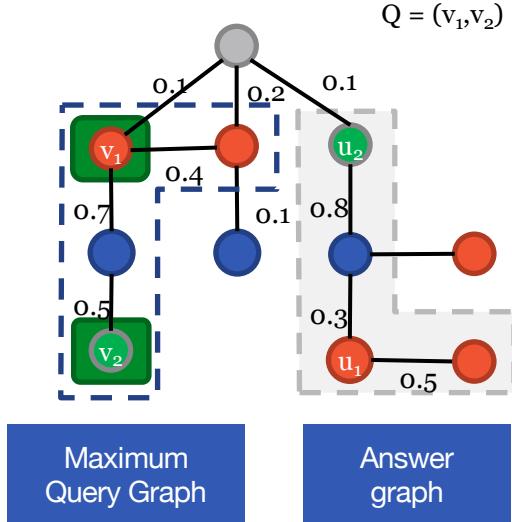
Case: Known Entities+Unknown Connections → Find Complex Connections



# GQBE: Maximum Query Graph

Jayaram et al. [2015]

Understand the connections implied by the tuples



1. Find the maximum query graph
  - Graph with M edges having the maximum weight

2. Answers subgraph-isomorphic to the query graph

NP-hard

3. Return top-k

Answer score:

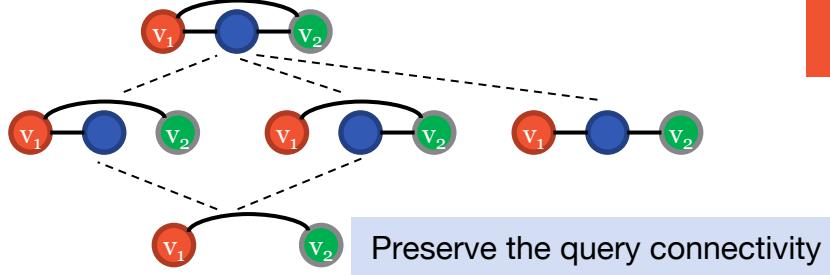
- Sum of query graph weights
- Similarity match between edges in the answer and the query (shared nodes take extra credit)

$$\text{match}(e, e') = \begin{cases} \frac{w(e)}{|E(u)|} & \text{if } u=f(u) \\ \frac{w(e)}{|E(v)|} & \text{if } v=f(v) \\ \frac{w(e)}{\min(|E(u)|, |E(v)|)} & \text{if } u=f(u), v=f(v) \\ 0 & \text{otherwise} \end{cases}$$

# GQBE: Multiple Query Tuples

Understand the connections implied by the tuples

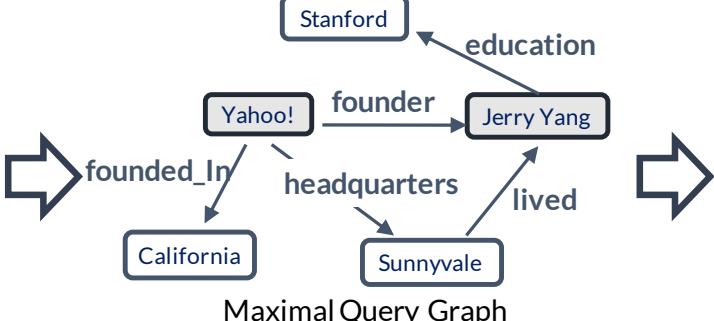
Subgraphs of  
Maximum  
Query graph



Maximum  
Query Graph  
is Very Large

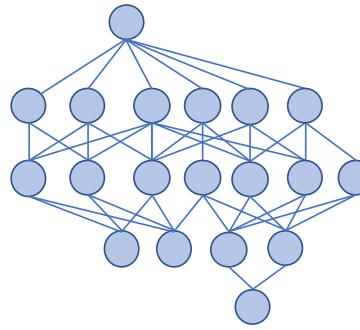
Preserve the query connectivity

(Jerry Yang, Yahoo!)



Entity Tuple

Maximal Query Graph



Query Lattice

Jayaram et al. [2015]

Find answers using a lattice obtained removing edges from the union graph

GQBE finds answers for multiple query tuples

Compute a re-weighted union graph of the individual query graphs

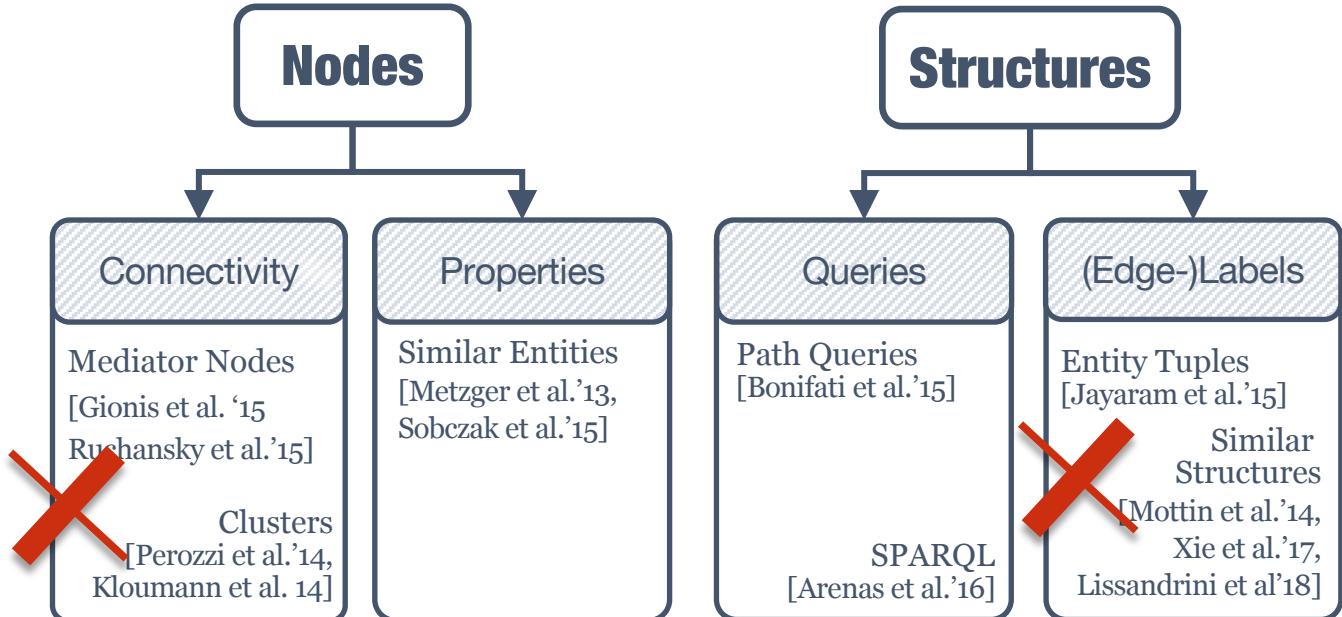
$\langle \text{David Fillo, Yahoo!} \rangle$   
 $\langle \text{Bill Gates, Microsoft} \rangle$   
 $\langle \text{Sergei Brin, Google} \rangle$



Top-k Answers



**SEARCHING FOR**  
**BY LOOKING AT**  
**PRODUCES**



**Few Approaches accept User Feedback**



<https://j.mp/ExploreSIGMOD>

# Where we are

Relational databases

Textual data

Graphs and networks

Challenges and Remarks

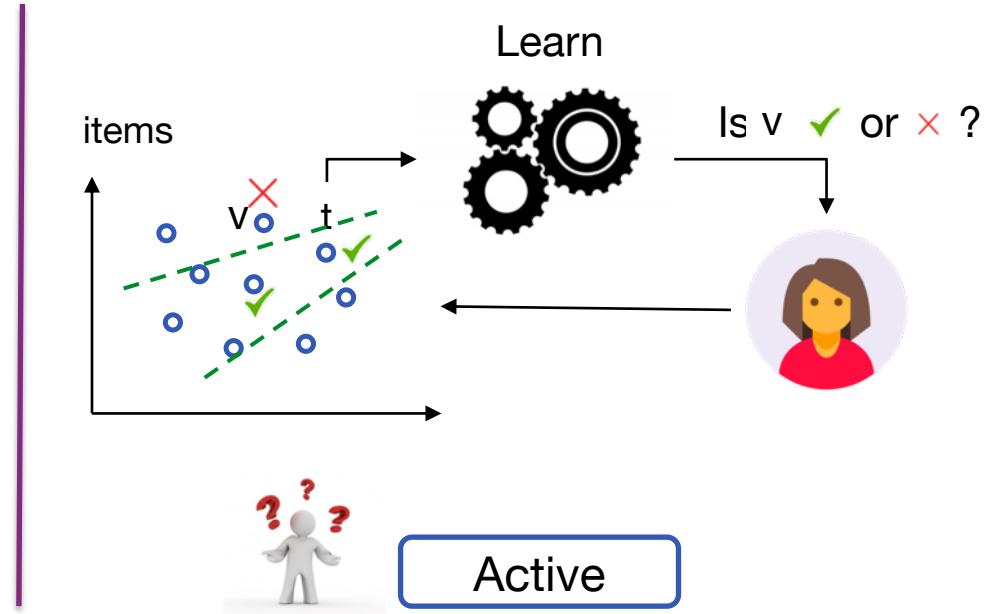
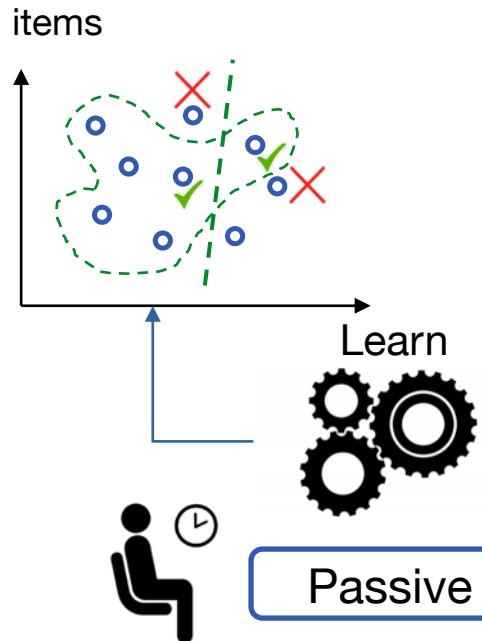


Machine  
learning

# Online exploration of datasets

Main idea: Learn the items to show online as more points are acquired

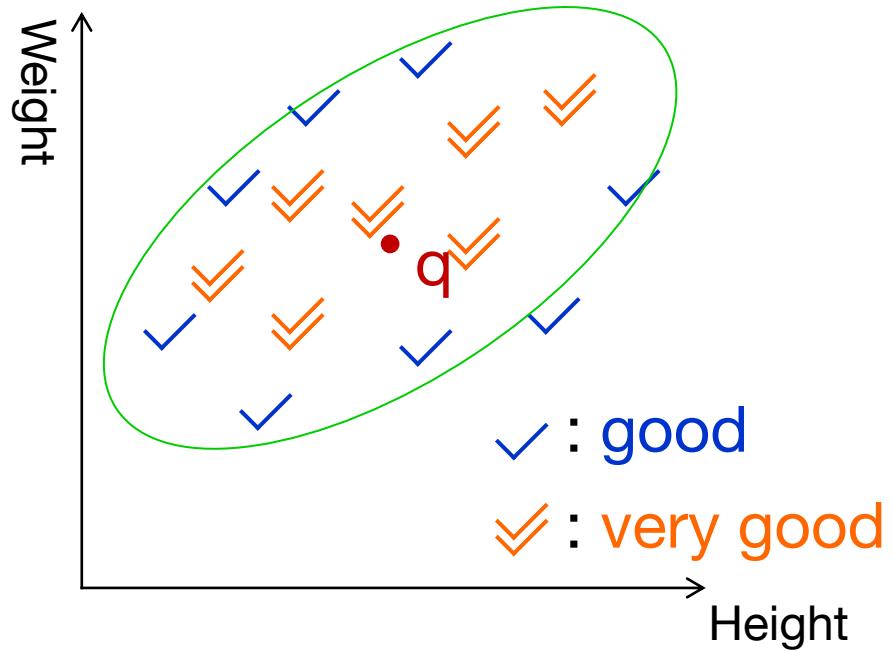
Two ways of learning: passive and active



Main idea: learn an **implicit query** from user examples and optional scores

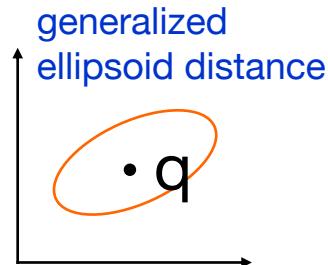
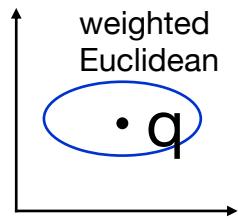
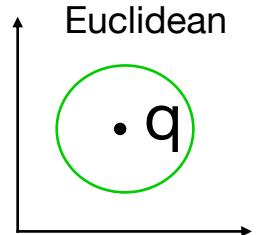
Searching “**mildly overweighted**” patients

- The doctor selects examples by **browsing** patient database
- The examples have “**oblique**” correlation
- We can “**guess**” the **implied query**



# Learning an ellipsoid distance

[Ishikawa et al., 1999]



$$D(x, q) = (x - q)^T M (x - q)$$

Implicit query

$$D(x, q) = \sum_j^n \sum_k^n m_{jk} (x_j - q_j)(x_k - q_k)$$

Learn the query minimizing the penalty = weighted sum of distances between query point and sample vectors

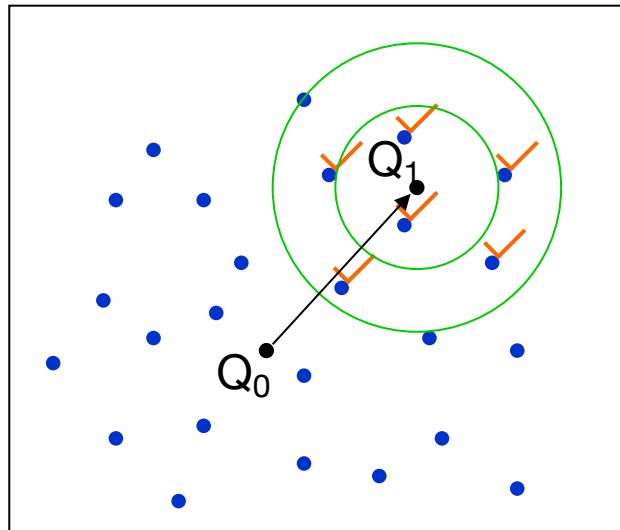
$$\begin{aligned} & \text{minimize} \quad \sum_i (x_i - q)^T M (x_i - q) \\ & \text{subject to} \quad \det(M) = 1 \end{aligned}$$



# Learning the distance

[Ishikawa et al., 1999]

Query point is moved towards “**good**” examples — **Rocchio formula** in IR



$Q_0$ : query point

● : retrieved data

✓ : relevance judgments

$Q_1$ : new query point

Learning can be done online!!!



# Explore-by-Example

[Dimitriadou et al., 2015]

Human-in-the-loop application

Explore big datasets to discover interesting data

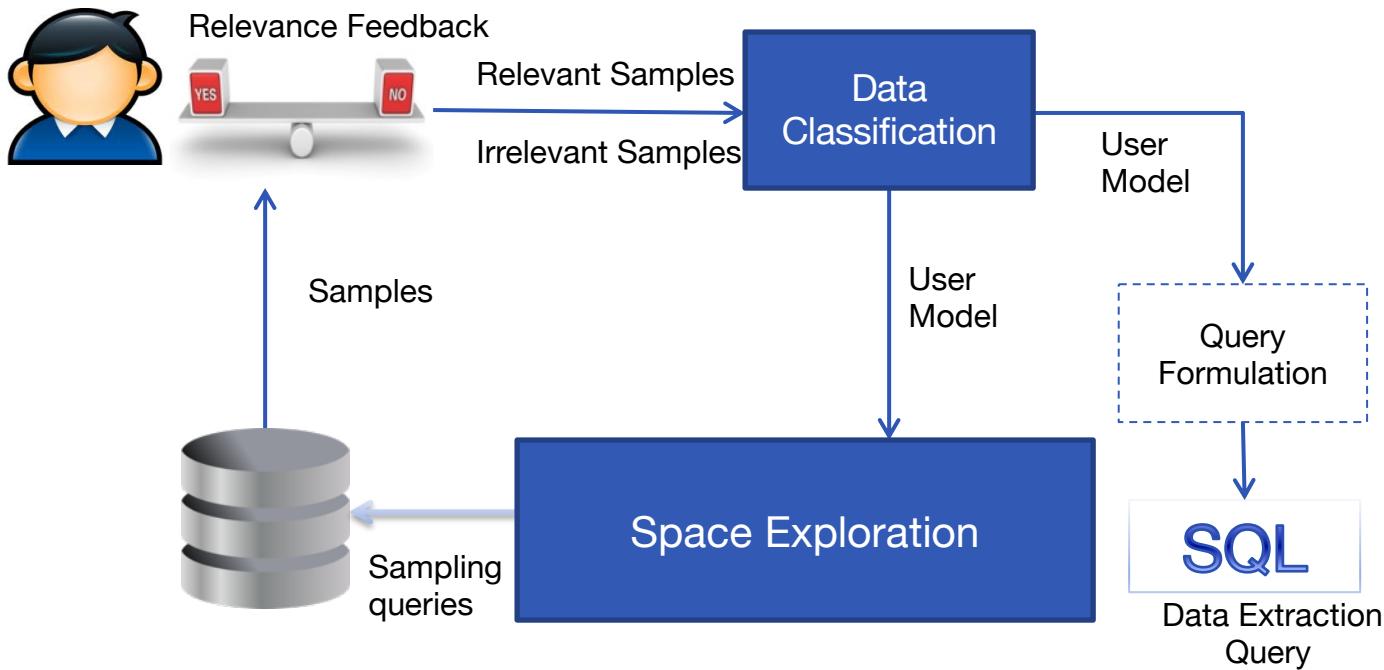
Challenges:

- Ad-hoc queries: “correct” predicates are unknown a priori
- Labor intensive: thousands of objects to review
- Resource intensive: execution of long query sequences on big data



# Explore-by-Example: AIDE

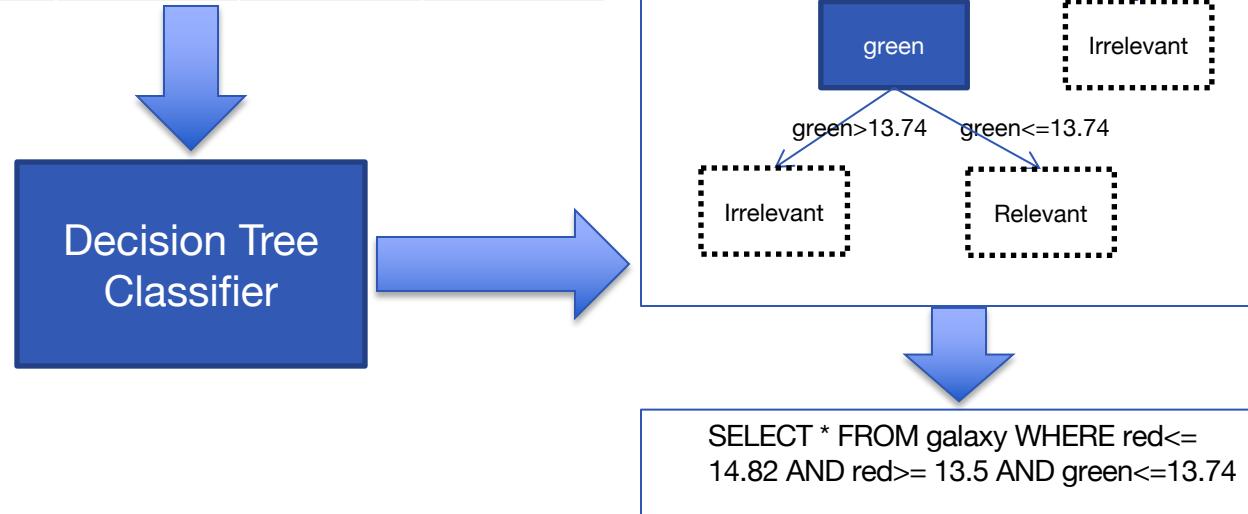
[Dimitriadou et al., 2015]



# Classification & Query Formulation

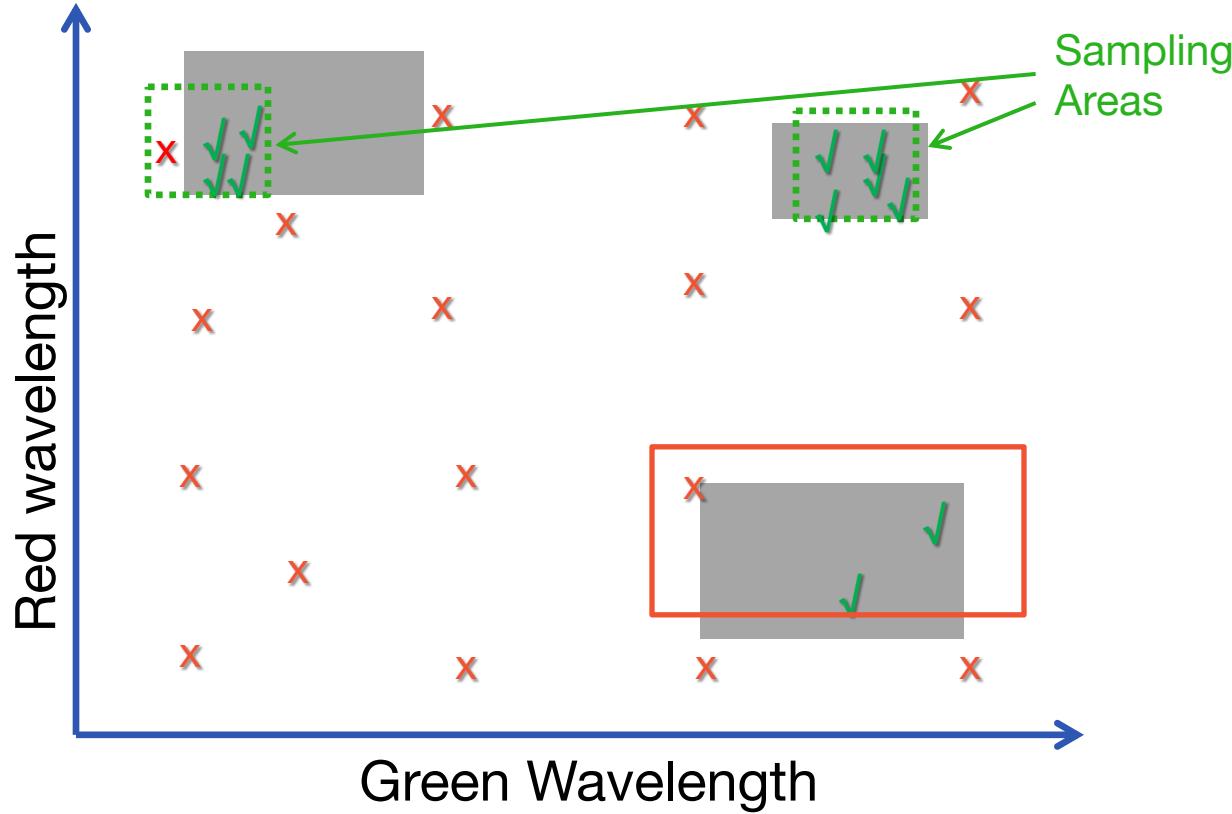
[Dimitriadou et al., 2015]

Sample	Red	Green	Relevant
Object A	13.67	12.34	Yes
Object B	15.32	14.50	No
..	..	..	...
Object X	14.21	13.57	Yes



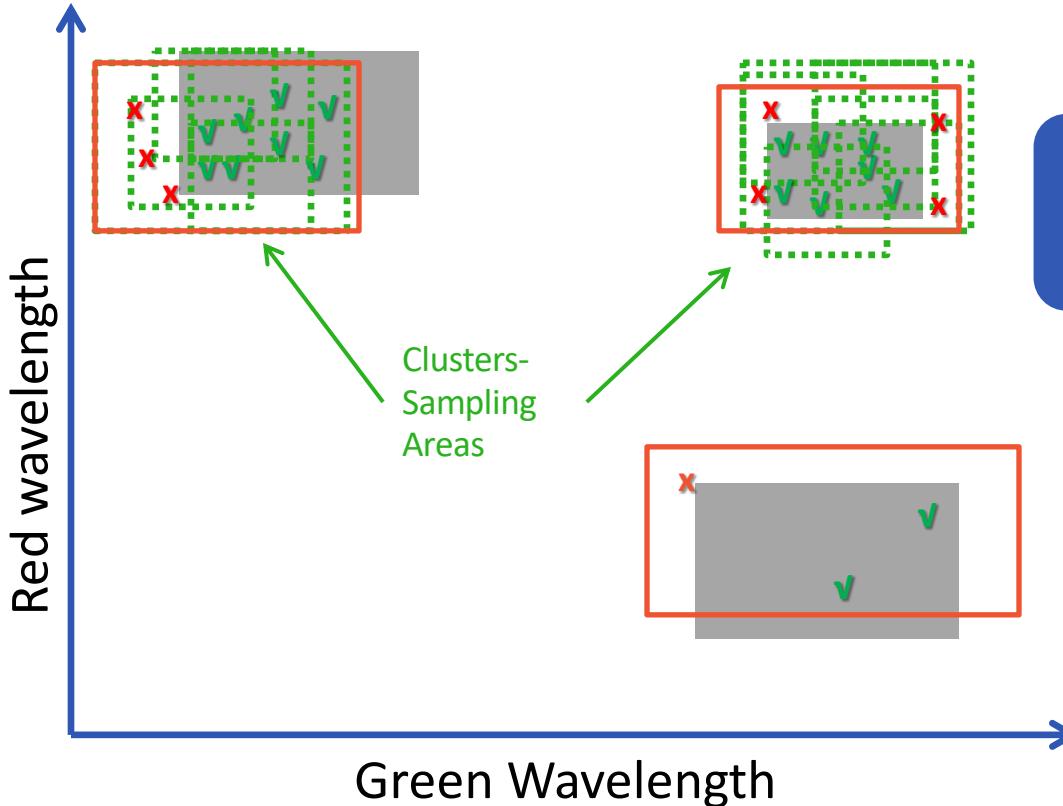
# Misclassified Sample Exploitation

[Dimitriadou et al., 2015]



# Clustering-based Sampling

[Dimitriadou et al., 2015]

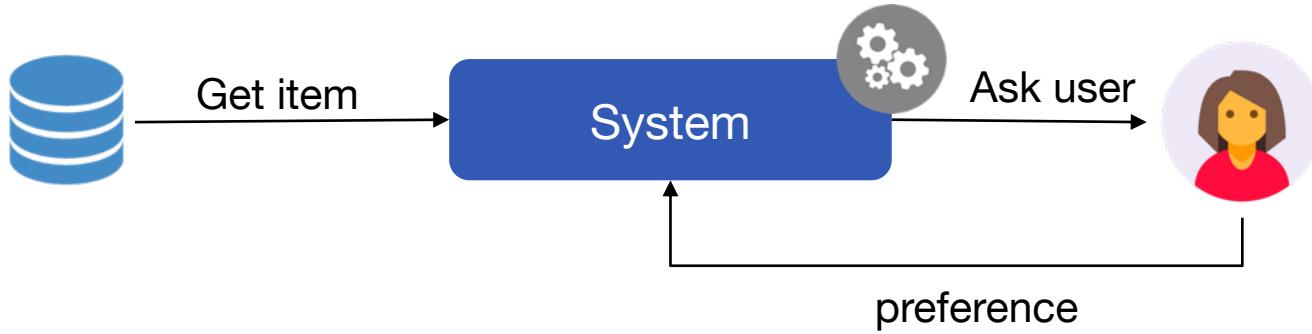


Idea: Use a k-medoid approach to find sampling areas

# Active learning for online query systems

[Vanchinathan et al., 2015]

Main idea: the system “query” the user to understand her preferences

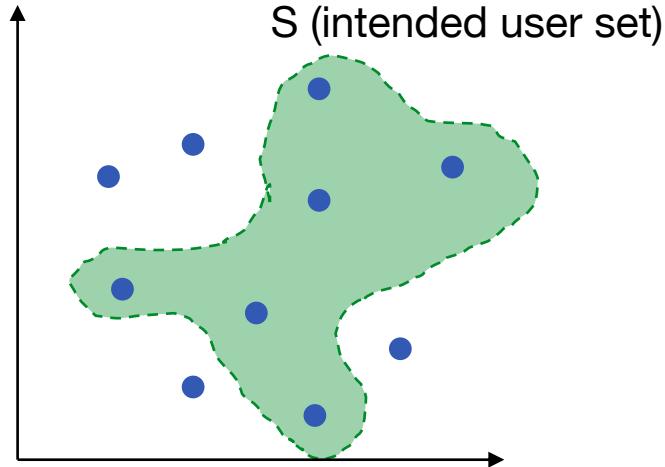


Learn unknown preferences and minimize the number of questions to the user

# Learning unknown preferences

[Vanchinathan et al., 2015]

**Problem:** Find a set  $S$  that maximize the unknown user preference within a budget (e.g., number of interactions)



$$\arg \max \sum_{v \in S} \text{pref}(v)$$

subject to  $\text{Cost}(S) \leq \text{budget}$

Cost for the set  $S$

# A step back ...

*Learning from an unknown environment ...*



# Multi-armed bandits

- Maximize the **reward** by successively playing gamble machines (the ‘arms’ of the bandits)
- Invented in **early 1950s** by Robbins for decision making under uncertainty when the environment is unknown
- The reward is unknown ahead of time



Reward  $X_1$



Reward  $X_2$



Reward  $X_3$

...

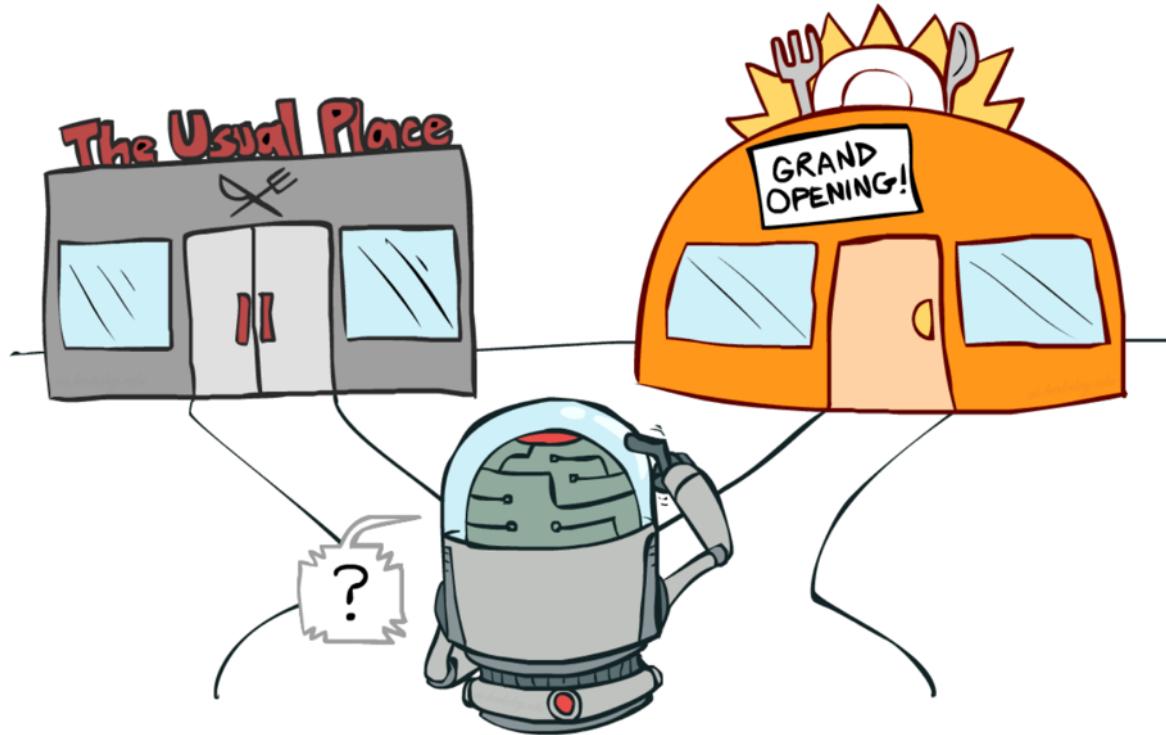
# Multi-armed bandits

- Reward = random variable  $X_{i,n}$ ;  $1 \leq i \leq K, n \geq 1$
- $i$  = index of the gambling machine
- $n$  = number of plays
- $\mu_i$  = expected reward of machine  $i$ .

A policy, or allocation strategy  $A$  is an algorithm that chooses the next machine to play based on the sequence of past plays and obtained rewards.



# Exploration vs Exploitation



<https://lilianweng.github.io/lil-log/2018/01/23/the-multi-armed-bandit-problem-and-its-solutions.html>

# A pure exploitation algorithm

Choose the machine with current best expected reward

- **Exploitation vs exploration dilemma:** Should you **exploit** the information you've learned or **explore** new options in the hope of greater payoff?
- In the greedy case, the balance is completely towards exploitation



# Quality measure - Regret

Total expected regret (after T plays):

$$R_T = \mu^* \cdot T - \sum_{i=1}^K \mu_j \cdot \mathbb{E}[N_{i,T}]$$

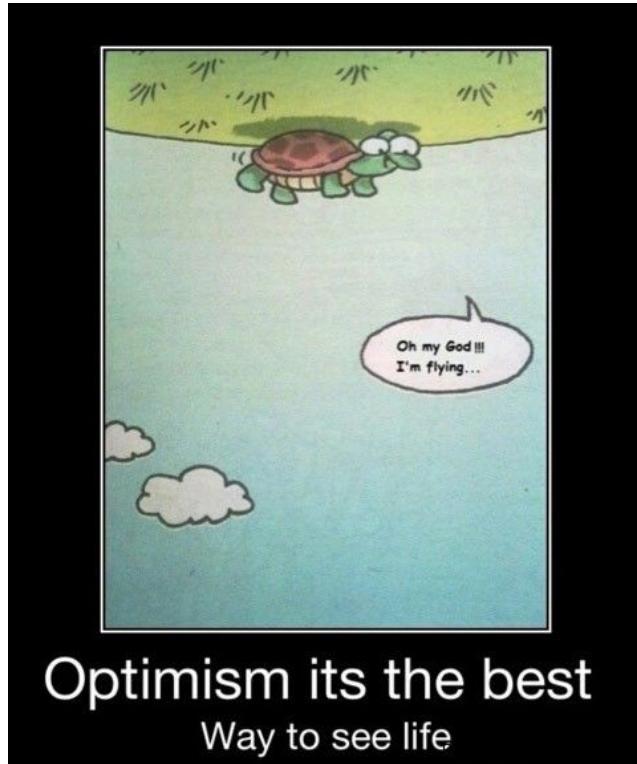
$\mu^*$ : highest expected reward

$\mathbb{E}[N_{i,T}]$ : expected number of times machine  $i$  is played

An algorithm solve the multi-armed problem if  
it matches the lower bound  $R_T = O(\log T)$  [Think about binary search]



# An optimistic view



# Upper confidence bound (UCB)

- Formalise the intuition using **confidence intervals**
- Optimistic estimate of the mean of arm = ‘largest value it could plausibly be’
- Suggests

$$\text{Optimistic estimate} = \frac{1}{n_j} \sum_{s=1}^{n_j} X_{j,s} + \sqrt{\frac{2 \log(1/t)}{n_j}}$$



# UCB algorithm

1. Pull at each time  $t$  the arm with the maximum probability of being the best

$$\frac{1}{n_j} \sum_{s=1}^{n_j} X_{j,s} + \sqrt{\frac{2 \log(1/t)}{n_j}}$$

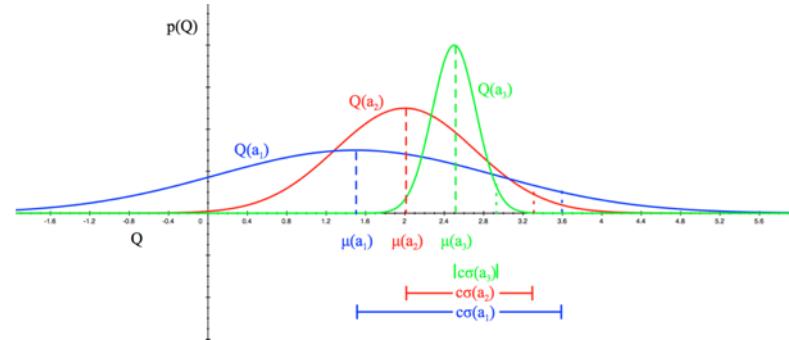
2. Repeat until the budget (number of steps  $T$ ) is depleted

$n_j$ : number of times the arm  $j$  has been pulled

**Balance exploration and exploitation:** The uncertainty diminishes as the time passes



# Bayesian UCB



- In UCB no prior on the reward
  - Hence, we use Hoeffding's Inequality for a very generalize estimation.
- If we consider a prior we can get a better estimate
- For example, if we expect the mean reward of every slot machine to be Gaussian
  - set the upper bound as 95% confidence interval by setting the exploration to be twice the standard deviation.

# Back to our problem



# Background: Gaussian processes

[Bishop et al., 2006]

Idea: Model the user preferences as a Gaussian Process

A Gaussian Process (GP) is an infinite set of variables, any subset of this is Gaussian

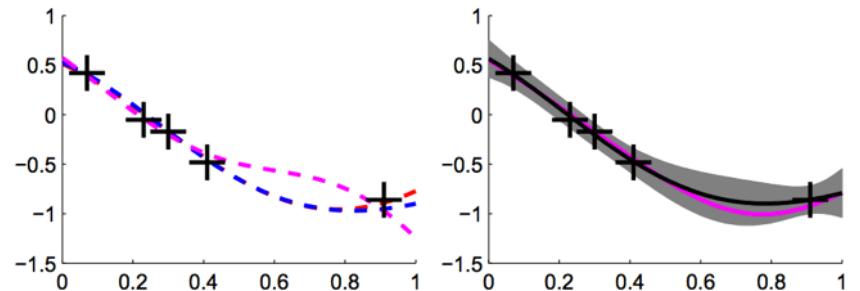
$$P(\mathbf{f}|\Sigma, \mu) = |2\pi\Sigma|^{\frac{1}{2}} \exp(-\frac{1}{2}(\mathbf{f} - \mu)^T \Sigma^{-1} (\mathbf{f} - \mu))$$

Gaussian prior

Specified only by mean and covariance

Given observations  $\{x, y\}_{i=1}^n$  over an unknown function  $f$  drawn from a Gaussian prior, the posterior is Gaussian

$$P(\mathbf{f}|\mathbf{y}) \propto \int dx P(\mathbf{f}, \mathbf{x}, \mathbf{y})$$



# GP-Select

[Vanchinathan et al., 2015]

---

## Algorithm 1 GP-SELECT

---

**Input:** Ground Set  $\mathbf{V}$ , kernel  $\kappa$  and budget  $B$

Initialize selection set  $S$

**for**  $t = 1, 2, \dots, B$  **do**

**Model Update:**

$[\mu_{t-1}(\cdot), \sigma_{t-1}^2(\cdot)] \leftarrow \text{GP-Inference}(\kappa, (S, y_{\{1:t-1\}}))$

**Item Selection:**

        Set  $v_t \leftarrow \underset{v \in \mathbf{V} / \{v_{1:t-1}\}}{\operatorname{argmax}} \mu_{t-1}(v) + \beta_t^{1/2} \sigma_{t-1}(v)$

$S \leftarrow S \cup \{v_t\}$

        Receive feedback  $y_t = f(v_t) + \epsilon_t$

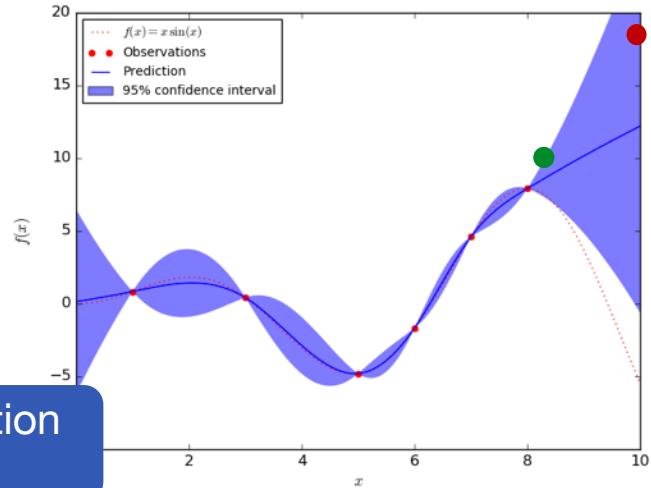
**end for**

Learn posterior

Trades off exploration  
exploitation

Ask user feedback

- Exploration: select items with high-variance
- Exploitation: select items with high-value

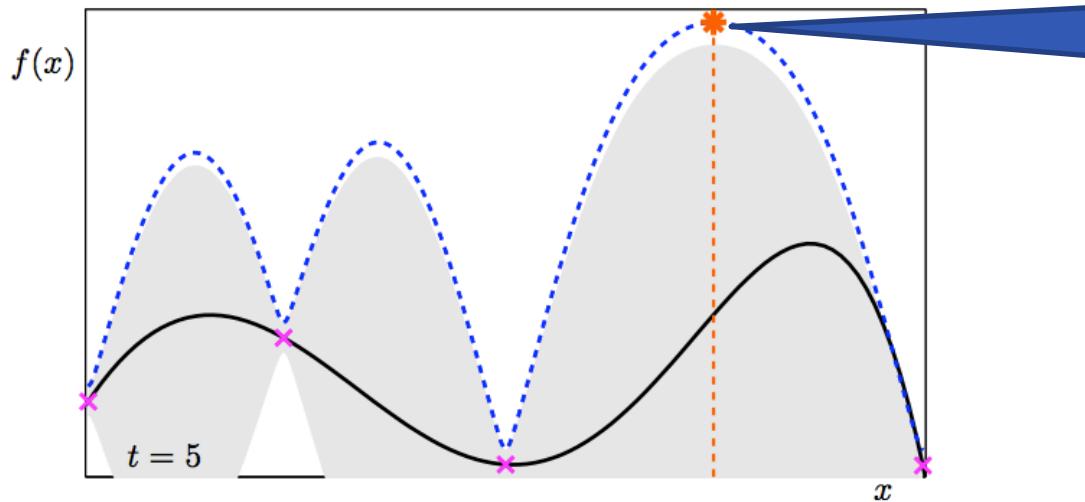


# GP-Select bound

[Srinivas et al, 2019]

- Multi-armed bandits with "infinite" arms
- Trades off exploration-exploitation
- $\Psi_T(X)$ : Maximum information gain

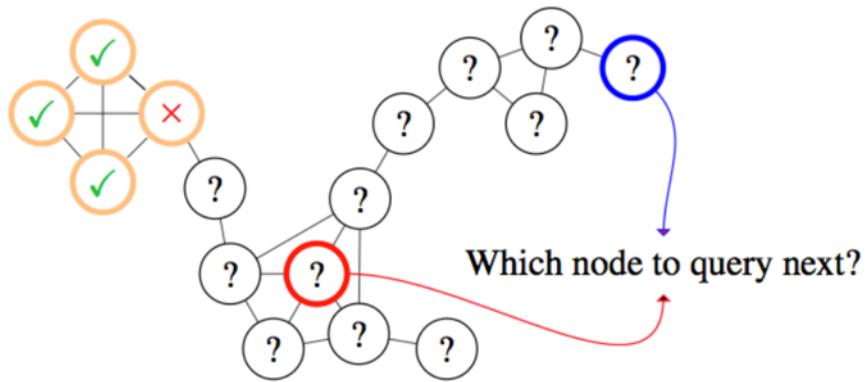
$$R_T = f(x^*) - \max_{t=1 \dots T} f(x_t) \leq \sqrt{\frac{\Psi_T(X)}{n}}$$



# Active learning on graphs – which prior?

[Ma et al., 2015]

Idea: Use the graph structure to infer the node classes



Use graph Laplacian as prior  
 $L = D - A$ ,  $A$  is the adjacency matrix

$$p(\mathbf{f}) \sim \mathcal{N}(0, L^{-1})$$

Laplacian: higher probability of having the same class if two nodes are connected



# Where could Active learning help?

## Reverse engineering queries and rules

- Interactive Refinement of example tuples
- Learning the most probable queries from their results



## Graph exploration

- Summarization of knowledge graphs with preferences
- Seed set expansion
- Recommendation of relevant nodes

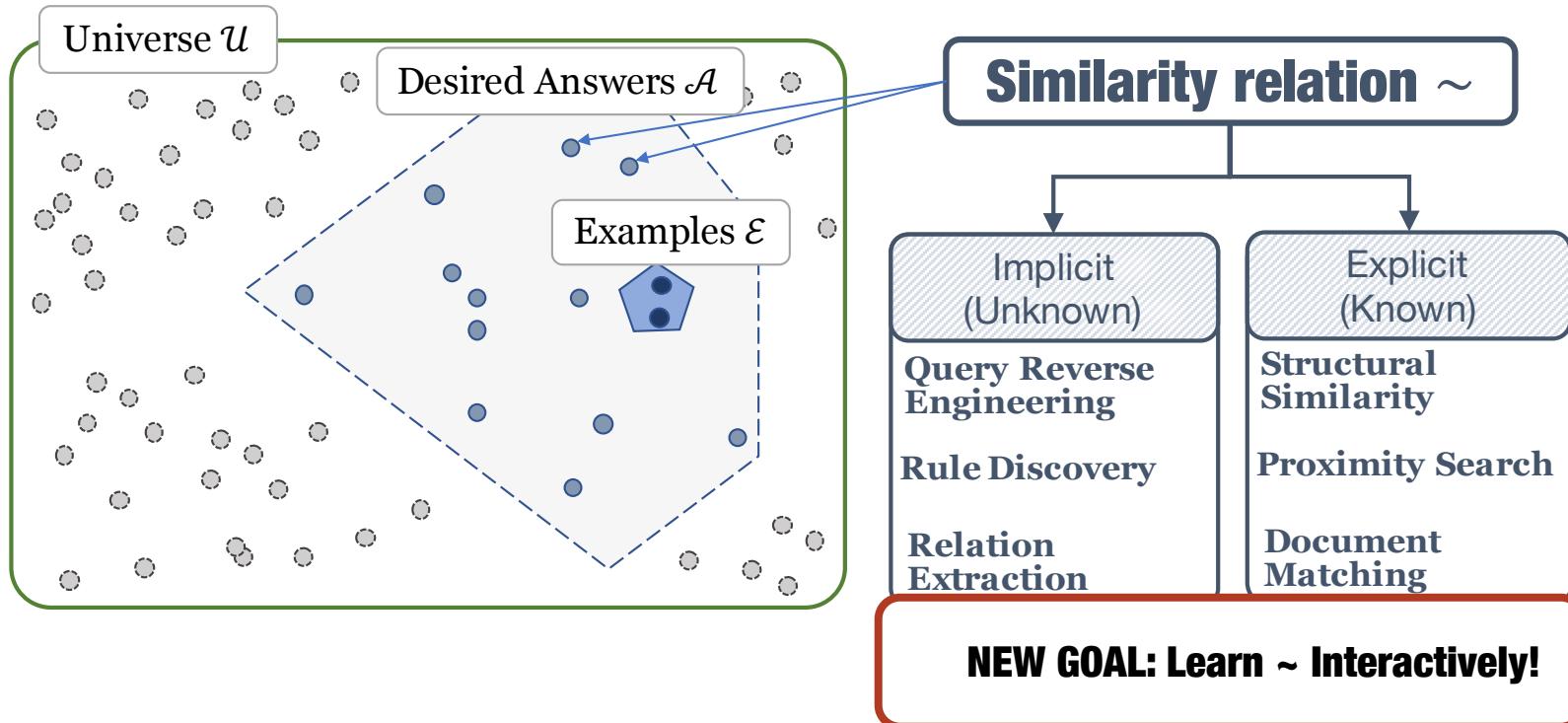


## Text processing

- Fast entity matching
- Advertising based on documents search



# Example-based methods



# MAB: good resources

## Books and surveys

- <http://slivkins.com/work/MAB-book.pdf>
- <http://downloads.tor-lattimore.com/book.pdf>
- <http://sbubeck.com/SurveyBCB12.pdf>

## Tutorials

- Lattimore - AAAI 2018: [part 1](#) - [part 2](#)
- [Tutorial on bayesian optimization of expensive cost functions](#)
- Blog on bandits: <http://banditalgs.com/>



# Where we are

Relational databases

Textual data

Graphs and networks

Machine learning

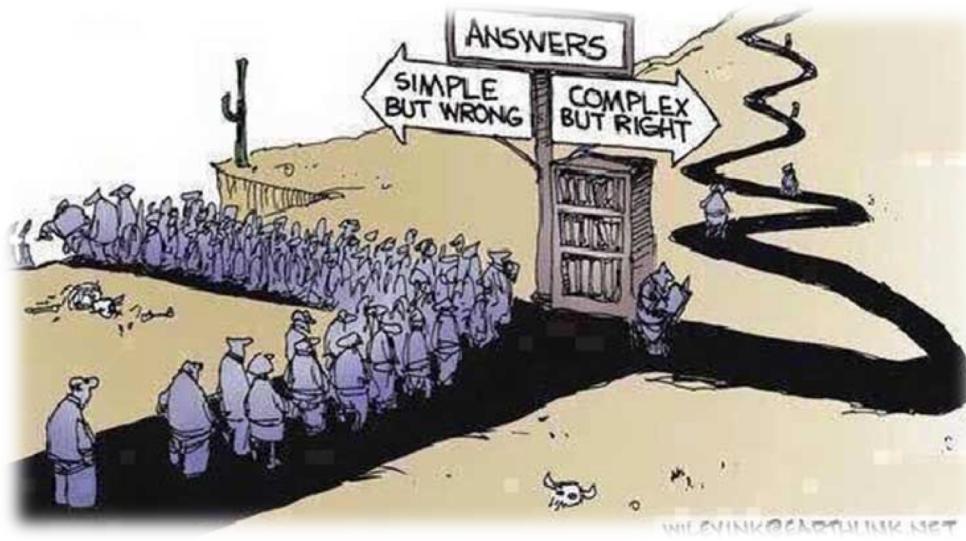
Challenges and Remarks



# Big data – Easy value?



© marktoonist.com



# Exploration

*We know where we start  
we don't know what we'll find*



# Traditional Search Methods are not Enough

## We need Specialized Methods for Data Exploration

**From broad views**

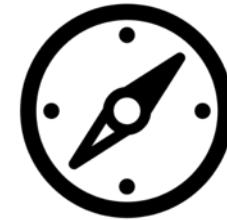
**to Detailed view**

**From exploration as  
select count(\*)**

**to find what is interesting**

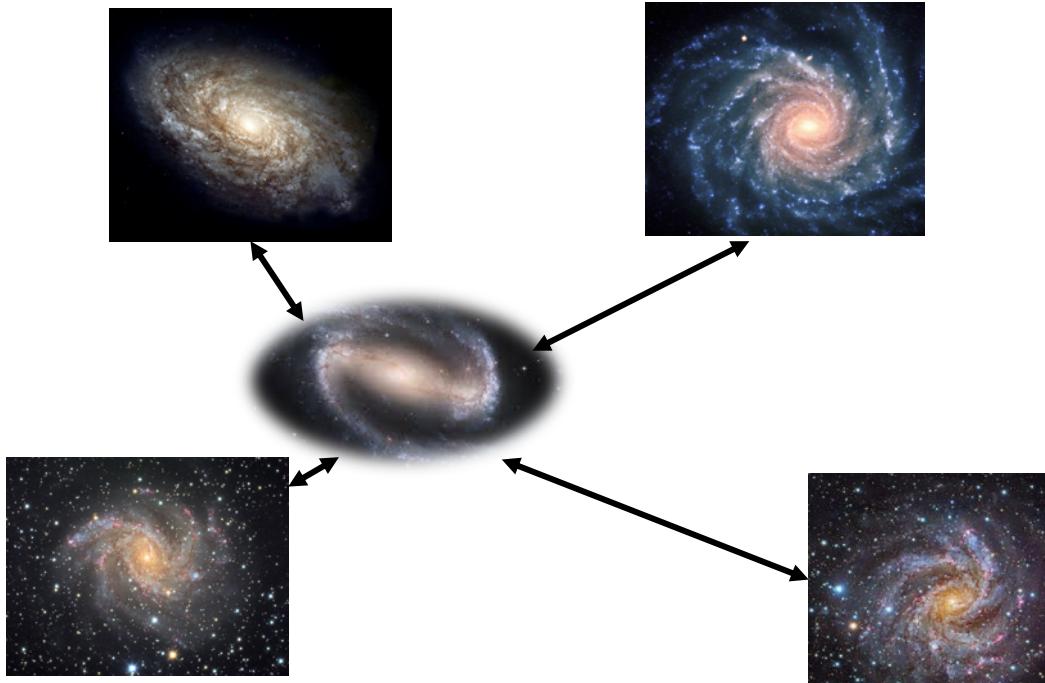
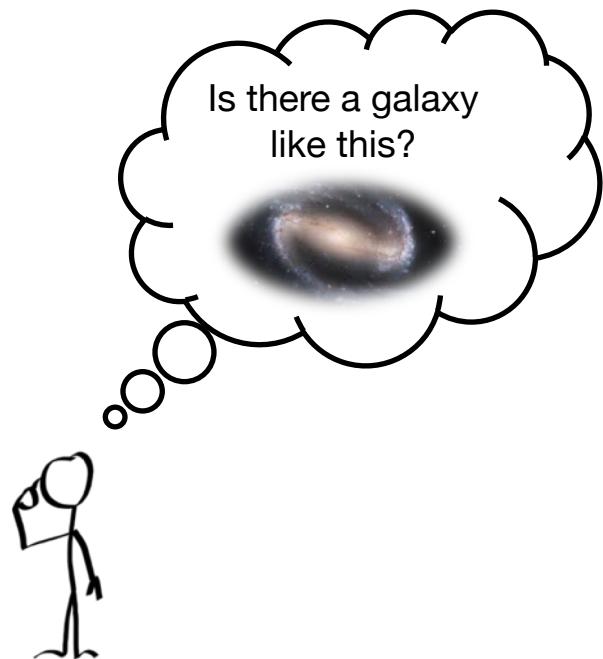


**From Exact Search  
based on explicit conditions**

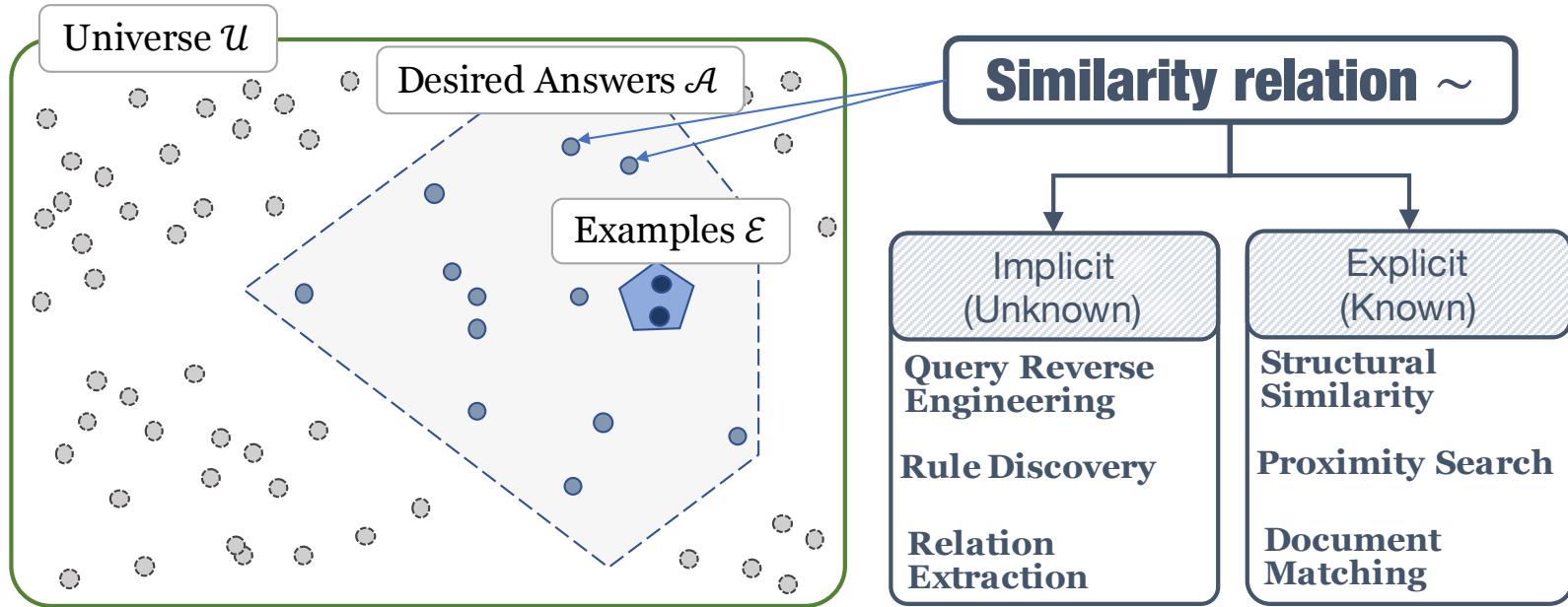


**to Exploratory Search  
based on Implicit needs**

# Similarities are the key ...



# Example-based methods: All You Need is ...



# Example-based methods

## Relational

- Reverse engineering queries
- Example-driven schema mapping
- Interactive data repairing



## Textual

- Search documents by example
- Entity extraction by example text
- Web table completion using examples



## Graph

- Community-based Node-retrieval
- Entity Search
- Path and SPARQL queries
- Graph structures as Examples



# Example-based methods: takeaways

## Relational

- Complex search space
- Exact and approximate
- Interactivity can improve the quality
- Limited to query inference



## Textual

- Allows serendipitous search
- Easier document finding
- Speed up entity matching
- Extract semi-structure data



## Graph

- Heterogenous Structures
- Exploit locality
- Entity attributes are expressive
- Large result-sets require ranking

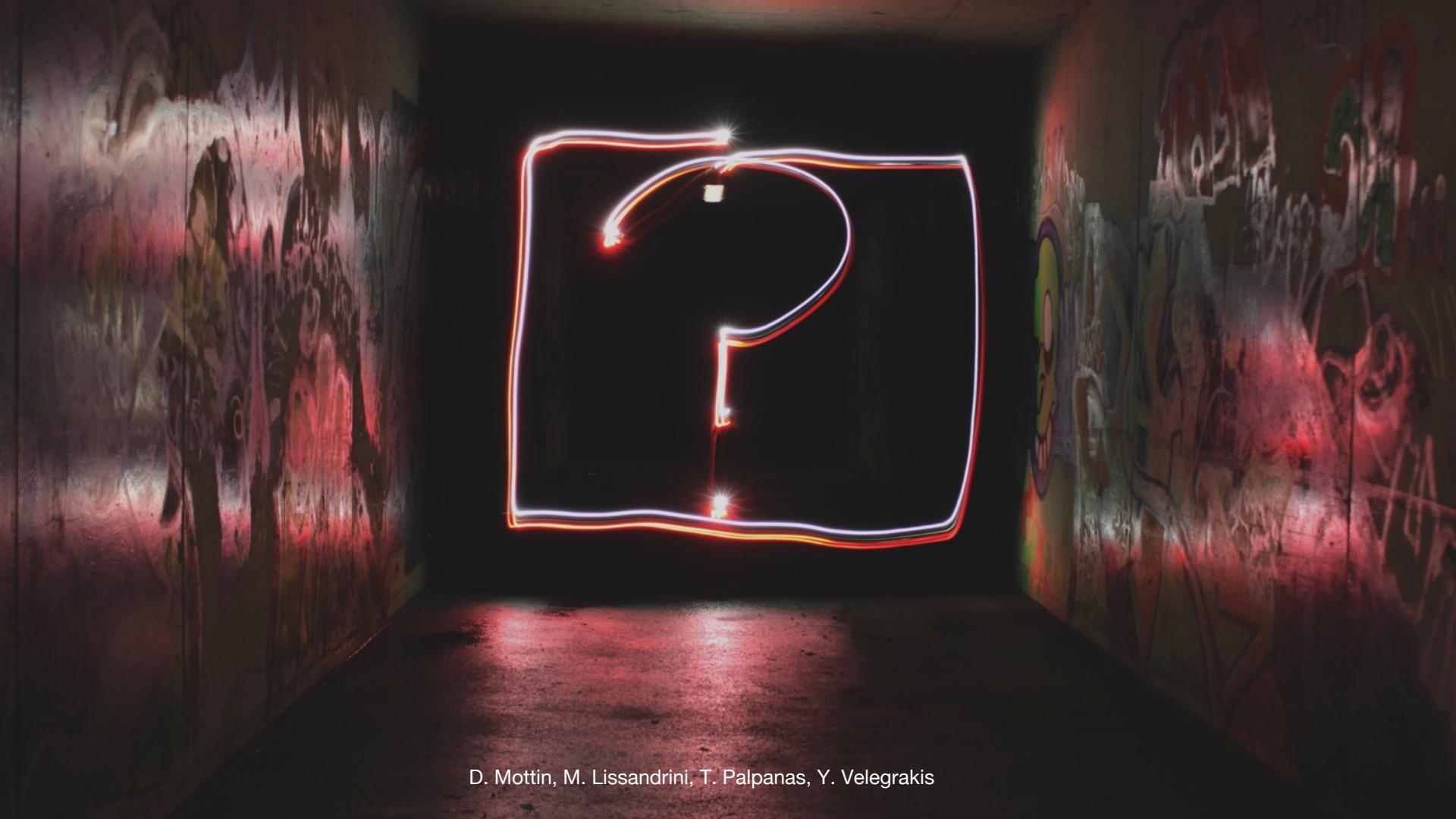


# The use of examples

## Examples can ease data exploration

- ... reduce need for complex queries / simplify user input
- ... require no schema knowledge
- ... allow uncertainty in search conditions
- ... require little data analytics expertise





D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# Acknowledgments

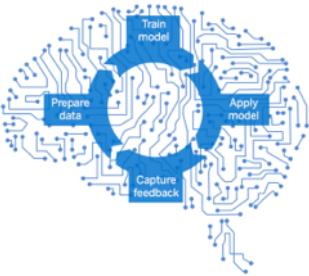
We would like to thank the authors of the papers  
who kindly provided us the slides

Angela Bonifati, Radu Ciucianu, Marcelo Arenas, Gonzalo Diaz, Egor Kostylev, Yaakov Weiss, Sarah Cohen, Fotis Psallidas, Li Hao, Chan Chee Yong, Ilaria Bordino, Mohamed Yakout, Kris Ganjam, Kaushik Chakrabati, Thibault Sellam, Rohit Singh, Maeda Hanafi, Dmitri Kalashnikov, Marcin Sydow, Mingzhu Zhu, Yoshiharu Ishikawa, Daniel Deutch, Nandish Jayaram, Paolo Papotti, Bryan Perozzi, Kiriaki Dimitriadou, Yifei Ma, Natali Ruchansky, Quoc Trung Tran, Hastagiri Prakash Vanchinathan

*... and many others (see references)*



# Where should we invest time?



Machine  
learning

Approximate  
Methods

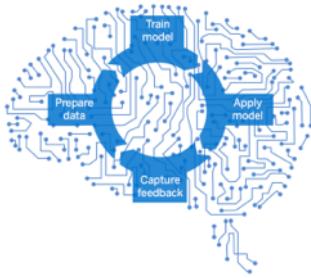


User models

Scalability



# Where should we invest time?



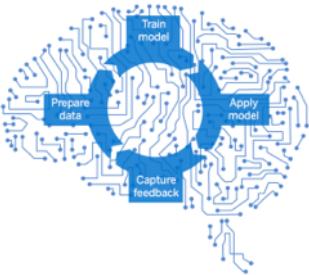
Machine  
learning

## Learn from Examples

- ... Similarity Measures: are often “fuzzy” and “implicit”
- ... New representations of the search space
- Challenge: Scale! Exploration of large search spaces



# Where should we invest time?



Machine  
learning



User models



## Learn from Examples

- ... Similarity Measures to represent User Interests
- ... User-centric, dynamic, Exploration-strategies: learn as you go
- Challenge: Distinct User have Different Goals! Explore in different ways

We need more data!

# Where should we invest time?

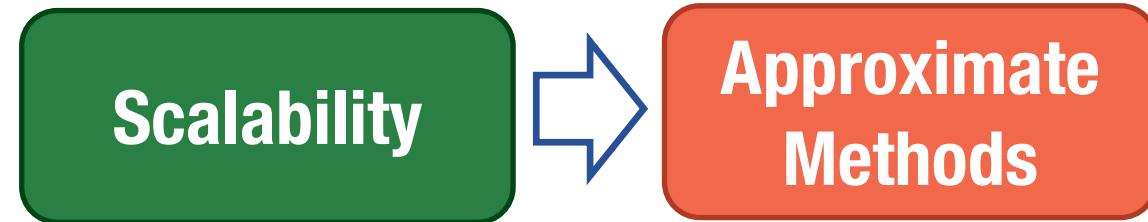


Scalability

## Scale Example-based search

- ... Huge search space, dynamic data, variety of data models
- ... Exploration is Interactive, requires Interactive response time
- Adaptive Data-structures, localized access, flexible schema, incremental index

# Where should we invest time?



## Scale Example-based search

- ... An approximate answer now is better than a precise answer in 1hour
- ... Approximate answers can provide insights without being accurate

Exploratory queries retrieve large resultsets: the user needs only a glimpse to figure out  
if they are moving in the right direction!

# Features of Exploratory Search Systems

[White and Roth, 2009]

## Support querying and rapid query refinement:

- Offer facets and metadata-based result filtering
- Leverage search context

## Example-driven

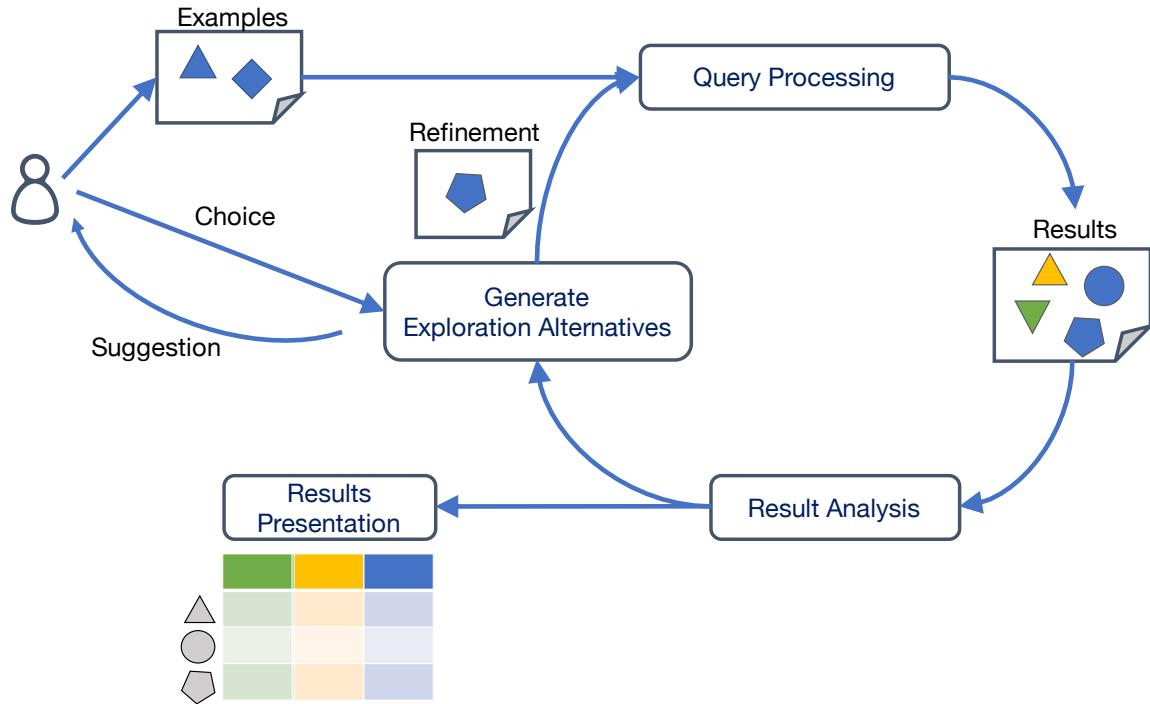
- visualizations, summarizations, and explanations
- paired with methods to suggest further example-based explorations.

Via Interactivity  
& Personalization!

## Support learning and understanding



# Interactive Example Based Exploration System?



## Requires:

### Fast Query Processing

Avoid the full recomputation of a query  
Limit the computation to only a sample  
Adaptive query executions  
Adaptive data-structures and indexes,

### Automatic Result Analysis

Automatically identify peculiar characteristics  
Data-summarization techniques  
Learn user interests automatically





# ADOPT HETEROGENEITY

Need for solutions that  
**operate across different models**

**operate on heterogeneous  
datastores**

**dataset search**

*Data Lakes??*





# **DEMOCRATIZATION**

## **easy access to data**

Tools that work on  
**commodity**  
**hardware, mobile devices**

Data-exploration for  
**everyday use-cases**

Users want back  
**the control on their data**





D. Mottin, M. Lissandrini, T. Palpanas, Y. Velegrakis

# NATURAL LANGUAGE INTERFACE

*flexible, vague,  
imprecise input*

**Exploration through  
conversation**



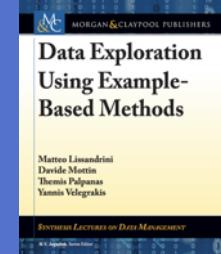
# Example is always more efficacious than precept

*Samuel Johnson, Rasselas (1759), Chapter 29.*

“New Trends on Exploratory Methods for Data Analytics.” *PVLDB*, 2017.

“Data Exploration Using Example-Based Methods.” *M&C*, 2018.

“Exploring the Data Wilderness through Examples.” *SIGMOD*, 2019.



**Slides:** <https://data-exploration.ml/>



# References

- M. Arenas, G. I. Diaz, and E. V. Kostylev. Reverse engineering sparql queries. WWW, 2016.
- Agichtein, E. and Gravano, L. Snowball: Extracting relations from large plain-text collections. ICDL, 2000.
- A.Bonifati, R.Ciucanu, and A.Lemay. Learning path queries on graph databases. EDBT, 2015.
- A. Bonifati, R. Ciucanu, and S. Staworko. Learning join queries from user examples. TODS, 2016.
- A. Bonifati, U. Comignani, E. Coquery, and R. Thion. Interactive mapping specification with exemplar tuples. SIGMOD, 2017.
- I. Bordino, G. De Francisci Morales, I. Weber, and F. Bonchi. From machu picchu to rafting the urubamba river: anticipating information needs via the entity-query graph. WSDM, 2013.
- D. Deutch and A. Gilad. Qplain: Query by explanation. ICDE, 2016.



# References

- G. Diaz, M. Arenas, and M. Benedikt. Sparqlbye: Querying rdf data by example. PVLDB, 2016.
- K. Dimitriadou, O. Papaemmanouil, and Y. Diao. Explore-by-example: An automatic query steering framework for interactive data exploration. In SIGMOD, 2014.
- B. Eravci and H. Ferhatoğlu. Diversity based relevance feedback for time series search. PVLDB, 2013.
- A. Gionis, M. Mathioudakis, and A. Ukkonen. Bump hunting in the dark: Local discrepancy maximization on graphs. ICDE, 2015.
- M. F. Hanafi, A. Abouzied, L. Chiticariu, and Y. Li. Synthesizing extraction rules from user examples with seer. SIGMOD, 2017.
- He, J., Veltri, E., Santoro, D., Li, G., Mecca, G., Papotti, P. and Tang, N. Interactive and deterministic data cleaning. SIGMOD, 2016.
- Y. Ishikawa, R. Subramanya, and C. Faloutsos. Mindreader: Querying databases through multiple examples. VLDB, 1998.



# References

- N. Jayaram, A. Khan, C. Li, X. Yan, and R. Elmasri. Querying knowledge graphs by example entity tuples. TKDE, 2015.
- H. Li, C.-Y. Chan, and D. Maier. Query from examples: An iterative, data-driven approach to query construction. PVLDB, 2015.
- M. Lissandrini, D. Mottin, Y. Velegrakis, T. Palpanas. Multi-Example Search in Rich Information Graphs ICDE 2018
- Y. Ma, T.-K. Huang, and J. G. Schneider. Active search and bandits on graphs using sigma-optimality. UAI, 2015.
- S. Metzger, R. Schenkel, and M. Sydow. Qbees: query by entity examples. CIKM, 2013.
- D. Mottin, M. Lissandrini, Y. Velegrakis, and T. Palpanas. Searching with xq: the exemplar query search engine. SIGMOD, 2014.
- D. Mottin, M. Lissandrini, Y. Velegrakis, and T. Palpanas. Exemplar queries: a new way of searching. VLDB J., 2016.
- B. Perozzi, L. Akoglu, P. Iglesias Sanchez, and E. Müller. Focused clustering and outlier detection in large attributed graphs. KDD, 2014.



# References

- F. Psallidas, B. Ding, K. Chakrabarti, and S. Chaudhuri. S4: Top-k spreadsheet-style search for query discovery. SIGMOD, 2015.
- R. Rolim, G. Soares, L. D'Antoni, O. Polozov, S. Gulwani, R. Gheyi, R. Suzuki, and B. Hartmann. Learning syntactic program transformations from examples. ICSE, 2017.
- N. Ruchansky, F. Bonchi, D. García-Soriano, F. Gullo, and N. Kourtellis. The minimum wiener connector problem. SIGMOD, 2015.
- T. Sellam and M. Kersten. Cluster-driven navigation of the query space. TKDE, 2016.
- Y. Shen, K. Chakrabarti, S. Chaudhuri, B. Ding, and L. Novik. Discovering queries based on example tuples. SIGMOD, 2014.
- R. Singh. Blinkfill: Semi-supervised programming by example for syntactic string transformations. PVLDB, 2016.
- G. Sobczak, M. Chochół, R. Schenkel, and M. Sydow. iqbees: Towards interactive semantic entity search based on maximal aspects. Foundations of Intelligent Systems, 2015.



# References

- Y. Su, S. Yang, H. Sun, M. Srivatsa, S. Kase, M. Vanni, and X. Yan. Exploiting relevance feedback in knowledge graph search. KDD, 2015.
- Q. T. Tran, C.-Y. Chan, and S. Parthasarathy. Query reverse engineering. VLDB J., 2014.
- H. P. Vanchinathan, A. Marfurt, C.-A. Robelin, D. Koss- mann, and A. Krause. Discovering valuable items from massive data. In KDD, 2015.
- C. Wang, A. Cheung, and R. Bodik. Interactive query synthesis from input-output examples. In SIGMOD, 2017.
- C. Wang, A. Cheung, and R. Bodik. Synthesizing highly expressive sql queries from input-output examples. In PLDI, 2017.
- Y. Y. Weiss and S. Cohen. Reverse engineering spj-queries from examples. SIGMOD, 2017.
- M. Yakout, K. Ganjam, K. Chakrabarti, and S. Chaudhuri. Infogather: Entity augmentation and attribute discovery by holistic matching with web tables. SIGMOD, 2012.
- M. Zhu and Y.-F. B. Wu. Search by multiple examples. WSDM, 2014.
- M. M. Zloof. Query by example. AFIPS NCC, 1975.

