

Dependable Data Repairing with Fixing Rules



مختبر قطر لبحوث الحوسبة
Qatar Computing Research Institute

Member of Qatar Foundation *جامعة قطر*

Jiannan Wang
Nan Tang

Data is Dirty

incomplete
inconsistent
inaccurate

...

Data is Dirty

incomplete
inconsistent
inaccurate

...

25% companies: flawed data
3+ trillion \$: US economy
20%: labor productivity

....

Data is Dirty

incomplete
inconsistent
inaccurate

...

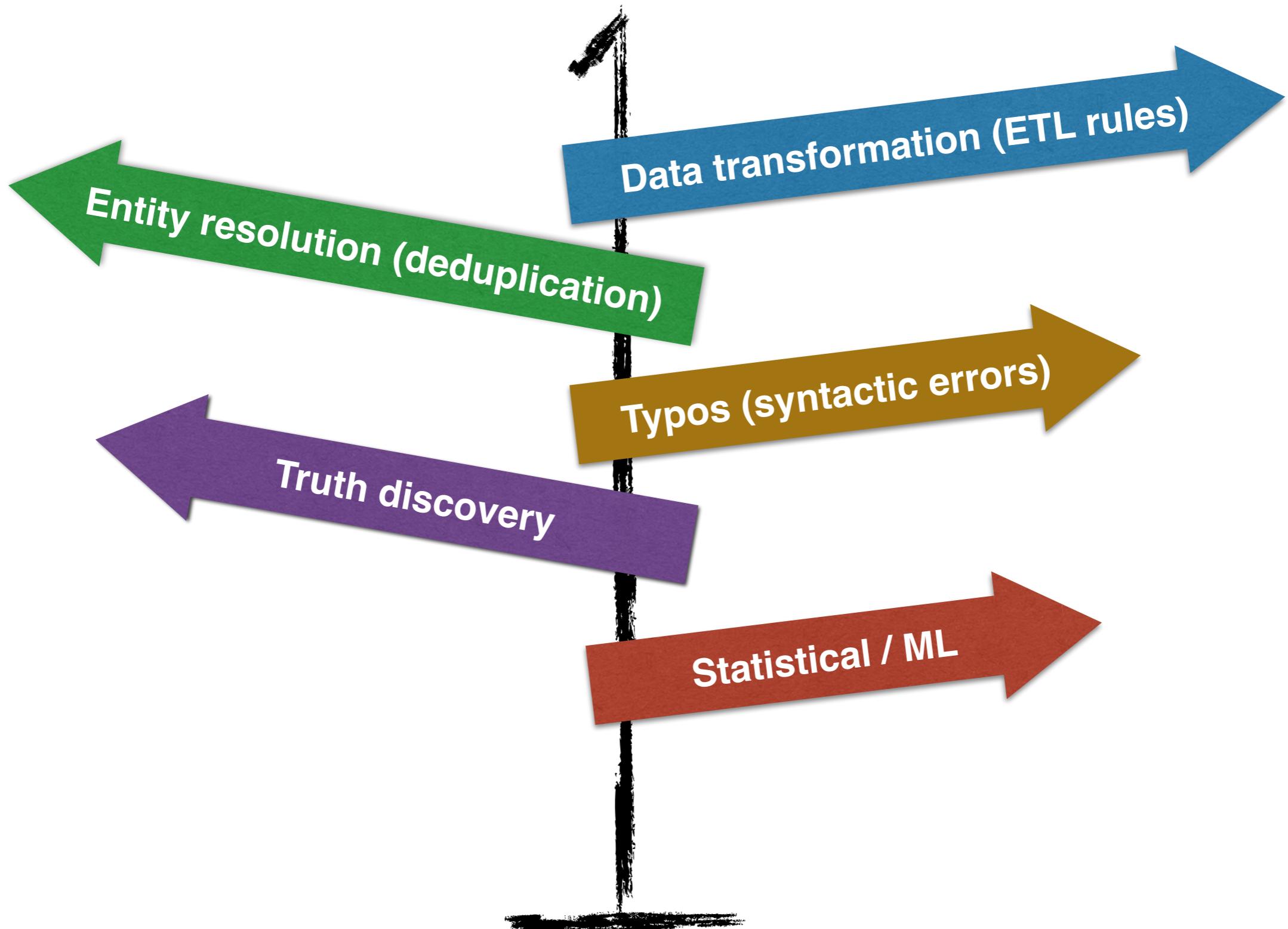
25% companies: flawed data
3+ trillion \$: US economy
20%: labor productivity

....

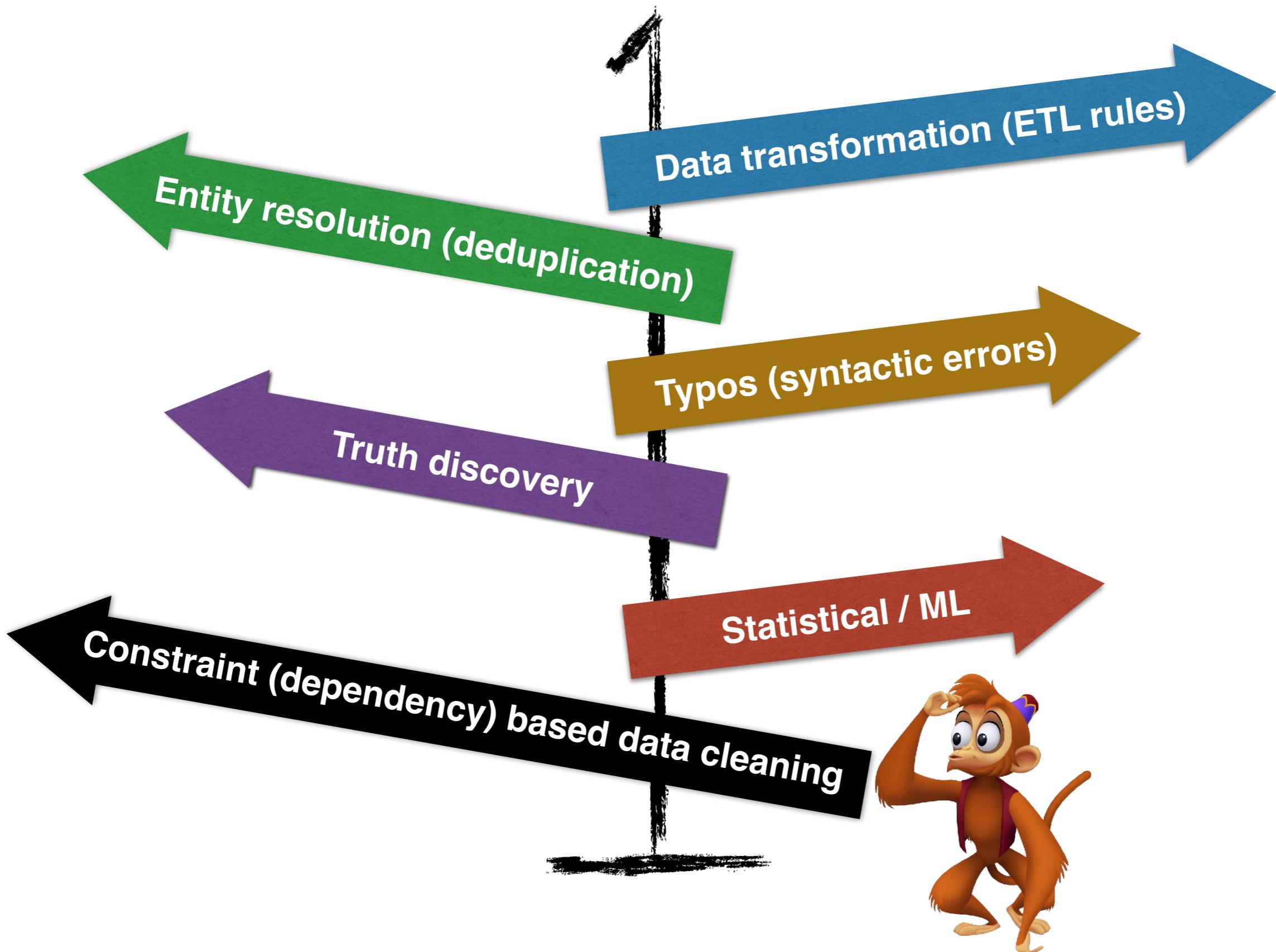
Data is Dirty

Big (clean) data: new oil

State-of-the-art



State-of-the-art



Dependency Theory

- Data dependencies (*a.k.a.* integrity constraints)

Dependency Theory

- Data dependencies (a.k.a. integrity constraints)

	name	country	capital	city	conf
r1	George	China	Beijing	Beijing	SIGMOD
r2	Ian	China	Shanghai	Hongkong	ICDE
r3	Peter	China	Tokyo	Tokyo	ICDE
r4	Mike	Canada	Toronto	Toronto	VLDB

Dependency Theory

- Data dependencies (a.k.a. integrity constraints)

FD: [country] \rightarrow [capital]

	name	country	capital	city	conf
r1	George	China	Beijing	Beijing	SIGMOD
r2	Ian	China	Shanghai	Hongkong	ICDE
r3	Peter	China	Tokyo	Tokyo	ICDE
r4	Mike	Canada	Toronto	Toronto	VLDB

Dependency Theory

- Data dependencies (a.k.a. integrity constraints)

FD: [country] \rightarrow [capital]



	name	country	capital	city	conf
r1	George	China	Beijing	Beijing	SIGMOD
r2	Ian	China	Shanghai	Hongkong	ICDE
r3	Peter	China	Tokyo	Tokyo	ICDE
r4	Mike	Canada	Toronto	Toronto	VLDB

Dependency Theory

- Data dependencies (a.k.a. integrity constraints)

FD: [country] \rightarrow [capital]



	name	country	capital	city	conf
r1	George	China	Beijing	Beijing	SIGMOD
r2	Ian	China	Shanghai	Hongkong	ICDE
r3	Peter	China	Tokyo	Tokyo	ICDE
r4	Mike	Canada	Toronto	Toronto	VLDB

Data dependencies are not sufficient to guide dependable data repairing

User Guidance

	name	country	capital	city	conf
r1	George	China	Beijing	Beijing	SIGMOD
r2	Ian	China	Shanghai	Hongkong	ICDE
r3	Peter	China	Tokyo	Tokyo	ICDE
r4	Mike	Canada	Toronto	Toronto	VLDB

	country	capital
	China	Beijing
	Canada	Ottawa
	Japan	Tokyo

User Guidance

editing rule: ((country, country) -> (capital, capital))

	name	country	capital	city	conf
r1	George	China	Beijing	Beijing	SIGMOD
r2	Ian	China	Shanghai	Hongkong	ICDE
r3	Peter	China	Tokyo	Tokyo	ICDE
r4	Mike	Canada	Toronto	Toronto	VLDB

	country	capital
	China	Beijing
	Canada	Ottawa
	Japan	Tokyo

User Guidance

editing rule: ((country, country) -> (capital, capital))

	name	country	capital	city	conf		country	capital
r1	George	China	Beijing	Beijing	SIGMOD		China	Beijing
r2	Ian	China	Shanghai	Hongkong	ICDE		Canada	Ottawa
r3	Peter	China	Tokyo	Tokyo	ICDE		Japan	Tokyo
r4	Mike	Canada	Toronto	Toronto	VLDB			

User Guidance

editing rule: ((country, country) -> (capital, capital))

	name	country	capital	city	conf		country	capital
r1	George	China	Beijing	Beijing	SIGMOD		China	Beijing
r2	Ian	China	Shanghai	Hongkong	ICDE		Canada	Ottawa
r3	Peter	China	Tokyo	Tokyo	ICDE		Japan	Tokyo
r4	Mike	Canada	Toronto	Toronto	VLDB			

Is r2[country] China?
YES.



User Guidance

editing rule: ((country, country) -> (capital, capital))

	name	country	capital	city	conf		country	capital
r1	George	China	Beijing	Beijing	SIGMOD		China	Beijing
r2	Ian	China	Beijing	Hongkong	ICDE		Canada	Ottawa
r3	Peter	China	Tokyo	Tokyo	ICDE		Japan	Tokyo
r4	Mike	Canada	Toronto	Toronto	VLDB			

Is r2[country] China?
YES.



User Guidance

editing rule: ((country, country) -> (capital, capital))

	name	country	capital	city	conf		country	capital
r1	George	China	Beijing	Beijing	SIGMOD		China	Beijing
r2	Ian	China	Beijing	Hongkong	ICDE		Canada	Ottawa
r3	Peter	China	Tokyo	Tokyo	ICDE		Japan	Tokyo
r4	Mike	Canada	Toronto	Toronto	VLDB			

Is r2[country] China?
YES.

Is r1[country] China?

Is r3[country] China?

Is r4[country] Canada?

.....



User Guidance

editing rule: ((country, country) -> (capital, capital))

	name	country	capital	city	conf		country	capital
r1	George	China	Beijing	Beijing	SIGMOD		China	Beijing
r2	Ian	China	Beijing	Hongkong	ICDE		Canada	Ottawa
r3	Peter	China	Tokyo	Tokyo	ICDE		Japan	Tokyo
r4	Mike	Canada	Toronto	Toronto	VLDB			

Is r2[country] China?
YES.

Is r1[country] China?

Is r3[country] China?

Is r4[country] Canada?

.....



check **each** tuple: not cheap !!

Heuristic (Automated)

precision: +
recall: ++

precision: ++
recall: ++

Certain (User guided)

precision: +
recall: ++

Heuristic **(Automated)**

precision: ++
recall: +

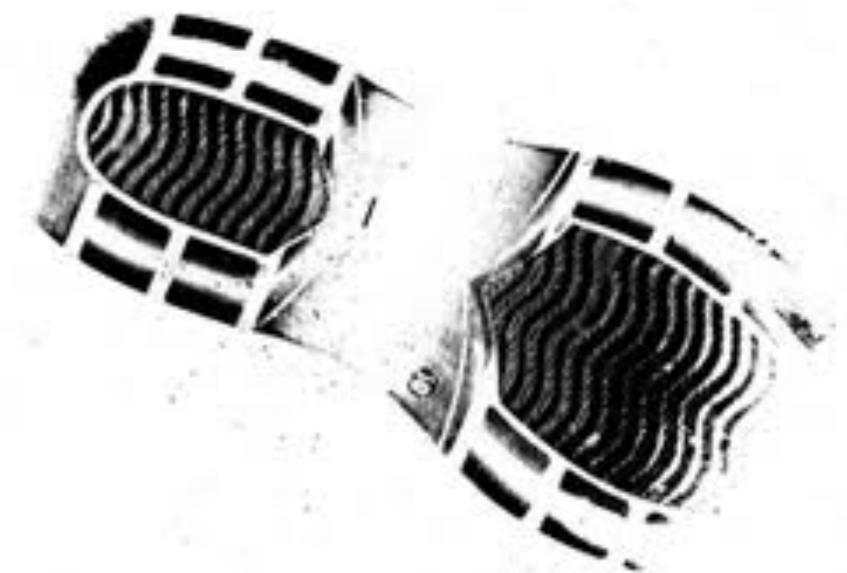
Fixing Rules **(Automated)**

precision: ++
recall: ++

Certain **(User guided)**



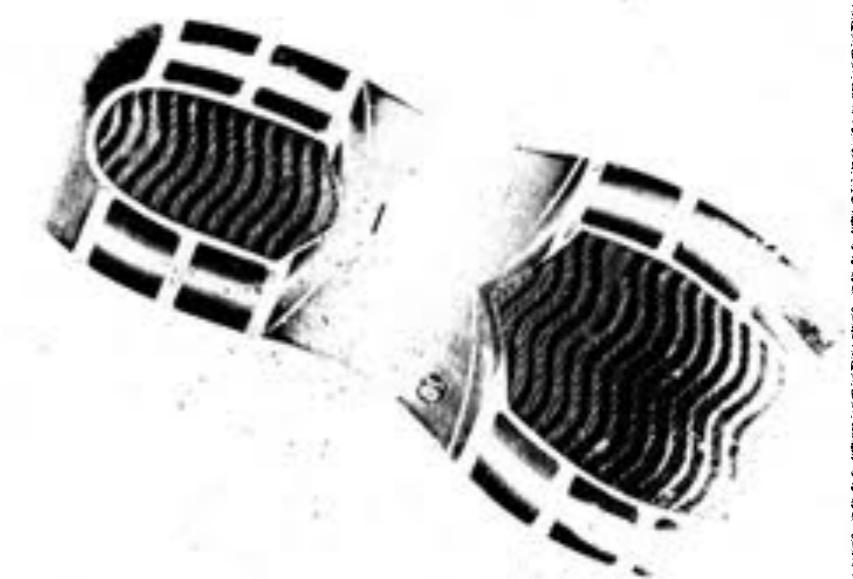






negative

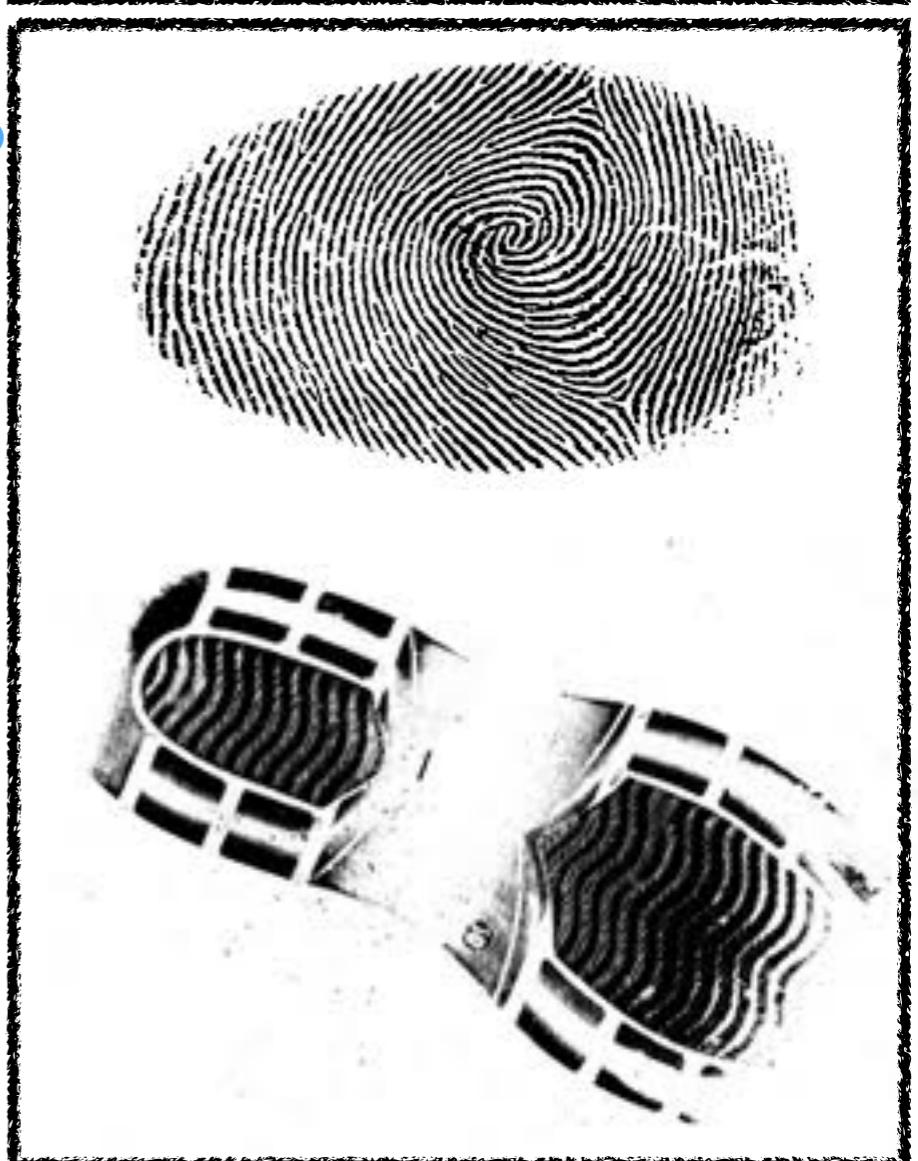
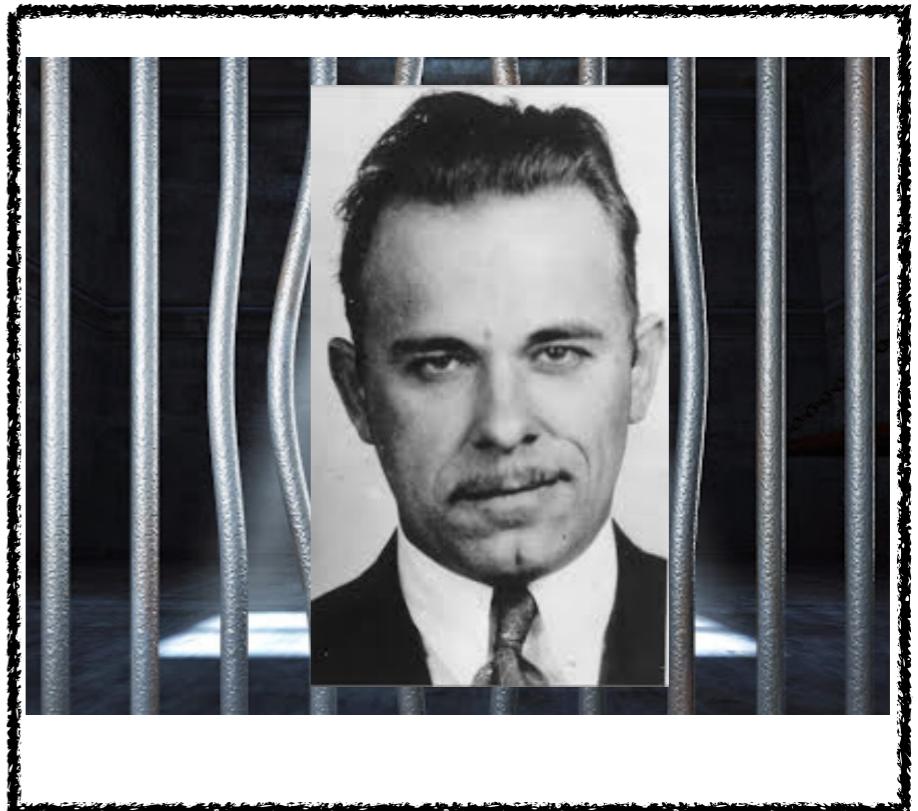
evidence





negative

evidence



Data patterns

country
capital

China
Shanghai

Data patterns

country
capital

China
Shanghai

evidence
negative

Data patterns

country
capital

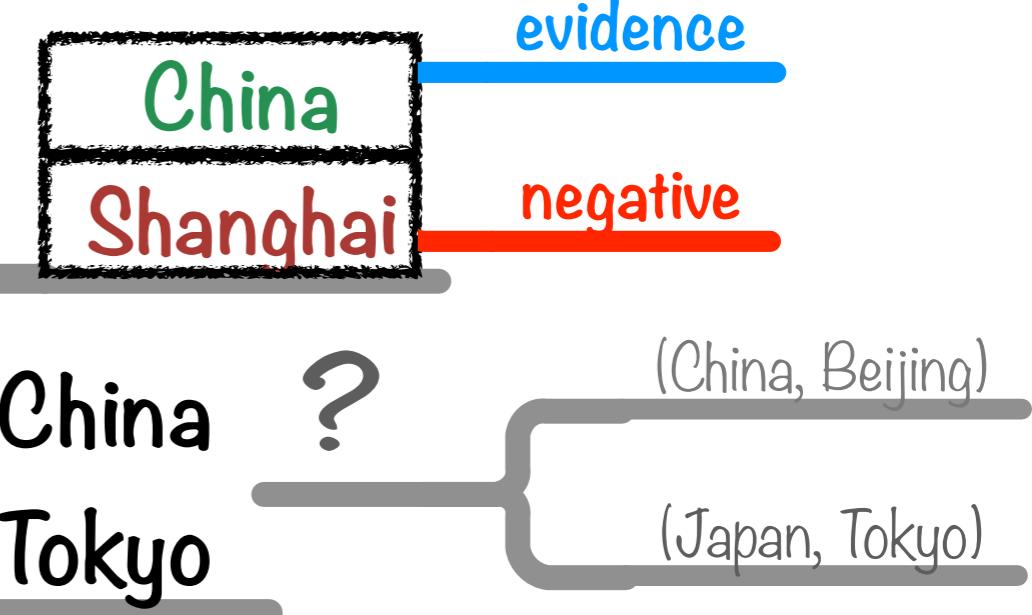
China
Shanghai

China
Tokyo

evidence
negative

Data patterns

country
capital



Data patterns

country
capital

name
work mail

China
Shanghai

China ?
Tokyo

ian

ian@gmail.com

evidence

negative

(China, Beijing)

(Japan, Tokyo)

Data patterns

country
capital

name
work mail

China
Shanghai

China ?
Tokyo

ian
ian@gmail.com

evidence

negative

(China, Beijing)

(Japan, Tokyo)

evidence

negative

Data patterns

country
capital

name
work mail

city
area code



Data patterns

country
capital

name
work mail

city
area code

China
Shanghai

China ?
Tokyo

ian
ian@gmail.com

Beijing
110002

evidence

negative

(China, Beijing)

(Japan, Tokyo)

evidence

negative

evidence

negative

Fixing Rules

• Syntax

fR1: (([country], [China]), (capital, {Shanghai, Hongkong})) -> Beijing

evidence	negative	fact
country	{capital	capital
China	Shanghai	Beijing
	Hongkong	

Fixing Rules

• Syntax

fR1: (([country], [China]), (capital, {Shanghai, Hongkong})) -> Beijing

evidence	negative	fact
country	{capital}	capital
China	Shanghai	Beijing
	Hongkong	

	name	country	capital	city	conf
r1	George	China	Beijing	Beijing	SIGMOD
r2	Ian	China	Shanghai	Hongkong	ICDE
r3	Peter	China	Tokyo	Tokyo	ICDE
r4	Mike	Canada	Toronto	Toronto	VLDB

Fixing Rules

• Syntax

fR1: (([country], [China]), (capital, {Shanghai, Hongkong})) -> Beijing

evidence	negative	fact
country	{capital}	capital
China	Shanghai	Beijing
	Hongkong	

	name	country	capital	city	conf
r1	George	China	Beijing	Beijing	SIGMOD
r2	Ian	China	Shanghai	Hongkong	ICDE
r3	Peter	China	Tokyo	Tokyo	ICDE
r4	Mike	Canada	Toronto	Toronto	VLDB

Fixing Rules

• Syntax

fR1: (([country], [China]), (capital, {Shanghai, Hongkong})) -> Beijing

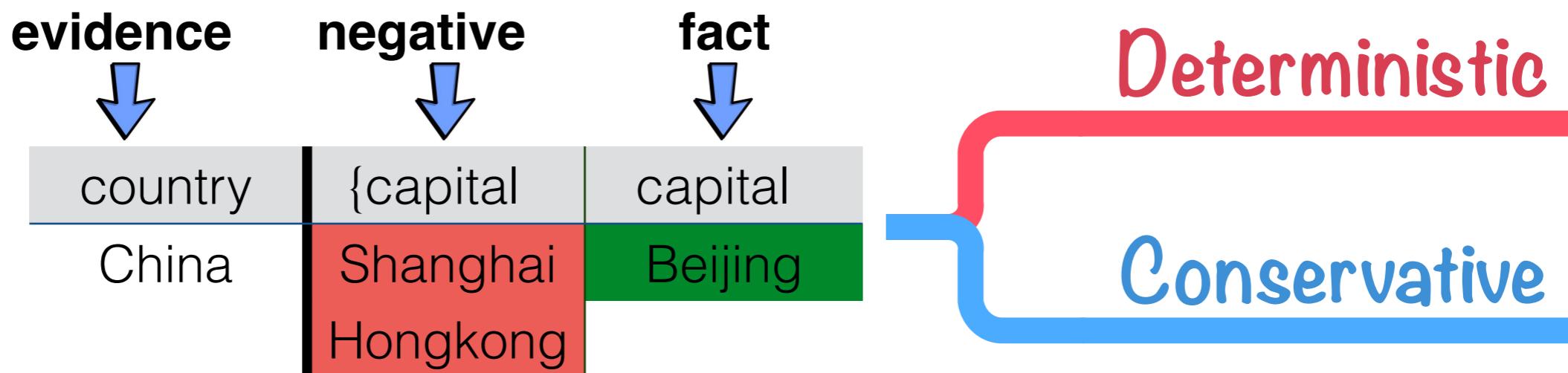
evidence	negative	fact
country	{capital}	capital
China	Shanghai	Beijing
	Hongkong	

	name	country	capital	city	conf
r1	George	China	Beijing	Beijing	SIGMOD
r2	Ian	China	Beijing	Hongkong	ICDE
r3	Peter	China	Tokyo	Tokyo	ICDE
r4	Mike	Canada	Toronto	Toronto	VLDB

Fixing Rules

• Syntax

fR1: (([country], [China]), (capital, {Shanghai, Hongkong})) -> Beijing



	name	country	capital	city	conf
r1	George	China	Beijing	Beijing	SIGMOD
r2	Ian	China	Beijing	Hongkong	ICDE
r3	Peter	China	Tokyo	Tokyo	ICDE
r4	Mike	Canada	Toronto	Toronto	VLDB

Applying One Fixing Rule

country	{capital	capital
China	Shanghai Hongkong	Beijing

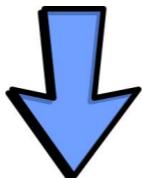
r2	Ian	China	Shanghai	Hongkong	ICDE
-----------	-----	-------	----------	----------	------

Applying One Fixing Rule

country	{capital	capital
China	Shanghai	Beijing
	Hongkong	

r2

Ian	China	Shanghai	Hongkong	ICDE
-----	-------	----------	----------	------



r2'

Ian	China	Beijing	Hongkong	ICDE
-----	-------	---------	----------	------



Applying Multiple Fixing Rules

- **Fixes**

fR₁'

country	{capital	capital
China	Shanghai	Beijing
	Hongkong	
	Tokyo	

fR₃

capital	city	conf	{country	country
Tokyo	Tokyo	ICDE	China	Japan

Applying Multiple Fixing Rules

- **Fixes**

fR_{1'}

country	{capital	capital
China	Shanghai	Beijing
	Hongkong	
	Tokyo	

fR₃

capital	city	conf	{country	country
Tokyo	Tokyo	ICDE	China	Japan

r2

Ian	China	Shanghai	Hongkong	ICDE
-----	-------	----------	----------	------

fR_{1'}



r2'

Ian	China	Beijing	Hongkong	ICDE
-----	-------	---------	----------	------



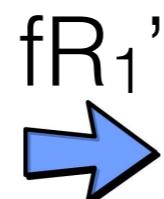
Applying Multiple Fixing Rules

- **Fixes**

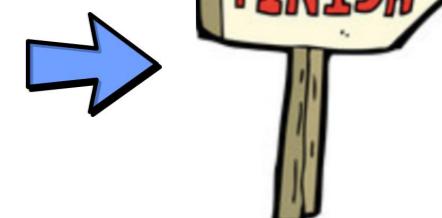
fR ₁ '		
country	{capital	capital
China	Shanghai	Beijing
	Hongkong	
	Tokyo	

fR ₃				
capital	city	conf	{country	country
Tokyo	Tokyo	ICDE	China	Japan

r2	Ian	China	Shanghai	Hongkong	ICDE
----	-----	-------	----------	----------	------

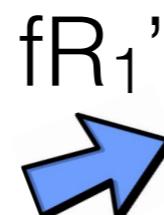


r2'	Ian	China	Beijing	Hongkong	ICDE
-----	-----	-------	---------	----------	------



- **Unique fixes**

r3	Peter	China	Tokyo	Tokyo	ICDE
----	-------	-------	-------	-------	------



r3'	Peter	China	Beijing	Tokyo	ICDE
-----	-------	-------	---------	-------	------



Applying Multiple Fixing Rules

- **Fixes**

fR _{1'}		
country	{capital}	capital
China	Shanghai	Beijing
	Hongkong	
	Tokyo	

fR ₃				
capital	city	conf	{country}	country
Tokyo	Tokyo	ICDE	China	Japan

- **Unique fixes**

r3	Peter	China	Tokyo	Tokyo	ICDE
----	-------	-------	-------	-------	------

fR_{1'}



r2'

Ian	China	Beijing	Hongkong	ICDE
-----	-------	---------	----------	------



fR_{1'}



r3'

Peter	China	Beijing	Tokyo	ICDE
-------	-------	---------	-------	------



fR₃



r3''

Peter	Japan	Tokyo	Tokyo	ICDE
-------	-------	-------	-------	------



Fundamentals

Fundamental Problems

Termination

Yes

Consistency

PTIME

Implication

coNP-complete

Determinism

Yes

Ensuring Consistency

- **Tuple enumeration**

fR ₁		
country	{capital}	capital
China	Shanghai	Beijing
	Hongkong	

fR ₂		
country	{capital}	capital
Canada	Toronto	Ottawa

Ensuring Consistency

- **Tuple enumeration**

fR ₁		
country	{capital}	capital
China	Shanghai	Beijing
	Hongkong	

fR ₂		
country	{capital}	capital
Canada	Toronto	Ottawa

(o, China, Shanghai, o, o)

(o, China, Hongkong, o, o)

(o, China, Toronto, o, o)

(o, Canada, Shanghai, o, o)

(o, Canada, Hongkong, o, o)

(o, Canada, Toronto, o, o)

Ensuring Consistency

- **Tuple enumeration**

fR ₁		
country	{capital}	capital
China	Shanghai Hongkong	Beijing

fR ₂		
country	{capital}	capital
Canada	Toronto	Ottawa

(o, China, Shanghai, o, o)

(o, China, Hongkong, o, o)

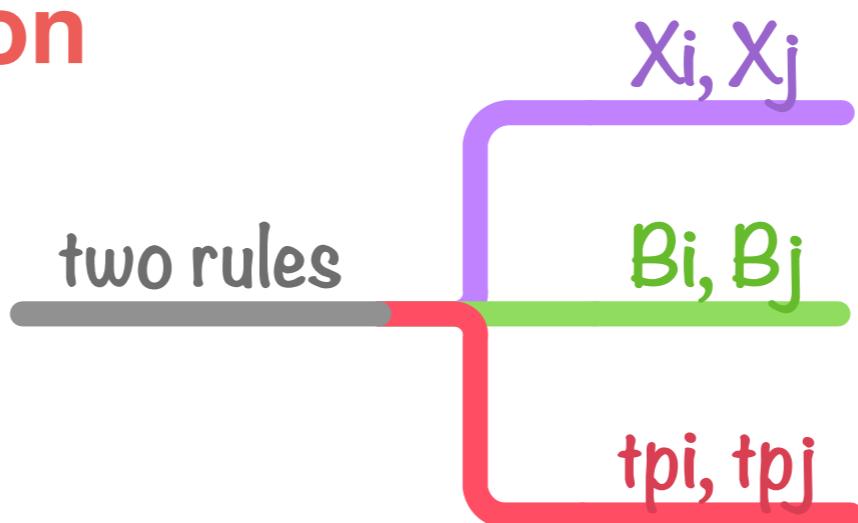
(o, China, Toronto, o, o)

(o, Canada, Shanghai, o, o)

(o, Canada, Hongkong, o, o)

(o, Canada, Toronto, o, o)

- **Rule characterization**



Repair

Repairing with Fixing Rules

	name	country	capital	city	conf
r1	George	China	Beijing	Beijing	SIGMOD
r2	Ian	China	Shanghai	Hongkong	ICDE
r3	Peter	China	Tokyo	Tokyo	ICDE
r4	Mike	Canada	Toronto	Toronto	VLDB

fR₁

country	{capital}	capital
China	Shanghai	Beijing
	Hongkong	

fR₃

capital	city	conf	{country}	country
Tokyo	Tokyo	ICDE	China	Japan

fR₂

country	{capital}	capital
Canada	Toronto	Ottawa

fR₄

capital	conf	{city}	city
Beijing	ICDE	Hongkong	Shanghai

Repairing with Fixing Rules

	name	country	capital	city	conf
r1	George	China	Beijing	Beijing	SIGMOD
r2	Ian	China	Shanghai	Hongkong	ICDE
r3	Peter	China	Tokyo	Tokyo	ICDE
r4	Mike	Canada	Toronto	Toronto	VLDB

fR₁

country	{capital}	capital
China	Shanghai	Beijing
	Hongkong	

fR₂

country	{capital}	capital
Canada	Toronto	Ottawa

fR₃

capital	city	conf	{country}	country
Tokyo	Tokyo	ICDE	China	Japan

fR₄

capital	conf	{city}	city
Beijing	ICDE	Hongkong	Shanghai

Key

List

- country, China → fR1
- country, Canada → fR2
- conf, ICDE → fR3, fR4
- capital, Tokyo → fR3
- city, Tokyo → fR3
- capital, Beijing → fR4

Repairing with Fixing Rules

	name	country	capital	city	conf
r1	George	China	Beijing	Beijing	SIGMOD
r2	Ian	China	Shanghai	Hongkong	ICDE
r3	Peter	China	Tokyo	Tokyo	ICDE
r4	Mike	Canada	Toronto	Toronto	VLDB

country	{capital}	capital
China	Shanghai	Beijing
	Hongkong	

capital	city	conf	{country}	country
Tokyo	Tokyo	ICDE	China	Japan

country	{capital}	capital
Canada	Toronto	Ottawa

capital	conf	{city}	city
Beijing	ICDE	Hongkong	Shanghai

Key

List

r1:
itr1:

cnt(fR1) = 1, cnt(fR4) = 1, rules = {fR1}
r1' = r1, rules = {}

country, China

→ fR1

country, Canada

→ fR2

conf, ICDE

→ fR3, fR4

capital, Tokyo

→ fR3

city, Tokyo

→ fR3

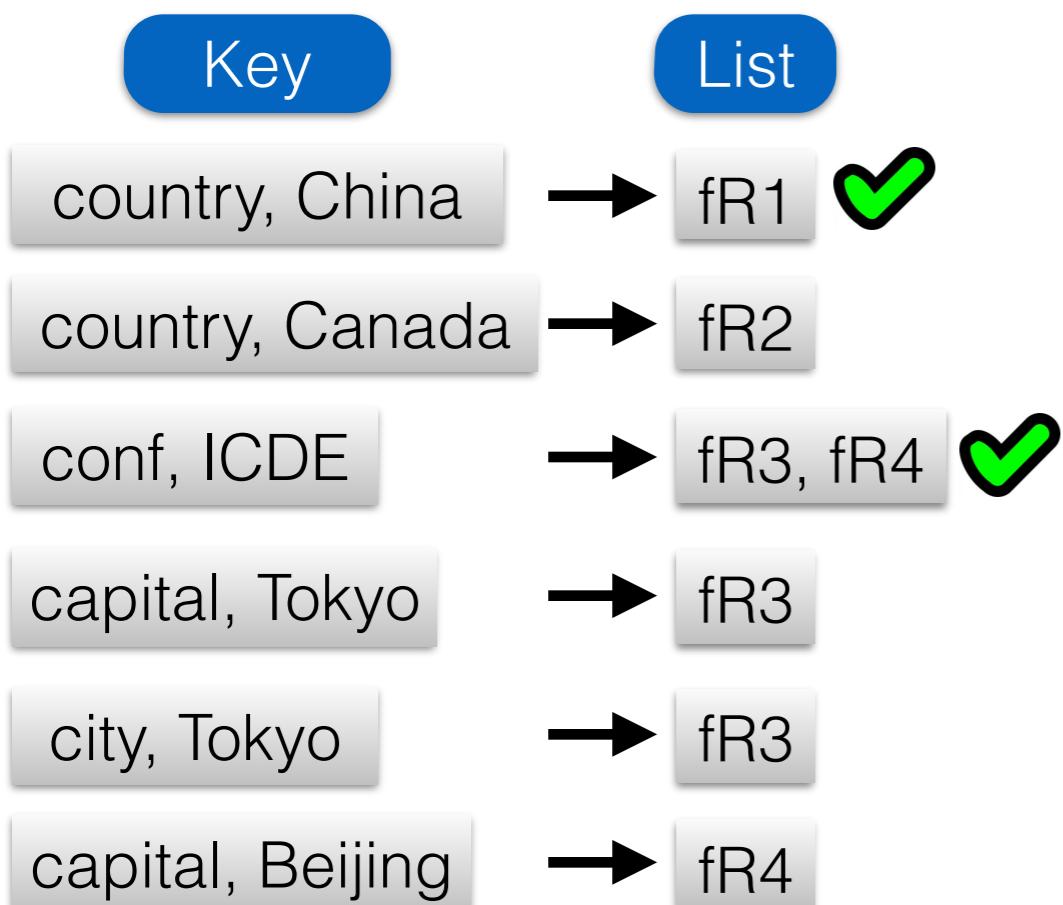
capital, Beijing

→ fR4

Repairing with Fixing Rules

	name	country	capital	city	conf
r1	George	China	Beijing	Beijing	SIGMOD
r2	Ian	China	Shanghai	Hongkong	ICDE
r3	Peter	China	Tokyo	Tokyo	ICDE
r4	Mike	Canada	Toronto	Toronto	VLDB

fR ₁		fR ₃	
country	{capital}	capital	country
China	Shanghai	Beijing	
	Hongkong		
fR ₂		fR ₄	
country	{capital}	capital	city
Canada	Toronto	Ottawa	
capital	conf	{city}	city
Beijing	ICDE	Hongkong	Shanghai

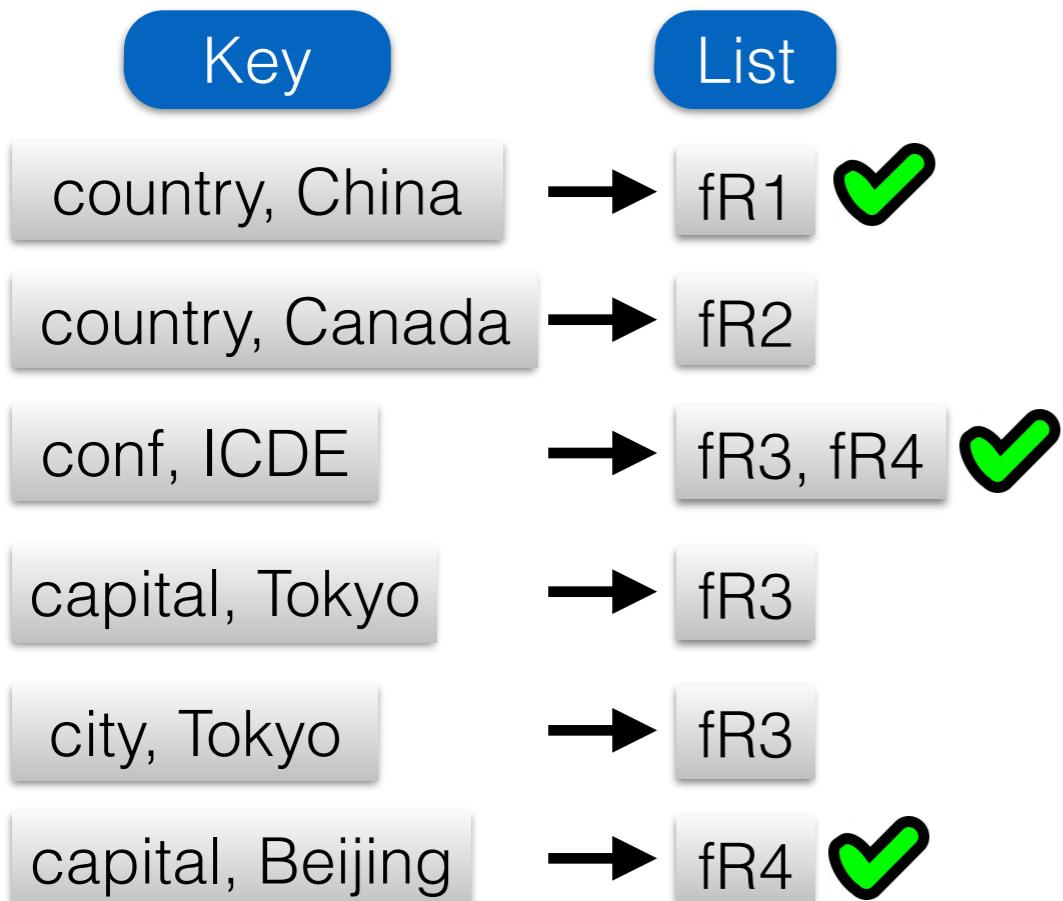


r1: itr1:	cnt(fR1) = 1, cnt(fR4) = 1, rules = {fR1} r1' = r1, rules = {}
r2: itr1: itr2:	cnt(fR1, fR3, fR4) = 1, rules = {fR1} r2'[capital] = Beijing, cnt(fR3) = 1, cnt(fR4) = 2, rules = {fR4} r2'[city] = Shanghai, rules = {}

Repairing with Fixing Rules

	name	country	capital	city	conf
r1	George	China	Beijing	Beijing	SIGMOD
r2	Ian	China	Shanghai	Hongkong	ICDE
r3	Peter	China	Tokyo	Tokyo	ICDE
r4	Mike	Canada	Toronto	Toronto	VLDB

fR ₁		fR ₃	
country	{capital}	capital	country
China	Shanghai	Beijing	
	Hongkong		
fR ₂		fR ₄	
country	{capital}	capital	city
Canada	Toronto	Ottawa	
capital	conf	{city}	city
Beijing	ICDE	Hongkong	Shanghai



r1: itr1:	cnt(fR1) = 1, cnt(fR4) = 1, rules = {fR1} r1' = r1, rules = {}
r2: itr1: itr2:	cnt(fR1, fR3, fR4) = 1, rules = {fR1} r2'[capital] = Beijing, cnt(fR3) = 1, cnt(fR4) = 2, rules = {fR4} r2'[city] = Shanghai, rules = {}

Repairing with Fixing Rules

	name	country	capital	city	conf
r1	George	China	Beijing	Beijing	SIGMOD
r2	Ian	China	Beijing	Shanghai	ICDE
r3	Peter	China	Tokyo	Tokyo	ICDE
r4	Mike	Canada	Toronto	Toronto	VLDB

country	{capital}	capital
China	Shanghai	Beijing
	Hongkong	

capital	city	conf	{country}	country
Tokyo	Tokyo	ICDE	China	Japan

country	{capital}	capital
Canada	Toronto	Ottawa

capital	conf	{city}	city
Beijing	ICDE	Hongkong	Shanghai

Key

List

country, China

→ fR1

country, Canada

→ fR2

conf, ICDE

→ fR3, fR4

capital, Tokyo

→ fR3

city, Tokyo

→ fR3

capital, Beijing

→ fR4

r1:
itr1:

cnt(fR1) = 1, cnt(fR4) = 1, rules = {fR1}
r1' = r1, rules = {}

r2:
itr1:

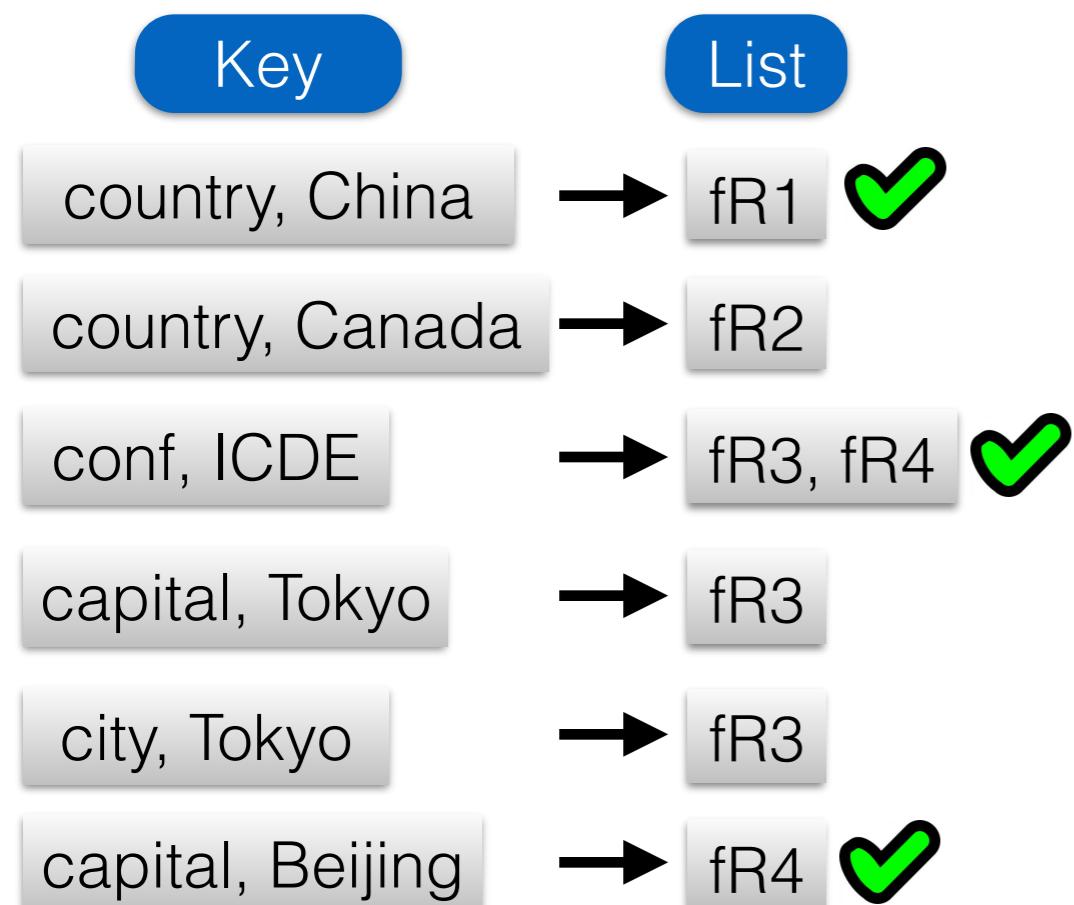
cnt(fR1, fR3, fR4) = 1, rules = {fR1}
r2'[capital] = Beijing, cnt(fR3) = 1,
cnt(fR4) = 2, rules = {fR4}
r2'[city] = Shanghai, rules = {}

itr2:

Repairing with Fixing Rules

	name	country	capital	city	conf
r1	George	China	Beijing	Beijing	SIGMOD
r2	Ian	China	Beijing	Shanghai	ICDE
r3	Peter	Japan	Tokyo	Tokyo	ICDE
r4	Mike	Canada	Toronto	Toronto	VLDB

fR ₁		fR ₃	
country	{capital}	capital	country
China	Shanghai	Beijing	
	Hongkong		
fR ₂		fR ₄	
country	{capital}	capital	city
Canada	Toronto	Ottawa	
capital	conf	{city}	city
Beijing	ICDE	Hongkong	Shanghai

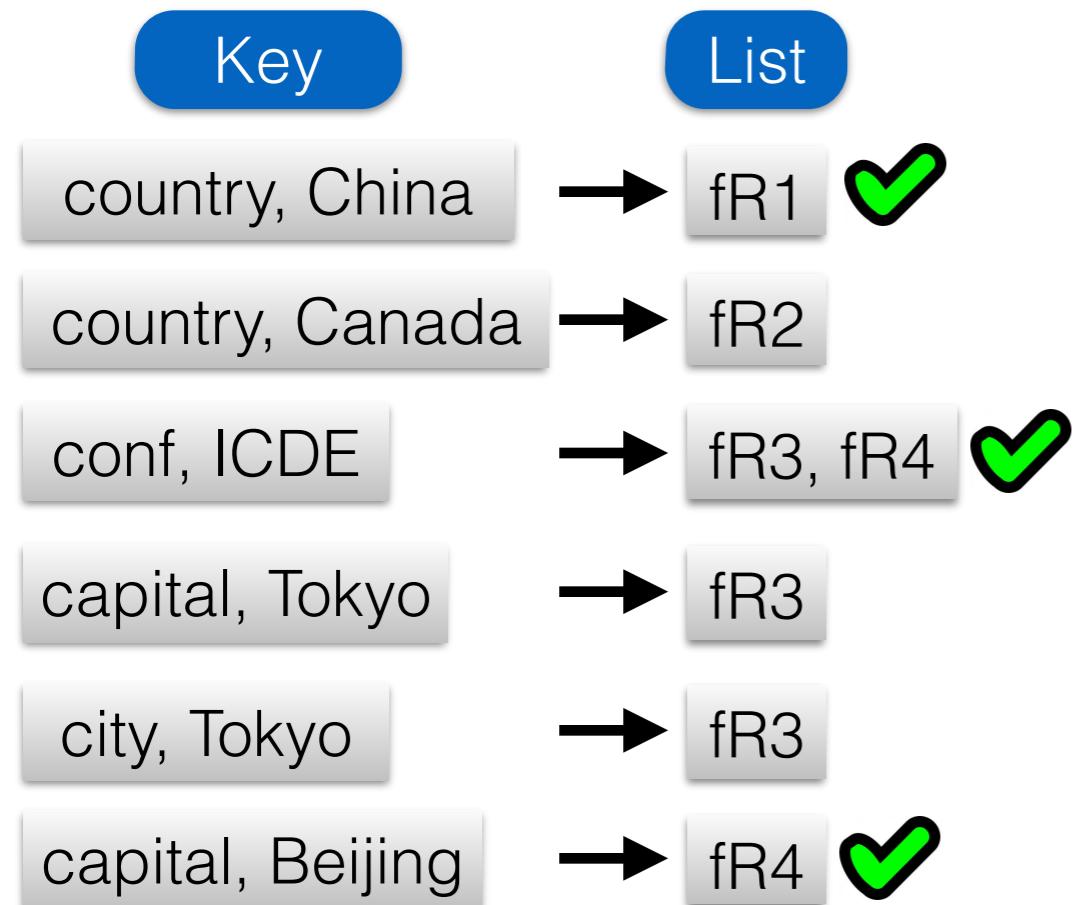


r1: itr1:	cnt(fR1) = 1, cnt(fR4) = 1, rules = {fR1} r1' = r1, rules = {}
r2: itr1: itr2:	cnt(fR1, fR3, fR4) = 1, rules = {fR1} r2'[capital] = Beijing, cnt(fR3) = 1, cnt(fR4) = 2, rules = {fR4} r2'[city] = Shanghai, rules = {}
r3: itr1:	cnt(fR3) = 3, cnt(fR4) = 1, rules = {fR3} r3'[country] = Japan, rules = {}

Repairing with Fixing Rules

	name	country	capital	city	conf
r1	George	China	Beijing	Beijing	SIGMOD
r2	Ian	China	Beijing	Shanghai	ICDE
r3	Peter	Japan	Tokyo	Tokyo	ICDE
r4	Mike	Canada	Ottawa	Toronto	VLDB

fR ₁		fR ₃	
country	{capital}	capital	country
China	Shanghai	Beijing	
	Hongkong		
fR ₂		fR ₄	
country	{capital}	capital	city
Canada	Toronto	Ottawa	
capital	conf	{city}	city
Beijing	ICDE	Hongkong	Shanghai



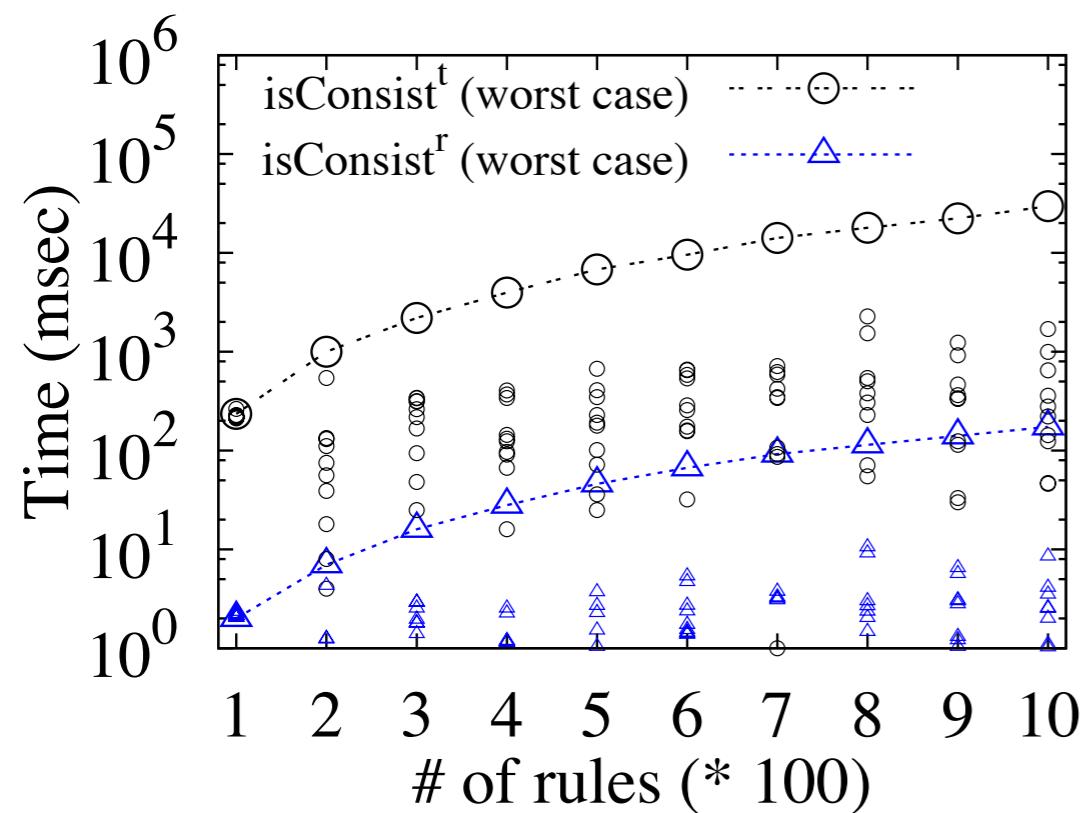
r1: itr1:	cnt(fR1) = 1, cnt(fR4) = 1, rules = {fR1} r1' = r1, rules = {}
r2: itr1:	cnt(fR1, fR3, fR4) = 1, rules = {fR1} r2'[capital] = Beijing, cnt(fR3) = 1, cnt(fR4) = 2, rules = {fR4}
itr2:	r2'[city] = Shanghai, rules = {}
r3: itr1:	cnt(fR3) = 3, cnt(fR4) = 1, rules = {fR3} r3'[country] = Japan, rules = {}
r4: itr1:	cnt(fR3) = 1, rules = {fR2} r4'[capital] = Ottawa, rules = {}

Experiment

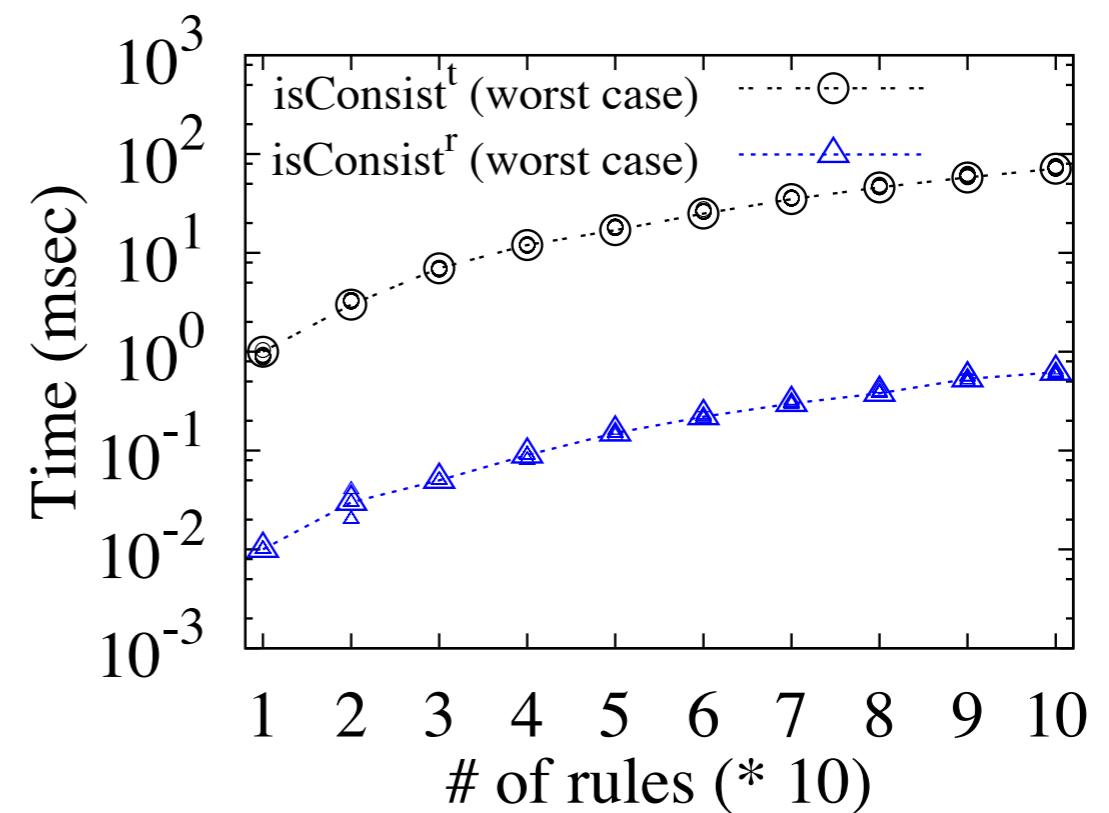
Experimental Study

- Efficiency of checking consistency
- Accuracy
- Efficiency of repairing algorithms

Efficiency of Checking Consistency

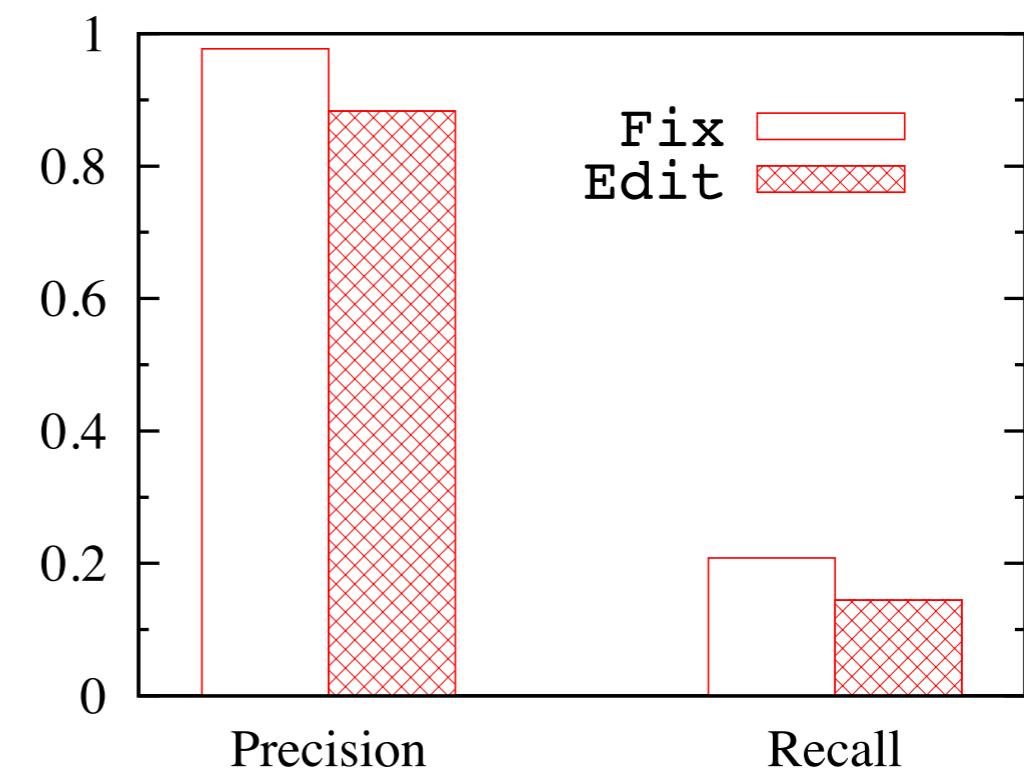
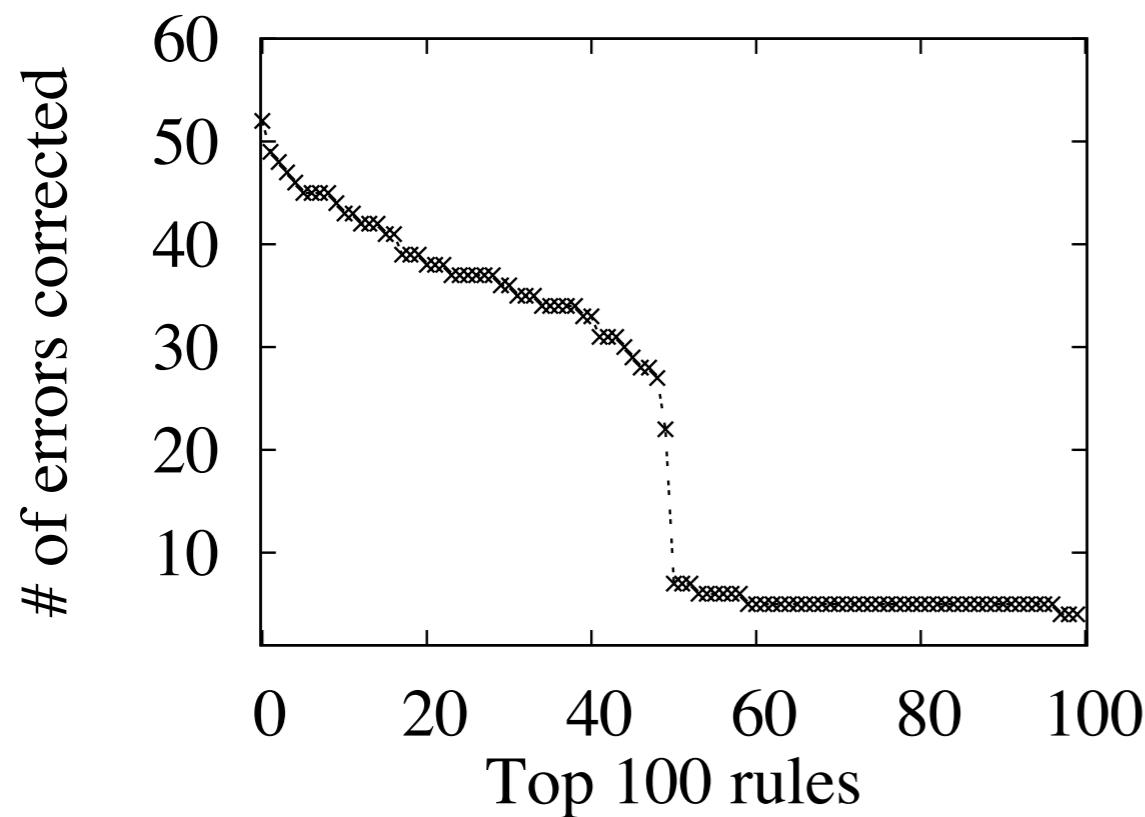


Hospital data



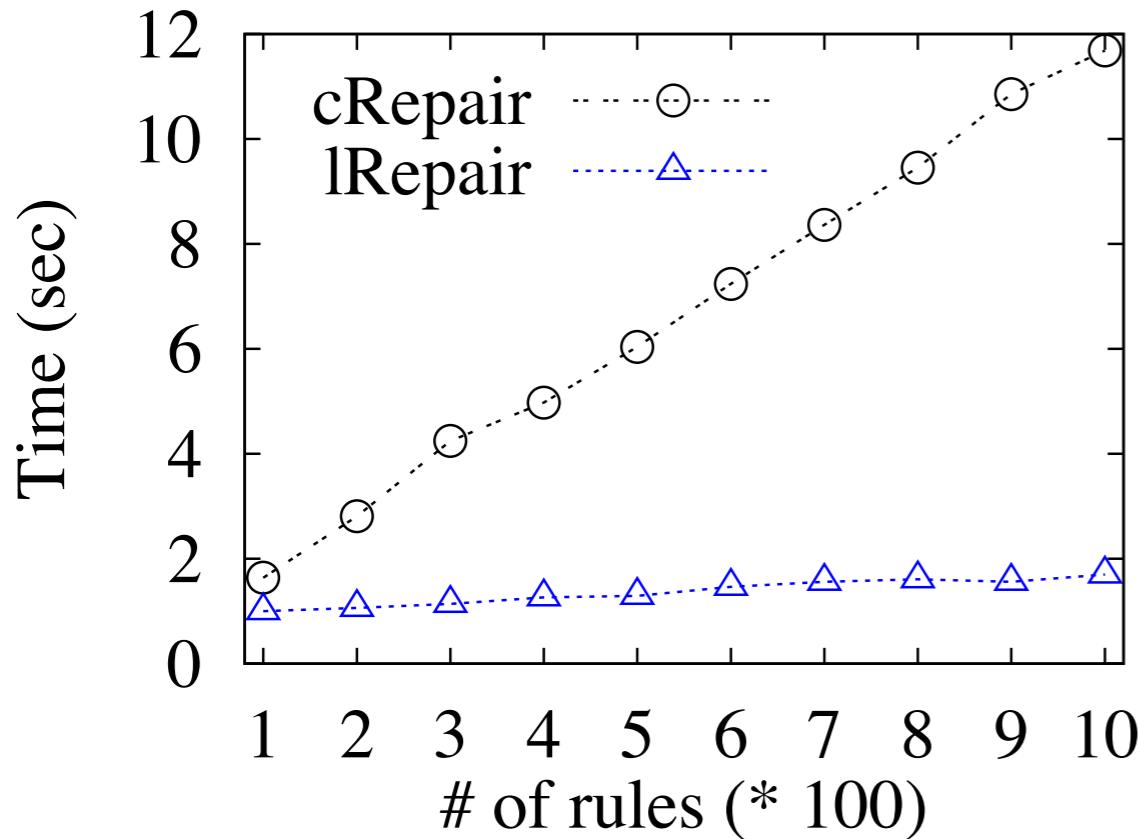
UIS

Accuracy

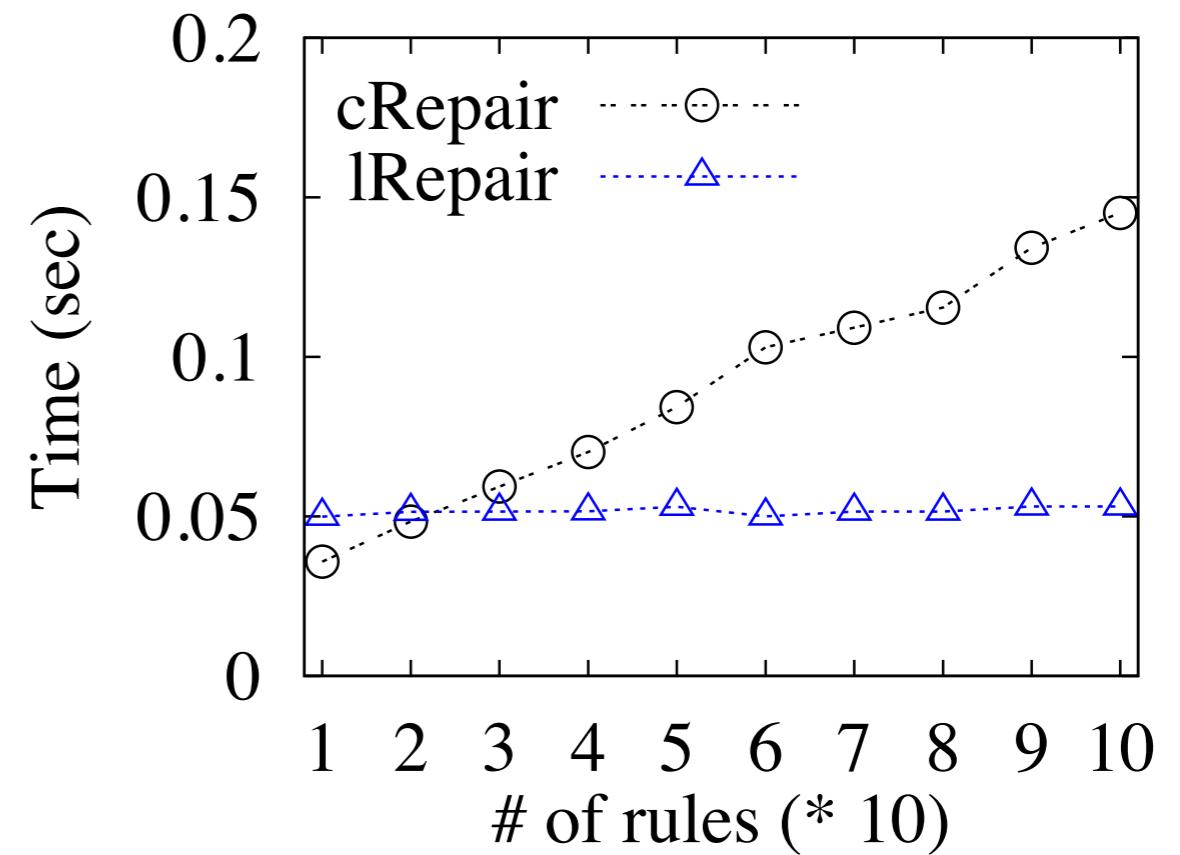


Hospital data

Efficiency of Repairing Algorithms



Hospital data



UIS

precision: +
recall: ++

**Heuristic
(Automated)**

precision: ++
recall: +

Fixing Rules

precision: ++
recall: ++

**Certain
(User guided)**

precision: +
recall: ++

**Heuristic
(Automated)**

precision: ++
recall: +

Conclusion:
Automated
Dependable
Fundamentals
Repair

Fixing Rules

precision: ++
recall: ++

**Certain
(User guided)**

precision: +
recall: ++

**Heuristic
(Automated)**

precision: ++
recall: +

Conclusion:
Automated
Dependable
Fundamentals
Repair

Fixing Rules

Future work:

Discovery
Generalized fixing rules

precision: ++
recall: ++

**Certain
(User guided)**

**How to Get
My Rules?**

Generating Fixing Rules

country	{capital	capital
China	Shanghai	Beijing
	Hongkong	

 Freebase®

Generating Fixing Rules

country	{capital	capital
China	Shanghai	Beijing
	Hongkong	



MQL1

```
[{"type": "/location/country",
"name": null,
"/location/country/capital": []
}]
```

Generating Fixing Rules

country	{capital	capital
China	Shanghai	Beijing
	Hongkong	



MQL1

```
[{"type": "/location/country",
"name": null,
"/location/country/capital": []}]
```

MQL2

```
[{"/location/country/iso3166_1_shortname": "CHINA",
"/location/location/contains": [
  {"name": null,
   "type": "/location/citytown"}]}]
```