

Data Cleaning

Nan Tang, QCRI



معهد قطر لبحوث الحوسبة
Qatar Computing Research Institute

عضو في مؤسسة قطر Member of Qatar Foundation

Big Data Cleaning

Nan Tang, QCRI



معهد قطر لبحوث الحوسبة
Qatar Computing Research Institute

عضو في مؤسسة قطر Member of Qatar Foundation

Big Data Cleaning

Nan Tang, QCRI



معهد قطر لبحوث الحوسبة
Qatar Computing Research Institute

عضو في مؤسسة قطر Member of Qatar Foundation

Data Cleaning?

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

Data Cleaning Market

INFORMATICA®

Data Explorer



Microsoft®
SQL Server®

Tamr

talend*
integration at any scale



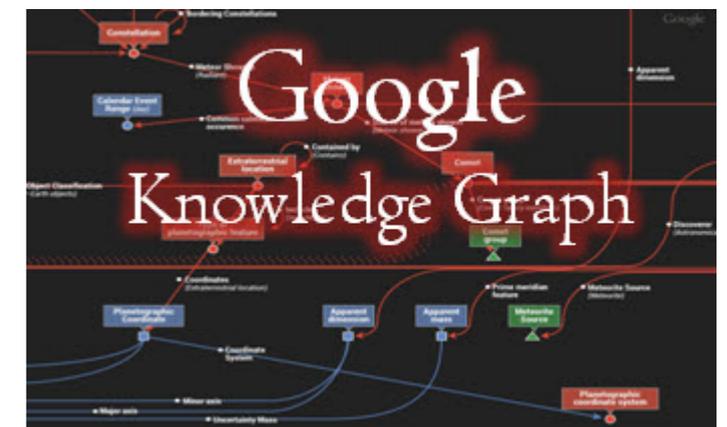
pentaho™
open source business intelligence™

ORACLE®
Data Quality 11g

IBM®

InfoSphere™
software

Trusted Information



Data Cleaning Problems

name	graduated	affiliation	country	capital	age
Nan Tang	CUHK	QCRI	Qatari	Doha	33
Xiaokui Xiao	CUHK	NTU	Singapore	Singapore	
Nan Tang	CUHK	University of Edinburgh	UK	Edinburgh	31
Gao Cong	NUS	University of Edinburgh	UK	London	36

Data Cleaning Problems

typo

name	graduated	affiliation	country	capital	age
Nan Tang	CUHK	QCRI	Qatari	Doha	33
Xiaokui Xiao	CUHK	NTU	Singapore	Singapore	
Nan Tang	CUHK	University of Edinburgh	UK	Edinburgh	31
Gao Cong	NUS	University of Edinburgh	UK	London	36

Data Cleaning Problems

typo

name	graduated	affiliation	country	capital	age
Nan Tang	CUHK	QCRI	Qatar	Doha	33
Xiaokui Xiao	CUHK	NTU	Singapore	Singapore	
Nan Tang	CUHK	University of Edinburgh	UK	Edinburgh	31
Gao Cong	NUS	University of Edinburgh	UK	London	36

Data Cleaning Problems

name	graduated	affiliation	country	capital	age
Nan Tang	CUHK	QCRI	Qatar	Doha	33
Xiaokui Xiao	CUHK	NTU	Singapore	Singapore	
Nan Tang	CUHK	University of Edinburgh	UK	Edinburgh	31
Gao Cong	NUS	University of Edinburgh	UK	London	36

Data Cleaning Problems

name	graduated	affiliation	country	capital	age
Nan Tang	CUHK	QCRI	Qatar	Doha	33
Xiaokui Xiao	CUHK	NTU	Singapore	Singapore	
Nan Tang	CUHK	University of Edinburgh	UK	Edinburgh	31
Gao Cong	NUS	University of Edinburgh	UK	London	36

typo

Currency

Completeness

Data Cleaning Problems

name	graduated	affiliation	country	capital	age
Nan Tang	CUHK	QCRI	Qatar	Doha	33
Xiaokui Xiao	CUHK	NTU	Singapore	Singapore	
Nan Tang	CUHK	University of Edinburgh	UK	Edinburgh	31
Gao Cong	NUS	University of Edinburgh	UK	London	36

typo

Currency

Completeness

Data Cleaning Problems

The table below illustrates data cleaning problems. Callouts highlight specific issues:

- typo**: Points to the 'country' cell 'Qatar' in the first row.
- Consistency**: Points to the 'country' cell 'UK' in the third row.
- Currency**: Points to the 'age' cell '31' in the third row.
- Completeness**: Points to the empty 'age' cell in the second row.

name	graduated	affiliation	country	capital	age
Nan Tang	CUHK	QCRI	Qatar	Doha	33
Xiaokui Xiao	CUHK	NTU	Singapore	Singapore	
Nan Tang	CUHK	University of Edinburgh	UK	Edinburgh	31
Gao Cong	NUS	University of Edinburgh	UK	London	36

Data Cleaning Problems

Duplicates

typo

Consistency

Currency

name	graduated	affiliation	country	capital	age
Nan Tang	CUHK	QCRI	Qatar	Doha	33
Xiaokui Xiao	CUHK	NTU	Singapore	Singapore	
Nan Tang	CUHK	University of Edinburgh	UK	Edinburgh	31
Gao Cong	NUS	University of Edinburgh	UK	London	36

Completeness

Data Cleaning Problems

Duplicates

typo

Consistency

Currency

name	graduated	affiliation	country	capital	age
Nan Tang	CUHK	QCRI	Qatar	Doha	33
Xiaokui Xiao	CUHK	NTU	Singapore	Singapore	
Nan Tang	CUHK	University of Edinburgh	UK	Edinburgh	31
Gao Cong	NUS	University of Edinburgh	UK	London	36

Completeness

source2

name	affiliation
Nan Tang	QCRI

source3

name	affiliation
Nan Tang	CWI

.....

Data Cleaning Problems

Duplicates

typo

Consistency

Currency

name	graduated	affiliation	country	capital	age
Nan Tang	CUHK	QCRI	Qatar	Doha	33
Xiaokui Xiao	CUHK	NTU	Singapore	Singapore	
Nan Tang	CUHK	University of Edinburgh	UK	Edinburgh	31
Gao Cong	NUS	University of Edinburgh	UK	London	36

Completeness

source2

name	affiliation
Nan Tang	QCRI

source3

name	affiliation
Nan Tang	CWI

.....

truth discovery

Data Cleaning Problems

Duplicates

typo

Consistency

Currency

name	graduated	affiliation	country	capital	age
Nan Tang	CUHK	QCRI	Qatar	Doha	33
Xiaokui Xiao	CUHK	NTU	Singapore	Singapore	
Nan Tang	CUHK	University of Edinburgh	UK	Edinburgh	31
Gao Cong	NUS	University of Edinburgh	UK	London	36

Completeness

source2

name	affiliation
Nan Tang	QCRI

source3

name	affiliation
Nan Tang	CWI

.....

truth discovery

name	full
NTU	Nanyang Technological University
NUS	National University of Singapore

Data Cleaning Problems

Duplicates

typo

Consistency

Currency

name	graduated	affiliation	country	capital	age
Nan Tang	CUHK	QCRI	Qatar	Doha	33
Xiaokui Xiao	CUHK	NTU	Singapore	Singapore	
Nan Tang	CUHK	University of Edinburgh	UK	Edinburgh	31
Gao Cong	NUS	University of Edinburgh	UK	London	36

Completeness

source2

name	affiliation
Nan Tang	QCRI

source3

name	affiliation
Nan Tang	CWI

.....

truth discovery

name

full

NTU	Nanyang Technological University
NUS	National University of Singapore

ETL (transformation)

Data Cleaning Problems

Duplicates

typo

Consistency

Currency

name	graduated	affiliation	country	capital	age
Nan Tang	CUHK	QCRI	Qatar	Doha	33
Xiaokui Xiao	CUHK	NTU	Singapore	Singapore	
Nan Tang	CUHK	University of Edinburgh	UK	Edinburgh	31
Gao Cong	NUS	University of Edinburgh	UK	London	36

Completeness

source2

name	affiliation
Nan Tang	QCRI

source3

name	affiliation
Nan Tang	CWI

.....

truth discovery

name

full

NTU	Nanyang Technological University
NUS	National University of Singapore

ETL (transformation)

(UK, hasCapital, London)
KBs (e.g., Yago)

Data Cleaning Problems

Duplicates

name	graduated	affiliation	country	capital	age
Nan Tang	CUHK	QCRI	Qatar	Doha	33
Xiaokui Xiao	CUHK	NTU	Singapore	Singapore	
Nan Tang	CUHK	University of Edinburgh	UK	Edinburgh	31
Gao Cong	NUS	University of Edinburgh	UK	London	36

typo

Consistency

Currency

Completeness

Heterogeneous sources

(UK, hasCapital, London)
KBs (e.g., Yago)

source2

name	affiliation
Nan Tang	QCRI

source3

name	affiliation
Nan Tang	CWI

.....

truth discovery

name	full
NTU	Nanyang Technological University
NUS	National University of Singapore

ETL (transformation)

Data Cleaning Problems

Duplicates

Volume

Velocity

typo

Consistency

Currency

name	graduated	affiliation	country	capital	age
	HK	QCRI	Qatar	Doha	33
	HK	NTU	Singapore	Singapore	
Nan Tang	CUHK	University of Edinburgh	UK	Edinburgh	31
Gao Cong	NUS	University of Edinburgh	UK	London	36

Completeness

source2

name	affiliation
Nan Tang	QCRI

source3

name	affiliation
Nan Tang	CWI

.....

truth discovery

V...

Heterogeneous sources

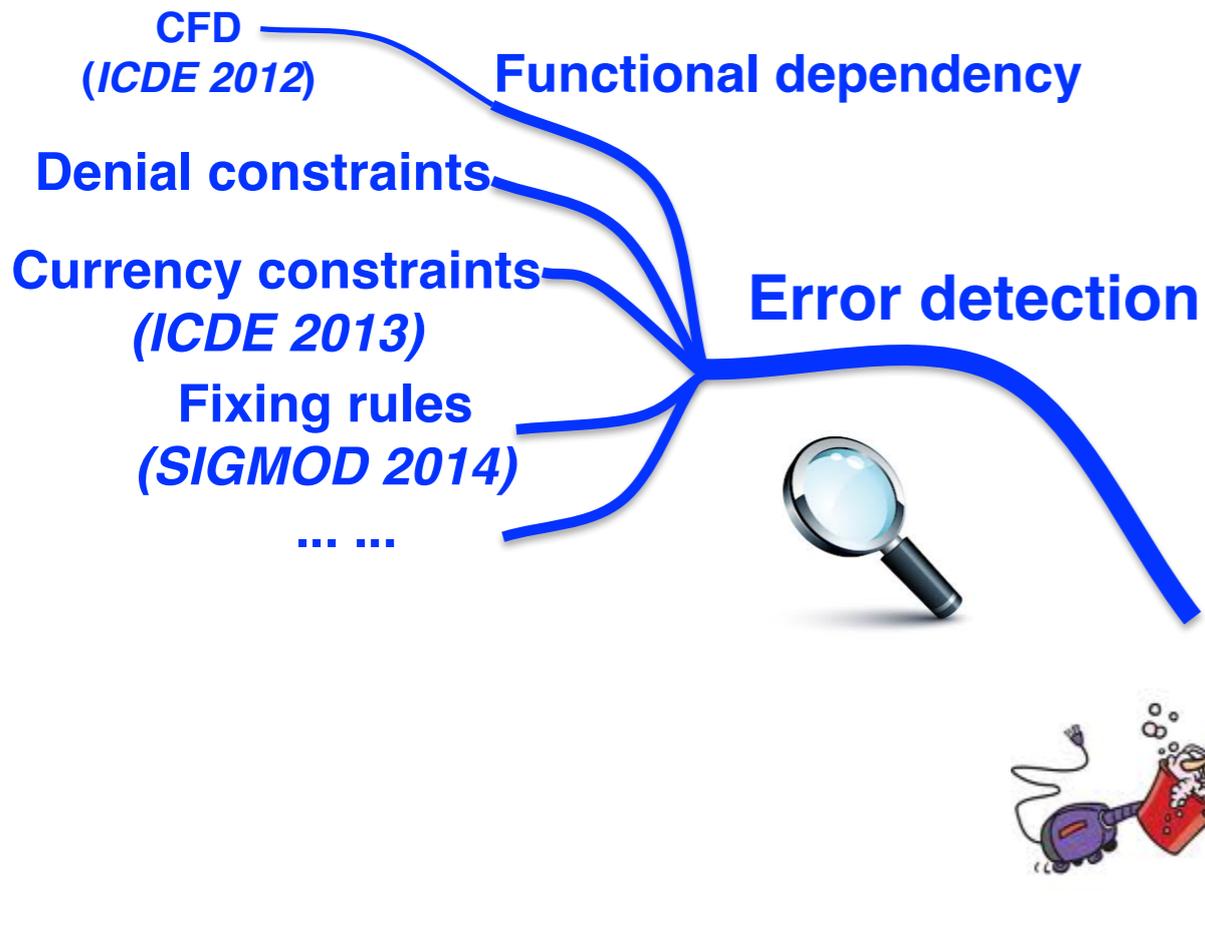
(UK, hasCapital, London)
KBs (e.g., Yago)

ETL (transformation)

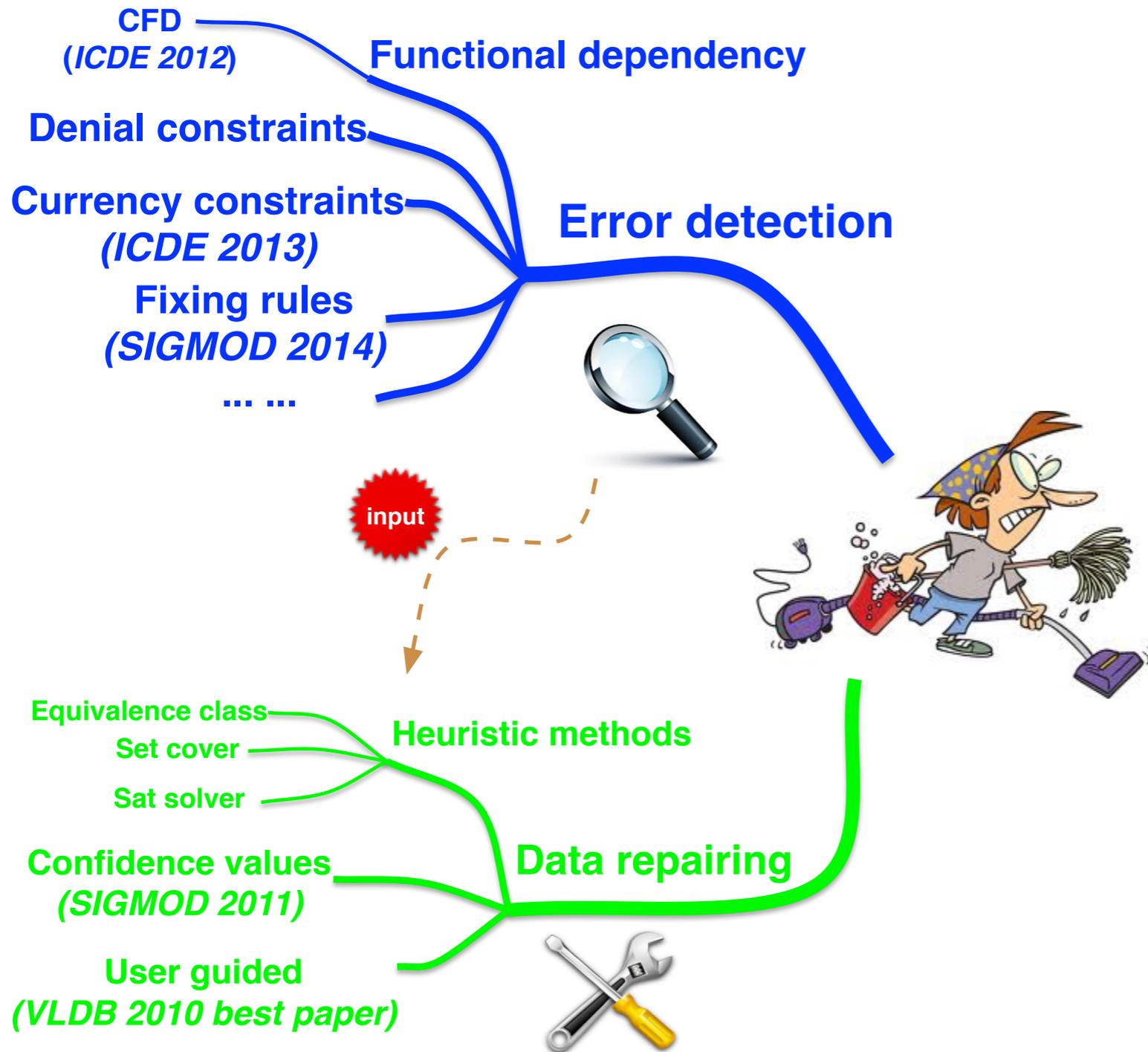
Data Cleaning Solutions



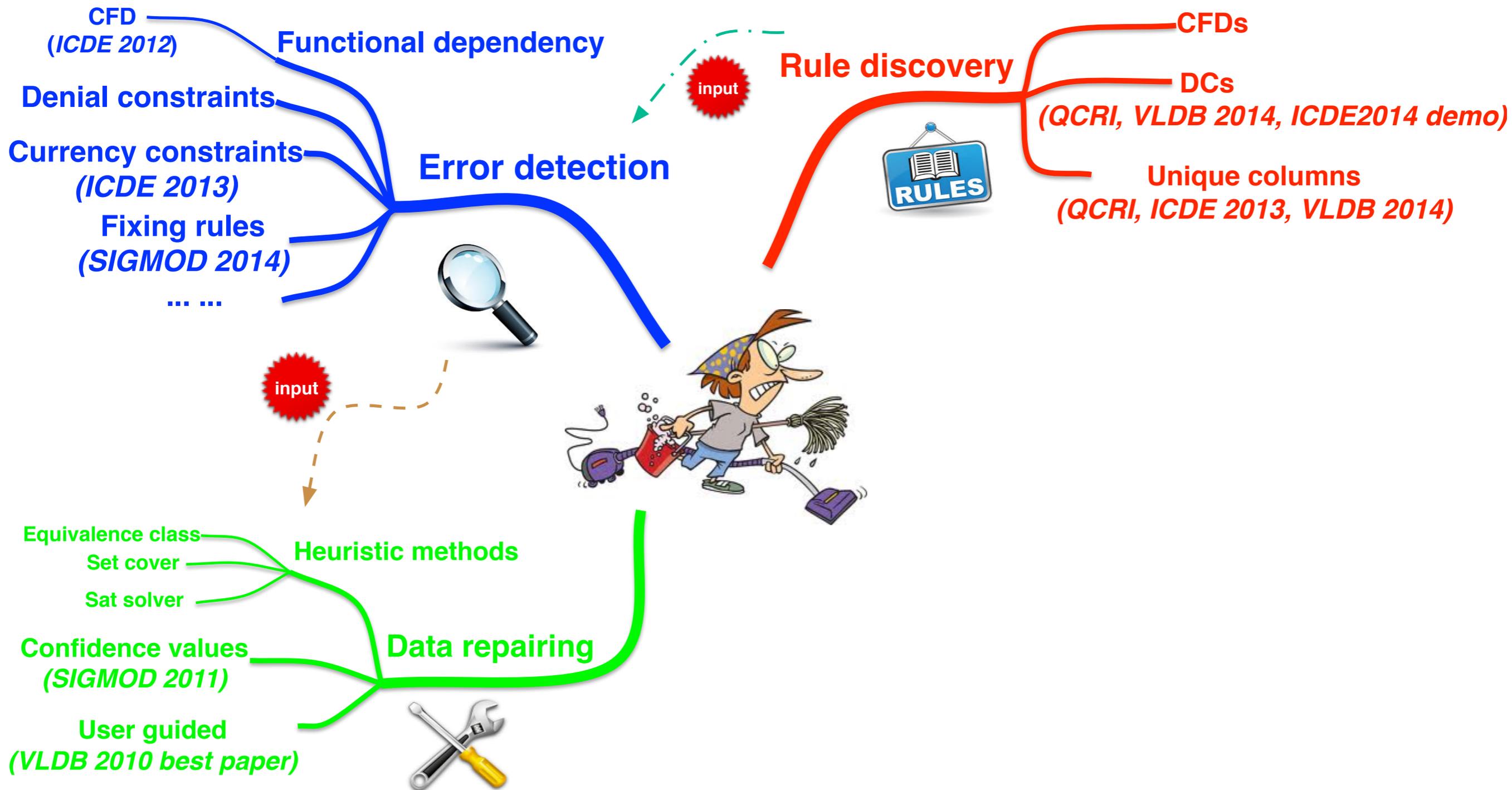
Data Cleaning Solutions



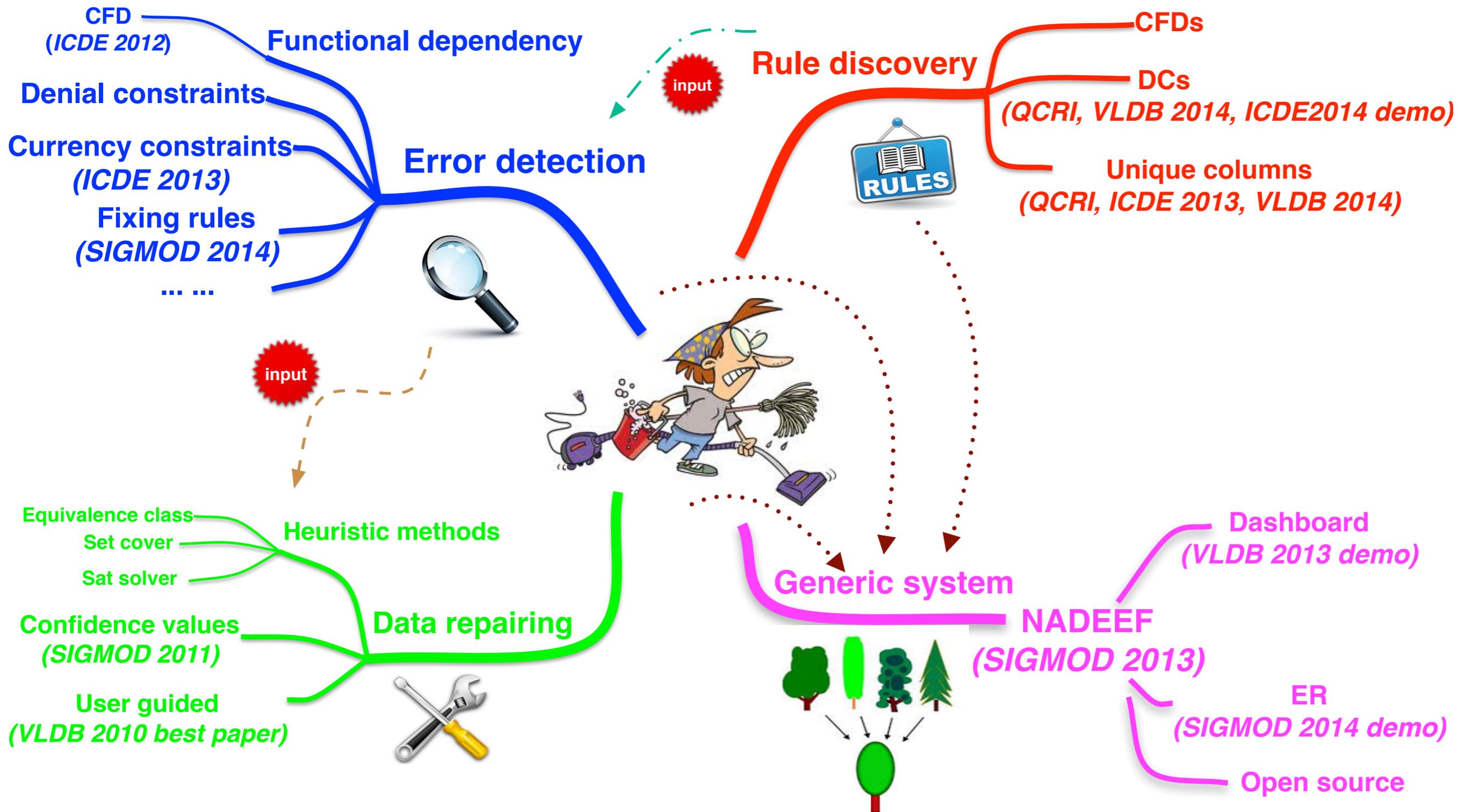
Data Cleaning Solutions



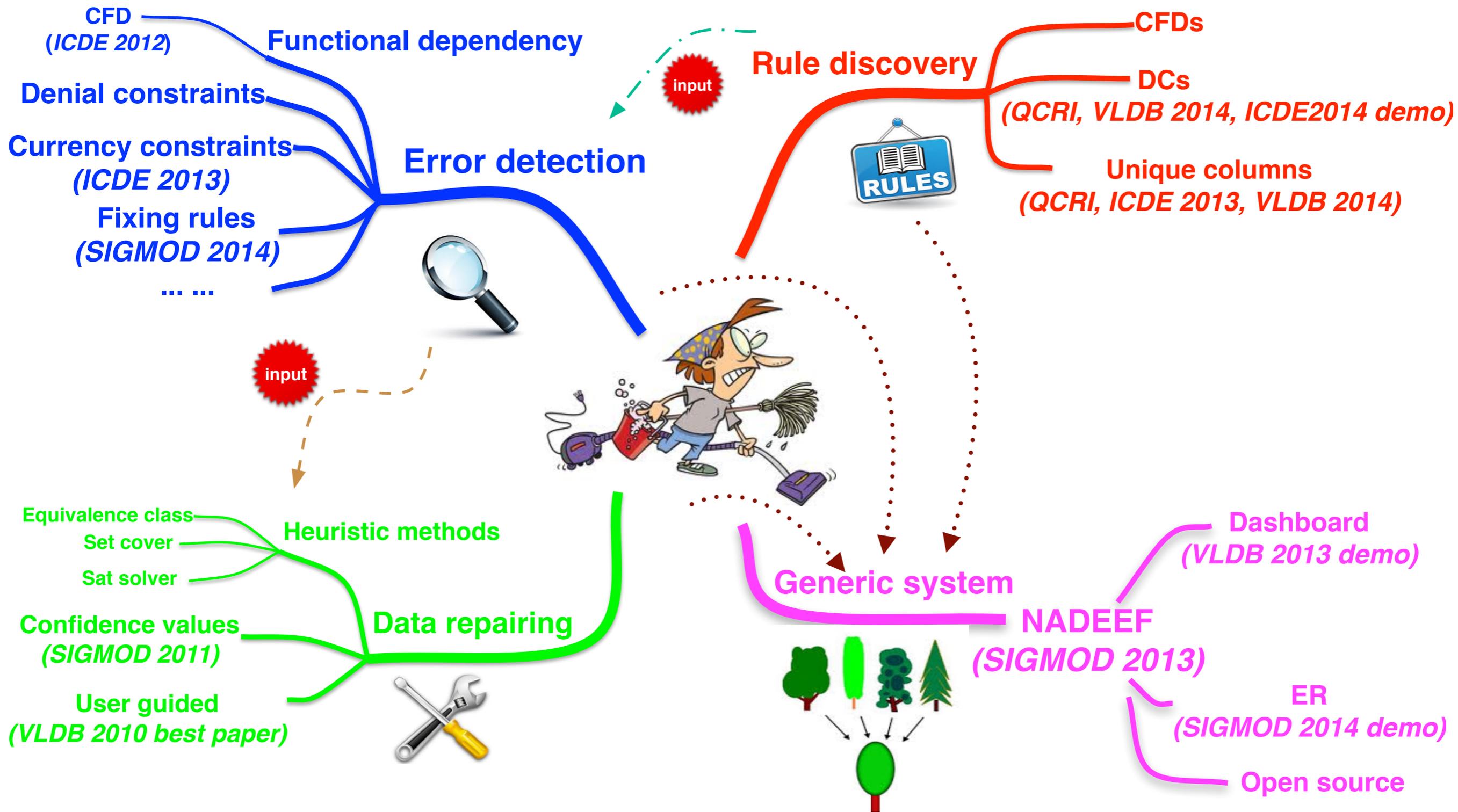
Data Cleaning Solutions



Data Cleaning Solutions



Data Cleaning Solutions



Error Detection

Error Detection

	name	country	capital	city	salary	tax
r1	Nan	China	Beijing	Beijing	50000	1000
r2	Yin	China	Shanghai	Hongkong	40000	1200
r3	Si	Netherlands	Den Hagg	Utrecht	60000	1400
r4	Lei	Netherlands	Amsterdam	Amsterdam	35000	800

emp

Error Detection

FD: [country] -> [capital]

	name	country	capital	city	salary	tax
r1	Nan	China	Beijing	Beijing	50000	1000
r2	Yin	China	Shanghai	Hongkong	40000	1200
r3	Si	Netherlands	Den Hagg	Utrecht	60000	1400
r4	Lei	Netherlands	Amsterdam	Amsterdam	35000	800

emp

Error Detection

FD: [country] -> [capital]

	name	country	capital	city	salary	tax
r1	Nan	China	Beijing	Beijing	50000	1000
r2	Yin	China	Shanghai	Hongkong	40000	1200
r3	Si	Netherlands	Den Hagg	Utrecht	60000	1400
r4	Lei	Netherlands	Amsterdam	Amsterdam	35000	800

emp

Error Detection

FD: [country] -> [capital]

CFD: [country = China] -> [capital = Beijing]

	name	country	capital	city	salary	tax
r1	Nan	China	Beijing	Beijing	50000	1000
r2	Yin	China	Shanghai	Hongkong	40000	1200
r3	Si	Netherlands	Den Hagg	Utrecht	60000	1400
r4	Lei	Netherlands	Amsterdam	Amsterdam	35000	800

emp

Error Detection

FD: [country] -> [capital]

CFD: [country = China] -> [capital = Beijing]

	name	country	capital	city	salary	tax
r1	Nan	China	Beijing	Beijing	50000	1000
r2	Yin	China	Shanghai	Hongkong	40000	1200
r3	Si	Netherlands	Den Hagg	Utrecht	60000	1400
r4	Lei	Netherlands	Amsterdam	Amsterdam	35000	800

emp

Error Detection

FD: [country] -> [capital]

CFD: [country = China] -> [capital = Beijing]

	name	country	capital	city	salary	tax
r1	Nan	China	Beijing	Beijing	50000	1000
r2	Yin	China	Shanghai	Hongkong	40000	1200
r3	Si	Netherlands	Den Hagg	Utrecht	60000	1400
r4	Lei	Netherlands	Amsterdam	Amsterdam	35000	800

emp

DC: $\neg t1, t2 (t1.salary > t2.salary \text{ and } t1.tax < t2.tax)$

Error Detection

FD: [country] -> [capital]

CFD: [country = China] -> [capital = Beijing]

	name	country	capital	city	salary	tax
r1	Nan	China	Beijing	Beijing	50000	1000
r2	Yin	China	Shanghai	Hongkong	40000	1200
r3	Si	Netherlands	Den Hagg	Utrecht	60000	1400
r4	Lei	Netherlands	Amsterdam	Amsterdam	35000	800

emp

DC: $\neg t1, t2 (t1.salary > t2.salary \text{ and } t1.tax < t2.tax)$

Error Detection

FD: [country] -> [capital]

CFD: [country = China] -> [capital = Beijing]

	name	country	capital	city	salary	tax
r1	Nan	China	Beijing	Beijing	50000	1000
r2	Yin	China	Shanghai	Hongkong	40000	1200
r3	Si	Netherlands	Den Hagg	Utrecht	60000	1400
r4	Lei	Netherlands	Amsterdam	Amsterdam	35000	800

emp

	country	capital
s1	China	Beijing
s2	Canada	Ottawa
s3

cap

DC: $\neg t1, t2 (t1.salary > t2.salary \text{ and } t1.tax < t2.tax)$

Error Detection

FD: [country] -> [capital]

CFD: [country = China] -> [capital = Beijing]

	name	country	capital	city	salary	tax
r1	Nan	China	Beijing	Beijing	50000	1000
r2	Yin	China	Shanghai	Hongkong	40000	1200
r3	Si	Netherlands	Den Hagg	Utrecht	60000	1400
r4	Lei	Netherlands	Amsterdam	Amsterdam	35000	800

emp

	country	capital
s1	China	Beijing
s2	Canada	Ottawa
s3

cap

DC: $\neg t1, t2 (t1.salary > t2.salary \text{ and } t1.tax < t2.tax)$

MD: $(emp[country] = cap[country]) \rightarrow (emp[capital] \Leftrightarrow cap[capital])$

Error Detection

FD: [country] -> [capital]

CFD: [country = China] -> [capital = Beijing]

	emp	cap	emp	cap
r1	Nan	China	Beijing	Beijing
r2
r3	Lei	Netherlands	Amsterdam	Amsterdam
r4

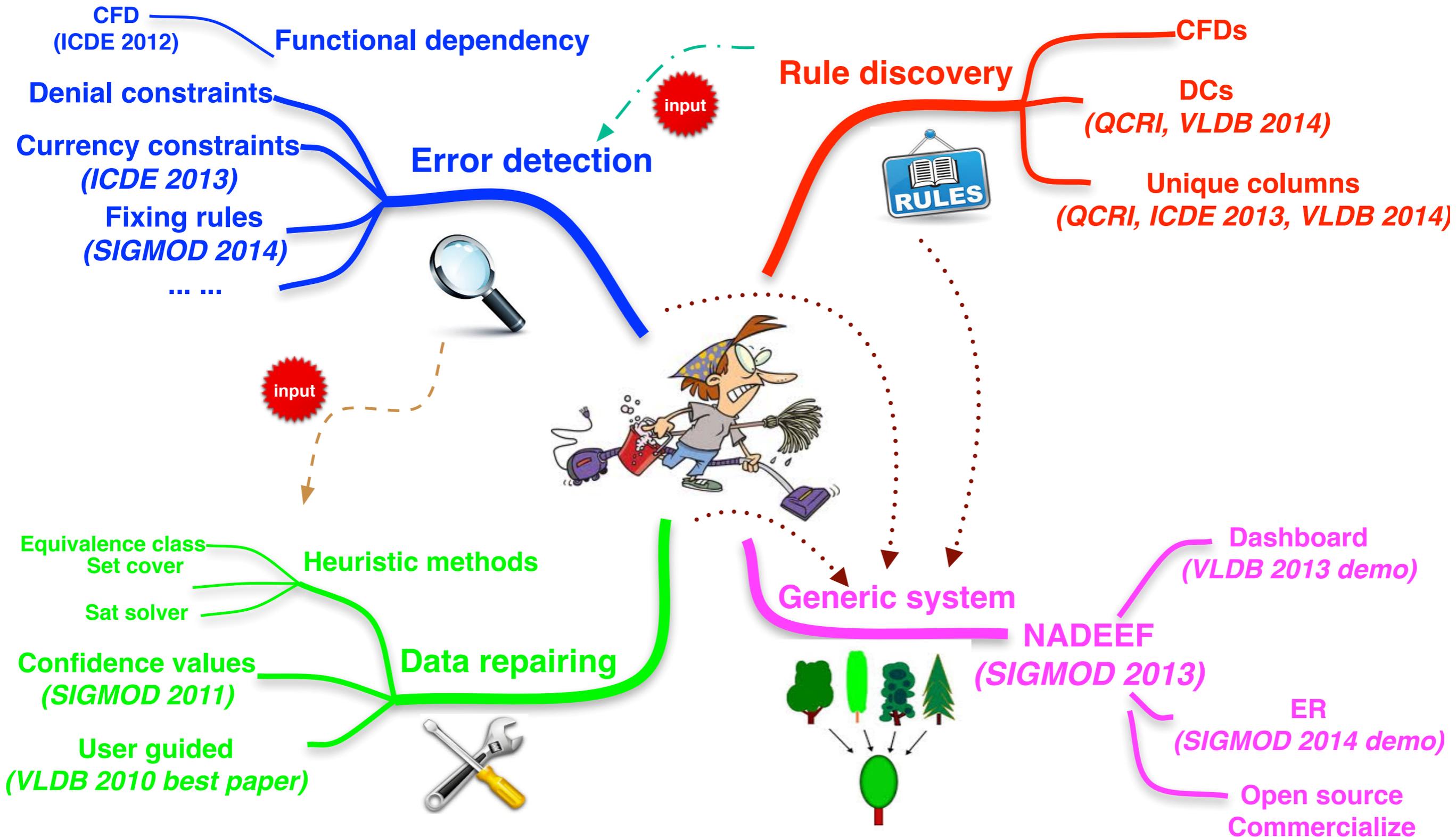
emp

cap

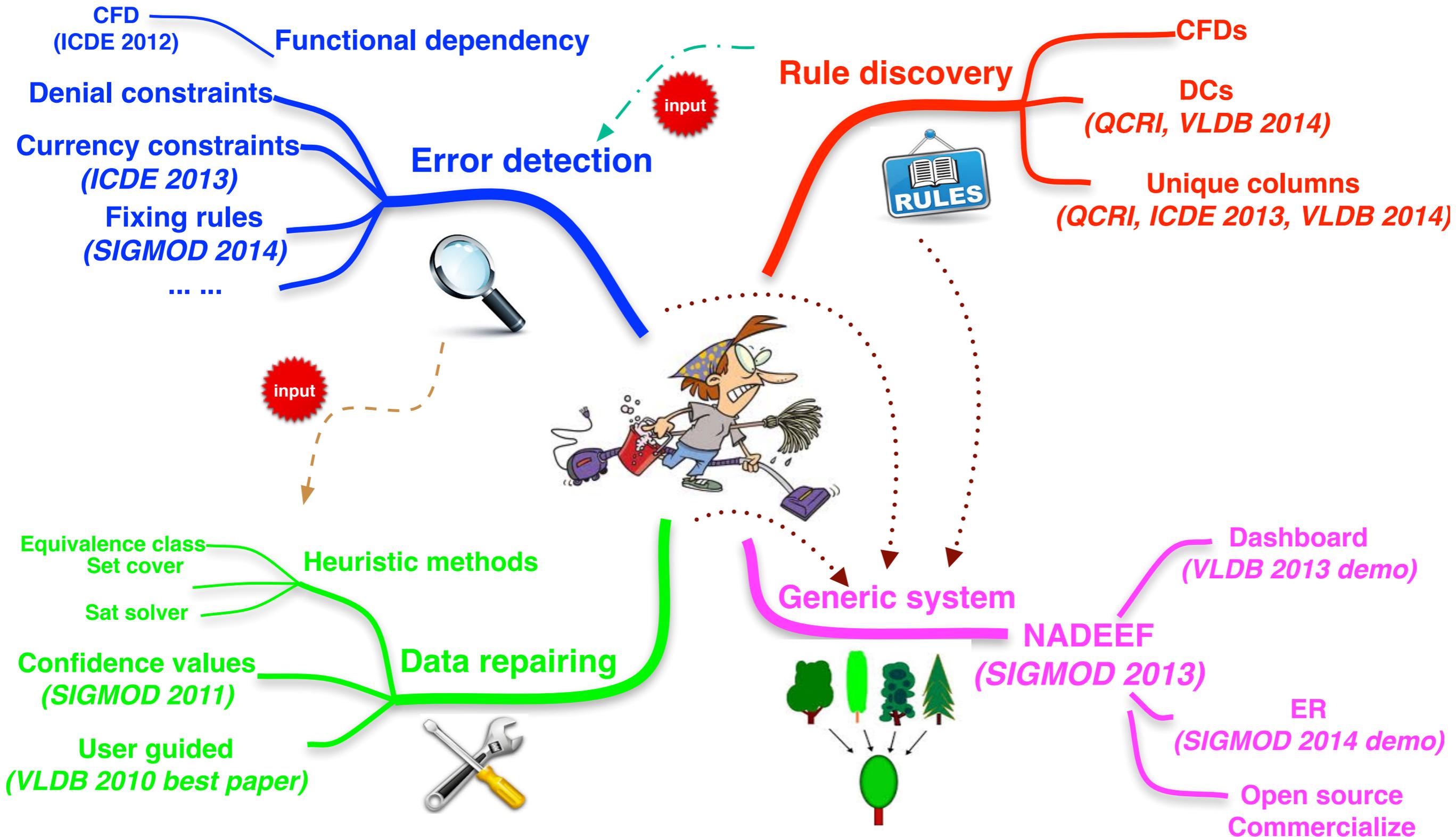
DC: $\neg t1, t2 (t1.salary > t2.salary \text{ and } t1.tax < t2.tax)$

MD: $(emp[country] = cap[country]) \rightarrow (emp[capital] \Leftrightarrow cap[capital])$

Data Repairing



Data Repairing



Automated



Computing a Consistent Database

D

Computing a Consistent Database

D

$D_g ?$

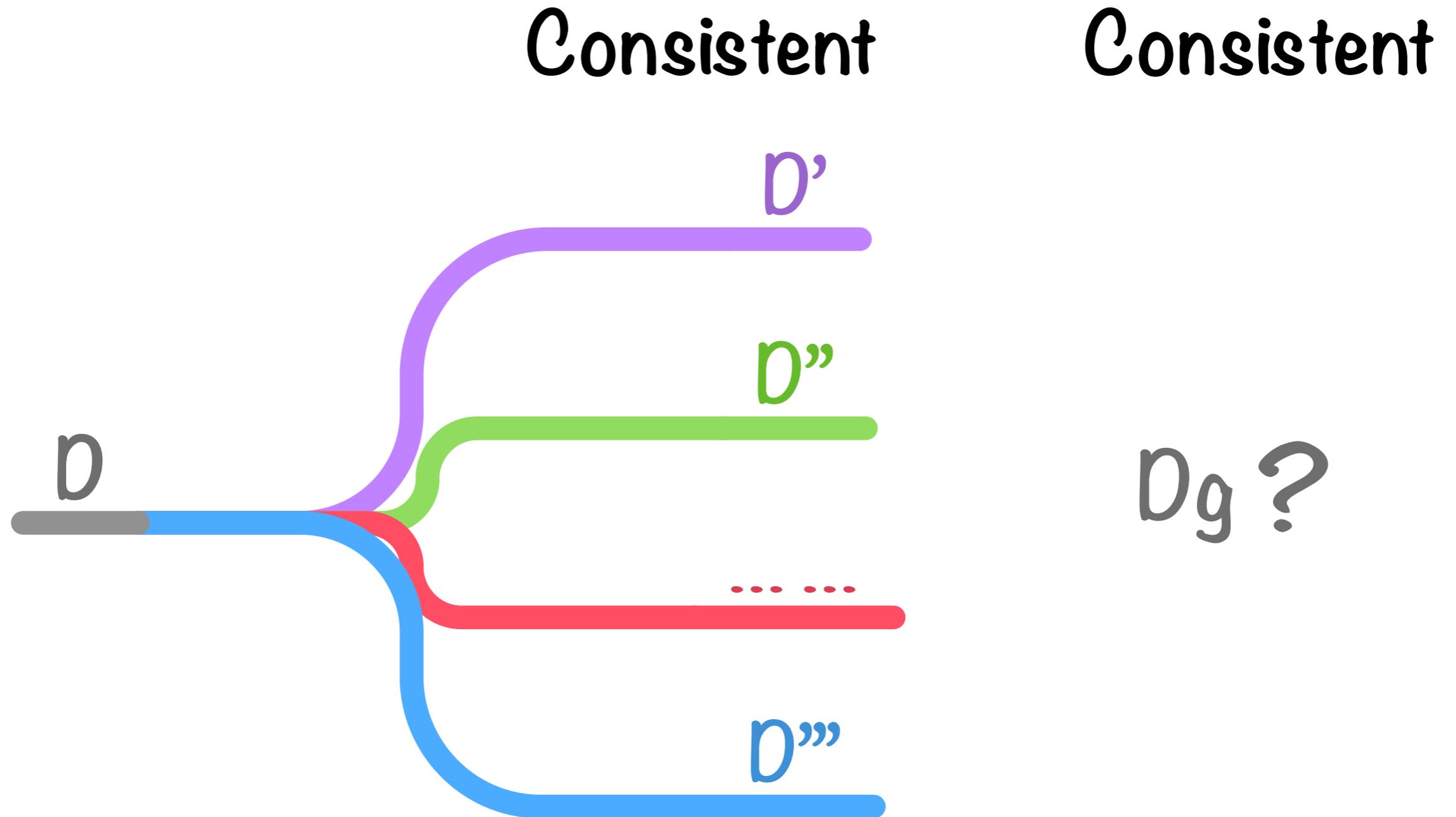
Computing a Consistent Database

Consistent

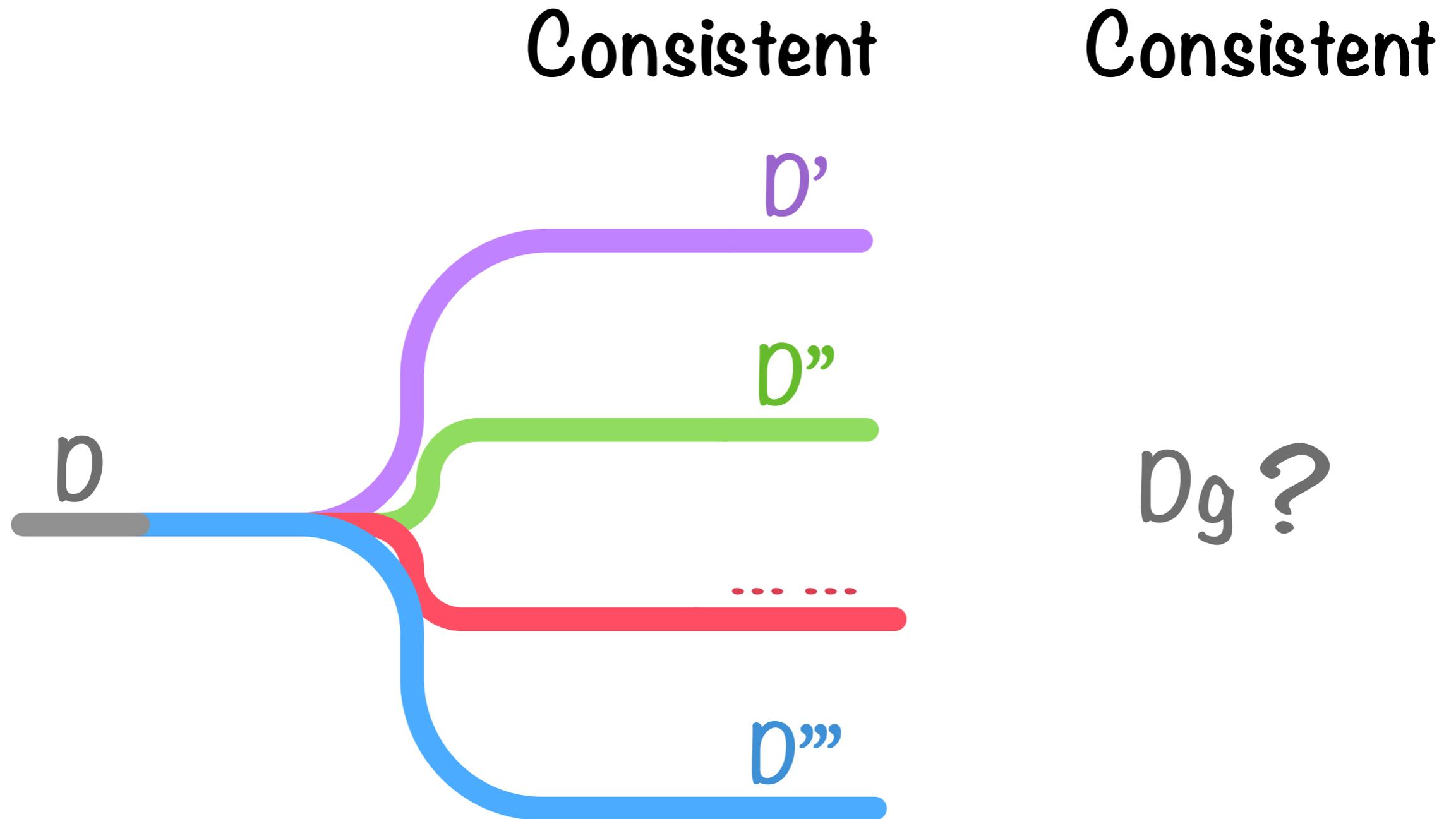
D

$D_g ?$

Computing a Consistent Database



Computing a Consistent Database



find a D' such that $\text{dist}(D, D')$ is minimum

Computing a Consistent Database

	name	nationality	capital	areacode	bornAt	salary	tax
r1	Nan	China	Beijing	10	Shenyang	50000	1000
r2	Yan	China	Shanghai	10	Hangzhou	40000	900
r3	Si	China	Beijing	10	Changsha	60000	1400
r4	Miura	China	Tokyo	3	Kyoto	35000	800

Computing a Consistent Database

FD1: [nationality] -> [capital]

FD2: [areacode] -> [capital]

	name	nationality	capital	areacode	bornAt	salary	tax
r1	Nan	China	Beijing	10	Shenyang	50000	1000
r2	Yan	China	Shanghai	10	Hangzhou	40000	900
r3	Si	China	Beijing	10	Changsha	60000	1400
r4	Miura	China	Tokyo	3	Kyoto	35000	800

Computing a Consistent Database

FD1: [nationality] -> [capital]

FD2: [areacode] -> [capital]

	name	nationality	capital	areacode	bornAt	salary	tax
r1	Nan	China	Beijing	10	Shenyang	50000	1000
r2	Yan	China	Shanghai	10	Hangzhou	40000	900
r3	Si	China	Beijing	10	Changsha	60000	1400
r4	Miura	China	Tokyo	3	Kyoto	35000	800

Computing a Consistent Database

FD1: [nationality] -> [capital]

FD2: [areacode] -> [capital]

	name	nationality	capital	areacode	bornAt	salary	tax
r1	Nan	China	Beijing	10	Shenyang	50000	1000
r2	Yan	China	Shanghai	10	Hangzhou	40000	900
			Beijing				
r3	Si	China	Beijing	10	Changsha	60000	1400
r4	Miura	China	Tokyo	3	Kyoto	35000	800
			Beijing				

Computing a Consistent Database

FD1: [nationality] -> [capital]

FD2: [areacode] -> [capital]

	name	nationality	capital	areacode	bornAt	salary	tax
r1	Nan	China	Beijing	10	Shenyang	50000	1000
r2	Yan	China	Shanghai	10	Hangzhou	40000	900
			Beijing				
r3	Si	China	Beijing	10	Changsha	60000	1400
r4	Miura	China	Tokyo	3	Kyoto	35000	800
			Beijing				

Equivalence
class

Vertex
cover

SAT
solver

...

Computing a Consistent Database

FD1: [nationality] -> [capital]

FD2: [areacode] -> [capital]

	name	nationality	capital	areacode	bornAt	salary	tax
r1	Nan	China	Beijing	10	Shenyang	50000	1000
r2	Yan	China	Shanghai	10	Hangzhou	40000	900
			Beijing				
r3	Si	China	Beijing	10	Changsha	60000	1400
r4	Miura	China	Tokyo	3	Kyoto	35000	800
			Beijing				



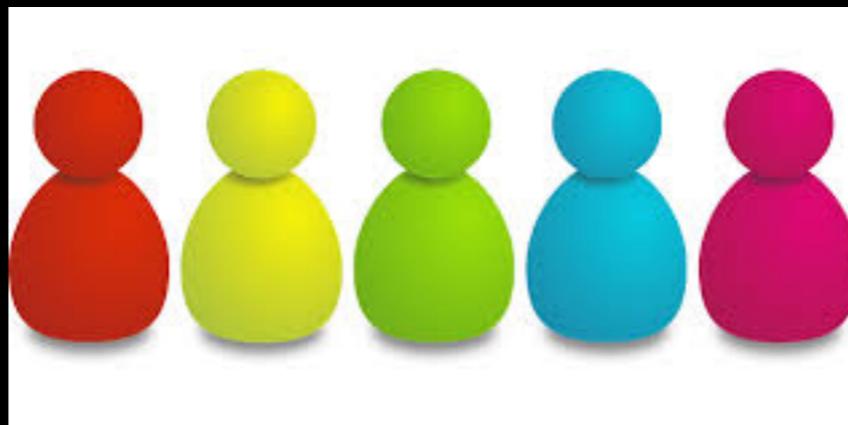
Equivalence class

Vertex cover

SAT solver

...

User Guided



Certain Fixes (VLDB 2010 Best Paper)

	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

Certain Fixes (VLDB 2010 Best Paper)

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

Certain Fixes (VLDB 2010 Best Paper)

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

Certain Fixes (VLDB 2010 Best Paper)

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

Is r2[country] China?
YES.



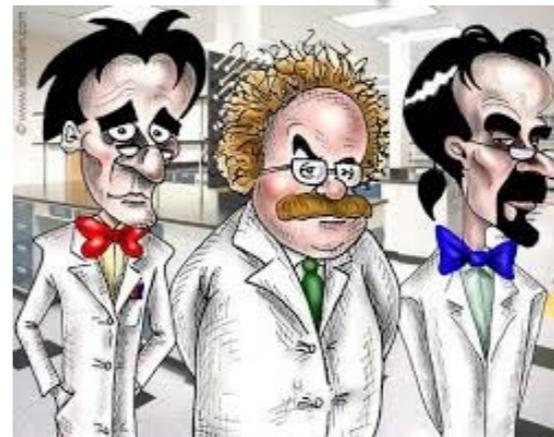
Certain Fixes (VLDB 2010 Best Paper)

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

	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

country	capital
China	Beijing
Canada	Ottawa
Japan	Tokyo

Is r2[country] China?
YES.



Certain Fixes (VLDB 2010 Best Paper)

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

	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

country	capital
China	Beijing
Canada	Ottawa
Japan	Tokyo

Is r2[country] China?
YES.

Is r1[country] China?

Is r3[country] China?

Is r4[country] Canada?

.....



Certain Fixes (VLDB 2010 Best Paper)

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

	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

country	capital
China	Beijing
Canada	Ottawa
Japan	Tokyo

Is r2[country] China?
YES.

Is r1[country] China?

Is r3[country] China?

Is r4[country] Canada?

.....



check *each* tuple: not cheap !!

precision: +
recall: ++

Heuristic
(Automated)

precision: ++
recall: ++

Certain
(User guided)

precision: +
recall: ++

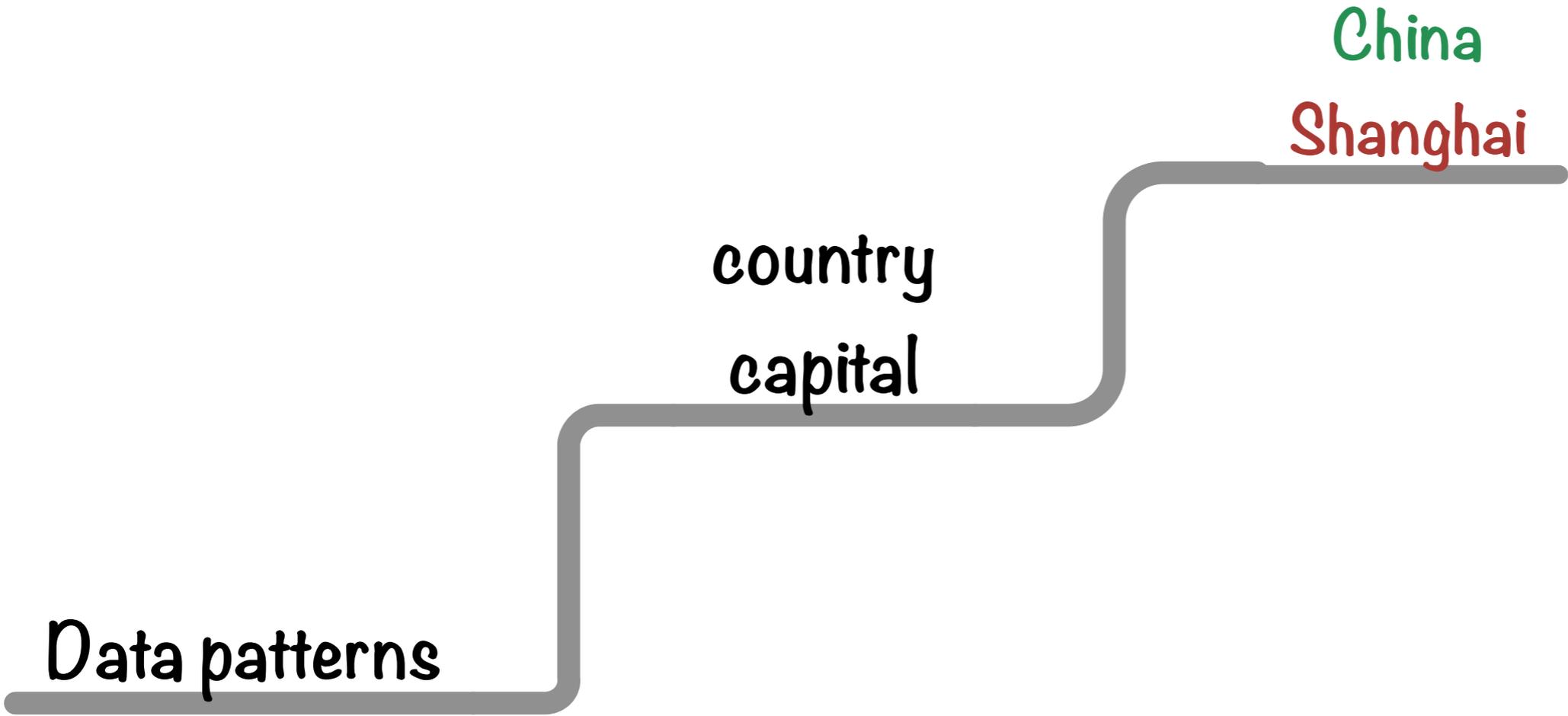
Heuristic
(Automated)

precision: ++
recall: +

Fixing Rules
(Automated)

precision: ++
recall: ++

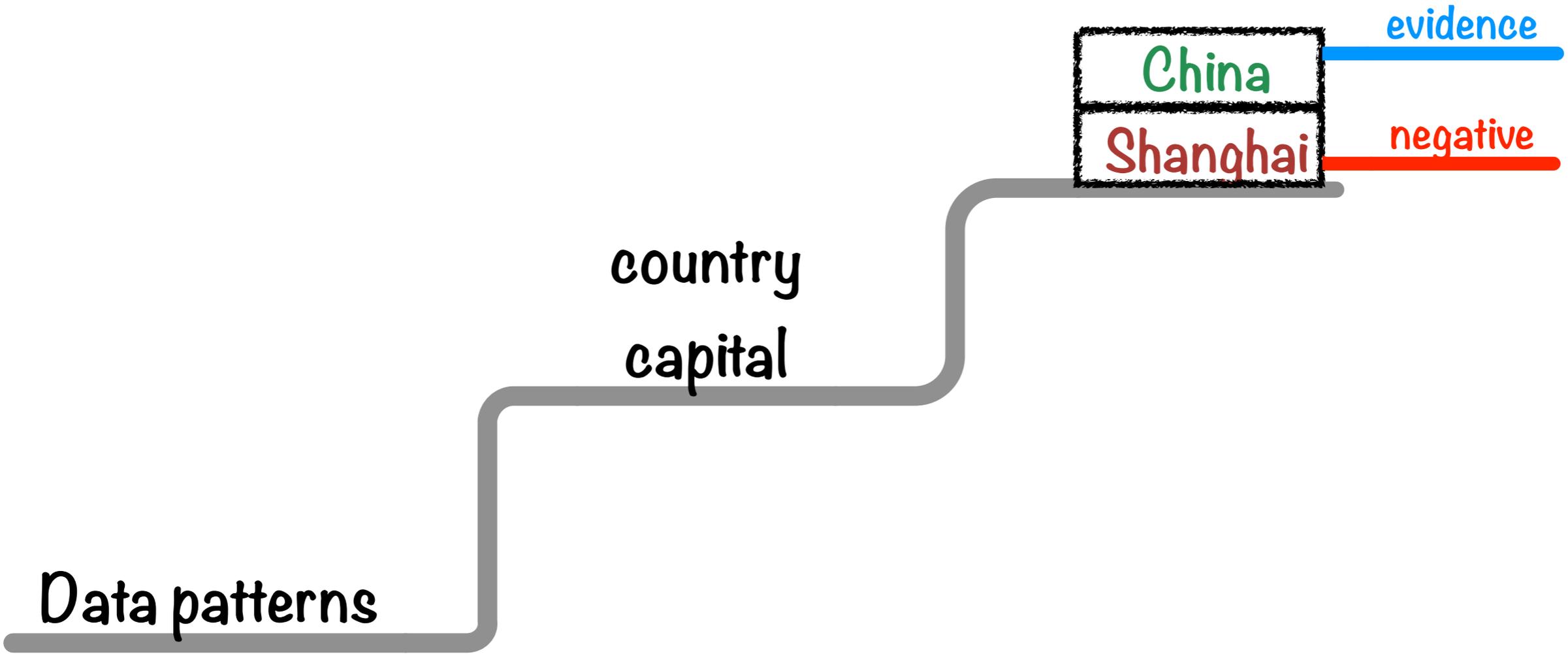
Certain
(User guided)

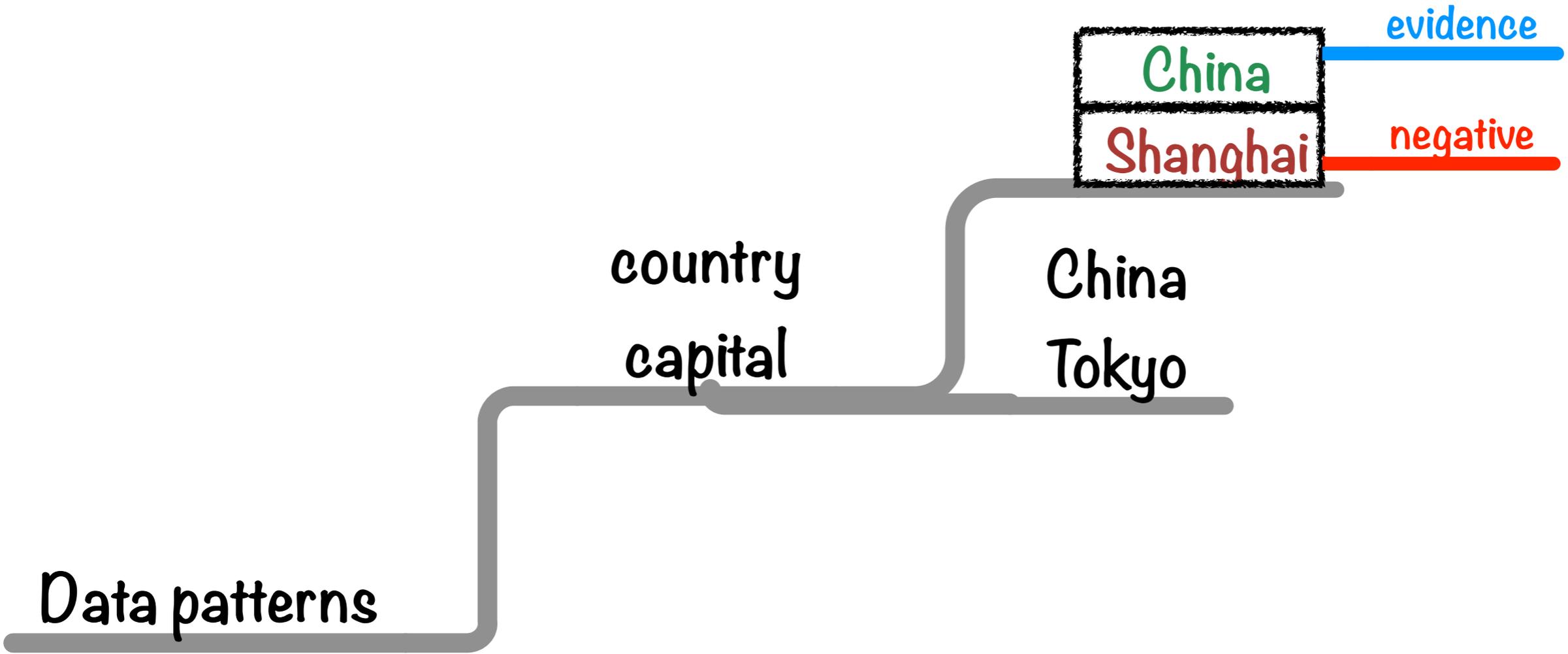


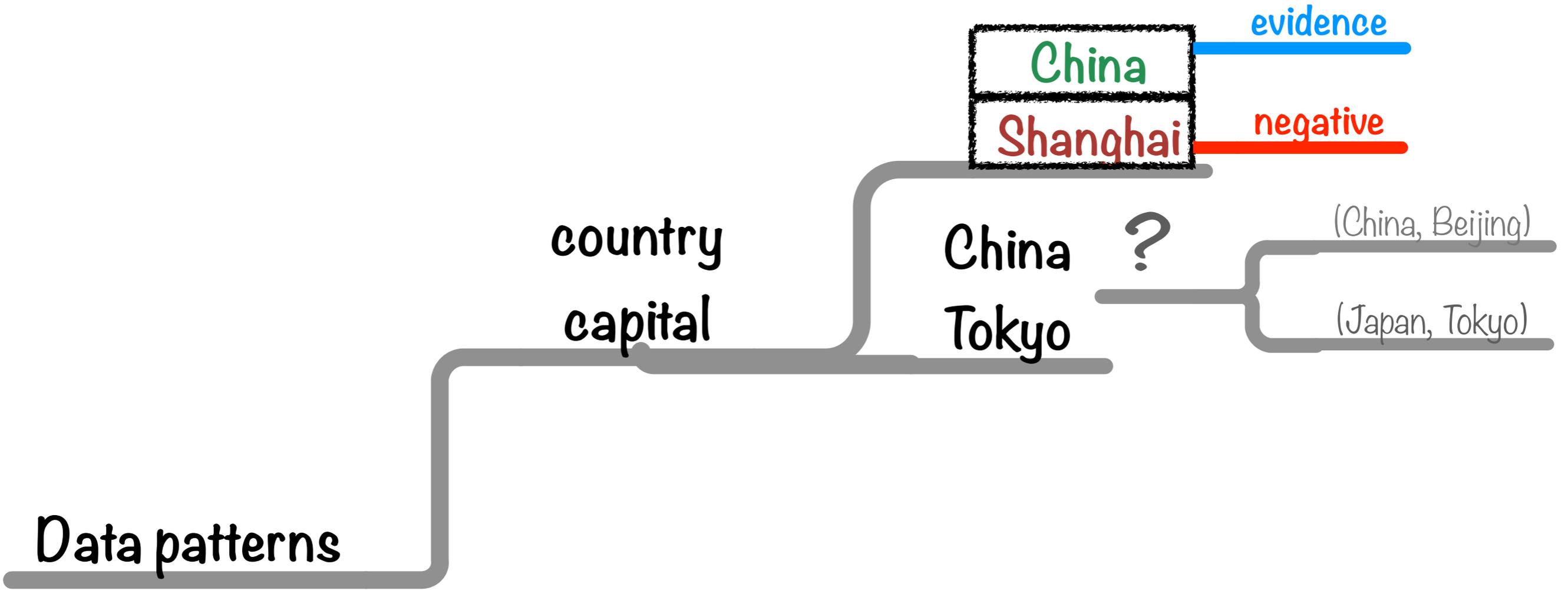
Data patterns

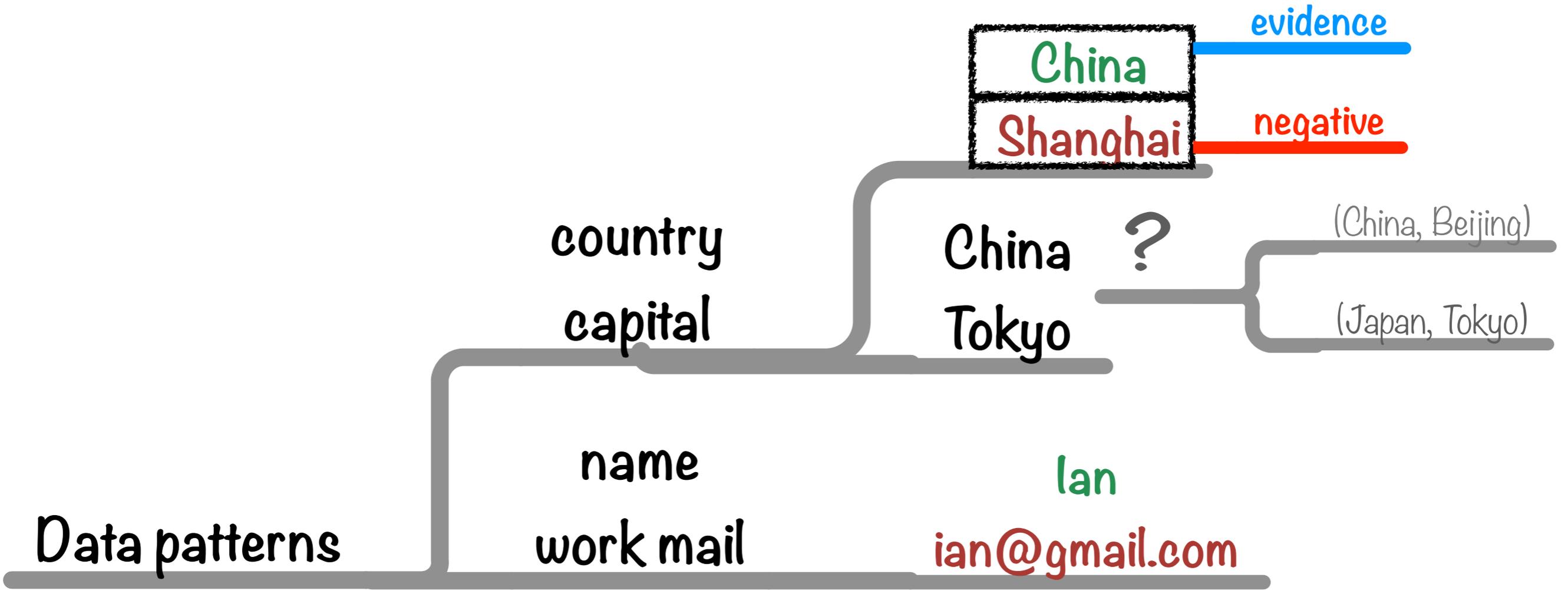
country
capital

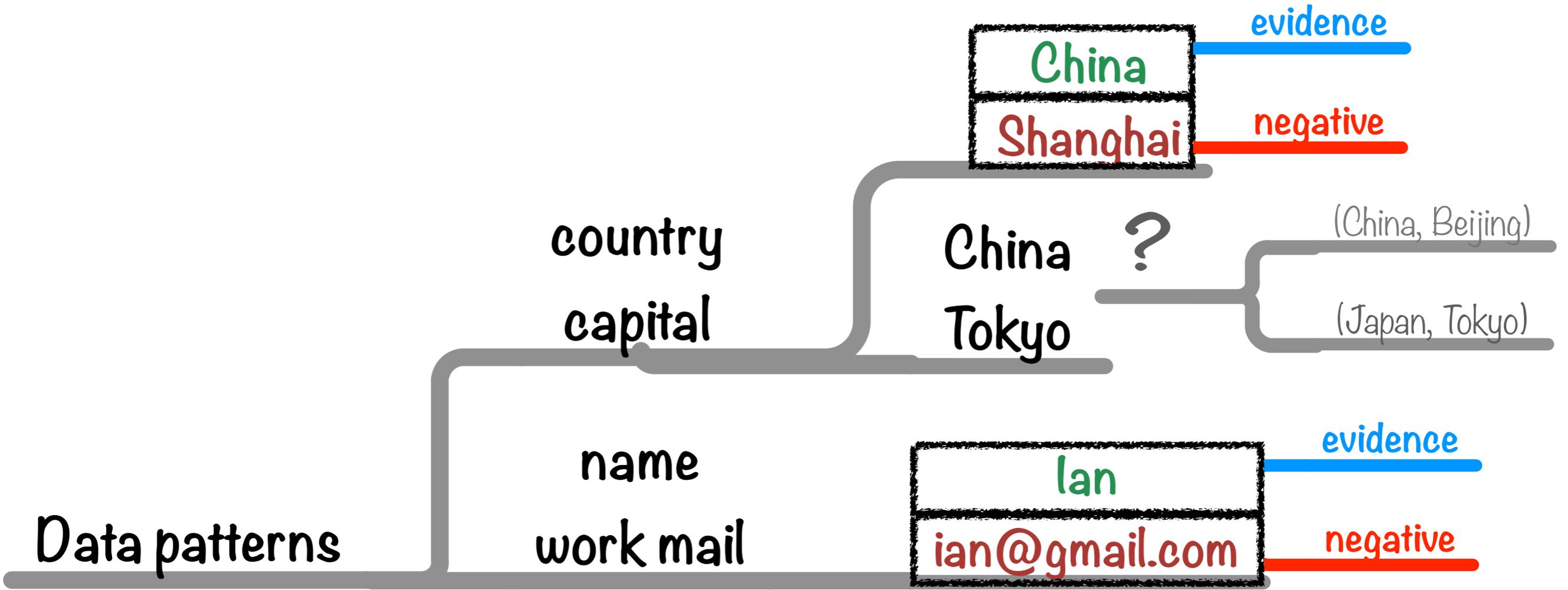
China
Shanghai

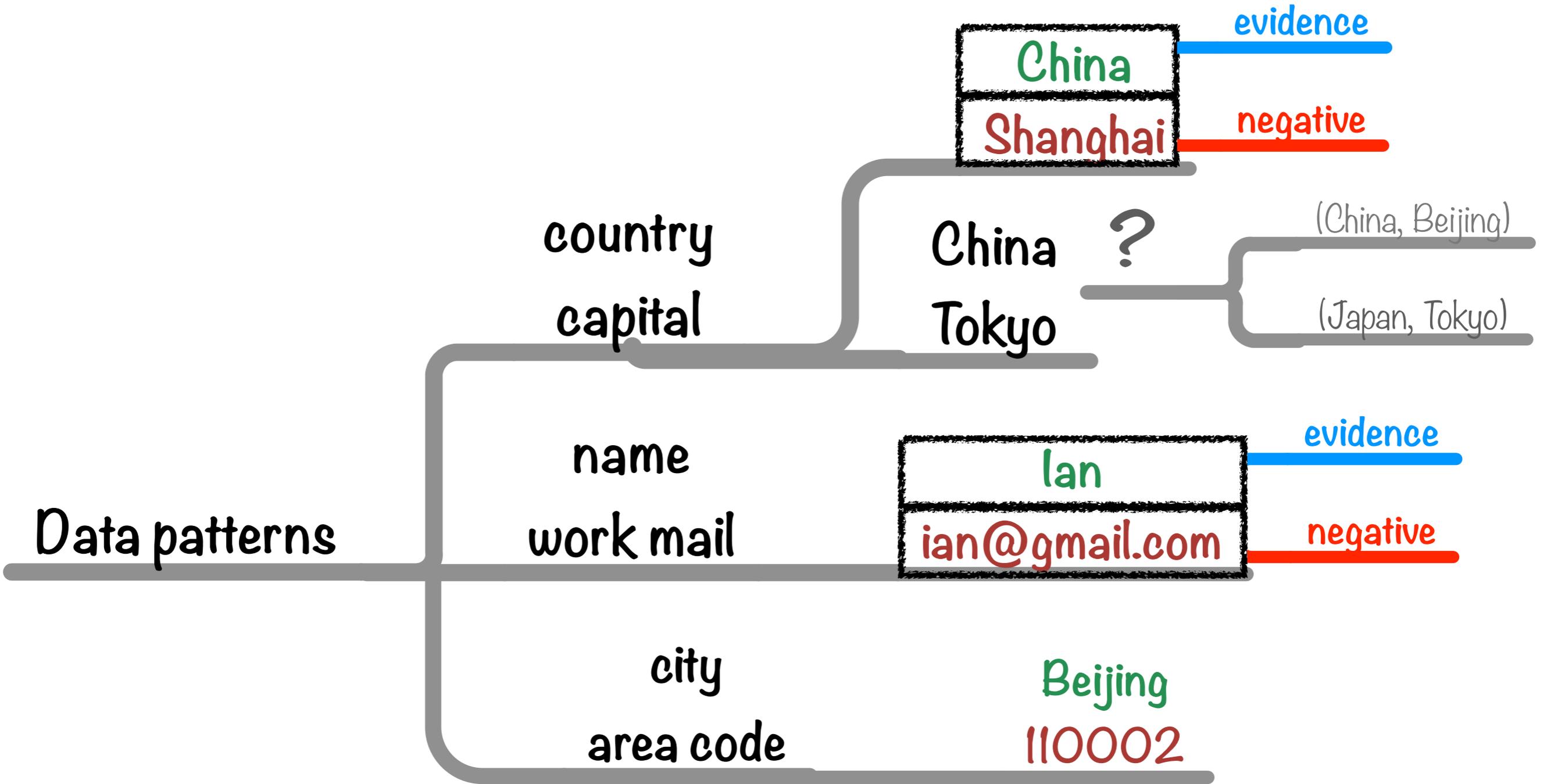


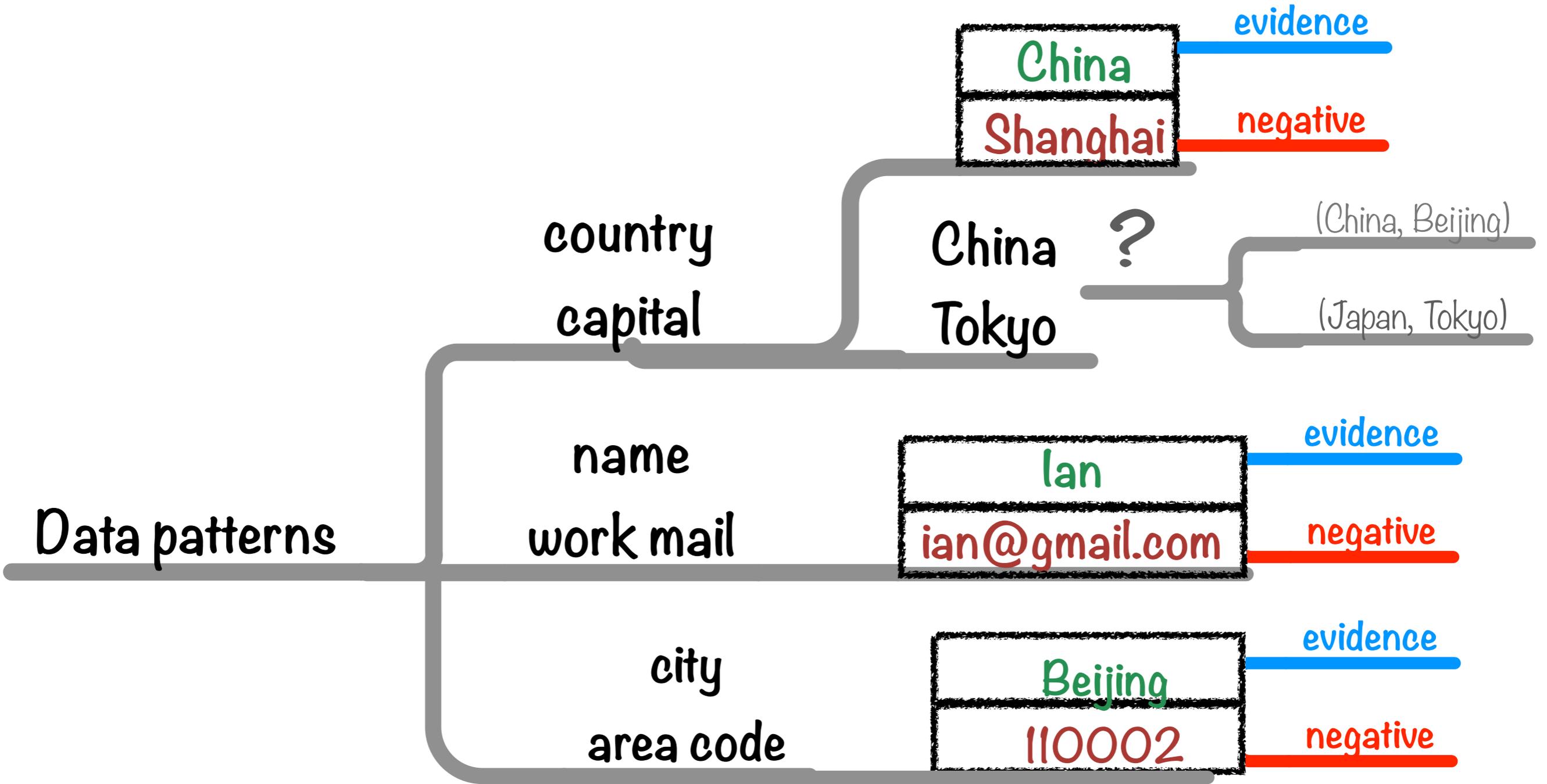












Data patterns

country
capital

China
Shanghai

evidence

negative

China ?
Tokyo

(China, Beijing)

(Japan, Tokyo)

name
work mail

Ian
ian@gmail.com

evidence

negative

city
area code

Beijing
110002

evidence

negative

Fixing Rules (SIGMOD 2014)

- **Syntax**

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

Fixing Rules (SIGMOD 2014)

- **Syntax**

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

country	{capital	capital
China	Shanghai	Beijing
	Hongkong	

Fixing Rules (SIGMOD 2014)

- **Syntax**

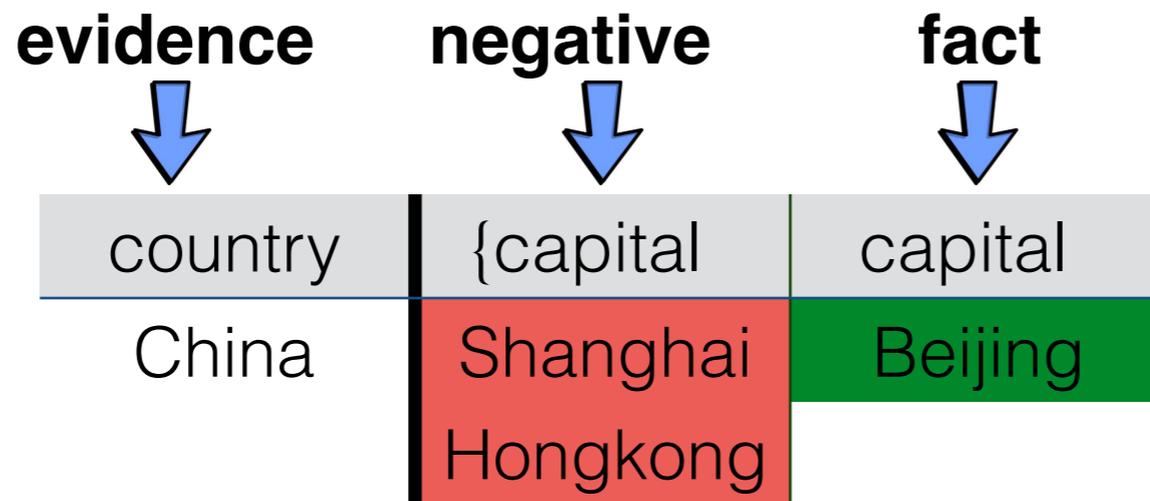
fR1: (([country], [China]), (capital, {Shanghai, Hongkong})) -> Beijing

evidence	negative	
country	{capital	capital
China	Shanghai	Beijing
	Hongkong	

Fixing Rules (SIGMOD 2014)

- **Syntax**

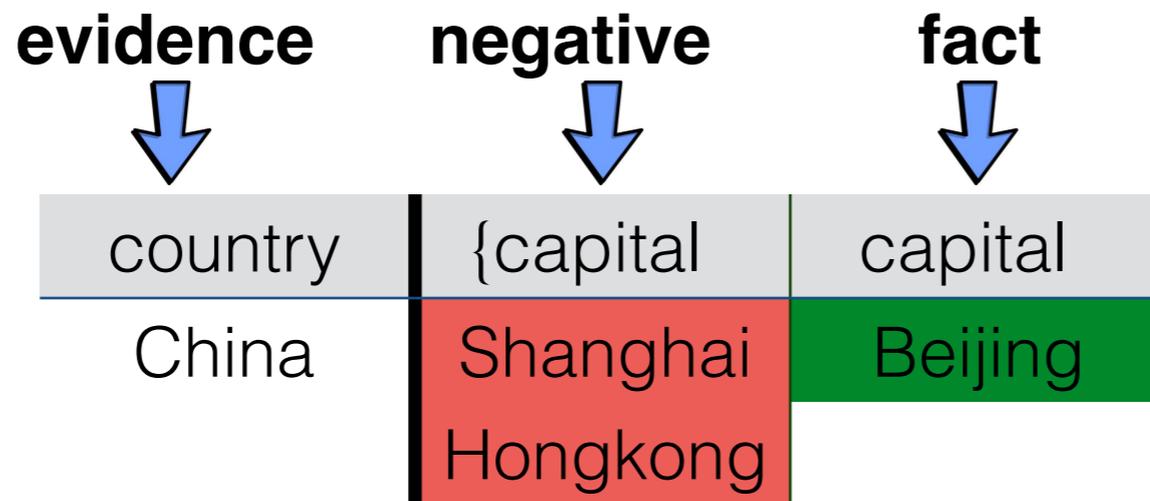
fR1: (([country], [China]), (capital, {Shanghai, Hongkong})) -> Beijing



Fixing Rules (SIGMOD 2014)

- Syntax**

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

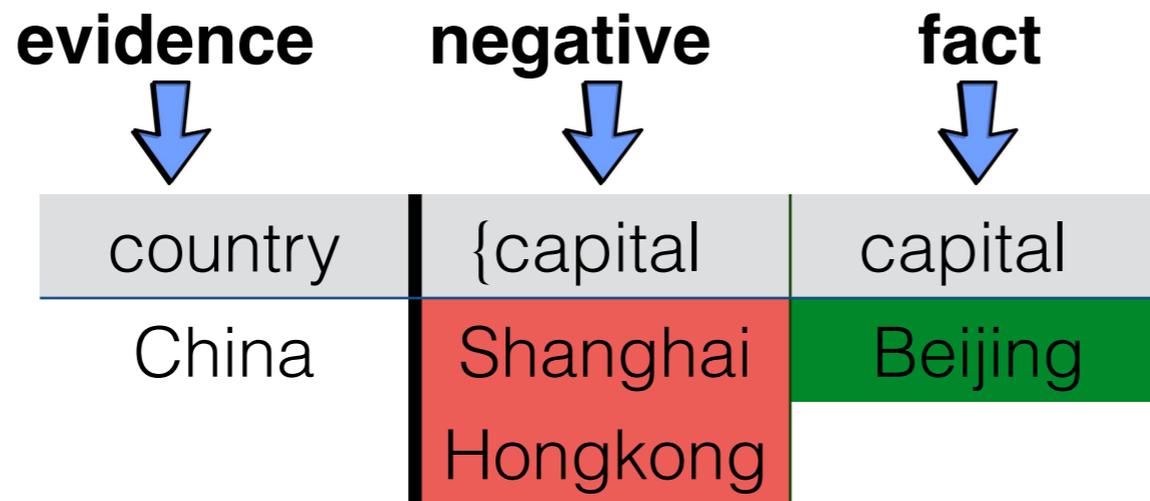


	name	nationality	capital	bornAt
r1	Nan	China	Beijing	Shenyang
r2	Yan	China	Shanghai	Hangzhou
r3	Si	China	Beijing	Changsha
r4	Miura	China	Tokyo	Kyoto

Fixing Rules (SIGMOD 2014)

- Syntax**

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

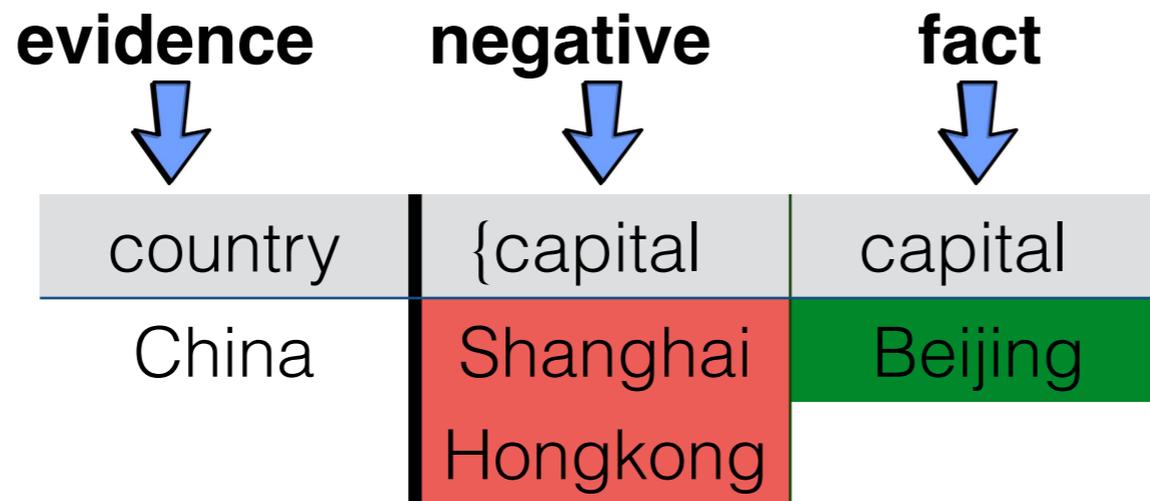


	name	nationality	capital	bornAt
r1	Nan	China	Beijing	Shenyang
r2	Yan	China	Shanghai	Hangzhou
r3	Si	China	Beijing	Changsha
r4	Miura	China	Tokyo	Kyoto

Fixing Rules (SIGMOD 2014)

- Syntax**

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

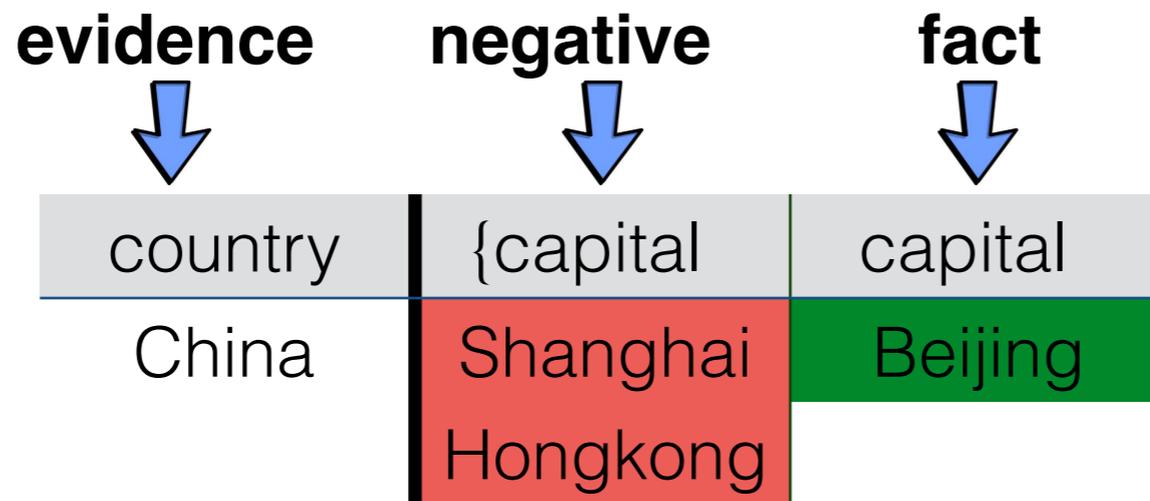


	name	nationality	capital	bornAt
r1	Nan	China	Beijing	Shenyang
r2	Yan	China	Beijing	Hangzhou
r3	Si	China	Beijing	Changsha
r4	Miura	China	Tokyo	Kyoto

Fixing Rules (SIGMOD 2014)

- Syntax**

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

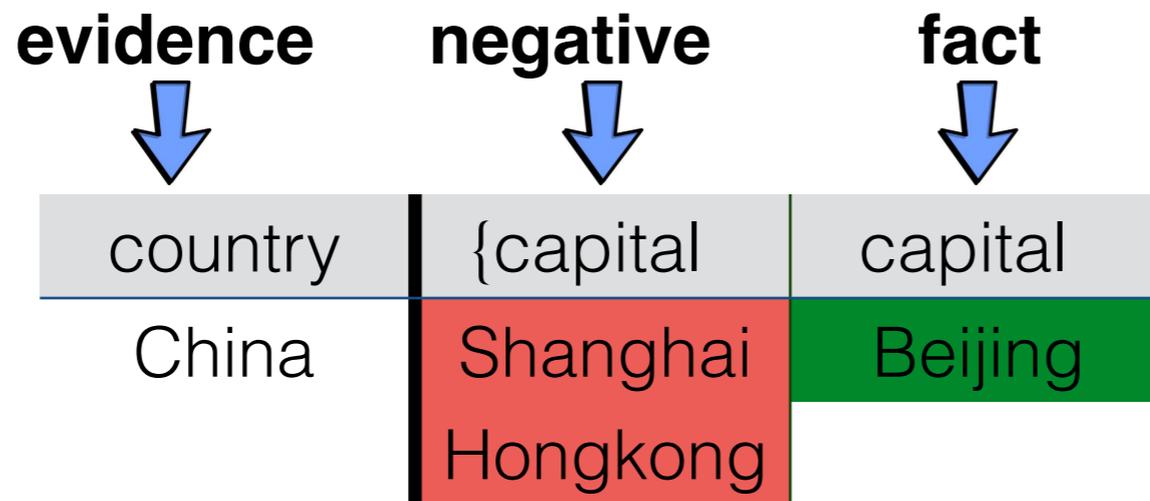


	name	nationality	capital	bornAt
r1	Nan	China	Beijing	Shenyang
r2	Yan	China	Beijing	Hangzhou
r3	Si	China	Beijing	Changsha
r4	Miura	China	Tokyo	Kyoto

Fixing Rules (SIGMOD 2014)

- Syntax**

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



	name	nationality	capital	bornAt
r1	Nan	China	Beijing	Shenyang
r2	Yan	China	Beijing	Hangzhou
r3	Si	China	Beijing	Changsha
r4	Miura	China	Tokyo	Kyoto



Confidence values Interaction

■ ■ ■ ■ ■ ■

Matching and Repairing (SIGMOD 2011)

FD: [nationality] -> [capital]

MD: ((nationality, country) -> (capital, capital))

	name	nationality	capital	bornAt
r1	Nan (0.9)	China (1.0)	Beijing (1.0)	Shenyang (0.9)
r2	Yan (0.8)	China (1.0)	Beijing (0.5)	Hangzhou (0.9)
r3	Si (0.9)	Canada (1.0)	Toronto (0.5)	Changsha (0.8)
r4	Miura (0.9)	Canada (0.9)	Vancouver (0.5)	Kyoto (1.0)

	country	capital
s1	China (1.0)	Beijing (1.0)
s2	Canada (1.0)	Ottawa (1.0)
s3	Japan (1.0)	Tokyo (1.0)

Matching and Repairing (SIGMOD 2011)



FD: [nationality] -> [capital]

MD: ((nationality, country) -> (capital, capital))

	name	nationality	capital	bornAt
r1	Nan (0.9)	China (1.0)	Beijing (1.0)	Shenyang (0.9)
r2	Yan (0.8)	China (1.0)	Beijing (0.5)	Hangzhou (0.9)
r3	Si (0.9)	Canada (1.0)	Toronto (0.5)	Changsha (0.8)
r4	Miura (0.9)	Canada (0.9)	Vancouver (0.5)	Kyoto (1.0)

	country	capital
s1	China (1.0)	Beijing (1.0)
s2	Canada (1.0)	Ottawa (1.0)
s3	Japan (1.0)	Tokyo (1.0)

Matching and Repairing (SIGMOD 2011)



FD: [nationality] -> [capital]

MD: ((nationality, country) -> (capital, capital))

	name	nationality	capital	bornAt
r1	Nan (0.9)	China (1.0)	Beijing (1.0)	Shenyang (0.9)
r2	Yan (0.8)	China (1.0)	Beijing (0.5)	Hangzhou (0.9)
r3	Si (0.9)	Canada (1.0)	Toronto (0.5)	Changsha (0.8)
r4	Miura (0.9)	Canada (0.9)	Vancouver (0.5)	Kyoto (1.0)

	country	capital
s1	China (1.0)	Beijing (1.0)
s2	Canada (1.0)	Ottawa (1.0)
s3	Japan (1.0)	Tokyo (1.0)

Matching and Repairing (SIGMOD 2011)



FD: [nationality] -> [capital]

MD: ((nationality, country) -> (capital, capital))

	name	nationality	capital	bornAt
r1	Nan (0.9)	China (1.0)	Beijing (1.0)	Shenyang (0.9)
r2	Yan (0.8)	China (1.0)	Beijing (0.5)	Hangzhou (0.9)
r3	Si (0.9)	Canada (1.0)	Ottawa (1.0)	Changsha (0.8)
r4	Miura (0.9)	Canada (0.9)	Vancouver (0.5)	Kyoto (1.0)

	country	capital
s1	China (1.0)	Beijing (1.0)
s2	Canada (1.0)	Ottawa (1.0)
s3	Japan (1.0)	Tokyo (1.0)

Matching and Repairing (SIGMOD 2011)



FD: [nationality] -> [capital]

MD: ((nationality, country) -> (capital, capital))

	name	nationality	capital	bornAt
r1	Nan (0.9)	China (1.0)	Beijing (1.0)	Shenyang (0.9)
r2	Yan (0.8)	China (1.0)	Beijing (0.5)	Hangzhou (0.9)
r3	Si (0.9)	Canada (1.0)	Ottawa (1.0)	Changsha (0.8)
r4	Miura (0.9)	Canada (0.9)	Vancouver (0.5)	Kyoto (1.0)

	country	capital
s1	China (1.0)	Beijing (1.0)
s2	Canada (1.0)	Ottawa (1.0)
s3	Japan (1.0)	Tokyo (1.0)

Summary of Data Repairing

Consistent database (heuristic)

Equivalence class
Vertex cover
Sat solver

Summary of Data Repairing

Consistent database
(heuristic)

Equivalence class
Vertex cover
Sat solver

improve accuracy

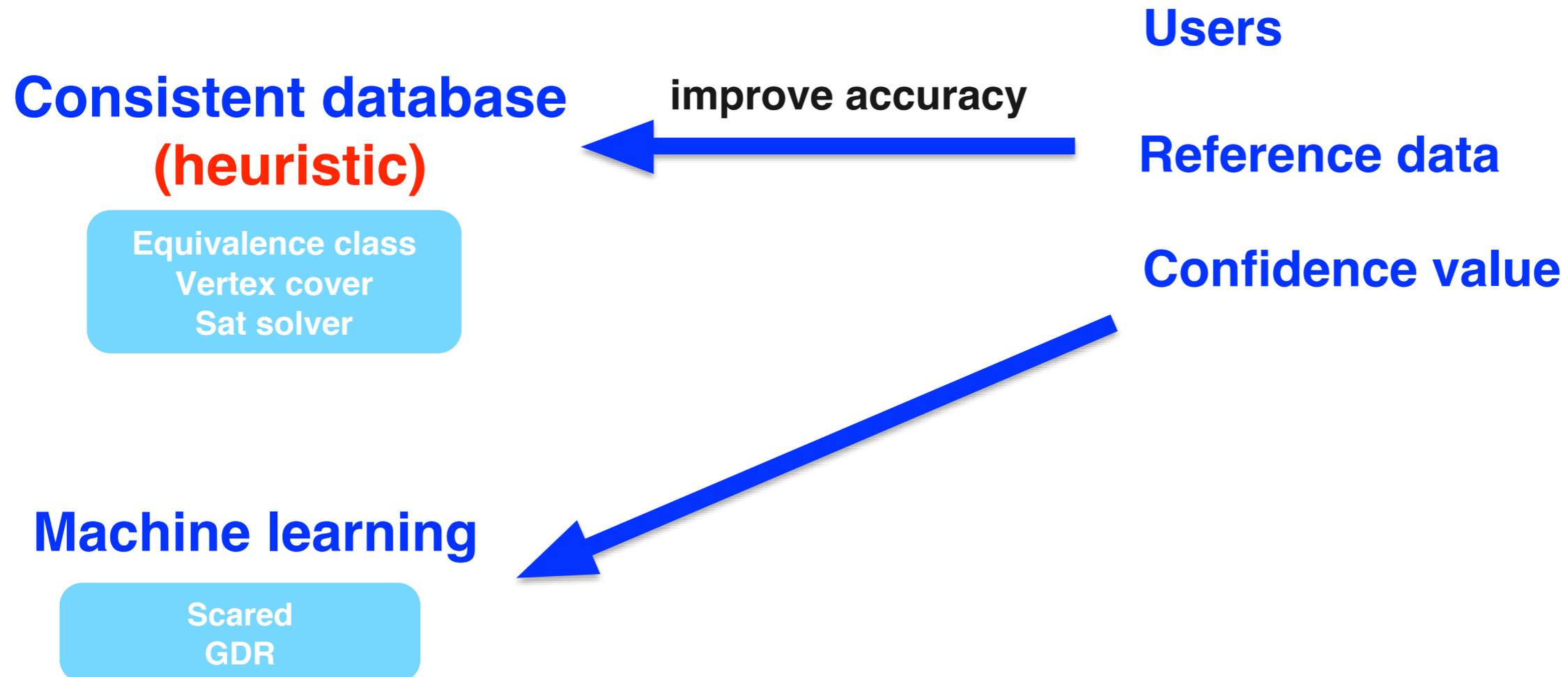


Users

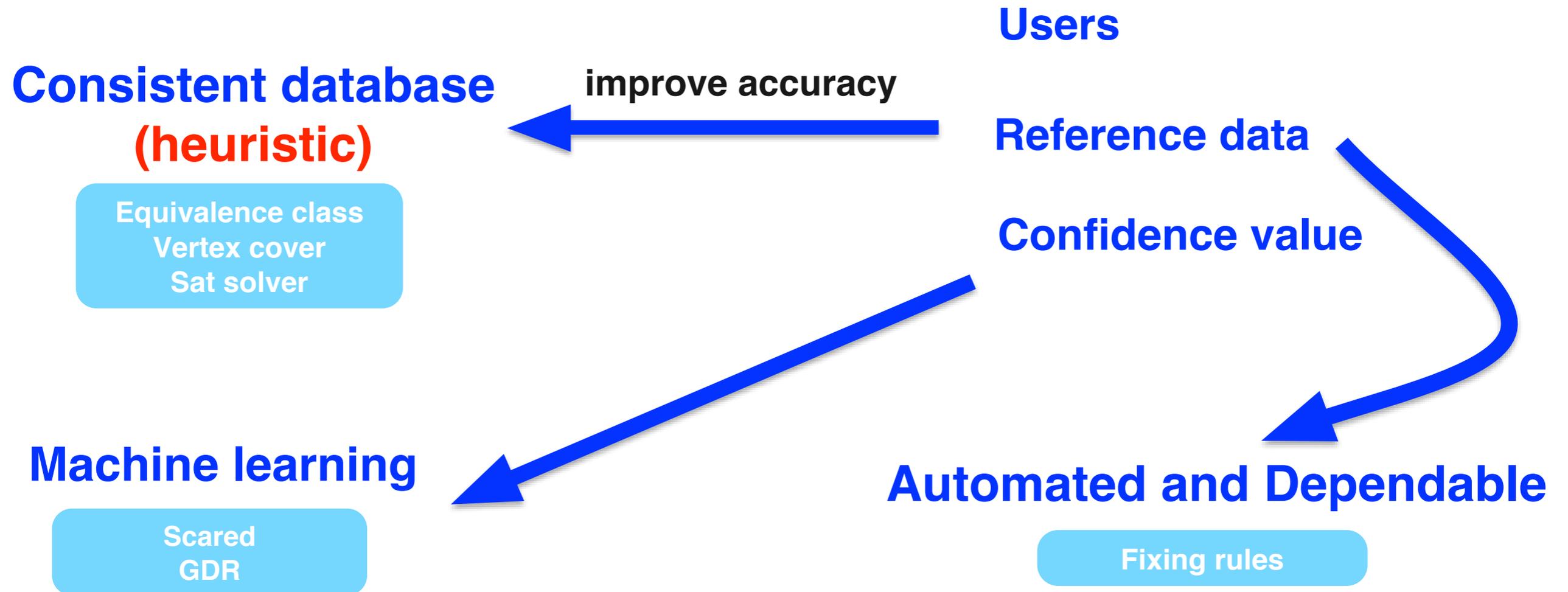
Reference data

Confidence value

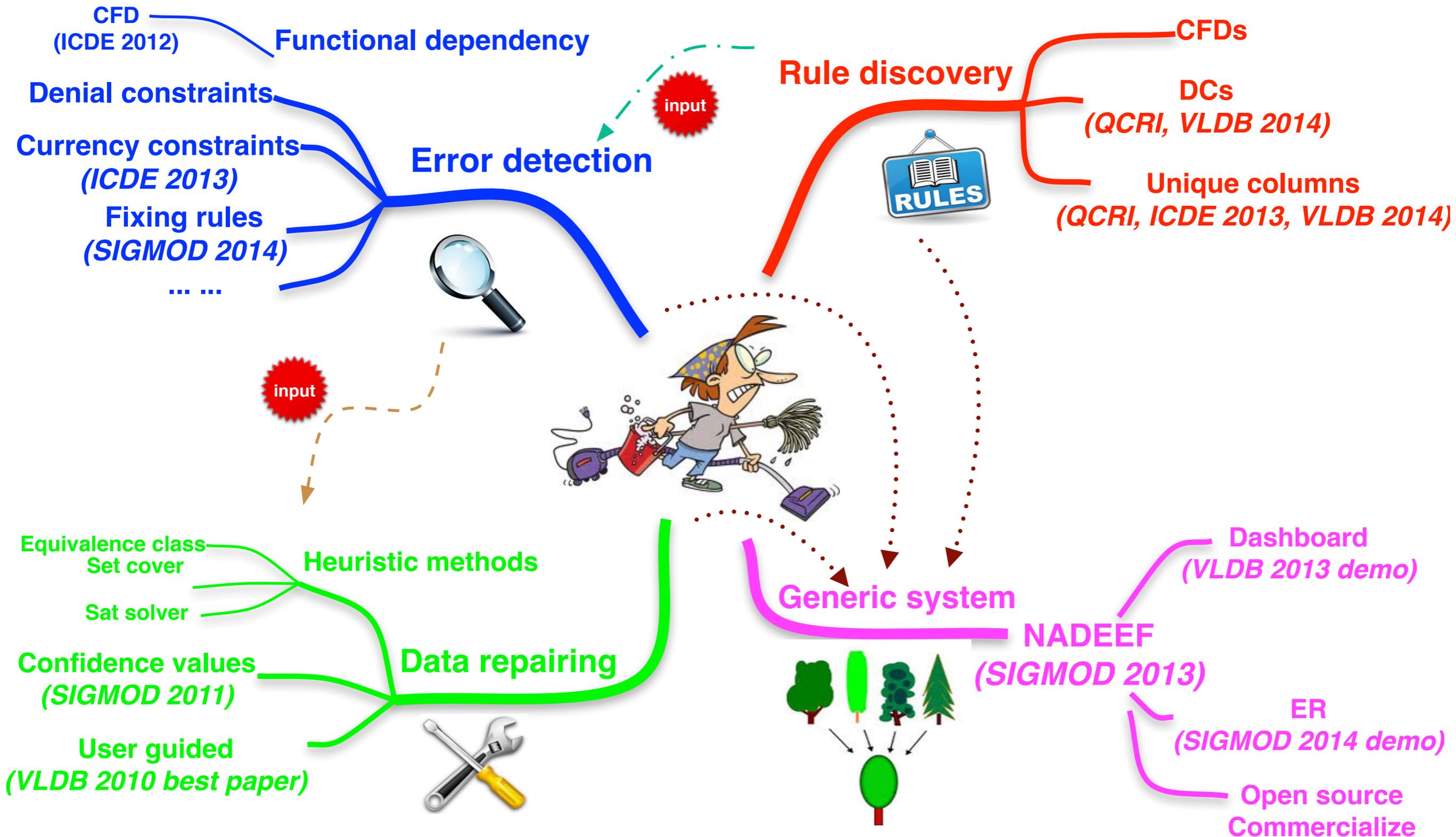
Summary of Data Repairing



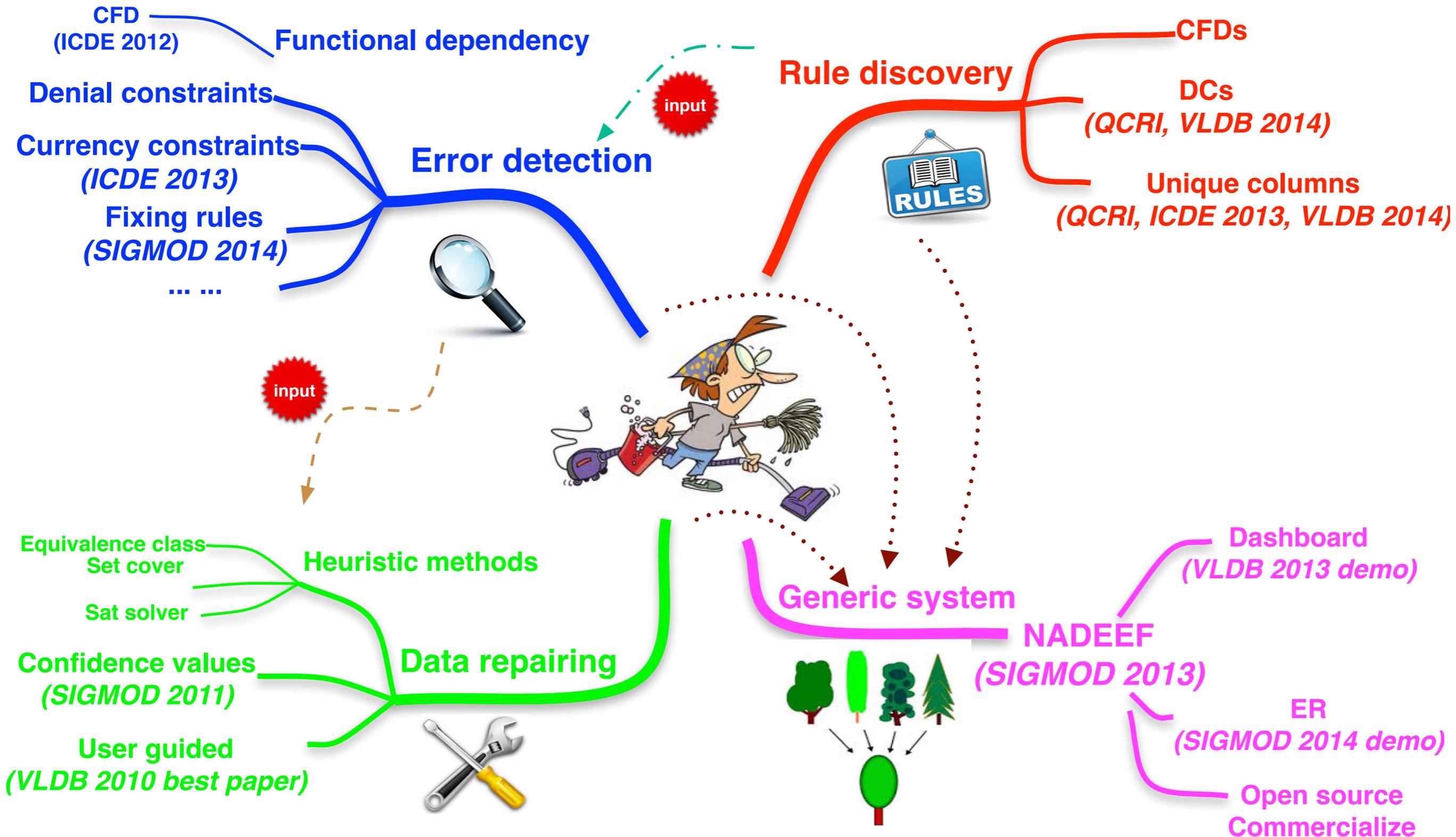
Summary of Data Repairing



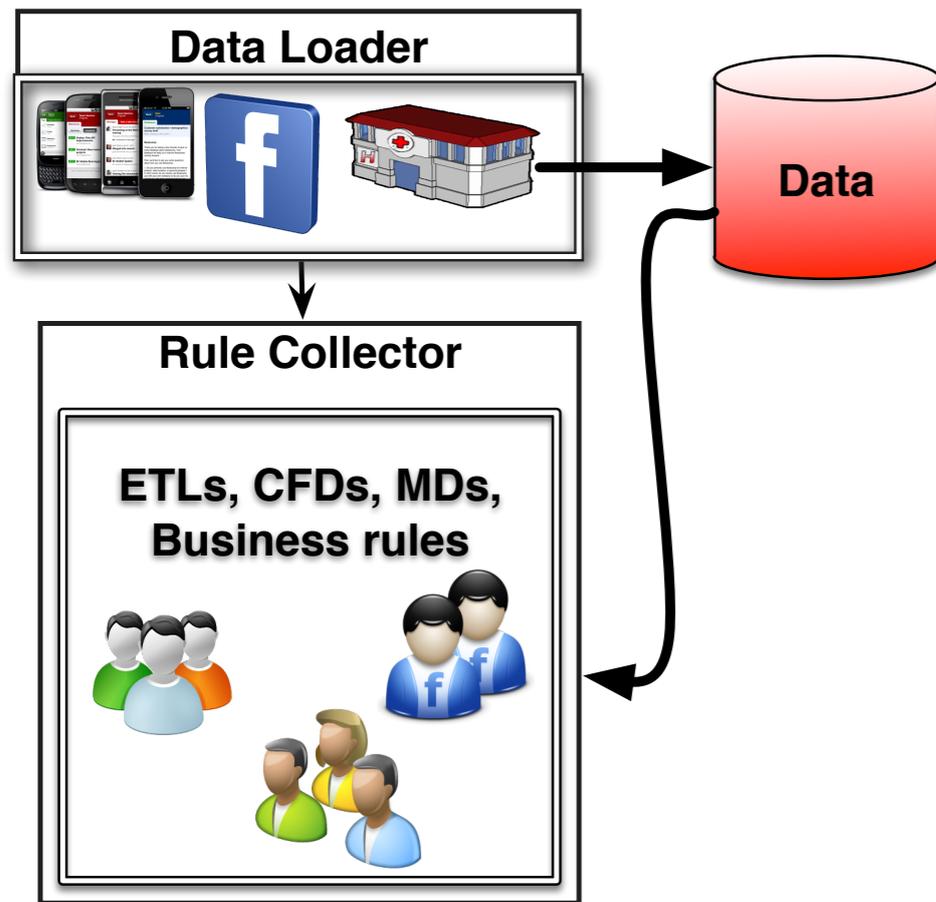
Generic Data Cleaning System



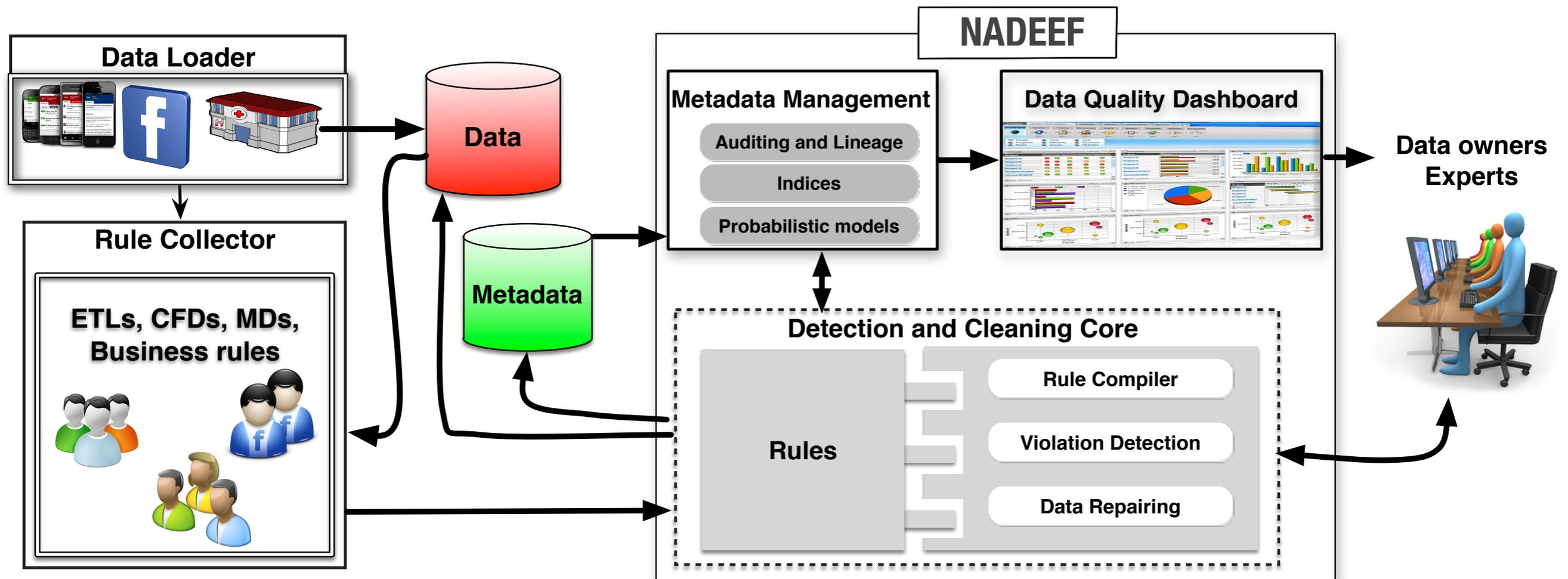
Generic Data Cleaning System



NADEEF (SIGMOD 2013)

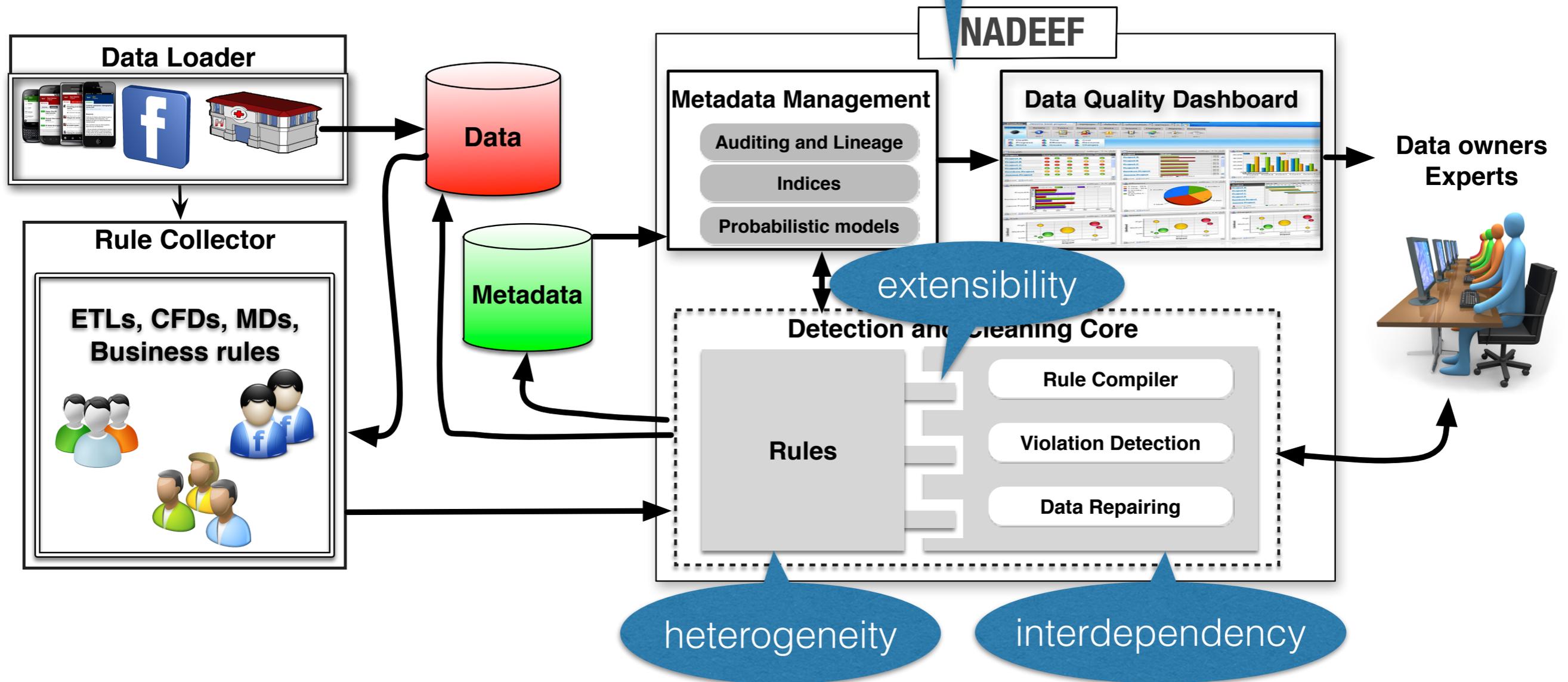


NADEEF (SIGMOD 2013)

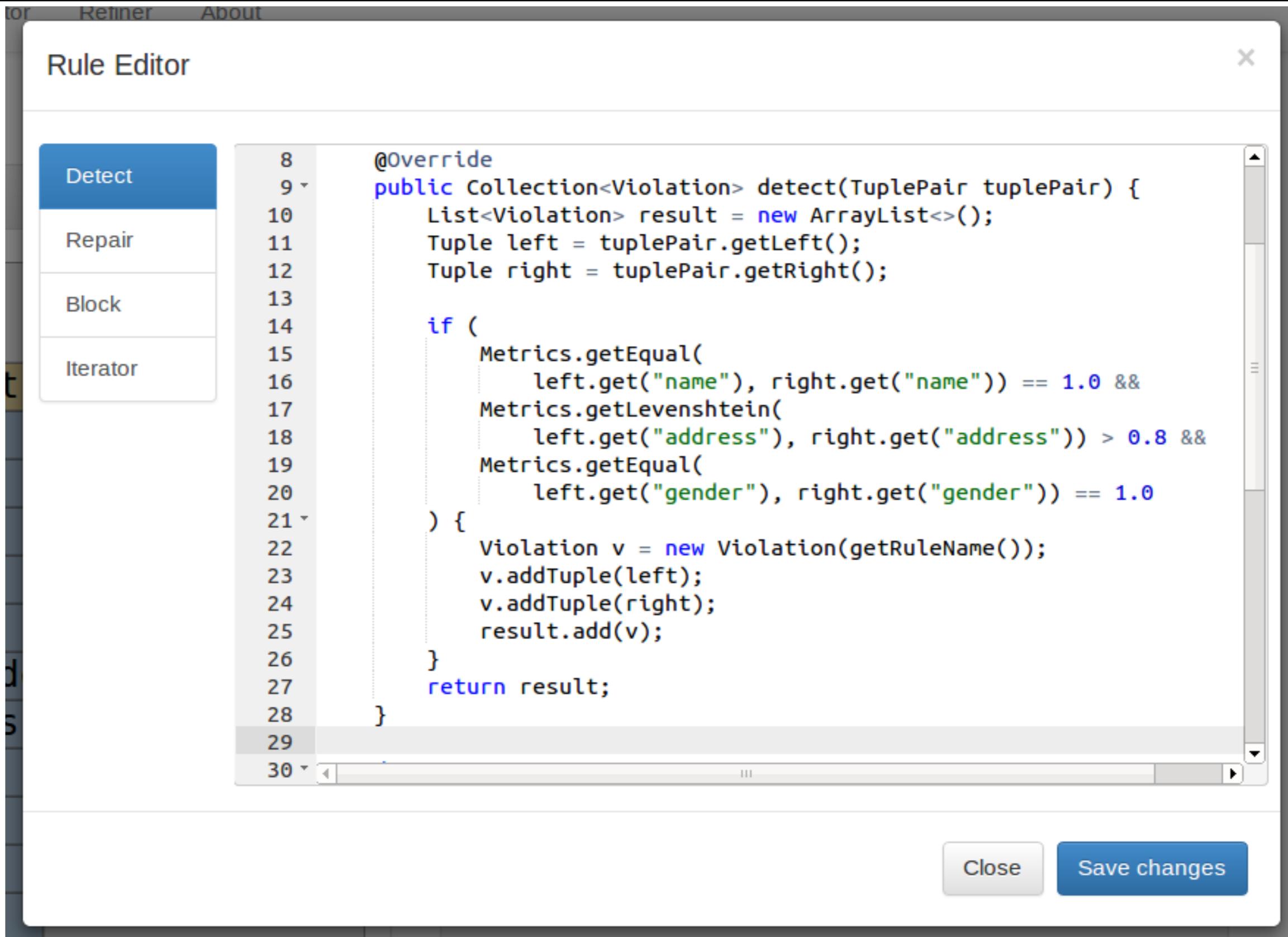


NADEEF (SIGMOD 2013)

metadata management and data custodians



NADEEF (SIGMOD 2013)

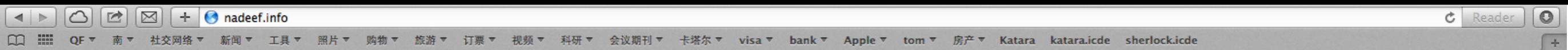


The screenshot shows a "Rule Editor" window with a sidebar on the left containing four buttons: "Detect" (highlighted in blue), "Repair", "Block", and "Iterator". The main area is a code editor with a line number column on the left (lines 8-30) and a code area on the right. The code is as follows:

```
8      @Override
9      public Collection<Violation> detect(TuplePair tuplePair) {
10         List<Violation> result = new ArrayList<>();
11         Tuple left = tuplePair.getLeft();
12         Tuple right = tuplePair.getRight();
13
14         if (
15             Metrics.getEqual(
16                 left.get("name"), right.get("name")) == 1.0 &&
17             Metrics.getLevenshtein(
18                 left.get("address"), right.get("address")) > 0.8 &&
19             Metrics.getEqual(
20                 left.get("gender"), right.get("gender")) == 1.0
21         ) {
22             Violation v = new Violation(getRuleName());
23             v.addTuple(left);
24             v.addTuple(right);
25             result.add(v);
26         }
27         return result;
28     }
29
30
```

At the bottom right of the window, there are two buttons: "Close" and "Save changes".

NADEEF Online



▶ QCRI DATA ANALYTICS

ABOUT FEATURES CONTACT

NADEEF

An extensible and user-friendly data cleaning system.

[GITHUB](#)

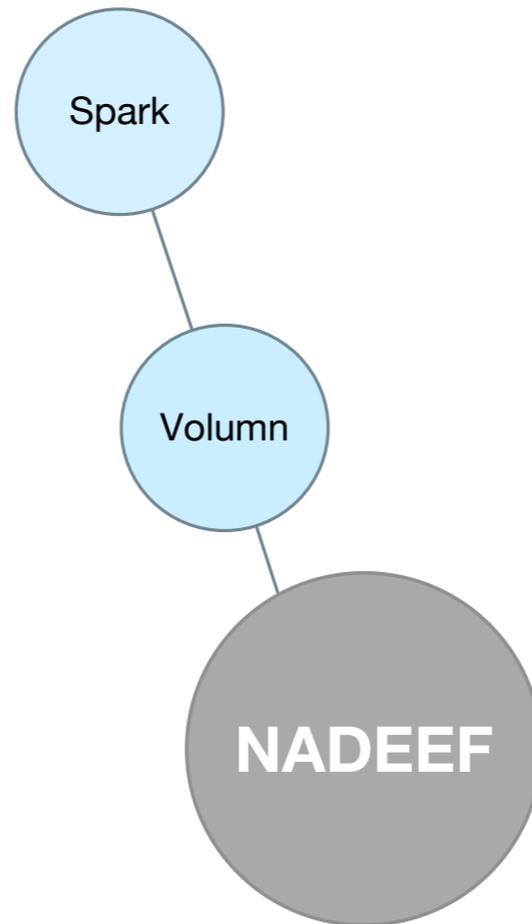
[LIVE DEMO](#)



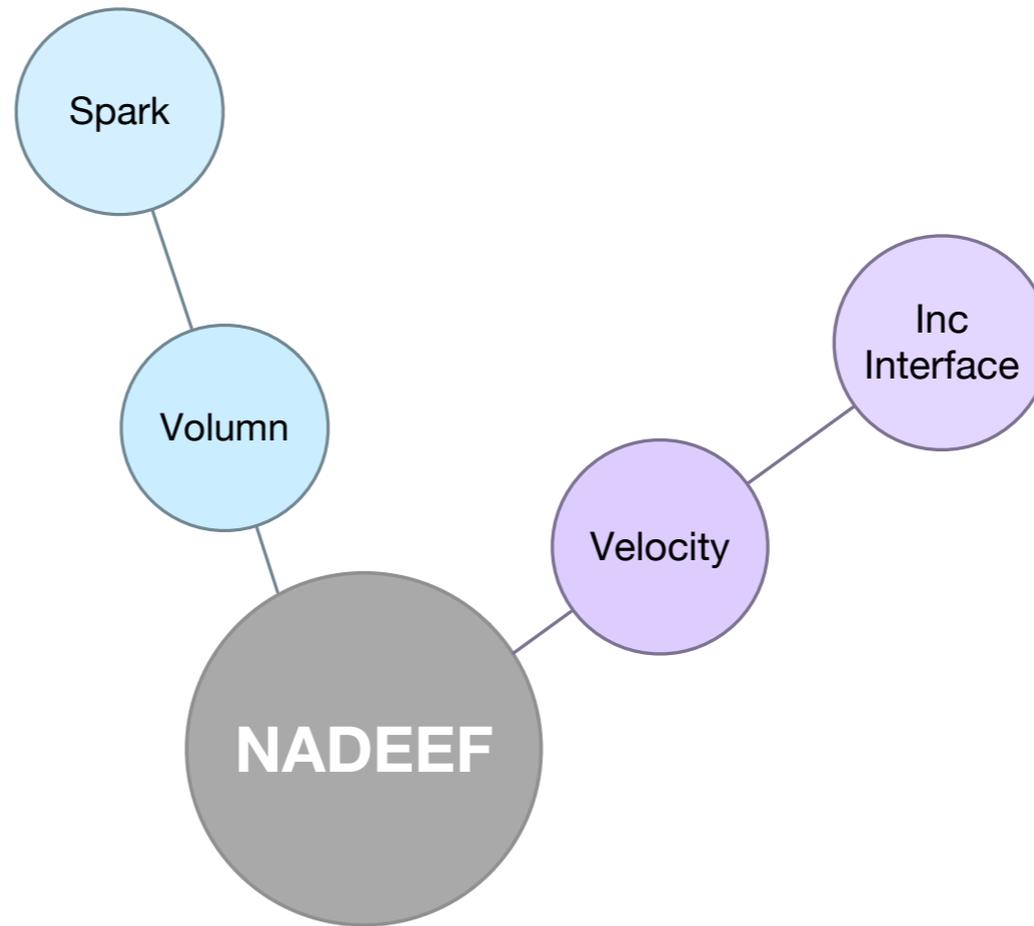
NADEEF for Big Data



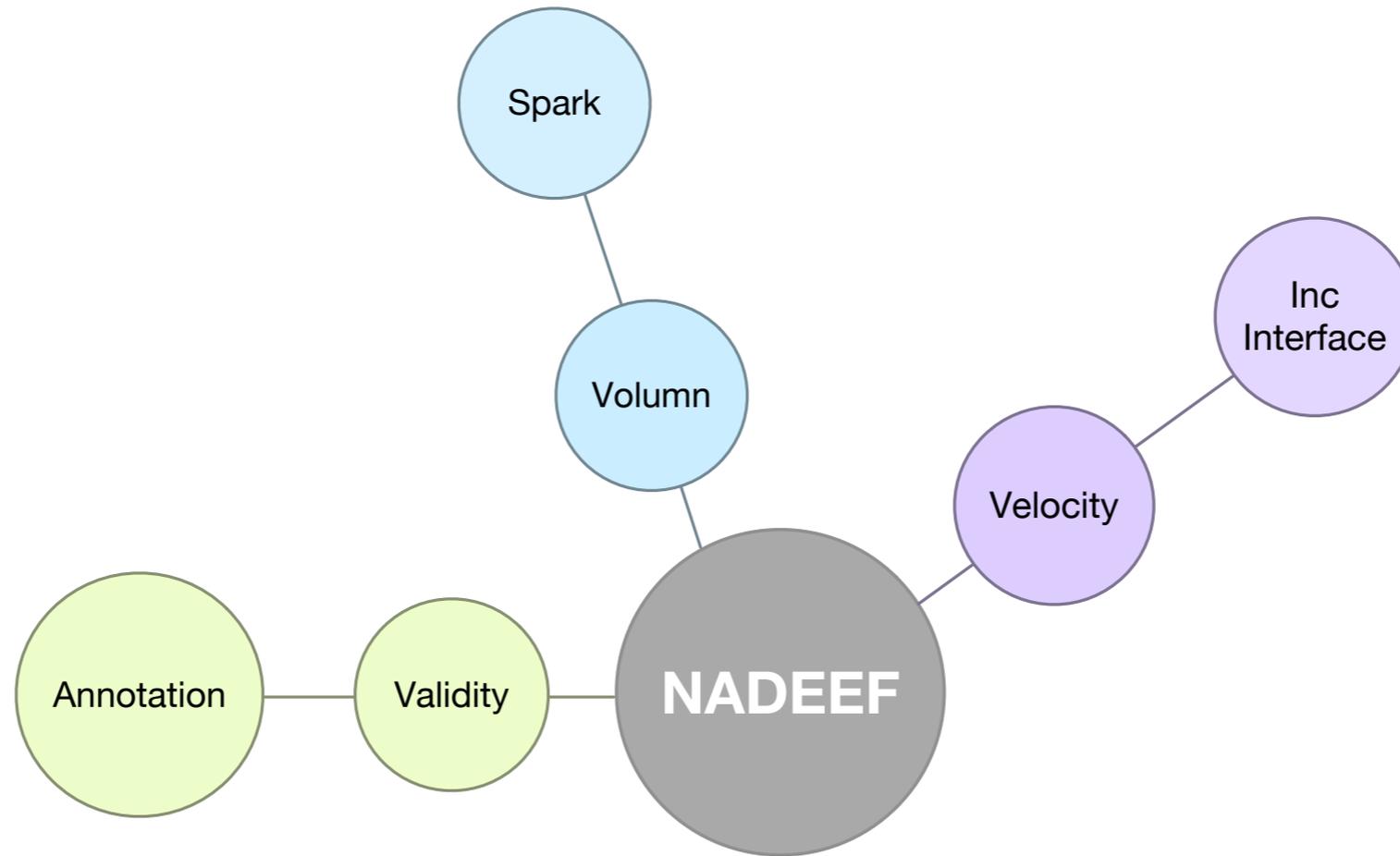
NADEEF for Big Data



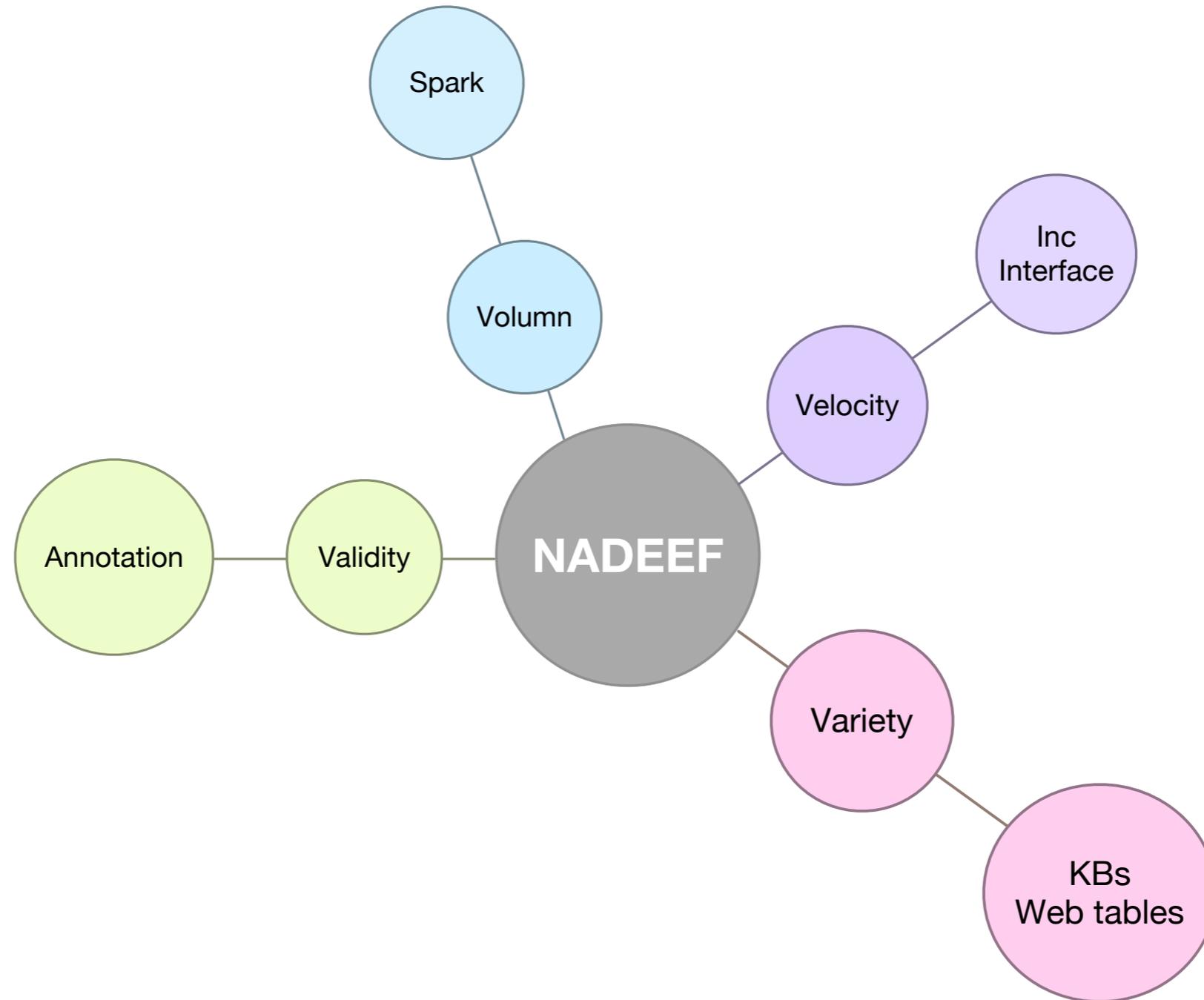
NADEEF for Big Data



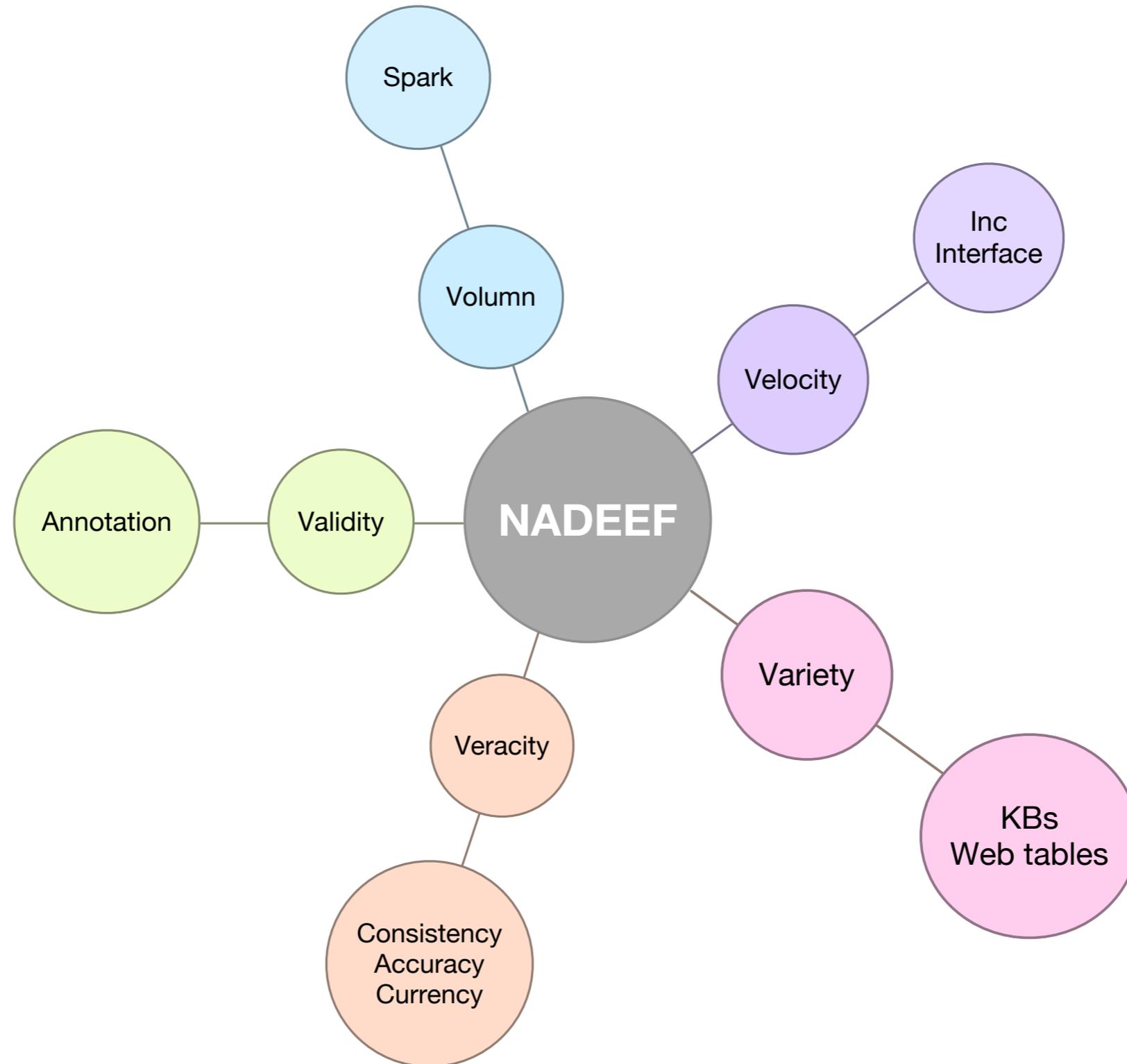
NADEEF for Big Data



NADEEF for Big Data



NADEEF for Big Data



Future Work

- Error detection
 - Rule discovery and validation
 - Combining different methods
- Explain errors to users
 - Summarization
 - Visualization
- Reliable data repairing
 - Effectively involve users as first-class citizens