

Detecting Inconsistencies in Distributed Data

Wenfei Fan
Floris Geerts
Shuai Ma
Heiko Müller

Agenda

Inconsistencies in Distributed Data

Problem Statement & Complexity

Detection Algorithms

Experimental Results

Conclusion & Outlook

Data Quality & Data Cleaning

Data quality tools support two methods for data cleaning:

- Detect data of poor quality (*detection*).
- Modify data of poor quality (*repair*).



Integrity constraints restrict valid database states, e.g., functional dependencies (FD).

Constraint violations highlight errors (*inconsistencies*).

How to express that $[CC, ZIP \rightarrow Street]$, but *only* for U.K. (CC=44) and Netherlands (CC=31)?

Conditional Functional Dependencies (CFDs)

CFDs **extend** FDs with conditions to **restrict their scope** and allow for **more expressive data quality rules**.

Embedded FD

Definition

A **CFD** φ is defined as $(X \rightarrow Y, T_p)$

Pattern Tableau

Example:

φ

CC	ZIP	\rightarrow	Street
44	-		-
31	-		-

Pattern Tuple

Bohannon, Fan, Geerts, Jia, Kementsietsidis, ICDE 2007.

Constraint Violations

\varnothing

CC	ZIP	→	Street
44	-		-
31	-		-

D

Name	Title	CC	AC	Phone	Street	City	ZIP
Sam	DMTS	44	131	2501984	Princess Str.	EDI	EH2 4HF
Philip	DMTS	44	131	4459011	Crichton Str.	EDI	EH4 8LE
Adam	VP	44	131	1326184	Mayfield Rd.	EDI	EH4 8LE
Joe	MTS	01	908	9075271	Queensway Drive	NYC	07974
Bob	DMTS	01	908	7732134	57 th St.	MH	07974
Jef	DMTS	31	20	7464774	Muntplein	AMS	1012 WR
Steven	MTS	31	20	4521633	Spuistraat	AMS	1012 WR
Bram	MTS	31	10	8974638	Kruisplein	ROT	3012 CC

Constraint Violations

\varnothing

CC	ZIP	→	Street
44	-		-
31	-		-

Note

Tuples that match the LHS of a pattern tuple are potential inconsistencies (*matching* or *relevant* tuples).

D

Name	Title	CC	AC				
Sam	DMTS	44	131				
Philip	DMTS	44	131	4459011	Crichton Str.	EDI	EH4 8LE
Adam	VP	44	131	1326184	Mayfield Rd.	EDI	EH4 8LE
Joe	MTS	01	908	9075271	Queensway Drive	NYC	07974
Bob	DMTS	01	908	7732134	57 th St.	MH	07974
Jef	DMTS	31	20	7464774	Muntplein	AMS	1012 WR
Steven	MTS	31	20	4521633	Spuistraat	AMS	1012 WR
Bram	MTS	31	10	8974638	Kruisplein	ROT	3012 CC

Constraint Violations

φ

CC	ZIP	→	Street
44	-		-
31	-		-

Violations

Matching tuples with *identical LHS* value, but *different RHS* value.

D

Name	Title	CC	AC				
Sam	DMTS	44	131	2501984	Princess Str.	EDI	EH2 4HF
Philip	DMTS	44	131	4459011	Crichton Str.	EDI	EH4 8LE
Adam	VP	44	131	1326184	Mayfield Rd.	EDI	EH4 8LE
Joe	MTS	01	908	9075271	Queensway Drive	NYC	07974
Bob	DMTS	01	908	7732134	57 th St.	MH	07974
Jef	DMTS	31	20	7464774	Muntplein	AMS	1012 WR
Steven	MTS	31	20	4521633	Spuistraat	AMS	1012 WR
Bram	MTS	31	10	8974638	Kruisplein	ROT	3012 CC

We distinguish two types of CFDs

Constant CFDs

 φ_1

CC	AC	→	City
44	131		EDI
01	908		MH

Variable CFDs

 φ_2

CC	ZIP	→	Street
44	-		-
31	-		-

 φ_3

ZIP	→	City
-		-

We distinguish two types of CFDs

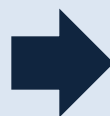
Constant CFDs

Variable CFDs

Decomposition of CFDs

φ_1

CC	ZIP	→	City
01	-		-
01	H0H 0H0		NORTH POLE



φ_2

CC	ZIP	→	City
01	-		-

φ_3

CC	ZIP	→	City
01	H0H 0H0		NORTH POLE

We distinguish two types of CFDs

Constant CFDs

φ_1

CC	AC	→	City
44	131		EDI
01	908		MH

Variable CFDs

φ_2

CC	ZIP	→	Street
44	-		-
31	-		-

φ_3

ZIP	→	City
-		-

Violations

Constant CFDs are violated by a **single tuple**.

Variable CFDs are violated by at least **two tuples**.

Detecting Constraint Violations

In a **centralized setting** violations of a set of CFDs can be detected using **two SQL queries**.

How to detect violations in a **distributed** setting?

Detecting Constraint Violations (cont.)

D_1

Name	Title	CC	AC	Phone	Street	City	ZIP
Sam	DMTS	44	131	2501984	Princess Str.	EDI	EH2 4HF
Philip	DMTS	44	131	4459011	Crichton Str.	EDI	EH4 8LE
Bob	DMTS	01	908	7732134	57 th St.	MH	07974
Jef	DMTS	31	20	7464774	Muntplein	AMS	1012 WR

D_2

Name	Title	CC	AC	Phone	Street	City	ZIP
Joe	MTS	01	908	9075271	Queensway Drive	NYC	07974
Steven	MTS	31	20	4521633	Spuistraat	AMS	1012 WR
Bram	MTS	31	10	8974638	Kruisplein	ROT	3012 CC

D_3

Name	Title	CC	AC	Phone	Street	City	ZIP
Adam	VP	44	131	1326184	Mayfield Rd.	EDI	EH4 8LE

Detecting Constraint Violations (cont.)

D_1	Name	Title	CC	AC	Violations			
	Sam	DMTS	44	131	Inconsistent tuples may be distributed across different sites.			
	Philip	DMTS	44	131				
	Bob	DMTS	01	908				
	Jef	DMTS	31	20	7464774	Muntplein	AMS	1012 WR

D_2	Name	Title	CC	AC	Phone	Street	City	ZIP
	Joe	MTS	01	908	9075271	Queensway Drive	NYC	07974
	Steven	MTS	31	20	4521633	Spuistraat	AMS	1012 WR
	Bram	MTS	31	10	8974638	Kruisplein	ROT	3012 CC

D_3	Name	Title	CC	AC	Phone	Street	City	ZIP
	Adam	VP	44	131	1326184	Mayfield Rd.	EDI	EH4 8LE

Detecting Constraint Violations (cont.)

In a **centralized setting** violations of a set of CFDs can be detected using **two** SQL queries.

How to detect violations in a **distributed** setting?

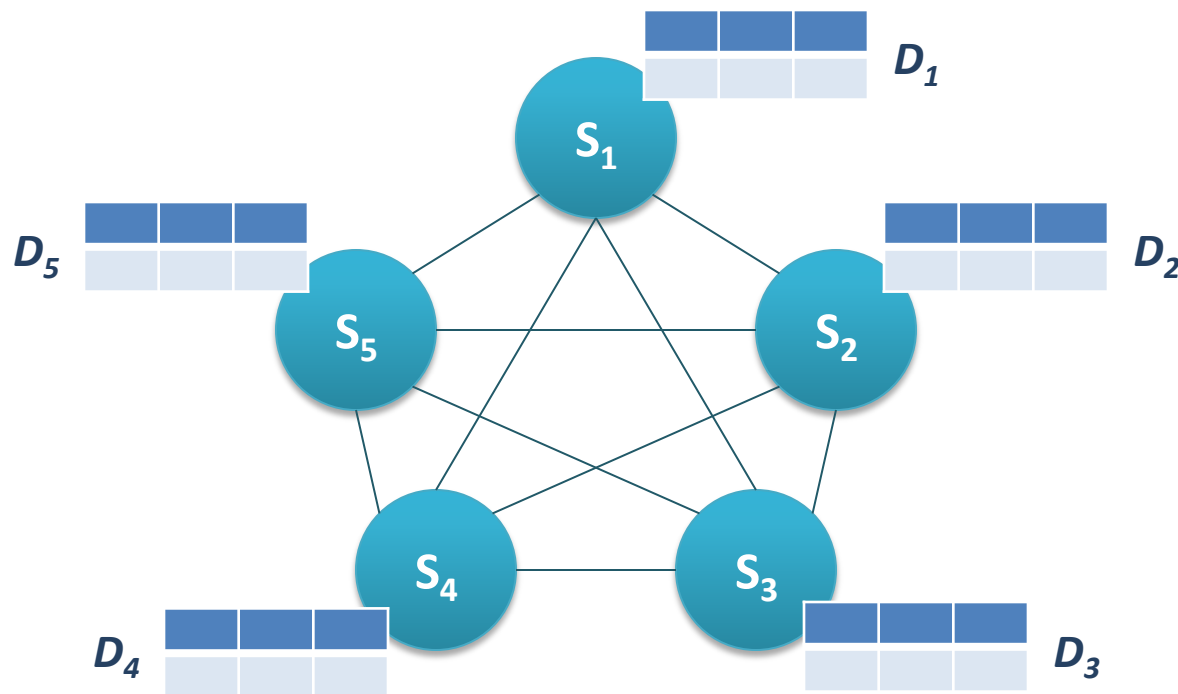
Data shipment required to detect violations.

Contributions

- 1) Violation detection as *optimization problem*.
- 2) Establish *complexity bounds*.
- 3) *Efficient algorithms* for violation detection.

Detecting Inconsistencies in Distributed Data

Given a set of CFDs $\Sigma = \{\varphi_1, \dots, \varphi_n\}$, and a fragmented relation $D = \{D_1, \dots, D_k\}$.



Detecting Inconsistencies in Distributed Data

Given a set of CFDs $\Sigma = \{\varphi_1, \dots, \varphi_n\}$, and a fragmented relation $D = \{D_1, \dots, D_k\}$.

Find all violations of CFDs $\varphi \in \Sigma$ in D with ...

- *minimal data shipment* (number of tuples), or
- *minimal overall response time* (shipment + local detection).

Complexity

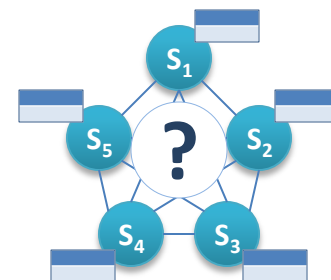
The problem is **NP-complete** in either setting!

Fragmentation of data

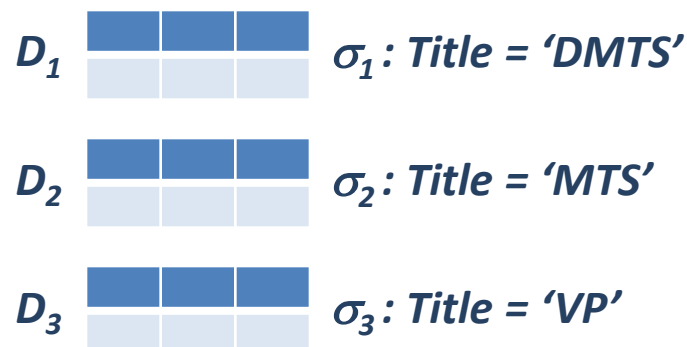
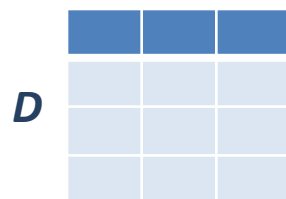
Horizontal fragmentation partitions D based on Boolean predicates $\sigma_1, \dots, \sigma_k$ such that:

- a) fragments are pair-wise disjoint, and
- b) original relation D results from $\sigma_1 \cup \dots \cup \sigma_k$.

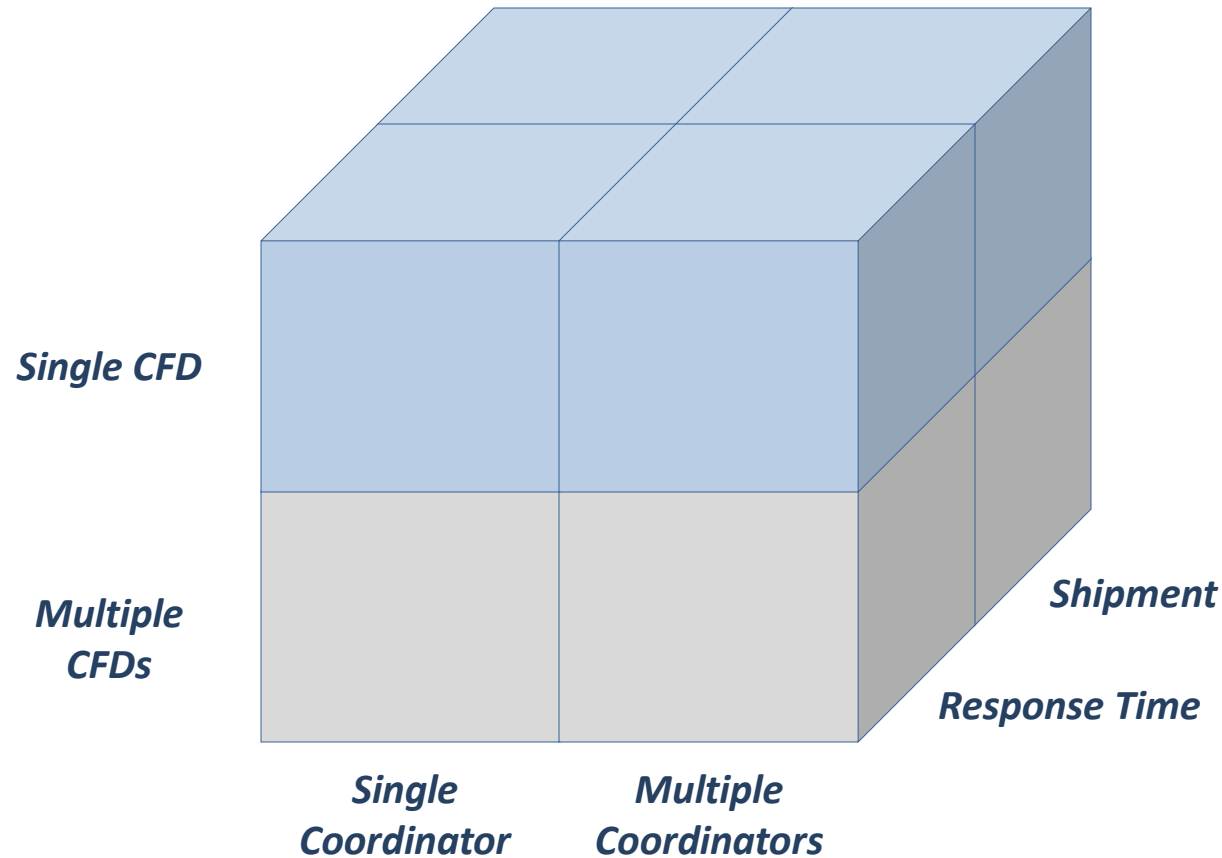
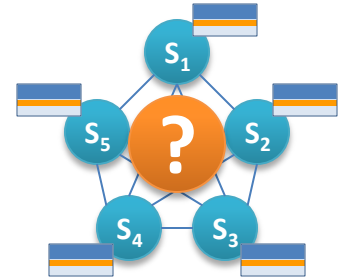
Each fragment resides at different network site.



Example:



Validating Horizontally Partitioned Data



Local Validation of CFDs

In two cases data shipping can be avoided.

1) Constant CFDs

\varnothing	CC	AC	→	City
	44	131		EDI
	01	908		MH

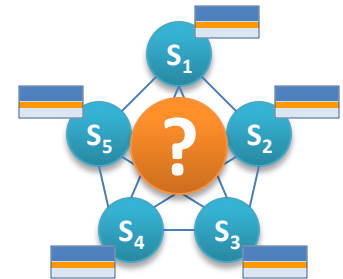
2) Partitioning Condition

\varnothing	CC	Title	→	Salary
	-	DMTS		-
	-	MTS		-
	-	VP		-

$\sigma_1 : \text{Title} = \text{'DMTS'}$

$\sigma_2 : \text{Title} = \text{'MTS'}$

$\sigma_3 : \text{Title} = \text{'VP'}$



Validation of CFDs involving Shipments

All algorithms are based on **four main steps**:



Algorithm Outline

- 1 Broadcast local statistics
- 2 Decide on coordinator(s)
- 3 Ship data to coordinator(s)
- 4 Violation detection at coordinator(s)

Validation of CFDs involving Shipments

All algorithms are based on **four main steps**:



Algorithm Outline

- 1 Broadcast local statistics
- 2 Decide on coordinator(s)
- 3 Ship data to coordinator(s)
- 4 Violation detection at coordinator(s)

Variation (1)

*Differ between **single coordinator** and **multiple coordinator** approach.*

Variation (2)

*Depends on **cost function**.*

Central approach

Ship all relevant tuples to a **single site** (*coordinator*).

Detect violations at coordinator site.



CentralDetect

- 1 Broadcast local statistics
- 2 Decide on coordinator
- 3 Ship data to coordinator
- 4 Violation detection at coordinator

Statistics

Total number of matching tuples.

Central approach (cont.)



	CC	ZIP	→	Street
\varnothing	44	-		-
	31	-		-

D_1	Name	Title	CC	AC	Phone	Street	City	ZIP
	Sam	DMTS	44	131	2501984	Princess Str.	EDI	EH2 4HF
	Philip	DMTS	44	131	4459011	Crichton Str.	EDI	EH4 8LE
	Bob	DMTS	01	908	7732134	57 th St.	MH	07974
	Jef	DMTS	31	20	7464774	Muntplein	AMS	1012 WR

3

D_2	Name	Title	CC	AC	Phone	Street	City	ZIP
	Joe	MTS	01	908	9075271	Queensway Drive	NYC	07974
	Steven	MTS	31	20	4521633	Spuistraat	AMS	1012 WR
	Bram	MTS	31	10	8974638	Kruisplein	ROT	3012 CC

2

D_3	Name	Title	CC	AC	Phone	Street	City	ZIP
	Adam	VP	44	131	1326184	Mayfield Rd.	EDI	EH4 8LE

1

Central approach (cont.)

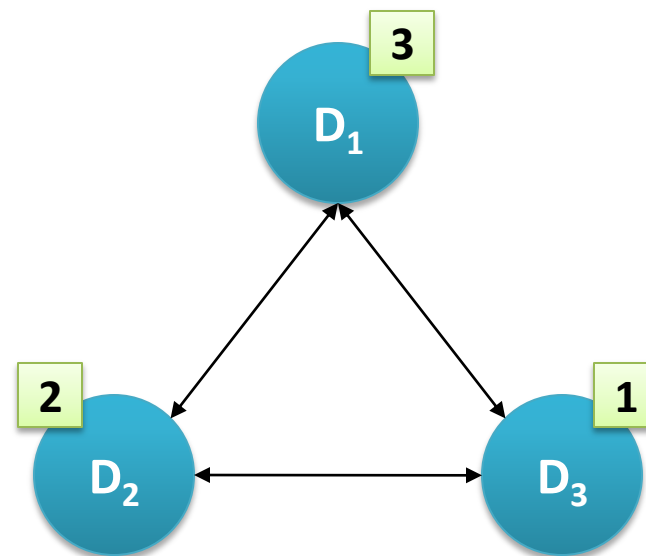
Ship all relevant tuples to a **single site** (*coordinator*).

Detect violations at **coordinator** site.



CentralDetect

- 1 Broadcast local statistics
- 2 Decide on coordinator
- 3 Ship data to coordinator
- 4 Violation detection at coordinator



Central approach (cont.)

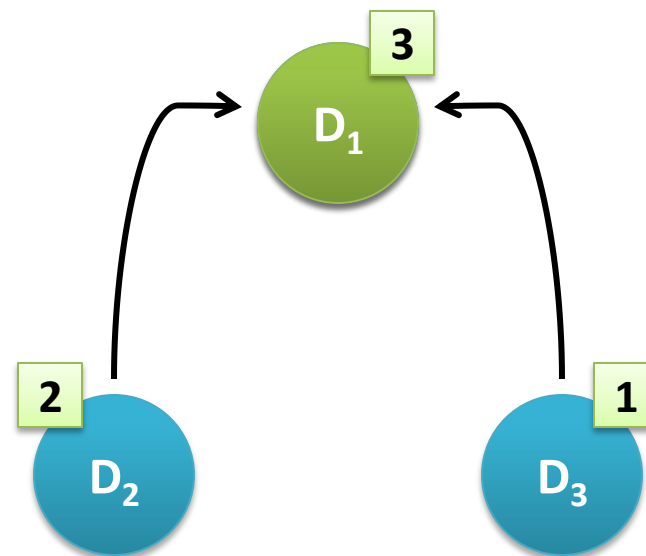
Ship all relevant tuples to a **single site** (*coordinator*).

Detect violations at **coordinator** site.



CentralDetect

- 1 Broadcast local statistics
- 2 Decide on coordinator
- 3 Ship data to coordinator
- 4 Violation detection at coordinator



Central approach (cont.)

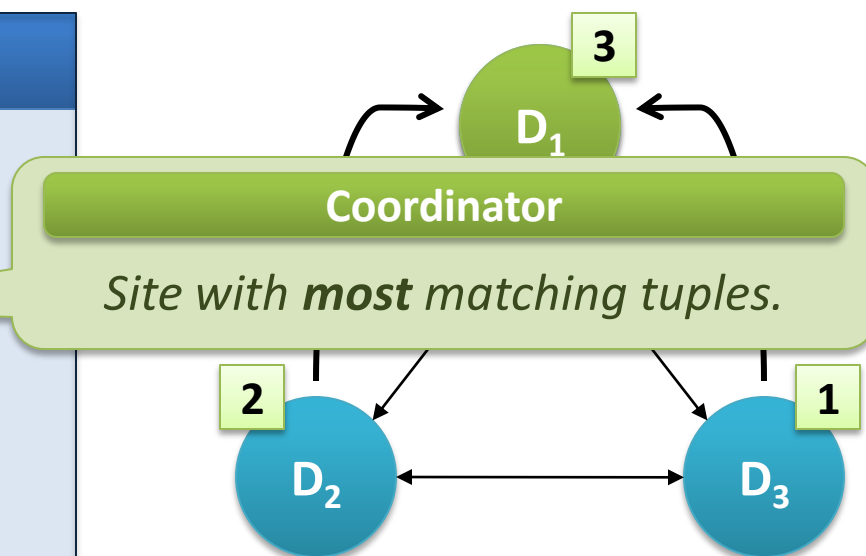
Ship all relevant tuples to a **single site (coordinator)**.

Detect violations at **coordinator** site.



CentralDetect

- 1 Broadcast local statistics
- 2 Decide on coordinator
- 3 Ship data to coordinator
- 4 Violation detection at coordinator



Pattern coordinator

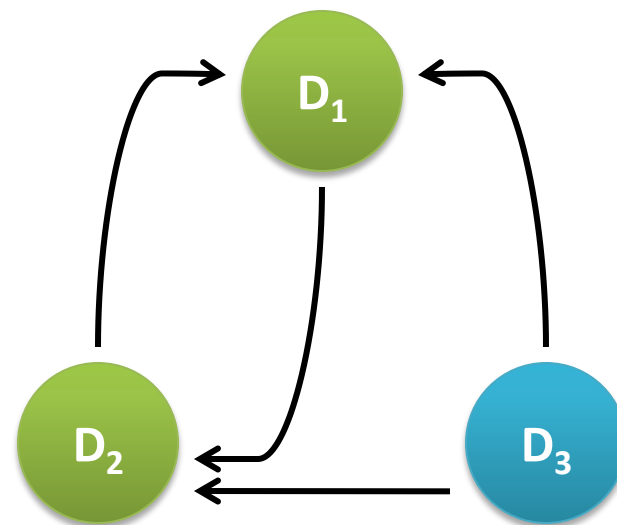
Leverage **structure of pattern tableau**.



- Assign **coordinator** for **individual pattern tuples**.
- Distribute detection process to **multiple coordinators**.

PatternDetect

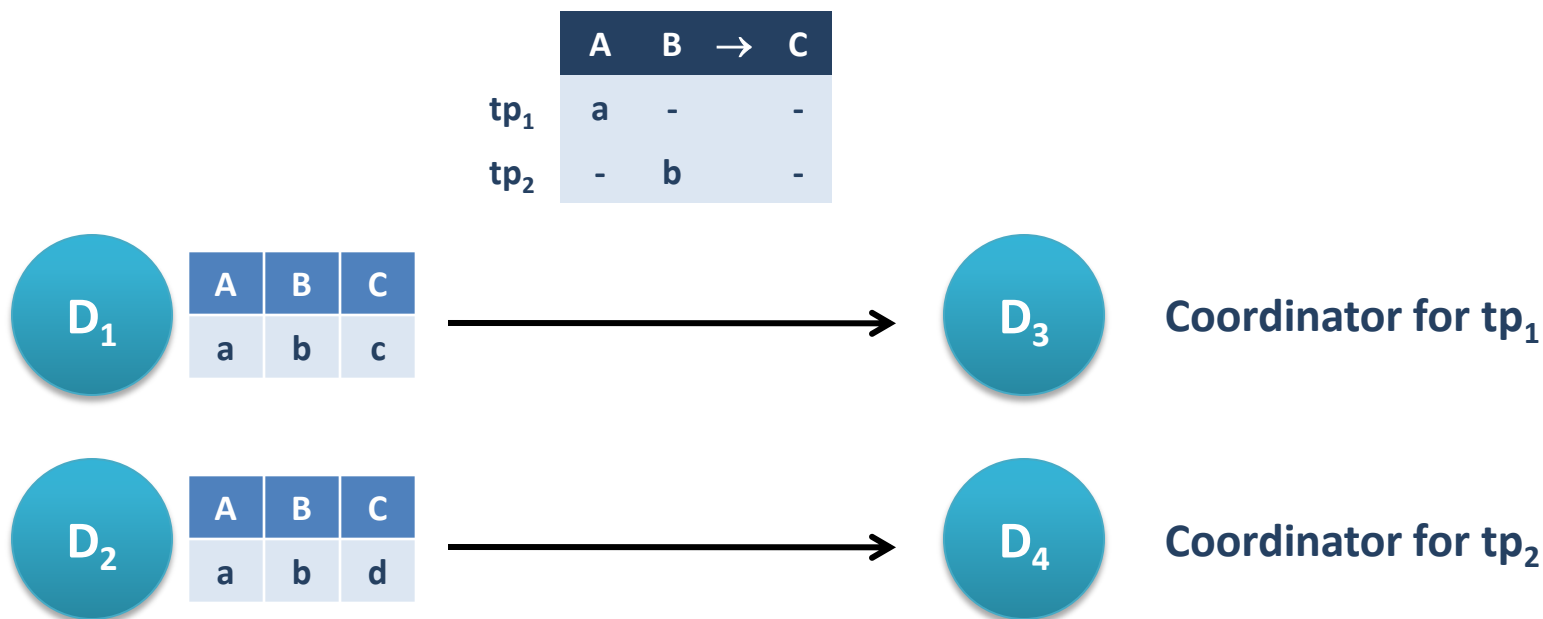
- 1 Broadcast local statistics
- 2 Decide on coordinator(s)
- 3 Ship data to coordinator(s)
- 4 Violation detection at coordinator(s)



Pattern coordinator (cont.)

How to ensure that tuples **matching multiple pattern tuples** are counted and send **only once**?

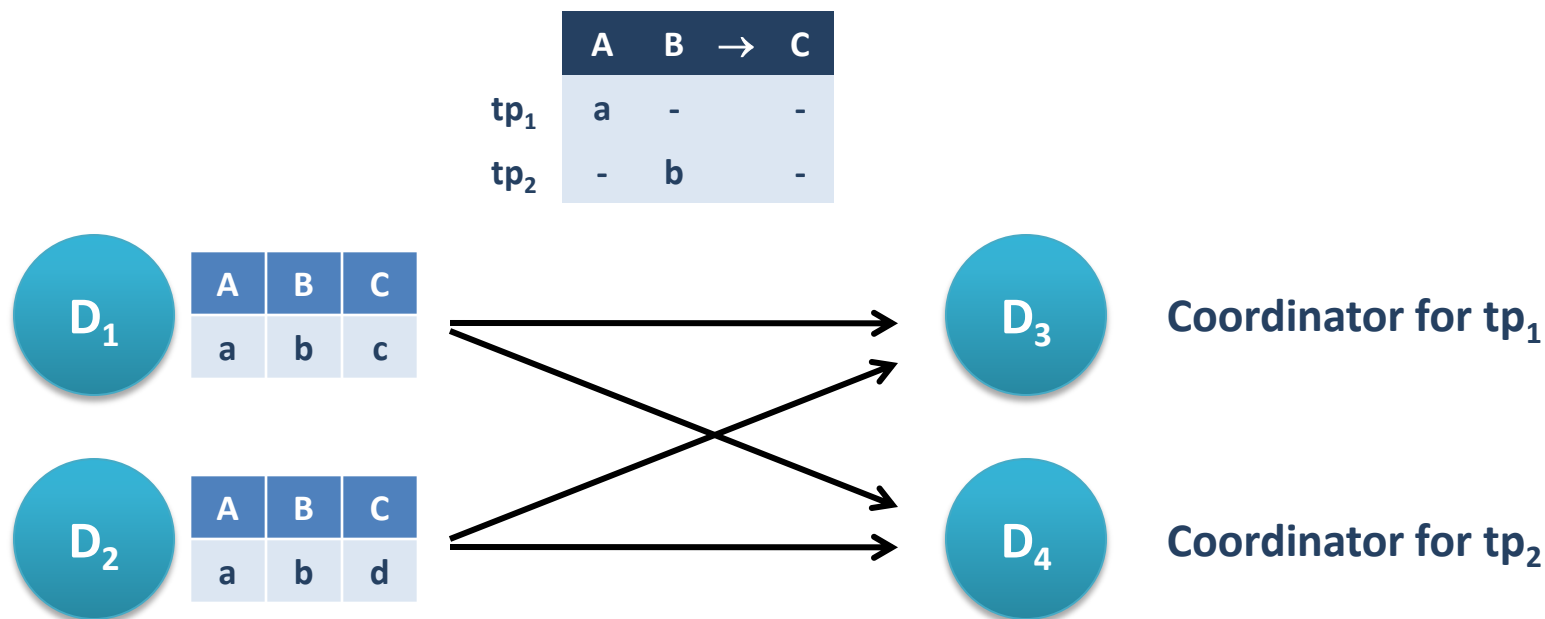
How to ensure that we **don't miss any violations**?



Pattern coordinator (cont.)

How to ensure that tuples **matching multiple pattern tuples** are counted and send **only once**?

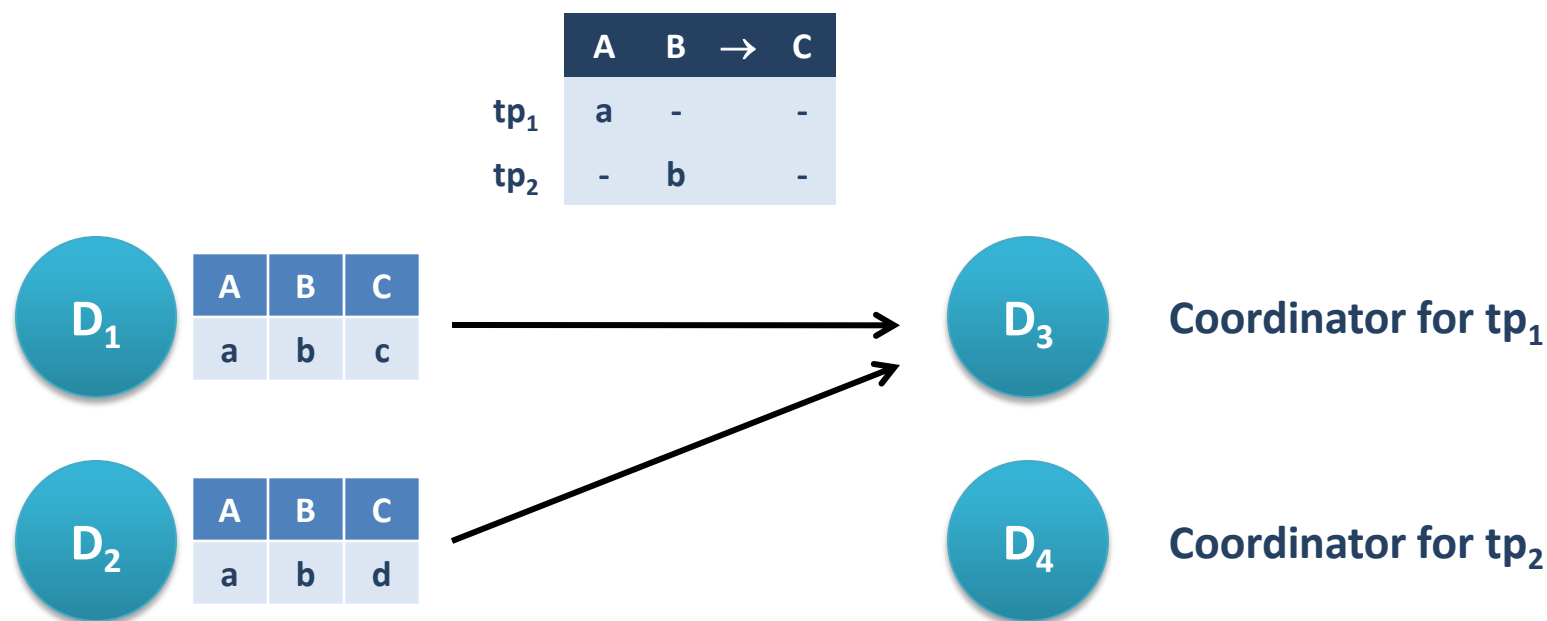
How to ensure that we **don't miss any violations**?



Pattern coordinator (cont.)

How to ensure that tuples **matching multiple pattern tuples** are counted and send **only once**?

How to ensure that we **don't miss any violations**?



Pattern coordinator (cont.)

Each tuples accounts for (and is send to the coordinator of) the **first matching pattern tuple**.

Requires a fixed **order on pattern tuples**!

Preferred order has **constants before wildcards** ('-').



Example:

ORDER BY A, B

A	B	→	C
a	b		-
-	b		-
-	-		-

A	B	→	C
-	b		-
a	b		-
-	-		-

A	B	→	C
-	-		-
a	b		-
-	b		-

Pattern coordinator (cont.)



\varnothing		CC	ZIP	→	Street
	tp ₁	44	-		-
	tp ₂	31	-		-

D_1	Name	Title	CC	AC	Phone	Street	City	ZIP
	Sam	DMTS	44	131	2501984	Princess Str.	EDI	EH2 4HF
	Philip	DMTS	44	131	4459011	Crichton Str.	EDI	EH4 8LE
	Bob	DMTS	01	908	7732134	57 th St.	MH	07974
	Jef	DMTS	31	20	7464774	Muntplein	AMS	1012 WR

tp ₁	2
tp ₂	1

D_2	Name	Title	CC	AC	Phone	Street	City	ZIP
	Joe	MTS	01	908	9075271	Queensway Drive	NYC	07974
	Steven	MTS	31	20	4521633	Spuistraat	AMS	1012 WR
	Bram	MTS	31	10	8974638	Kruisplein	ROT	3012 CC

tp ₁	0
tp ₂	2

D_3	Name	Title	CC	AC	Phone	Street	City	ZIP
	Adam	VP	44	131	1326184	Mayfield Rd.	EDI	EH4 8LE

tp ₁	1
tp ₂	0

Pattern coordinator (cont.)

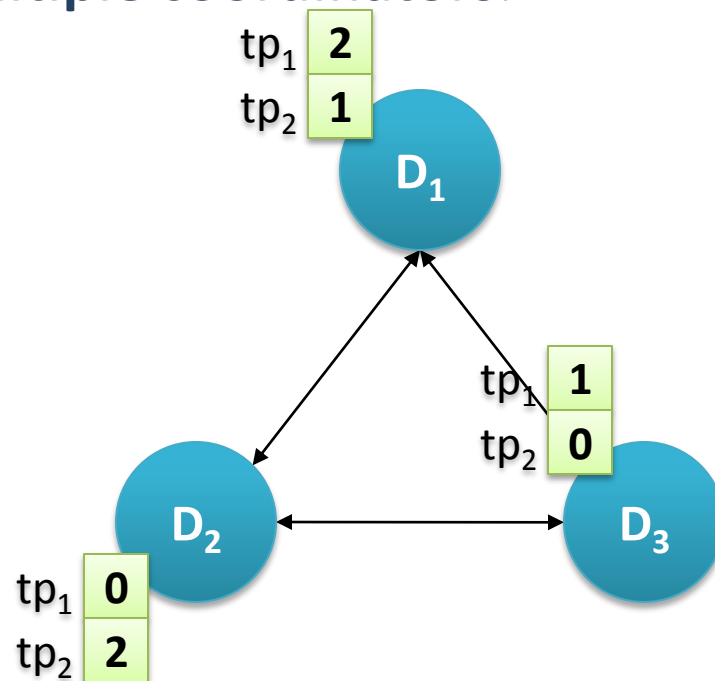
Leverage **structure of pattern tableau**.

- Assign **coordinator** for **individual pattern tuples**.
- Distribute detection process to **multiple coordinators**.



PatternDetect

- 1 Broadcast local statistics
- 2 Decide on coordinator(s)
- 3 Ship data to coordinator(s)
- 4 Violation detection at coordinator(s)



Pattern coordinator (cont.)

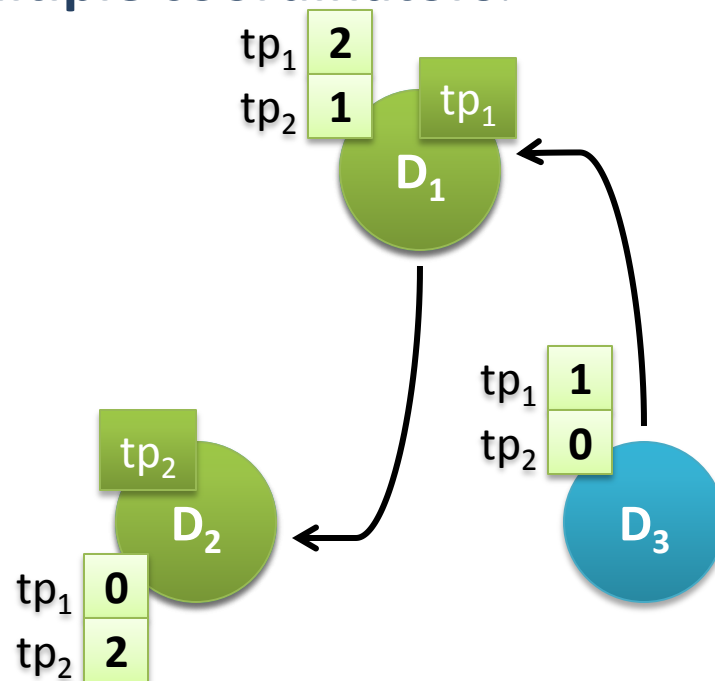
Leverage **structure of pattern tableau**.

- Assign **coordinator** for **individual pattern tuples**.
- Distribute detection process to **multiple coordinators**.



PatternDetect

- 1 Broadcast local statistics
- 2 Decide on coordinator(s)
- 3 Ship data to coordinator(s)
- 4 Violation detection at coordinator(s)



Pattern coordinator (cont.)

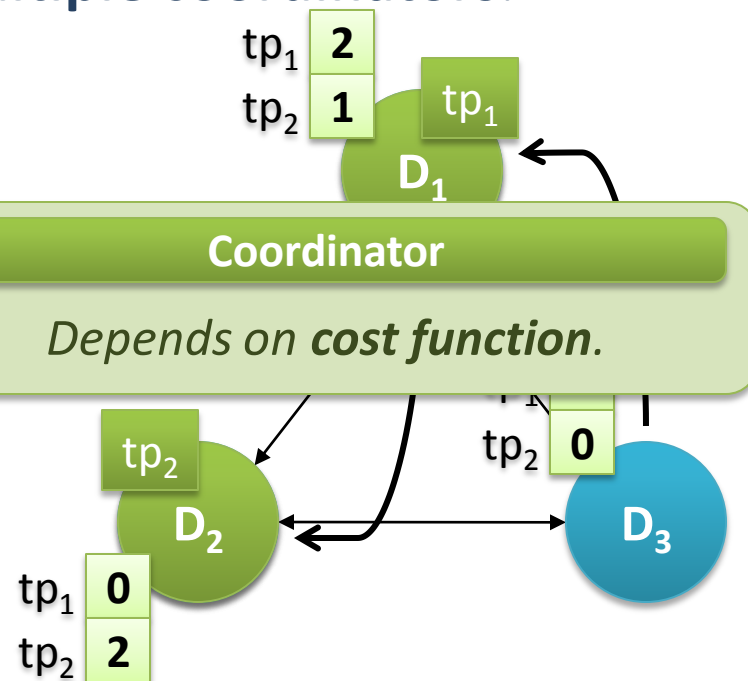
Leverage **structure of pattern tableau**.

- Assign **coordinator** for **individual pattern tuples**.
- Distribute detection process to **multiple coordinators**.



PatternDetect

- 1 Broadcast local statistics
- 2 Decide on coordinator(s)
- 3 Ship data to coordinator(s)
- 4 Violation detection at coordinator(s)



Impact of the presence of wildcards

For **sparse tableaus** (e.g., **FDs**) partitioning has minor (or no) impact.

Extend pattern tableau by **instantiating wildcards** with **frequent pattern tuples**.



MinePatternDetect

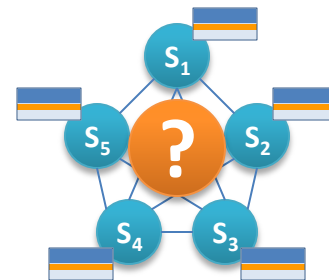
- 1a Mine locally for frequent patterns above given threshold
- 1b Exchange patterns and statistics
- 1c Construct extended pattern tableau
- 2 ...

Note

*Can significantly reduce **shipment**.*

Validating a set of CFDs

Given a set of CFDs $\Sigma = \{\varphi_1, \dots, \varphi_n\}$.



SeqDetect:

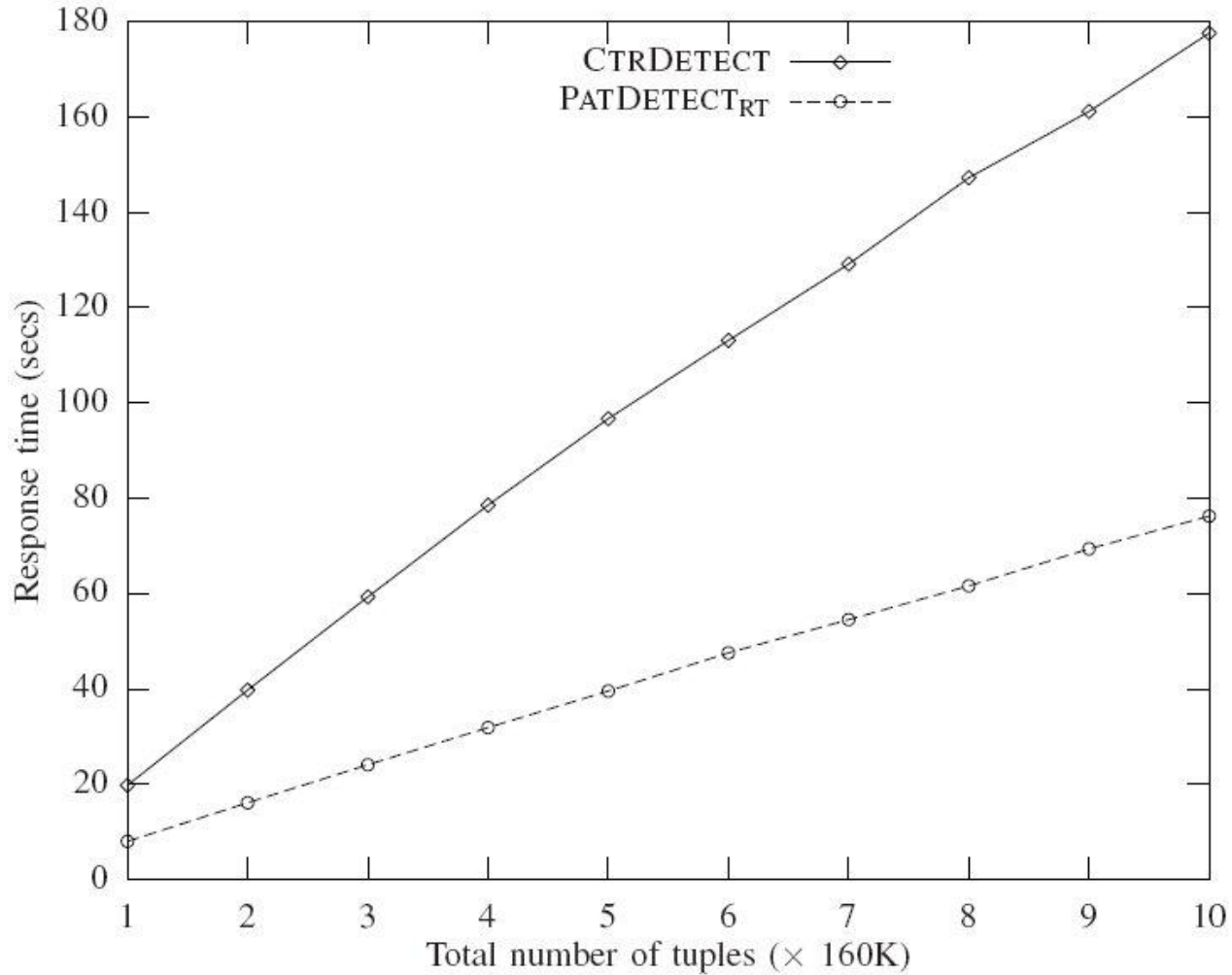
Sequentially executing one of the previous algorithms.

ClustDetect:

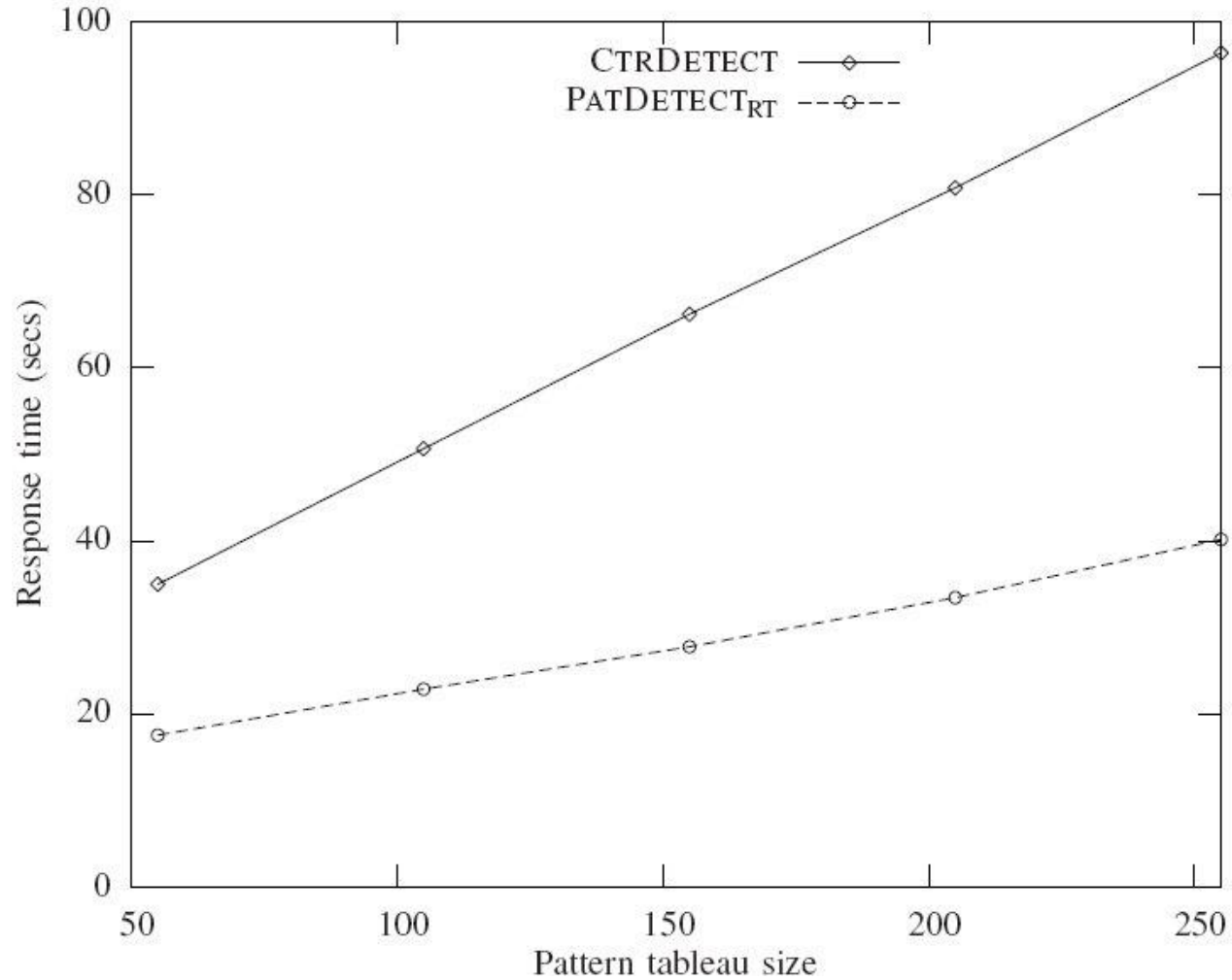
Reduce unnecessary data shipment by leveraging structure of **embedded FDs**.

- **Merge** CFDs that overlap on LHS.
- **Sequentially** validate merged CFDs.

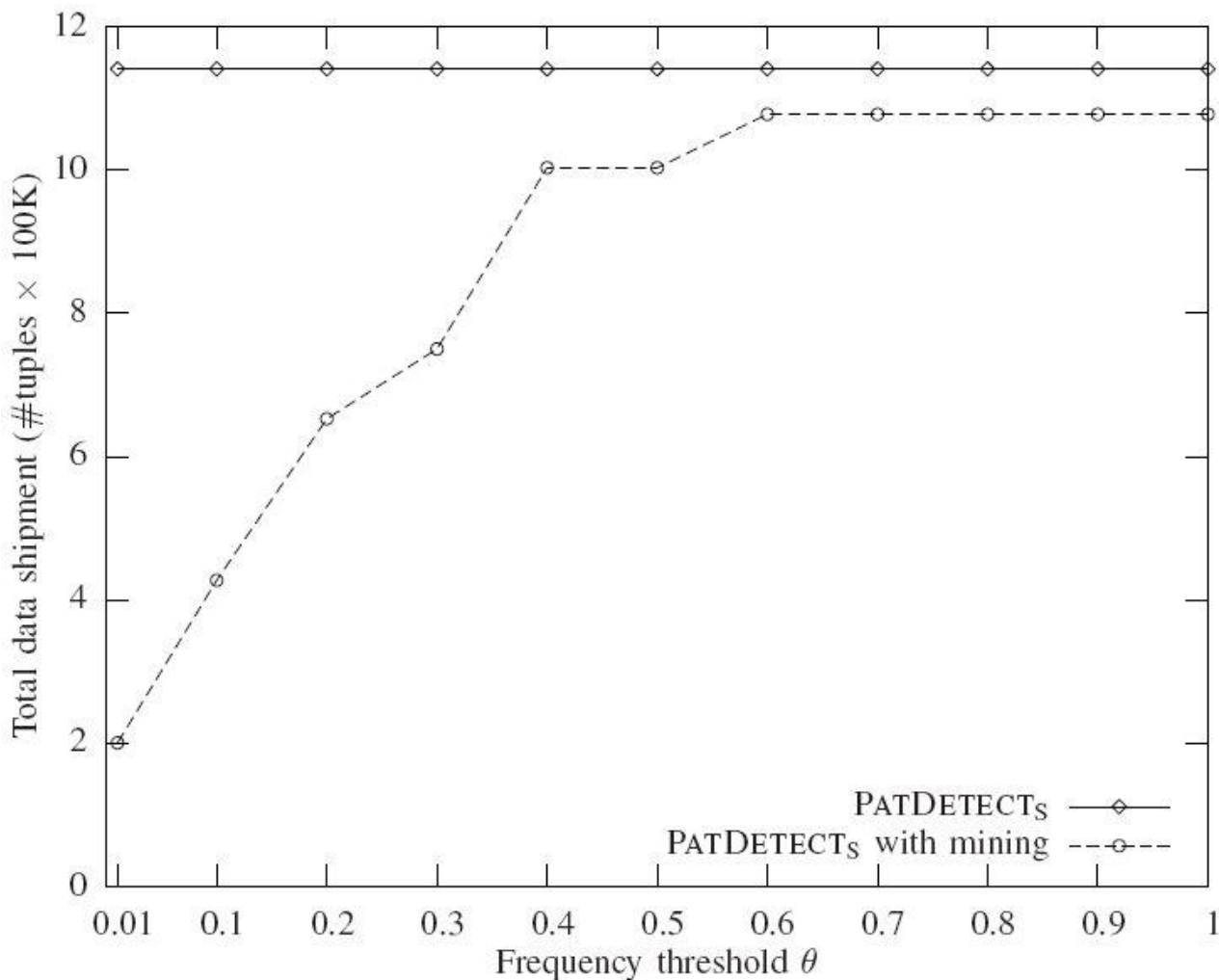
Scalability with $|D|$



Scalability with $|Tp|$



Impact of mining on shipment



Conclusion

Validating CFDs over distributed data requires data shipment.

Reduce shipment and response time by leveraging **pattern tableau** and structure of **embedded FD**.

Outlook

Develop validation algorithms for vertically (and horizontally) partitioned data.

Agenda

- ✓ **Inconsistencies in Distributed Data**
- ✓ **Problem Statement & Complexity**
- ✓ **Detection Algorithms**
- ✓ **Experimental Results**
- ✓ **Conclusion & Outlook**