

DART: a Data Acquisition and Repairing Tool

Bettina Fazzinga, Sergio Flesca, Filippo Furfaro and **Francesco Parisi**

D.E.I.S.
Università della Calabria

{bfazzinga, flesca, furfaro, fparisi}@deis.unical.it

Motivation

- Error-free acquisition of data is mandatory in several application scenarios
 - balance sheet analysis

BALANCE SHEET				
	Actual	1993	APR	MAY
ASSETS				
CURRENT ASSETS				
Cash and cash equivalents	\$401,000	\$12,500	\$120,518	
Accounts Receivable	250,000	\$12,328	1,082,575	
Inventory	400,000	900,398	1,085,899	
Other current assets	50,000	0	0	
Total current assets	1,101,000	1,025,226	2,289,192	
FIXED ASSETS				
Land	100,000	115,000	115,000	
Buildings	1,500,000	1,500,000	1,500,000	
Equipment	800,000	810,000	819,000	
Less-accumulated depreciation	2,400,000	2,474,000	2,474,000	
Total fixed assets	2,800,000	2,189,499	2,119,000	
INTANGIBLE ASSETS				
Cost	50,000	55,000	55,000	
Less-accumulated amortization	20,000	20,429	20,889	
Total intangible assets	30,000	34,571	34,112	
OTHER ASSETS	0	0	0	
Total Assets	\$3,266,000	\$4,078,902	\$4,454,816	
LIABILITIES AND STOCKHOLDERS' EQUITY				
CURRENT LIABILITIES				
Accounts payable	\$800,000	\$295,900	\$345,205	
Notes payable	100,000	700,000	730,000	
Current portion of long-term debt	100,000	100,000	100,000	
Income taxes	30,000	137,375	241,557	
Accrued expenses	50,000	111,380	121,200	
Other current liabilities	18,000	18,000	18,000	
Total current liabilities	800,000	1,354,645	1,586,002	
NON-CURRENT LIABILITIES				
Long-term debt	800,000	800,000	680,000	
Deferred income	100,000	100,000	100,000	
Deferred income taxes	30,000	30,000	30,000	
Other long-term liabilities	50,000	50,000	0	
Total liabilities	1,780,000	2,134,645	2,296,002	
STOCKHOLDERS' EQUITY				
Capital stock issued	100,000	125,000	125,000	
Additional paid in capital	50,000	75,000	75,000	
Retained earnings	1,436,000	1,754,317	1,988,814	
Total Liabilities and Equity	\$3,266,000	\$4,078,902	\$4,454,816	
[*] CORPORATION (Y/N)				
CASH BALANCE POSITIVE (NEGATIVE)	Y	POSITIVE	POSITIVE	POSITIVE
BALANCE SHEET Amount Out Of Balance		0	0	0
CASH FLOW STATEMENTS		0	0	0

electronic doc

Balance sheet analysis tool

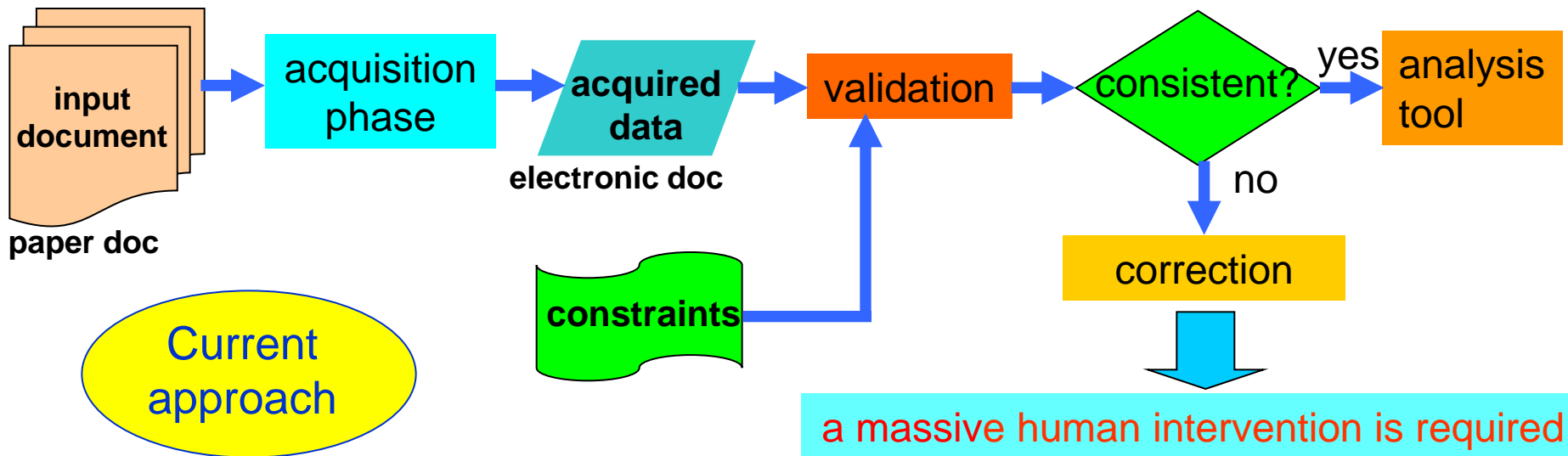


analysis report

- generally balance sheets are available as paper documents, thus they cannot be processed by balance analysis tools, since these work only on electronic data

Motivation

- Error-free acquisition of data is mandatory in several application scenarios
 - balance sheet analysis



- currently, **integrity constraints** defined on the input data are exploited **only for validating** acquired data
- if **data** are **inconsistent** all the document portions involved into unsatisfied constraint **must be checked for locating and correcting errors**

Motivation

- For instance

OCR tool

source document					acquired document	
cash sales	100	100 +		100	✓	
receivables	120	120 =		120	✓	
total cash receipts	220	220	220 -	250	✓	
payment of accounts	120	120 +		120	✓	
long-term financing	40	40 =		40	✓	
total disbursements	160	160	160 =	160	✓	
net cash inflow	60		60	60	✓	

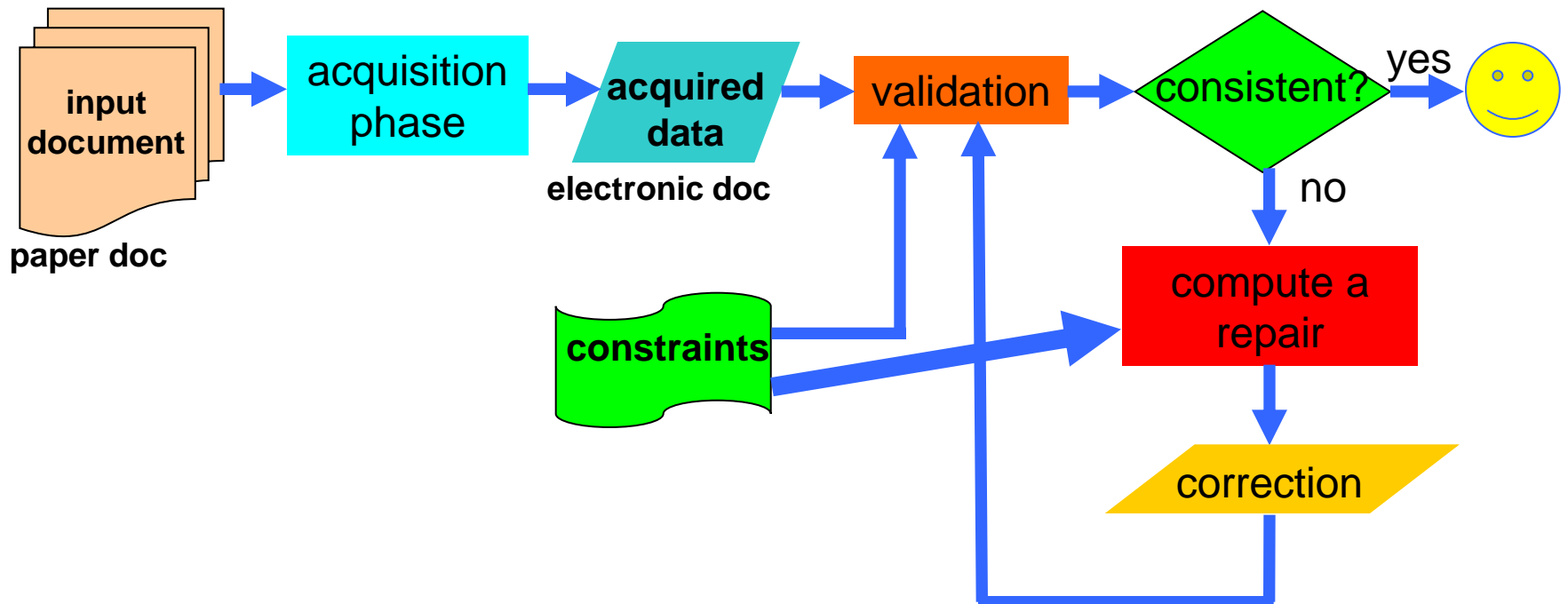
a massive human intervention is required for correcting errors

- constraints like those defined in the context of balance-sheet data can be express by **aggregate constraints**

Key Idea



exploit integrity constraints for **suggesting corrections**



the **human intervention** will be **limited** to verify only located suggestions

Key Idea

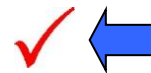


exploit integrity constraints for **suggesting corrections**

- For instance

acquired document

cash sales	100
receivables	120
total cash receipts	250
payment of accounts	120
long-term financing	40
total disbursement	160
net cash inflow	60



DART suggests decreasing the value down to 220

- in this case the operator will have to **verify a single value** instead of all the values in the table

Outline

- Repairing strategies
- DART architecture
- Aggregate constraints
- Steady aggregate constraints (SAC)
- Computing a card-minimal repair

Repairing strategy

- What is a **reasonable strategy** for repairing the acquired data?

~~Tuple deletion / insertion~~

The inconsistent cash budget

Receipts	cash sales	100
	receivables	120
	<i>total cash</i>	250

$$100 + 120 \neq 250$$

The **repaired** cash budget

Receipts	cash sales	100
	receivables	120
	XXXXXX	30
	<i>total cash</i>	250

$$\begin{array}{r} 100 + \\ 120 + \\ \underline{30} = \\ 250 \end{array}$$

Adding a new tuple means that the OCR tool skipped a whole row when acquiring ... **It's rather unrealistic!!!**

Repairing strategy

- What is a **reasonable strategy** for repairing the acquired data?
- The most natural approach is **updating directly the numerical data**
 - Work at attribute-level, rather than tuple-level

The inconsistent cash budget

Receipts	cash sales	100
	receivables	120
	<i>total cash</i>	250

$$100 + 120 \neq 250$$

The repaired cash budget

Receipts	cash sales	100
	receivables	120
	<i>total cash</i>	220

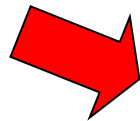
$$\begin{array}{r} 100 + \\ 120 = \\ \hline 220 \end{array}$$

- In our context, we can reasonably assume that **inconsistencies are due to symbol recognition errors**
- Thus, trying to re-construct the actual data values (without changing the number of tuples) is well founded

Card-minimal semantics

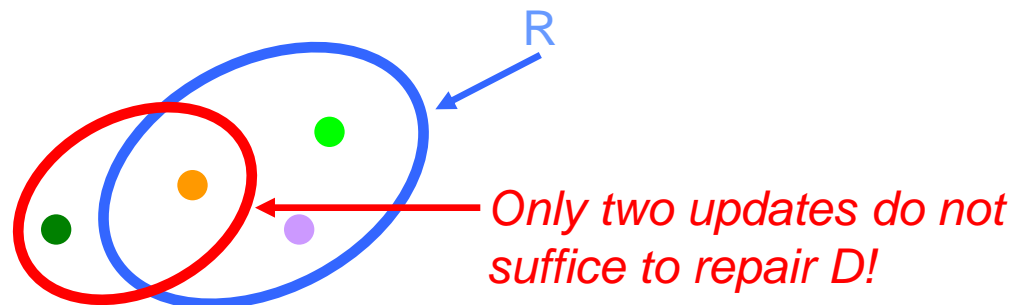
The most probable case is that the acquiring system made the **minimum number of errors**

Card-minimal semantics



It means assuming that the minimum number of errors occurred

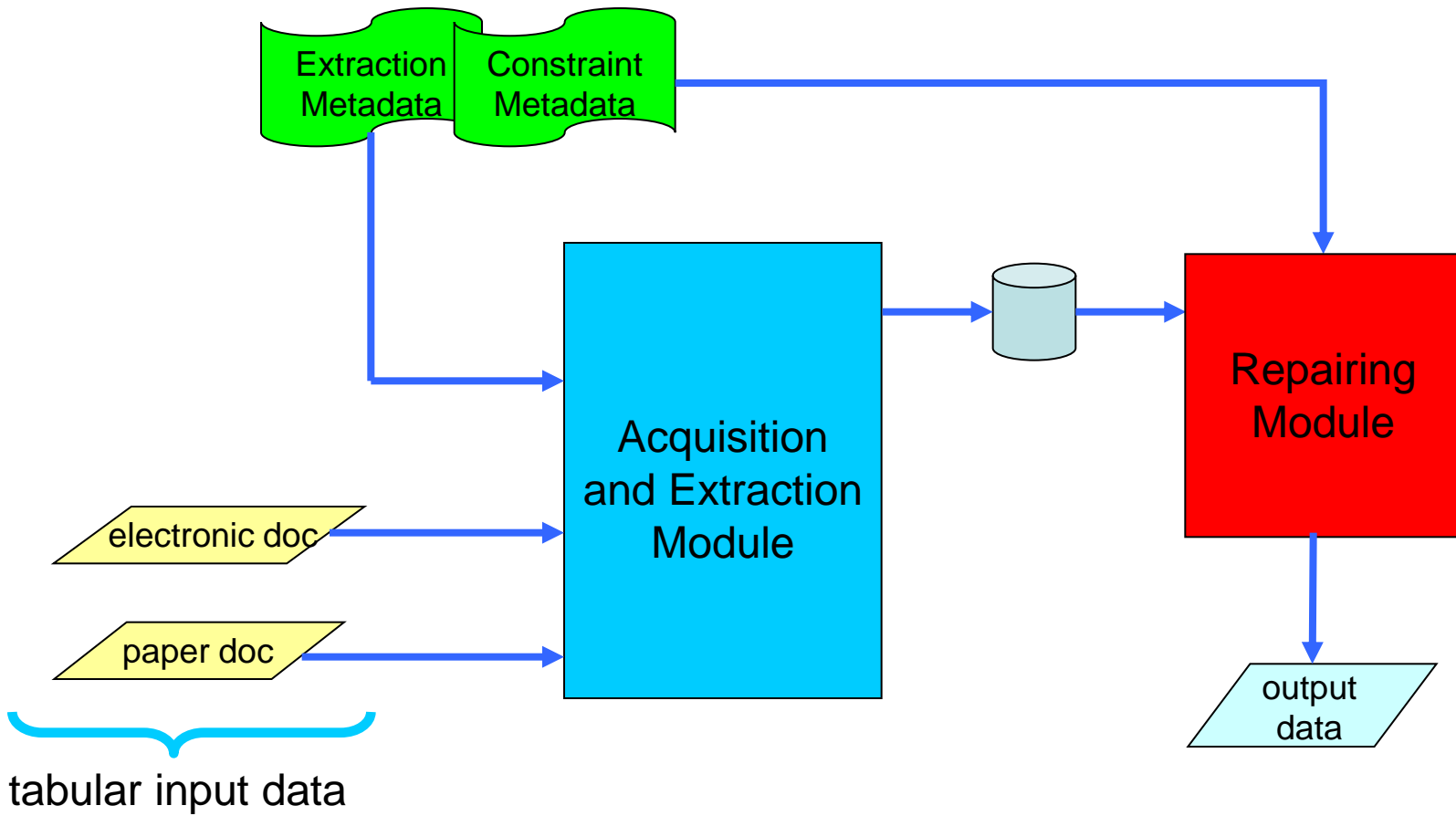
A repair R is **card-minimal** for D iff there is no repair R' for D consisting of fewer updates than R



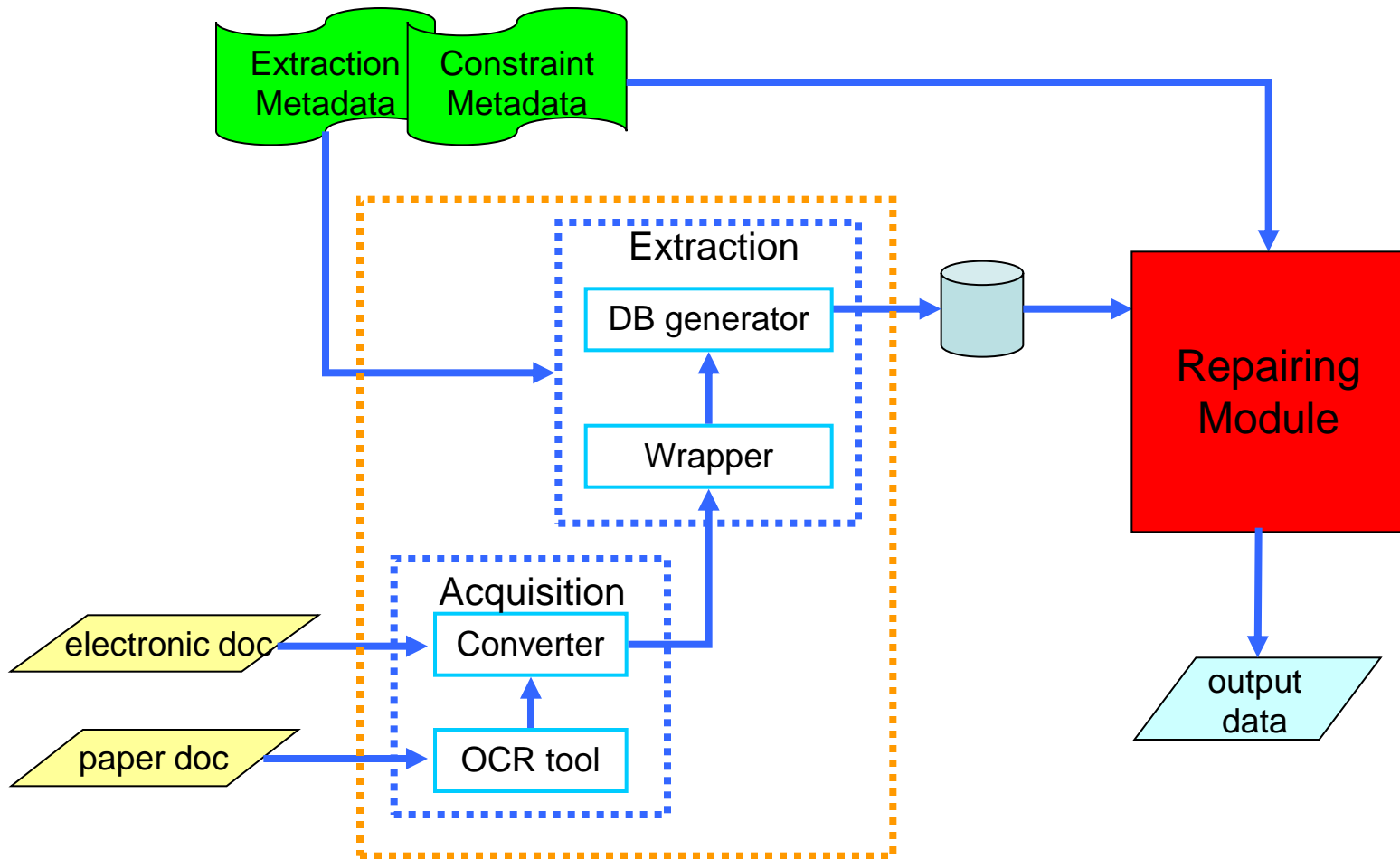
Outline

- Repairing strategies
- **DART architecture**
- Aggregate constraints
- Steady aggregate constraints (SAC)
- Computing a card-minimal repair

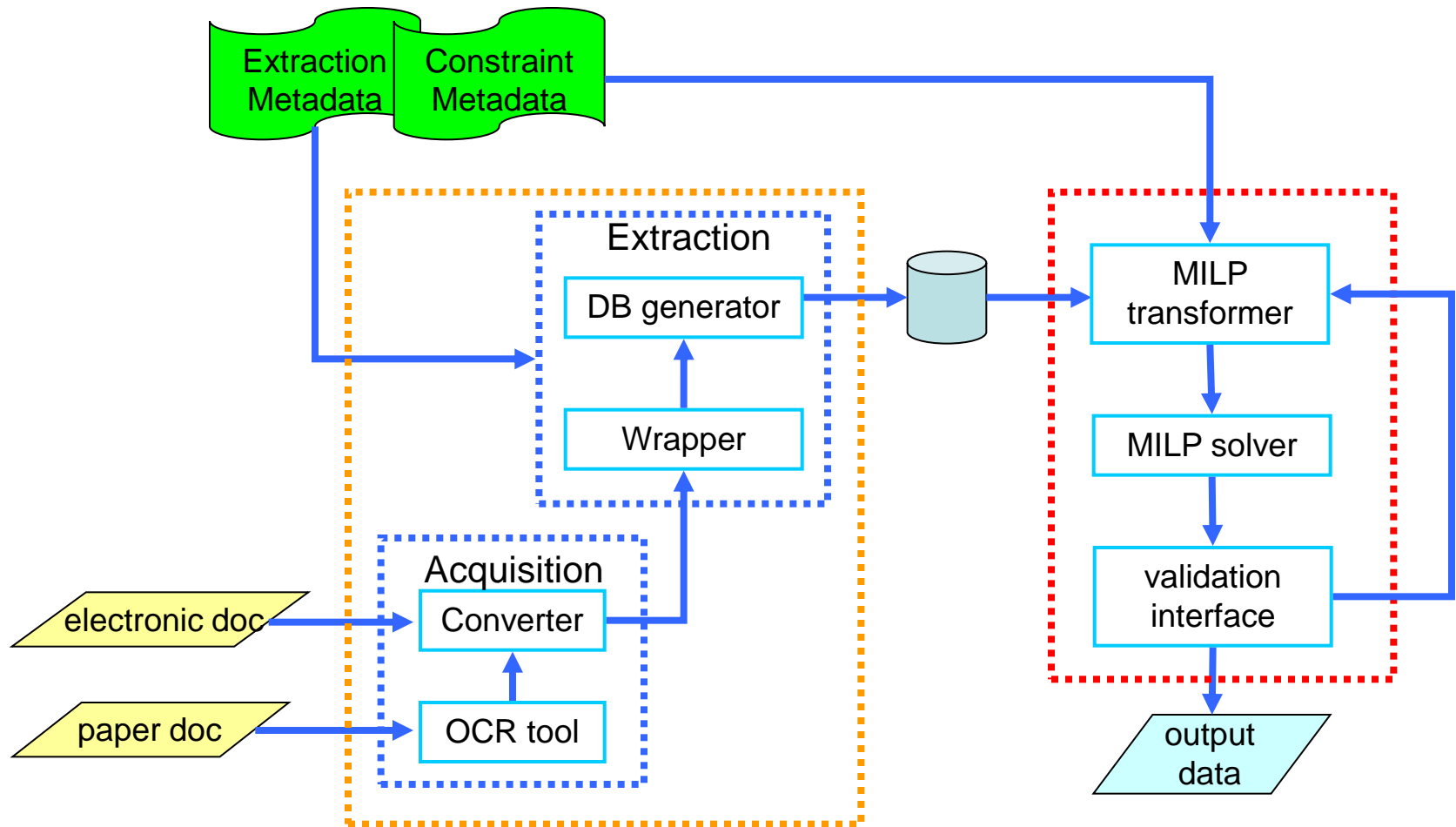
DART architecture



DART architecture - Acquisition and Extraction Module



DART architecture - Repairing Module



Outline

- Repairing strategies
- DART architecture
- **Aggregate constraints**
- Steady aggregate constraints (SAC)
- Computing a card-minimal repair

Aggregate constraints: the application context

- A cash budget for a firm:

Year	2004
Receipts	
beginning cash	20
cash sales	100
receivables	120
total cash receipts	220
Disbursements	
payment of accounts	120
capital expenditure	0
long-term financing	40
total disbursements	160
Balance	
net cash inflow	60
ending cash balance	80

Sections

Subsections

aggregate items

are obtained by aggregating
detail items of the same
section

Aggregate constraints: the application context

- A cash budget for a firm:

Year	2004
Receipts	
beginning cash	20
cash sales	100
receivables	120
total cash receipts	220
Disbursements	
payment of accounts	120
capital expenditure	0
long-term financing	40
total disbursements	160
Balance	
net cash inflow	60
ending cash balance	80

Sections

Subsections

derived items

are obtained using the value of other item of any type and belonging to any section

Aggregate constraints: the application context

- A cash budget satisfy some integrity constraints:

Year	2004	
Receipts		
beginning cash	20	
cash sales	100	100 +
receivables	120	120 =
total cash receipts	220	<u>220</u>
Disbursements		
payment of accounts	120	120 +
capital expenditure	0	0 +
long-term financing	40	40 =
total disbursements	160	<u>160</u>
Balance		
net cash inflow	60	
ending cash balance	80	

1)

for each section, the sum of all detail items must be equal to the value of the aggregate item

Aggregate constraints: the application context

- A cash budget satisfy some integrity constraints:

Year	2004
Receipts	
beginning cash	20
cash sales	100
receivables	120
total cash receipts	220
Disbursements	
payment of accounts	120
capital expenditure	0
long-term financing	40
total disbursements	160
Balance	
net cash inflow	60
ending cash balance	80

2)

the net cash inflow must be equal to the difference between total cash receipts and total disbursements

$$\begin{array}{r} 220 - \\ 160 = \\ \hline 60 \end{array}$$

From the paper document to its digitized version

Year	2004
Receipts	
beginning cash	20
cash sales	100
receivables	120
total cash receipts	220
Disbursements	
payment of accounts	120
capital expenditure	0
long-term financing	40
total disbursements	160
Balance	
net cash inflow	60
ending cash balance	80

CashBudget

Section	Subsection	Type	Value
Receipts	beginning cash	<i>drv</i>	20
Receipts	cash sales	<i>det</i>	100
Receipts	receivables	<i>det</i>	120
Receipts	total cash receipts	<i>aggr</i>	250
Disbursements	payment of accounts	<i>det</i>	120
Disbursements	capital expenditure	<i>det</i>	0
Disbursements	long-term financing	<i>det</i>	40
Disbursements	total disbursements	<i>aggr</i>	160
Balance	net cash inflow	<i>drv</i>	60
Balance	ending cash balance	<i>drv</i>	80



Acquisition and Extraction Module

Aggregate constraints

- can express constraints like those defined in the context of balance-sheet data

$$\forall \bar{X} \left(\phi(\bar{X}) \implies \text{Aggr } F(\bar{X}) \leq K \right)$$

where:

1. $\phi(\bar{X})$ is a conjunction of atoms
2. K is a constant
3. The aggregation formula $\text{Aggr } F(\bar{X})$ is the linear combination of aggregation functions

with $\bar{X}_i \subseteq \bar{X}$

$$\sum_{i=1}^n c_i \cdot \chi_i(\bar{X}_i)$$

Aggregation function

- Relational scheme $R(A_1, A_2, \dots, A_n)$
 - **Measure attributes:** numerical attributes representing measures
 - Such as weight, length, price, etc.

Linear combination of attributes

- Aggregation function

$\chi(x_1, \dots, x_k) = \text{SELECT } \textit{sum}(e)$
 $\text{FROM } R$
 $\text{WHERE } \alpha(x_1, \dots, x_k)$

Boolean formula on constants and attributes of R

Aggregate constraints

- CashBudget(Section, Subsection, Type, Value)

Section	Subsection	Type	Value
Receipts	beginning cash	<i>drv</i>	20
Receipts	cash sales	<i>det</i>	100
Receipts	receivables	<i>det</i>	120
Receipts	total cash receipts	<i>aggr</i>	250
Disbursements	payment of accounts	<i>det</i>	120
Disbursements	capital expenditure	<i>det</i>	0
Disbursements	long-term financing	<i>det</i>	40
Disbursements	total disbursements	<i>aggr</i>	160
Balance	net cash inflow	<i>drv</i>	60
Balance	ending cash balance	<i>drv</i>	80

1)

for each **section**, the sum of all **detail items** must be equal to the value of the **aggregate item**

Aggregation function:

$$\chi_1(s, t) = \text{SELECT } \text{sum}(\text{Value})$$

$$\text{FROM } \text{CashBudget}$$

$$\text{WHERE } \text{Section} = s$$

$$\text{AND } \text{Type} = t$$

Aggregate constraint:

$$\text{CashBudget}(s, -, -, -) \implies \chi_1(s, \text{det}) - \chi_1(s, \text{aggr}) = 0$$

Aggregate constraints

- CashBudget(Section, Subsection, Type, Value)

Section	Subsection	Type	Value
Receipts	beginning cash	<i>drv</i>	20
Receipts	cash sales	<i>det</i>	100
Receipts	receivables	<i>det</i>	120
Receipts	total cash receipts	<i>aggr</i>	250
Disbursements	payment of accounts	<i>det</i>	120
Disbursements	capital expenditure	<i>det</i>	0
Disbursements	long-term financing	<i>det</i>	40
Disbursements	total disbursements	<i>aggr</i>	160
Balance	net cash inflow	<i>drv</i>	60
Balance	ending cash balance	<i>drv</i>	80

2)

the net cash inflow must be equal to the difference between total cash receipts and total disbursements

Aggregation function:

```
 $\chi_2(ss) = SELECT \ sum(Value)
FROM \ CashBudget
WHERE \ Subsection = ss$ 
```

Aggregation constraint:

$$CashBudget(., ., ., .) \implies \chi_2(net\ cash\ in\ flow) - [\chi_2(total\ cash\ receipts) - \chi_2(total\ disbursements)] = 0$$

Outline

- Repairing strategies
- DART architecture
- Aggregate constraints
- **Steady aggregate constraints (SACs)**
- Computing a card-minimal repair

Steady aggregate constraints (SACs)

- a restricted form of aggregate constraints
- computing a card-minimal repair w.r.t. a set of SAC can be accomplished by solving an instance of MILP problem

CashBudget

Section	Subsection	Type	Value
Receipts	beginning cash	<i>drv</i>	20
Receipts	cash sales	<i>det</i>	100
Receipts	receivables	<i>det</i>	120
Receipts	total cash receipts	<i>aggr</i>	250
Disbursements	payment of accounts	<i>det</i>	120
Disbursements	capital expenditure	<i>det</i>	0
Disbursements	long-term financing	<i>det</i>	40
Disbursements	total disbursements	<i>aggr</i>	160
Balance	net cash inflow	<i>drv</i>	60
Balance	ending cash balance	<i>drv</i>	80

$\rightarrow Z_1$
 $\rightarrow Z_2$
 $\rightarrow Z_3$
 Z_4
 Z_5
 Z_6
 $\rightarrow Z_7$

$$\begin{cases} Z_1 + Z_2 = Z_3 \\ Z_4 + Z_5 + Z_6 = Z_7 \end{cases}$$

a system of inequalities can be associated if values “involved” in the constraints are independent on repairs

$$CashBudget(s, -, -, -) \implies \chi_1(s, det) - \chi_1(s, aggr) = 0$$

Steady aggregate constraints (SACs)

An aggregate constraint is an SAC if:

- 1) no attributes in the WHERE clause are measure attributes
 - 2) no attributes corresponding to variables in the WHERE clause are measure attributes
 - 3) no attributes corresponding to variables shared by two atoms are measure attributes
- CashBudget(Section, Subsection, Type, Value)

$$\text{CashBudget}(s, ss, t, v) \implies \chi_1(s, \text{det}) - \chi_1(s, \text{aggr}) = 0$$

where: $\chi_1(s, t) = \text{SELECT } \text{sum}(\text{Value})$
 $\text{FROM } \text{CashBudget}$
 $\text{WHERE } \text{Section} = s$
 $\text{AND } \text{Type} = t$

Steady aggregate constraints (SACs)

An aggregate constraint is an SAC if:

- 1) no attributes in the WHERE clause are **measure attributes**
- 2) no **attributes corresponding to variables in the WHERE clause** are **measure attributes**
- 3) no attributes corresponding to variables shared by two atoms are **measure attributes**

- CashBudget(**Section**, Subsection, **Type**, **Value**)

$$\text{CashBudget}(s, ss, t, v) \implies \chi_1(s, det) - \chi_1(s, aggr) = 0$$

where: $\chi_1(s, t) = \text{SELECT } \text{sum}(\text{Value})$
 $\text{FROM } \text{CashBudget}$
 $\text{WHERE } \text{Section} = s$
 $\text{AND } \text{Type} = t$

Steady aggregate constraints (SACs)

An aggregate constraint is an SAC if:

- 1) no attributes in the WHERE clause are **measure attributes**
 - 2) no attributes corresponding to variables in the WHERE clause are **measure attributes**
 - 3) no **attributes corresponding to variables shared by two atoms** are **measure attributes**
- CashBudget(Section, Subsection, Type, **Value**)

$$\text{CashBudget}(s, ss, t, v) \implies \chi_1(s, det) - \chi_1(s, aggr) = 0$$

where: $\chi_1(s, t) = \text{SELECT } \text{sum}(\text{Value})$
 $\text{FROM } \text{CashBudget}$
 $\text{WHERE } \text{Section} = s$
 $\text{AND } \text{Type} = t$

Complexity results under SACs

- even if SACs are a restricted form of (general) aggregate constraints, **results obtained for (general) aggregate constraints are still valid for SACs**
- the **repair existence problem**
 - deciding whether there is a repair for a database violating a given set of SACs is **NP-complete**
- the **minimal repair checking problem**
 - deciding whether a repair is minimal in **CoNP-complete**
- the **consistent query answer problem**
 - deciding whether a query is true in every card-minimal repair is $\Delta_2^p[\log n]$ – *complete*

Outline

- Repairing strategies
- DART architecture
- Aggregate constraints
- Steady aggregate constraints (SAC)
- **Computing a card-minimal repair**

Repairing Module – MILP transformer

- Under **SACs** a **card-minimal repair** can be computed solving an **MILP problem** instance
 - SACs are translated into a system of inequalities $A Z \leq B$
 - $Z=[z_1, z_2, \dots, z_N]$ is a vector of variables associated to database values v_1, v_2, \dots, v_N which are involved in a constraint

Section	Subsection	Type	Value
Receipts	beginning cash	drv	20
Receipts	cash sales	det	100
Receipts	receivables	det	120
Receipts	total cash receipts	aggr	250
Disbursements	payment of accounts	det	120
Disbursements	capital expenditure	det	0
Disbursements	long-term financing	det	40
Disbursements	total disbursements	aggr	160
Balance	net cash inflow	drv	60
Balance	ending cash balance	drv	80

Z_1
 Z_2
 Z_3
 Z_4
 Z_5
 Z_6
 Z_7

1) $\begin{cases} Z_1 + Z_2 = Z_3 \\ Z_4 + Z_5 + Z_6 = Z_7 \end{cases}$

1) $CashBudget(s, -, -, -) \implies \chi_1(s, det) - \chi_1(s, aggr) = 0$

Repairing Module – MILP transformer

- Under **SACs** a **card-minimal repair** can be computed solving an **MILP problem** instance
 - SACs are translated into a system of inequalities $A Z \leq B$
 - $Z=[z_1, z_2, \dots, z_N]$ is a vector of variables associated to database values v_1, v_2, \dots, v_N which are involved in a constraint

Section	Subsection	Type	Value
Receipts	beginning cash	drv	20
Receipts	cash sales	det	100
Receipts	receivables	det	120
Receipts	total cash receipts	aggr	250
Disbursements	payment of accounts	det	120
Disbursements	capital expenditure	det	0
Disbursements	long-term financing	det	40
Disbursements	total disbursements	aggr	160
Balance	net cash inflow	drv	60
Balance	ending cash balance	drv	80

$\rightarrow Z_1$
 $\rightarrow Z_2$
 $\rightarrow Z_3$
 Z_4
 Z_5
 Z_6
 Z_7
 $\rightarrow Z_8$

1) $\begin{cases} Z_1 + Z_2 = Z_3 \\ Z_4 + Z_5 + Z_6 = Z_7 \end{cases}$

2) $Z_3 - Z_7 = Z_8$

2) $CashBudget(-, -, -, -) \implies \chi_2(\text{net cash inflow}) - [\chi_2(\text{total cash receipts}) - \chi_2(\text{total disbursements})] = 0$

Repairing Module – MILP transformer

- Under SACs a card-minimal repair can be computed solving an MILP problem instance
 - SACs are translated into a system of inequalities $A Z \leq B$
 - $Z=[z_1, z_2, \dots, z_N]$ is a vector of variables associated to database values v_1, v_2, \dots, v_N which are involved in a constraint

Section	Subsection	Type	Value
Receipts	beginning cash	<i>drv</i>	20
Receipts	cash sales	<i>det</i>	100
Receipts	receivables	<i>det</i>	120
Receipts	total cash receipts	<i>aggr</i>	250
Disbursements	payment of accounts	<i>det</i>	120
Disbursements	capital expenditure	<i>det</i>	0
Disbursements	long-term financing	<i>det</i>	40
Disbursements	total disbursements	<i>aggr</i>	160
Balance	net cash inflow	<i>drv</i>	60
Balance	ending cash balance	<i>drv</i>	80

$\rightarrow Z_1$
 $\rightarrow Z_2$
 $\rightarrow Z_3$
 Z_4
 Z_5
 Z_6
 Z_7
 $\rightarrow Z_8$

$$\begin{cases}
 Z_1 + Z_2 = Z_3 \\
 Z_4 + Z_5 + Z_6 = Z_7 \\
 Z_3 - Z_7 = Z_8
 \end{cases}$$

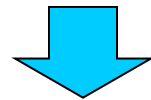
Repairing Module – MILP transformer

- Under SACs a card-minimal repair can be computed solving an MILP problem instance
 - SACs are translated into a system of inequalities $A Z \leq B$
 - $Z=[z_1, z_2, \dots, z_N]$ is a vector of variables associated to database values v_1, v_2, \dots, v_N which are involved in a constraint

Section	Subsection	Type	Value
Receipts	beginning cash	<i>drv</i>	20
Receipts	cash sales	<i>det</i>	100
Receipts	receivables	<i>det</i>	120
Receipts	total cash receipts	<i>aggr</i>	250
Disbursements	payment of accounts	<i>det</i>	120
Disbursements	capital expenditure	<i>det</i>	0
Disbursements	long-term financing	<i>det</i>	40
Disbursements	total disbursements	<i>aggr</i>	160
Balance	net cash inflow	<i>drv</i>	60
Balance	ending cash balance	<i>drv</i>	80

$z_1=130$
 $z_2=120$
 $z_3=250$
 $z_4=120$
 $z_5=0$
 $z_6=40$
 $z_7=160$
 $z_8=90$

$$\begin{cases}
 z_1 + z_2 = z_3 \\
 z_4 + z_5 + z_6 = z_7 \\
 z_3 - z_7 = z_8
 \end{cases}$$



each **solution**
 corresponds to a
 (possible not minimal)
repair

Repairing Module – MILP transformer

- In order to decide whether a solution corresponds to a **card-minimal repair**

– we define a variable $y_i = z_i \cdot v_i$

Section	Subsection	Type	Value
Receipts	beginning cash	<i>drv</i>	20
Receipts	cash sales	<i>det</i>	100
Receipts	receivables	<i>det</i>	120
Receipts	total cash receipts	<i>aggr</i>	250
Disbursements	payment of accounts	<i>det</i>	120
Disbursements	capital expenditure	<i>det</i>	0
Disbursements	long-term financing	<i>det</i>	40
Disbursements	total disbursements	<i>aggr</i>	160
Balance	net cash inflow	<i>drv</i>	60
Balance	ending cash balance	<i>drv</i>	80

→ z_1

→ z_2

→ z_3

z_4

z_5

z_6

z_7

→ z_8

$$\begin{cases} z_1 + z_2 = z_3 \\ z_4 + z_5 + z_6 = z_7 \\ z_3 - z_7 = z_8 \end{cases}$$

Repairing Module – MILP transformer

- In order to decide whether a solution corresponds to a **card-minimal repair**

– we define a variable $y_i = z_i - v_i$

Section	Subsection	Type	Value
Receipts	beginning cash	<i>drv</i>	20
Receipts	cash sales	<i>det</i>	100
Receipts	receivables	<i>det</i>	120
Receipts	total cash receipts	<i>aggr</i>	250
Disbursements	payment of accounts	<i>det</i>	120
Disbursements	capital expenditure	<i>det</i>	0
Disbursements	long-term financing	<i>det</i>	40
Disbursements	total disbursements	<i>aggr</i>	160
Balance	net cash inflow	<i>drv</i>	60
Balance	ending cash balance	<i>drv</i>	80

$\rightarrow z_1$
 $\rightarrow z_2$
 $\rightarrow z_3$
 z_4
 z_5
 z_6
 z_7
 $\rightarrow z_8$

$$\begin{aligned}
 z_1 + z_2 &= z_3 \\
 z_4 + z_5 + z_6 &= z_7 \\
 z_3 - z_7 &= z_8 \\
 y_1 &= z_1 - 100 \\
 y_2 &= z_2 - 120 \\
 y_3 &= z_3 - 250 \\
 y_4 &= z_4 - 120 \\
 y_5 &= z_5 - 0 \\
 y_6 &= z_6 - 40 \\
 y_7 &= z_7 - 160 \\
 y_8 &= z_8 - 60
 \end{aligned}$$

Repairing Module – MILP transformer

- In order to decide whether a solution corresponds to a **card-minimal repair**

– we define a variable $y_i = z_i - v_i$

Section	Subsection	Type	Value
Receipts	beginning cash	<i>drv</i>	20
Receipts	cash sales	<i>det</i>	100
Receipts	receivables	<i>det</i>	120
Receipts	total cash receipts	<i>aggr</i>	250
Disbursements	payment of accounts	<i>det</i>	120
Disbursements	capital expenditure	<i>det</i>	0
Disbursements	long-term financing	<i>det</i>	40
Disbursements	total disbursements	<i>aggr</i>	160
Balance	net cash inflow	<i>drv</i>	60
Balance	ending cash balance	<i>drv</i>	80

→ $z_1=130$
 → $z_2=120$
 → $z_3=250$
 $z_4=120$
 $z_5=0$
 $z_6=40$
 $z_7=160$
 → $z_8=90$

$$\begin{aligned}
 z_1 + z_2 &= z_3 \\
 z_4 + z_5 + z_6 &= z_7 \\
 z_3 - z_7 &= z_8 \\
 y_1 &= z_1 - 100 \\
 y_2 &= z_2 - 120 \\
 y_3 &= z_3 - 250 \\
 y_4 &= z_4 - 120 \\
 y_5 &= z_5 - 0 \\
 y_6 &= z_6 - 40 \\
 y_7 &= z_7 - 160 \\
 y_8 &= z_8 - 60
 \end{aligned}$$

Repairing Module – MILP transformer

- In order to decide whether a solution corresponds to a **card-minimal repair**

– we define a variable $y_i = z_i - v_i$

Section	Subsection	Type	Value
Receipts	beginning cash	<i>drv</i>	20
Receipts	cash sales	<i>det</i>	100
Receipts	receivables	<i>det</i>	120
Receipts	total cash receipts	<i>aggr</i>	250
Disbursements	payment of accounts	<i>det</i>	120
Disbursements	capital expenditure	<i>det</i>	0
Disbursements	long-term financing	<i>det</i>	40
Disbursements	total disbursements	<i>aggr</i>	160
Balance	net cash inflow	<i>drv</i>	60
Balance	ending cash balance	<i>drv</i>	80

$$\begin{array}{ll}
 \rightarrow z_1=130 & y_1=30 \\
 \rightarrow z_2=120 & y_2=0 \\
 \rightarrow z_3=250 & y_3=0 \\
 & z_4=120 & y_4=0 \\
 & z_5=0 & y_5=0 \\
 & z_6=40 & y_6=0 \\
 & z_7=160 & y_7=0 \\
 \rightarrow z_8=90 & y_8=30
 \end{array}$$

$$\begin{array}{l}
 z_1 + z_2 = z_3 \\
 z_4 + z_5 + z_6 = z_7 \\
 z_3 - z_7 = z_8 \\
 y_1 = z_1 - 100 \\
 y_2 = z_2 - 120 \\
 y_3 = z_3 - 250 \\
 y_4 = z_4 - 120 \\
 y_5 = z_5 - 0 \\
 y_6 = z_6 - 40 \\
 y_7 = z_7 - 160 \\
 y_8 = z_8 - 60
 \end{array}$$

$$y_i \neq 0$$



atomic updated on database value v_i

Repairing Module – MILP transformer

- In order to decide whether a solution corresponds to a **card-minimal repair**

– we define a variable $y_i = z_i - v_i$

Section	Subsection	Type	Value
Receipts	beginning cash	<i>drv</i>	20
Receipts	cash sales	<i>det</i>	100
Receipts	receivables	<i>det</i>	120
Receipts	total cash receipts	<i>aggr</i>	250
Disbursements	payment of accounts	<i>det</i>	120
Disbursements	capital expenditure	<i>det</i>	0
Disbursements	long-term financing	<i>det</i>	40
Disbursements	total disbursements	<i>aggr</i>	160
Balance	net cash inflow	<i>drv</i>	60
Balance	ending cash balance	<i>drv</i>	80

$$\begin{aligned} \rightarrow z_1 &= 130 & y_1 &= 30 \\ \rightarrow z_2 &= 120 & y_2 &= 0 \\ \rightarrow z_3 &= 250 & y_3 &= 0 \\ z_4 &= 120 & y_4 &= 0 \\ z_5 &= 0 & y_5 &= 0 \\ z_6 &= 40 & y_6 &= 0 \\ z_7 &= 160 & y_7 &= 0 \\ \rightarrow z_8 &= 90 & y_8 &= 30 \end{aligned}$$

$$\begin{aligned} z_1 + z_2 &= z_3 \\ z_4 + z_5 + z_6 &= z_7 \\ z_3 - z_7 &= z_8 \\ y_1 &= z_1 - 100 \\ y_2 &= z_2 - 120 \\ y_3 &= z_3 - 250 \\ y_4 &= z_4 - 120 \\ y_5 &= z_5 - 0 \\ y_6 &= z_6 - 40 \\ y_7 &= z_7 - 160 \\ y_8 &= z_8 - 60 \end{aligned}$$

– we have to **count the number of variables** y_i such that $y_i \neq 0$

Repairing Module – MILP transformer

- In order to detect if a variable z_i is assigned a value different v_i , a **binary variable** δ_i is defined
- we add the following constraints entailing that $y_i \neq 0 \implies \delta_i = 1$

$y_i \leq M\delta_i$ $\implies y_i > 0$ implies $\delta_i = 1$

$-M\delta_i \leq y_i$ $\implies y_i < 0$ implies $\delta_i = 1$

If a system of equalities has a solution, it has also one where each variable takes a value in $[-M, M]$

Repairing Module – MILP transformer

- In order to detect if a variable z_i is assigned (for each M -bounded solution) a value different v_i , a **binary variable δ_i** is defined
- we add the following constraints entailing that $y_i \neq 0 \implies \delta_i = 1$

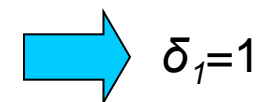
$$\begin{array}{l}
 y_i \leq M\delta_i \\
 -M\delta_i \leq
 \end{array}
 \begin{array}{l}
 \implies y_i > 0 \text{ implies } \delta_i = 1 \\
 \implies y_i < 0 \text{ implies } \delta_i = 1
 \end{array}$$

v.

$$\begin{array}{l}
 z_1 + z_2 = z_3 \\
 z_4 + z_5 + z_6 = z_7 \\
 z_3 - z_7 = z_8 \\
 y_1 = z_1 - 100 \\
 y_2 = z_2 - 120 \\
 y_3 = z_3 - 250 \\
 y_4 = z_4 - 120 \\
 y_5 = z_5 - 0 \\
 y_6 = z_6 - 40 \\
 y_7 = z_7 - 160 \\
 y_8 = z_8 - 60
 \end{array}$$

$$\begin{array}{l}
 y_1 \leq M\delta_1 \\
 -M\delta_1 \leq y_1 \\
 y_2 \leq M\delta_2 \\
 -M\delta_2 \leq y_2 \\
 \dots \\
 \dots \\
 y_8 \leq M\delta_8 \\
 -M\delta_8 \leq y_8
 \end{array}$$

$$\begin{array}{l}
 z_1=130 \\
 z_2=120 \\
 z_3=250 \\
 z_4=120 \\
 z_5=0 \\
 z_6=40 \\
 z_7=160 \\
 z_8=90
 \end{array}
 \begin{array}{l}
 y_1=30 \\
 y_2=0 \\
 y_3=0 \\
 y_4=0 \\
 y_5=0 \\
 y_6=0 \\
 y_7=0 \\
 y_8=30
 \end{array}$$



$$\delta_1 = 1$$

Repairing Module – MILP transformer

- In order to detect if a variable z_i is assigned (for each M -bounded solution) a value different v_i , a **binary variable δ_i** is defined
- we add the following constraints entailing that $y_i \neq 0 \implies \delta_i = 1$

$$\begin{array}{l}
 y_i \leq M\delta_i \\
 -M\delta_i \leq y_i
 \end{array}
 \implies
 \begin{array}{l}
 y_i > 0 \text{ implies } \delta_i = 1 \\
 y_i < 0 \text{ implies } \delta_i = 1
 \end{array}$$

v.

$$\begin{array}{l}
 z_1 + z_2 = z_3 \\
 z_4 + z_5 + z_6 = z_7 \\
 z_3 - z_7 = z_8 \\
 y_1 = z_1 - 100 \\
 y_2 = z_2 - 120 \\
 y_3 = z_3 - 250 \\
 y_4 = z_4 - 120 \\
 y_5 = z_5 - 0 \\
 y_6 = z_6 - 40 \\
 y_7 = z_7 - 160 \\
 y_8 = z_8 - 60
 \end{array}$$

$$\begin{array}{l}
 y_1 \leq M\delta_1 \\
 -M\delta_1 \leq y_1 \\
 y_2 \leq M\delta_2 \\
 -M\delta_2 \leq y_2 \\
 \dots \\
 \dots \\
 y_8 \leq M\delta_8 \\
 -M\delta_8 \leq y_8
 \end{array}$$

$$\begin{array}{ll}
 z_1=130 & v_1=30 \\
 z_2=120 & y_2=0 \\
 z_3=250 & y_3=0 \\
 z_4=120 & y_4=0 \\
 z_5=0 & y_5=0 \\
 z_6=40 & y_6=0 \\
 z_7=160 & y_7=0 \\
 z_8=90 & y_8=30
 \end{array}$$

$y_i = 0$ entails that either $\delta_i = 1$ or $\delta_i = 0$

Repairing Module – MILP transformer

- In order to consider solutions where each $\delta_i=0$ if $y_i=0$, we minimize the sum of values assigned to binary variables δ_i

$$\min \delta_1 + \delta_2 + \dots + \delta_8$$

$$z_1 + z_2 = z_3$$

$$z_4 + z_5 + z_6 = z_7$$

$$z_3 - z_7 = z_8$$

$$y_1 = z_1 - 100$$

...

$$y_8 = z_8 - 60$$

$$y_1 \leq M\delta_1$$

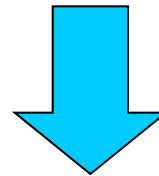
$$-M\delta_1 \leq y_1$$

...

$$y_8 \leq M\delta_8$$

$$-M\delta_8 \leq y_8$$

- any solution corresponds to an M -bounded repair having minimum cardinality w.r.t. all M -bounded repairs
- It can be shown that if a repair exists then there is a card-minimal repair that is M -bounded



any solution corresponds to a card-minimal repair

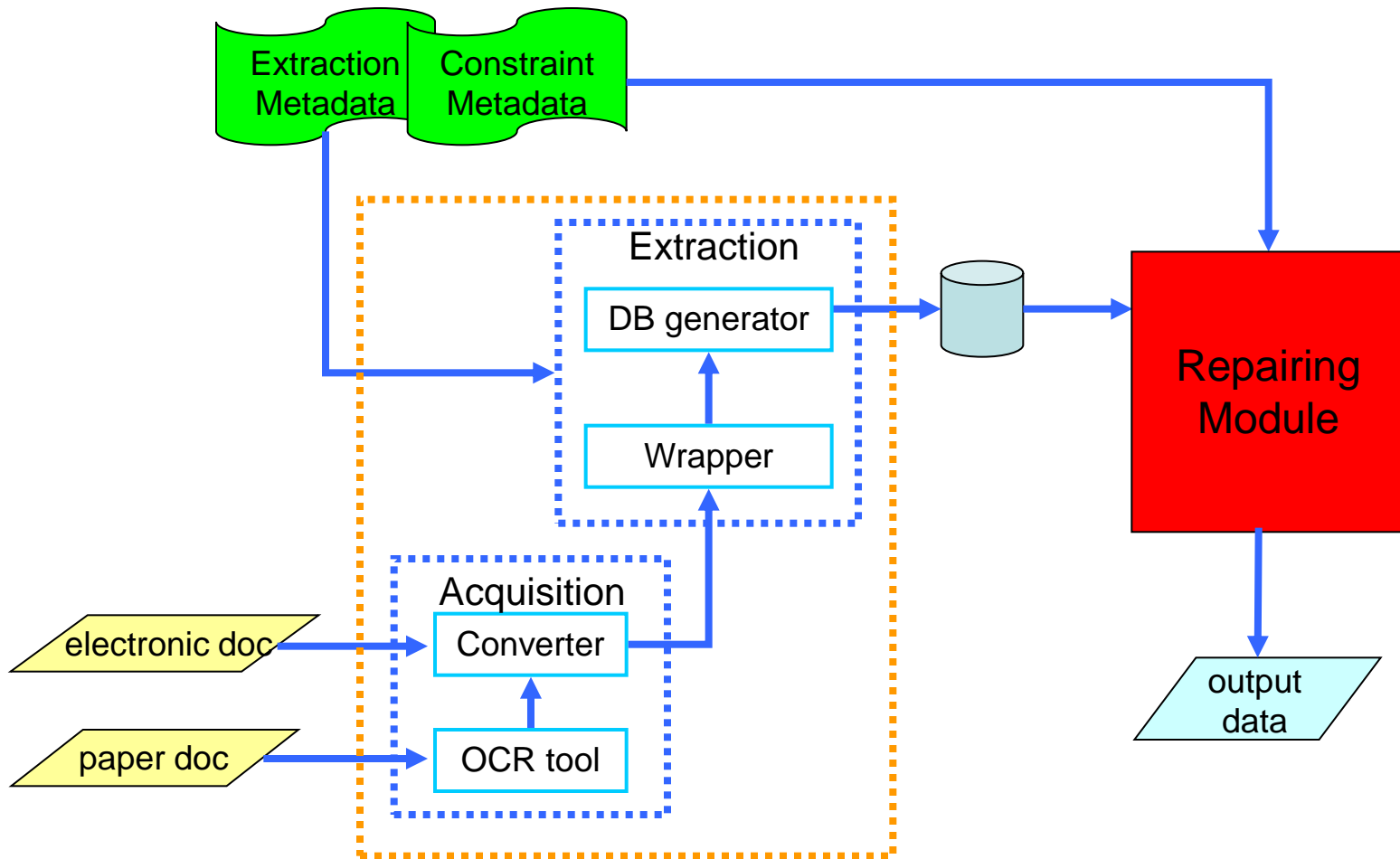
Conclusions and future work

- An architecture providing robust data acquisition facilities has been proposed
- A restricted, but useful in many real-life scenario, class of aggregate constraints has been located
- An approach for computing a card-minimal repair in presence of SACs has been provided
 - standard techniques addressing MILP problem can be re-used for computing a repair
- Experimental evaluation of the system effectiveness on large data sets (working with real databases) will be accomplished

Thank you!

...any questions?

DART architecture - Acquisition and Extraction Module

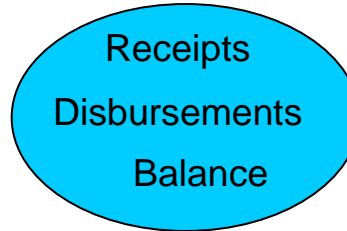


Data Extraction Sub-Module - Wrapper

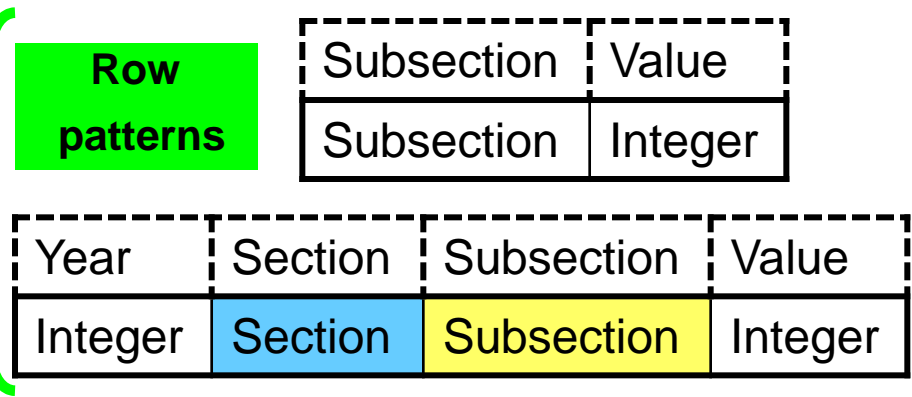
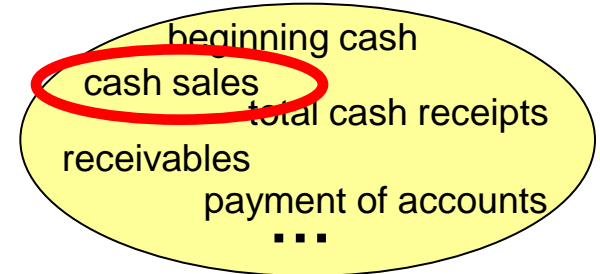
digitized document

2003	Receipts		
		beginning cash	20
		cash sales	100
		receivables	120
		total cash receipts	250
	Disbursements		
		payment of accounts	120
		capital expenditure	0
		long-term financing	40
		total disbursements	160
	Balance		
		net cash inflow	60
		ending cash balance	80

domain Section



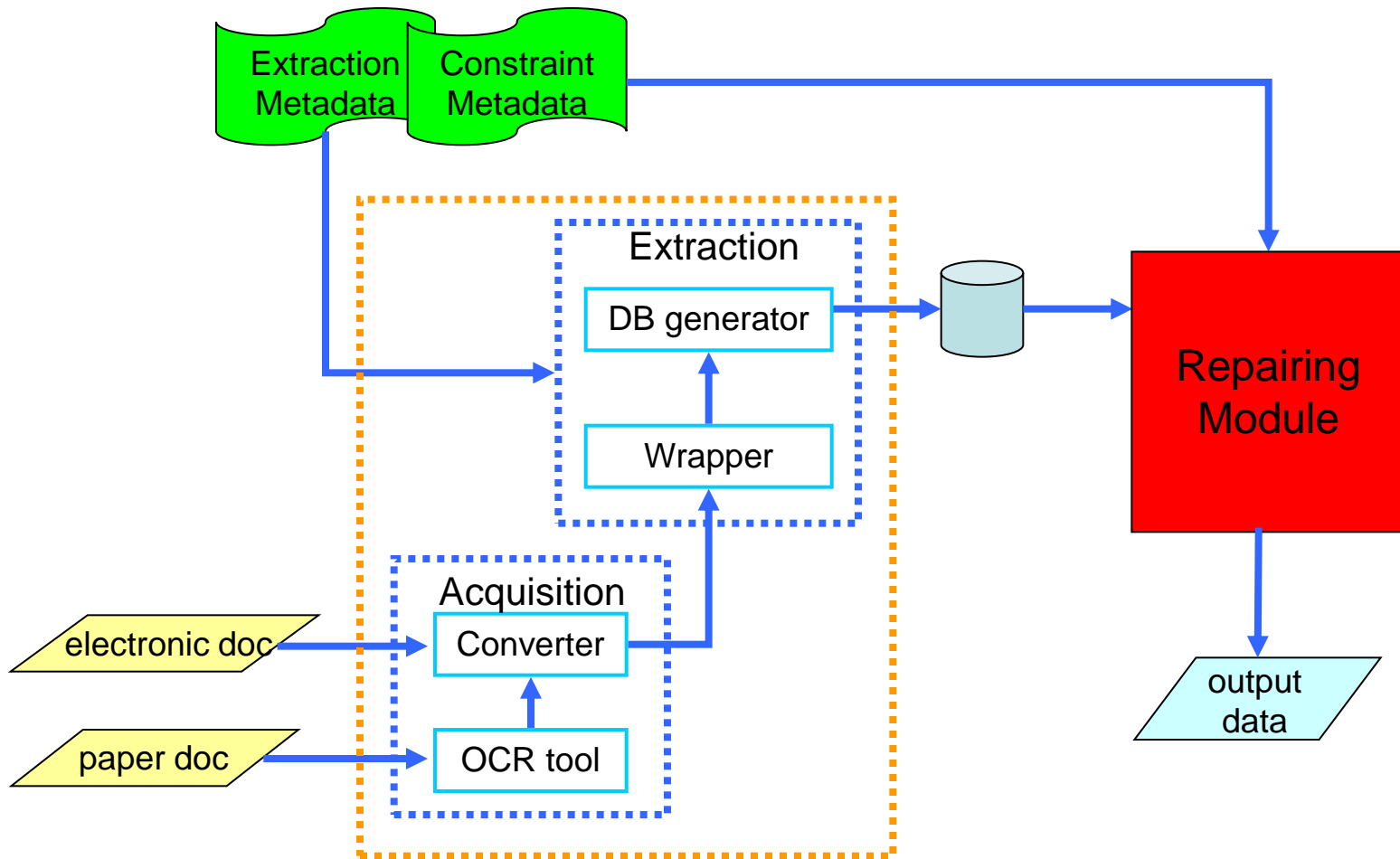
domain Subsection



Row pattern instance

2003	Receipts	cash sales	100
------	----------	------------	-----

DART architecture - Acquisition and Extraction Module



Data Extraction Sub-Module – DB generator

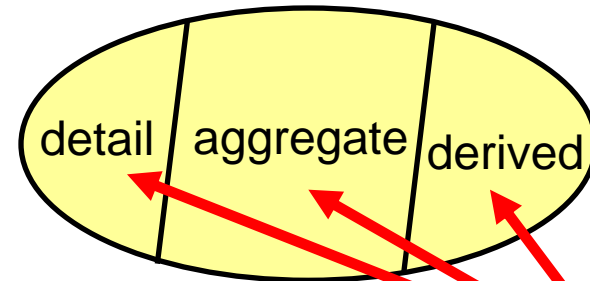
Row pattern

Year	Section	Subsection	Value
Integer	Section	Subsection	Integer

Row pattern instances

2003	Receipts	beginning cash	20
2003	Receipts	cash sales	100
2003	Receipts	receivables	120
2003	Receipts	total cash receipts	250

Subsection



CashBudget(Year, Section, Subsection, Type, Value)

CashBudget

Year	Section	Subsection	Type	Value
2003	Receipts	beginning cash	<i>drv</i>	20
2003	Receipts	cash sales	<i>det</i>	100
2003	Receipts	receivables	<i>det</i>	120
2003	Receipts	total cash receipts	<i>aggr</i>	250
...