

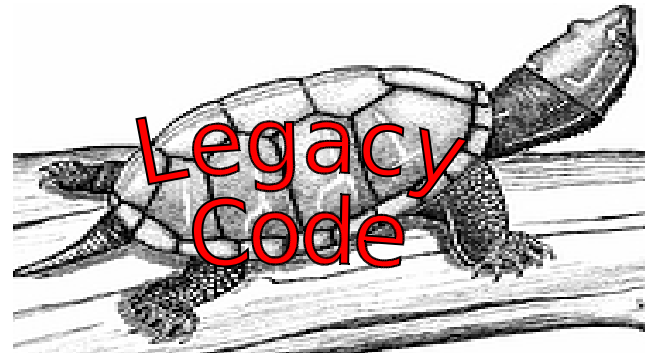
Semantics-based reverse engineering of data models from programs

Komondoor V Raghavan
IBM India Research Lab

(with G. Ramalingam, J. Field, et al)

Understanding legacy software

- Common scenario
 - huge existing legacy code base
 - building on top of existing code
 - transforming existing code
 - integrating legacy systems



- Legacy code can be surprisingly hard to work with
 - lack of documentation and understanding of existing code
- Need tools to help understand legacy code

Reverse engineering data models

- Goal: Reverse engineer a *logical data model* of a given (legacy) program
 - or Type Inference
 - focused on weakly-typed languages like Cobol
- Understanding logical structure of data is key to program understanding
- A logical data model can assist in common legacy transformation and maintenance tasks

An example Cobol program – *Data declarations*

```
01 CARD-TRANSACTION-REC  
  05 LOCATION-TYPE PIC X.  
  05 LOCATION-DETAILS PIC X(20).  
  05 CARD-INFO PIC X(19).  
  05 AMT PIC X(4).
```

Picture clauses

```
01 ATM-DETAILS.  
  05 ATM-ID PIC X(5).  
  05 ATM-ADDRESS X(12).  
  05 ATM-OWNER-ID PIC X(3).
```

Outermost variables

```
01 MERC-DETAILS.  
  05 MERCHANT-ID PIC X(8).  
  05 MERCHANT-ADDRESS PIC X(12).
```

Inner variables (fields)

```
01 CARD-NUM PIC X(16).  
01 CASHBACK-RATE X(2).  
01 CASHBACK X(3).
```

Example program -- code

```
/1/ READ CARD-TRANSACTION-REC.  
/2/ IF LOCATION-TYPE = 'M'  
/3/   MOVE LOCATION-DETAILS TO MERC-DETAILS  
/4/ ELSE  
/5/   MOVE LOCATION-DETAILS TO ATM-DETAILS  
/6/ ENDIF  
/7/ IF CARD-INFO[1:1] = 'C'  
/8/   MOVE CARD-INFO[2:3] TO CASHBACK-RATE  
/9/   MOVE AMT*CASHBACK-RATE/100 TO CASHB  
/10/  MOVE CARD-INFO[4:19] TO CARD-NUM  
/11/  WRITE CARD-NUM, CASHBACK TO CASHBACK-FILE  
/12/ ELSE  
/13/   MOVE CARD-INFO[2:17] TO CARD-NUM  
/14/ ENDIF  
/15/ IF LOCATION-TYPE = 'M'  
/16/   WRITE MERCHANT-ID, AMT, CARD-NUM TO M-FILE  
/17/ ELSE  
/18/   WRITE ATM-ID, ATM-OWNER-ID, AMT, CARD-NUM TO A-FILE.  
/19/ ENDIF
```

Types not obvious!

CreditCdNum

Disjoint union not obvious!

DebitCdNum

CreditCdNum | DebitCdNum

An example Cobol program – *Data declarations*

Implicit aggregate structure!

```
01 CARD-TRANSACTION-REC.  
  05 LOCATION-TYPE PIC X.  
  05 LOCATION-DETAILS PIC X(20).  
  05 CARD-INFO PIC X(19).  
  05 AMT PIC X(4).
```

AtmID ; OwnerID

MerchantID

'C':CreditTag ; CashBkRate ; CreditCdNum

!{'C'}:DebitTag ; DebitCdNum ; Unused

```
01 ATM-DETAILS.  
  05 ATM-ID PIC X(5).  
  05 ATM-ADDRESS X(12).  
  05 ATM-OWNER-ID PIC X(3).  
  
01 MERC-DETAILS.  
  05 MERCHANT-ID PIC X(8).  
  05 MERCHANT-ADDRESS PIC X(12).
```

```
01 CARD-NUM PIC X(16).  
01 CASHBACK-RATE X(2).  
01 CASHBACK X(3).
```

Algorithm 1 [TACAS '05]

- A “guarded” (dependent) type system, involving guarded type variables, records (concatenation), and unions
 - Example: ($'E': \alpha^1 ; \beta^7 ; \gamma^4 ; \delta^2$) |
(! $\{ 'E' \}$): $\varepsilon^1 ; \phi^9 ; \eta^4$)

Algorithm 1 [TACAS '05]

- A “guarded” (dependent) type system, involving guarded type variables, records (concatenation), and unions
 - Example: (`'E':Emp1 ; Eld7 ; Salary4 ; Unused2`) |
(`!{'E'}:Vis1 ; SSN9 ; Stipend4`)
- Formal characterization of a correct typing solution for a program
- Path-sensitive type inference algorithm
 - Improved accuracy; program-point specific types
 - Computed solution helps in constructing class diagram

Meaningful
names for
clarity

Applications of guarded type system

- Program understanding
- Understanding impact of changes
- Program transformation
 - Field expansion (e.g., Y2K expansion)
 - Porting from weakly-typed languages to object-oriented languages
 - Refactoring data declarations to make them better reflect logical structure

Key features of algorithm

- Based on dataflow analysis
 - Dataflow fact at each point is a *type* for the *entire memory*
 - Each *origin statement* (READ, MOVE literal TO var) gets a unique type variable
- Interprets predicates of the form
var == literal, var != literal
- Two key operations:
 - **Split:** Replace α^i by concatenation $\beta^j; \gamma^k$, $i = j + k$.
 - **Specialize:** Replace α^i by union $\beta^i \mid \gamma^i$.

CARD-TRANSACTION-REC

ATM-

MERC-

CARD-NUM

CASHBACK-RATE

CASHBACK

DETAILS

DETAILS

/1/ READ CARD-TRANSACTION-REC.

b^1

a^{44}

c^{43}

b^1

a^{44}

c^{43}

/2/ IF LOCATION-TYPE = 'M'

Split

$a^{44} \rightarrow b^1 ; c^{43}$

Specialize

$b^1 \rightarrow 'M':d^1 \mid !\{'M'\}:e^1$

CARD-TRANSACTION-REC

ATM-
DETAILS

MERC-
DETAILS

CARD-
NUM

CASHBACK
-
RATE

CASHBACK

/1/ READ CARD-TRANSACTION-REC.

'M':d¹

c⁴³

!{'M'}:e¹

f⁴³

/2/ IF LOCATION-TYPE = 'M'

CARD-TRANSACTION-REC

ATM-
DETAILS

MERC-
DETAILS

CARD-
NUM

CASHBACK
-
RATE

CASHBACK

/1/ READ CARD-TRANSACTION-REC.

'M':d¹

c⁴³

!{'M'}:e¹

f⁴³

'M':d¹

c⁴³

!{'M'}:e¹

f⁴³

/2/ IF LOCATION-TYPE = 'M'

'M':d¹

c⁴³

/3/ MOVE LOCATION-DETAILS TO MERC-DETAILS

/4/ ELSE

!{'M'}:e¹

f⁴³

/5/ MOVE LOCATION-DETAILS TO ATM-DETAILS

CARD-TRANSACTION-REC

ATM-
DETAILS

MERC-
DETAILS

CARD-
NUM

CASHBACK
-
RATE

CASHBACK

/1/ READ CARD-TRANSACTION-REC.

'M':d¹

h²⁰

j²³

!{'M'}:e¹

f⁴³

'M':d¹

h²⁰

j²³

!{'M'}:e¹

f⁴³

/2/ IF LOCATION-TYPE = 'M'

'M':d¹

h²⁰

j²³

/3/ MOVE LOCATION-DETAILS TO MERC-DETAILS

'M':d¹

h²⁰

j²³

h²⁰

/4/ ELSE

!{'M'}:e¹

f⁴³

/5/ MOVE LOCATION-DETAILS TO ATM-DETAILS

CARD-TRANSACTION-REC

ATM-
DETAILS

MERC-
DETAILS

CARD-
NUM

CASHBACK
-
RATE

CASHBACK

/1/ READ CARD-TRANSACTION-REC.

'M':d ¹	h ²⁰	j ²³
!{'M'}:e ¹	j ²⁰	k ²³

'M':d ¹	h ²⁰	j ²³				
!{'M'}:e ¹	j ²⁰	k ²³				

/2/ IF LOCATION-TYPE = 'M'

'M':d ¹	h ²⁰	j ²³				
--------------------	-----------------	-----------------	--	--	--	--

/3/ MOVE LOCATION-DETAILS TO MERC-DETAILS

'M':d ¹	h ²⁰	j ²³		h ²⁰		
--------------------	-----------------	-----------------	--	-----------------	--	--

/4/ ELSE

!{'M'}:e ¹	j ²⁰	k ²³				
-----------------------	-----------------	-----------------	--	--	--	--

/5/ MOVE LOCATION-DETAILS TO ATM-DETAILS

!{'M'}:e ¹	j ²⁰	k ²³		j ²⁰		
-----------------------	-----------------	-----------------	--	-----------------	--	--



CARD-TRANSACTION-REC

ATM-
DETAILS

MERC-
DETAILS

CARD-
NUM

CASHBACK
RATE

CASHBACK

/1/ READ CARD-TRANSACTION-REC.

'M':d ¹	h ²⁰	i ²³
!{'M'}:e ¹	j ²⁰	k ²³

'M':d ¹	h ²⁰	i ²³		h ²⁰		
!{'M'}:e ¹	j ²⁰	k ²³		j ²⁰		

/7/ IF CARD-INFO[1:1] = 'C'

CARD-TRANSACTION-REC

ATM-
DETAILS

MERC-
DETAILS

CARD-
NUM

CASHBACK
-
RATE

CASHBACK

/1/ READ CARD-TRANSACTION-REC.

'M':d ¹	h ²⁰	l ¹	m ²²
!{'M'}:e ¹	j ²⁰	n ¹	o ²²

'M':d ¹	h ²⁰	l ¹	m ²²		h ²⁰			
!{'M'}:e ¹	j ²⁰	n ¹	o ²²		j ²⁰			

/7/ IF CARD-INFO[1:1] = 'C'

Specialize

$l^1 \rightarrow 'C':p^1 \mid !{'C'}:q^1$
 $n^1 \rightarrow 'C':r^1 \mid !{'C'}:s^1$

CARD-TRANSACTION-REC

ATM-

MERC-

CARD-
NUM

CASHBACK
-
RATE

CASHBACK

--

DETAILS

DETAILS

--

--

--

/1/ READ CARD-TRANSACTION-REC.

'M':d ¹	h ²⁰	'C':p ¹	m ²²
!{'M'}:e ¹	j ²⁰	'C':r ¹	o ²²
'M':t ¹	u ²⁰	!{'C'}:q ¹	v ²²
!{'M'}:w ¹	x ²⁰	!{'C'}:s ¹	y ²²



'M':d ¹	h ²⁰	'C':p ¹	m ²²				
!{'M'}:e ¹	j ²⁰	'C':r ¹	o ²²				
'M':t ¹	u ²⁰	!{'C'}:q ¹	v ²²				
!{'M'}:w ^x	x ²⁰	!{'C'}:s ¹	y ²²				

/2/ IF LOCATION-TYPE = 'M'

CARD-TRANSACTION-REC

ATM-DETAILS MERC-DETAILS

'M':d ¹	h ²⁰	'C':p ¹	m ²²			CARD- NUM	CASHBACK -RATE	CASH BACK
!{'M'}:e ¹	j ²⁰	'C':r ¹	o ²²					
'M':t ¹	u ²⁰	!{'C'}:q ¹	v ²²					
!{'M'}:w ¹	x ²⁰	!{'C'}:s ¹	y ²²					

```

/2/ IF LOCATION-TYPE = 'M'
/3/ MOVE LOCATION-DETAILS TO MERC-DETAILS
/4/ ELSE
/5/ MOVE LOCATION-DETAILS TO ATM-DETAILS
/7/ IF CARD-INFO[1:1] = 'C'
/8/ MOVE CARD-INFO[2:3] TO CASHBACK-RATE
/9/ MOVE AMT*CASHBACK-RATE/100 TO CASHBACK
/10/ MOVE CARD-INFO[4:19] TO CARD-NUM
/11/ WRITE CARD-NUM, CASHBACK TO CASHBACK-FILE
/12/ ELSE
/13/ MOVE CARD-INFO[2:17] TO CARD-NUM
/15/ IF LOCATION-TYPE = 'M'
/16/ WRITE MERCHANT-ID, AMT, CARD-NUM TO M-FILE
/17/ ELSE
/18/ WRITE ATM-ID, ATM-OWNER-ID, AMT,CARD-NUM TO A-FILE.
    
```

'M':d ¹	h ₁ ⁸	h ₂ ¹²	'C':p ¹	m ₁ ²	m ₂ ¹⁶	m ₃ ⁴			h ₁ ⁸	h ₂ ¹²	m ₂ ¹⁶	m ₁ ²	z	
!{'M'}:e ¹	j ₁ ⁵	j ₂ ¹²	j ₃ ³	'C':r ¹	o ₁ ²	o ₂ ¹⁶	o ₃ ⁴	j ₁ ⁵	j ₂ ¹²	j ₃ ³		o ₂ ¹⁶	o ₁ ²	z
'M':t ¹	u ₁ ⁸	u ₂ ¹²	!{'C'}:q ¹	v ₁ ¹⁶	v ₂ ²	v ₃ ⁴			u ₁ ⁸	u ₂ ¹²	v ₁ ¹⁶			
!{'M'}:w ¹	x ₁ ⁵	x ₂ ¹²	x ₃ ³	!{'C'}:s ¹	y ₁ ¹⁶	y ₂ ²	y ₃ ⁴	x ₁ ⁵	x ₂ ¹²	x ₃ ³		y ₁ ¹⁶		

Correctness characterization

Input:

a	b	c	d	e	f...
---	---	---	---	---	------

 α β γ η

typing solution is **correct** because there exists an **atomization**...

... and a typing of each atom ...

... such that types completely describe runtime values

REPEAT ...

MOVE X TO ...

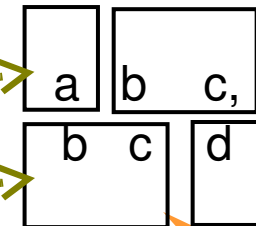
$\alpha; \beta$ |
 $\beta; \gamma$

Is type of

Is type of

Typing solution

Runtime values

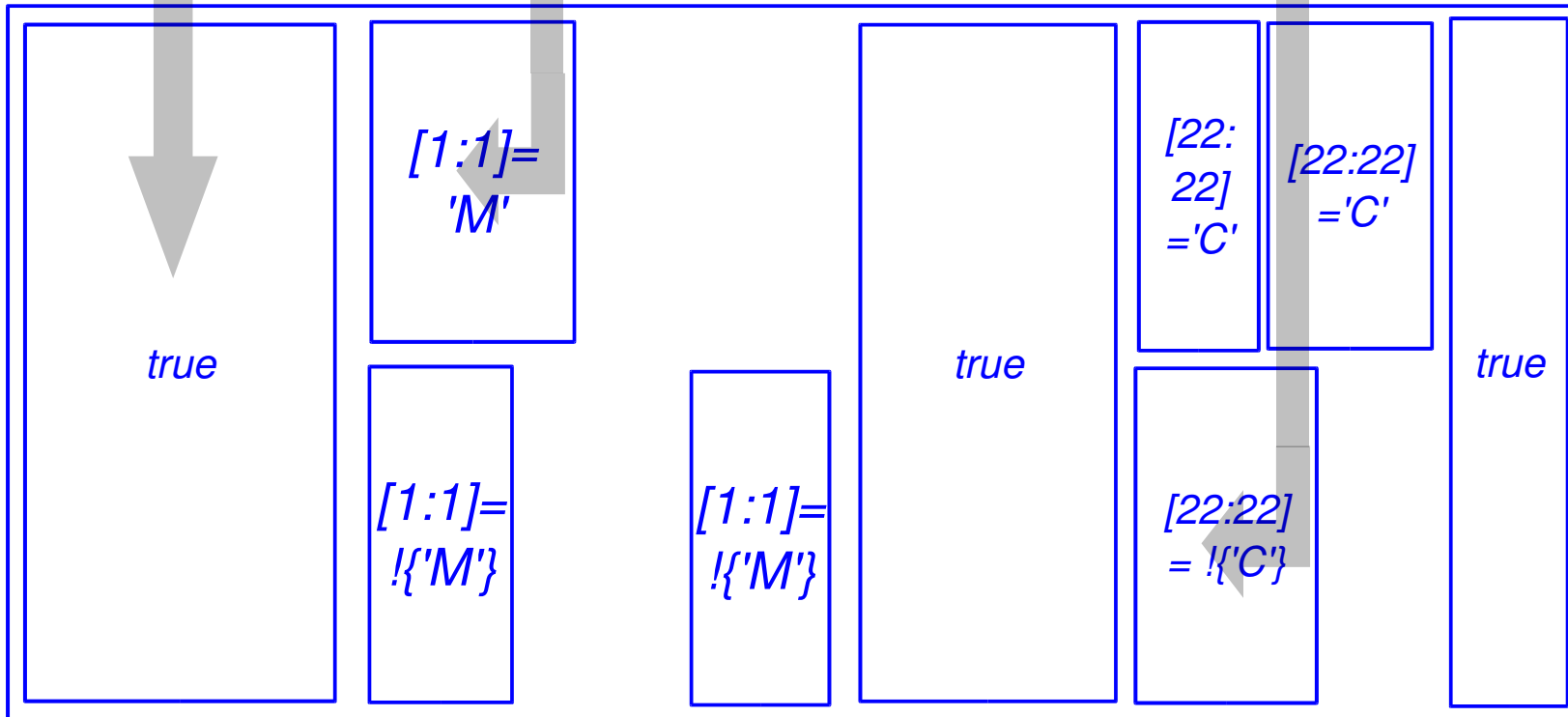


Characteristics of the solution

- Flow and path sensitive:
 - Each occurrence of a variable is assigned a type
 - Uses guards to ignore certain infeasible paths
- Determines variables of the same type, reveals record structure within variables, as well as disjoint unions
- Shortcomings:
 - Dataflow facts are “unfactored”, potentially of exponential size

/1/ READ CARD-TRANSACTION-REC.

'M':d ¹	h ₁ ⁸		h ₂ ¹²	'C':p ¹	m ₁ ²	m ₂ ¹⁶	m ₃ ⁴
!{'M'}:e ¹	j ₁ ⁵	j ₂ ¹²	j ₃ ³	'C':r ¹	o ₁ ²	o ₂ ¹⁶	o ₃ ⁴
'M':t ¹	u ₁ ⁸		u ₂ ¹²	!{'C'}:q ¹	v ₁ ¹⁶	v ₂ ²	v ₃ ⁴
!{'M'}:w ¹	x ₁ ⁵	x ₂ ¹²	x ₃ ³	!{'C'}:s ¹	y ₁ ¹⁶	y ₂ ²	y ₃ ⁴



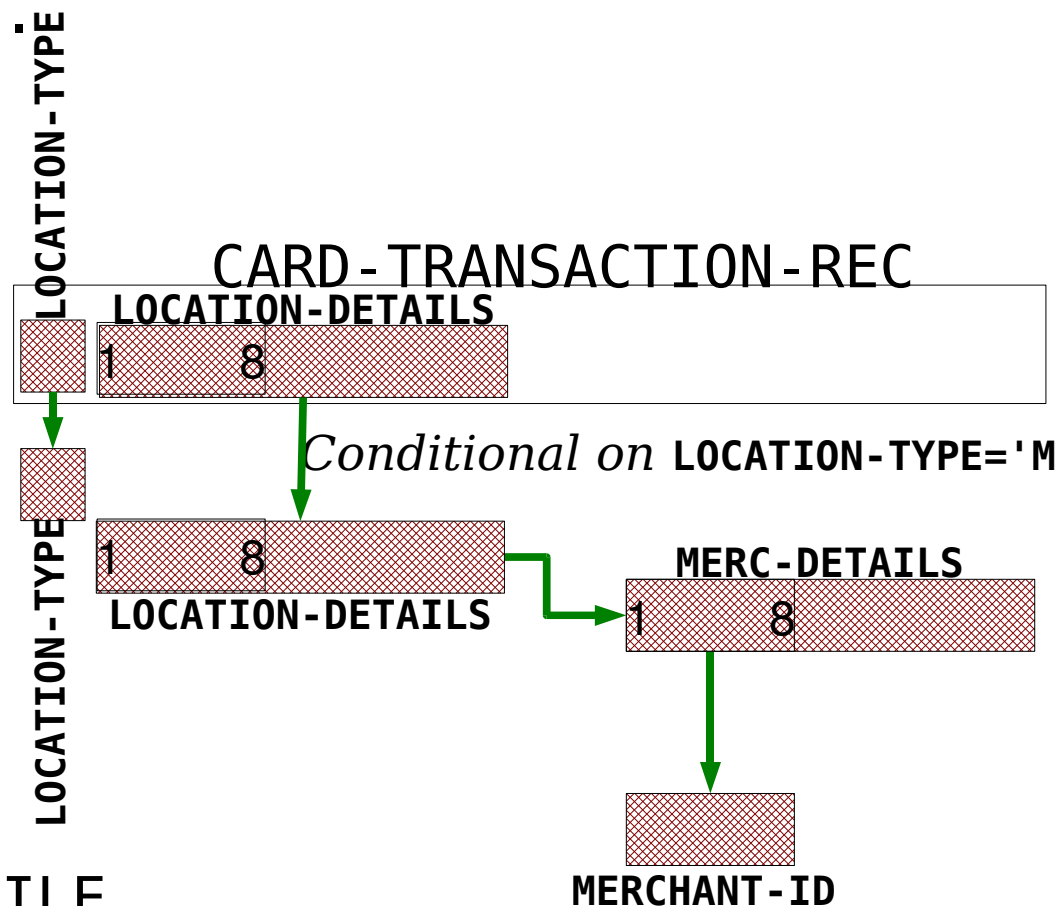
Algorithm 2 [ICSE '06, WCRE '07]

1. Compute guarded dependences
2. Compute *cuts* at each data-source statement (i.e., READ statement).
3. Organize the cuts as a cut-structure tree
 - *It is possible, but not desirable, to translate cut-structure tree directly into a class hierarchy*
4. *Factor* the cut-structure tree to capture better the grouping/structure of sibling cuts
5. Translate cut-structure tree into a class hierarchy

Step 1. Compute guarded dependences

```
01 CARD-TRANSACTION-REC.  
 05 LOCATION-TYPE PIC X.  
 05 LOCATION-DETAILS PIC X(20).  
   ... 23 more bytes ...  
01 MERC-DETAILS.  
 05 MERCHANT-ID PIC X(8).  
   ...  
/1/  READ CARD-TRANSACTION-REC  
  
/2/  IF LOCATION-TYPE = 'M'  
/3/  MOVE LOCATION-DETAILS  
      TO MERC-DETAILS  
  
/4/  ELSE ...  
/15/ IF LOCATION-TYPE = 'M'  
/16/ WRITE MERCHANT-ID,  
      AMOUNT, CARD-NUM TO M-FILE  
  
/17/ ELSE ...
```

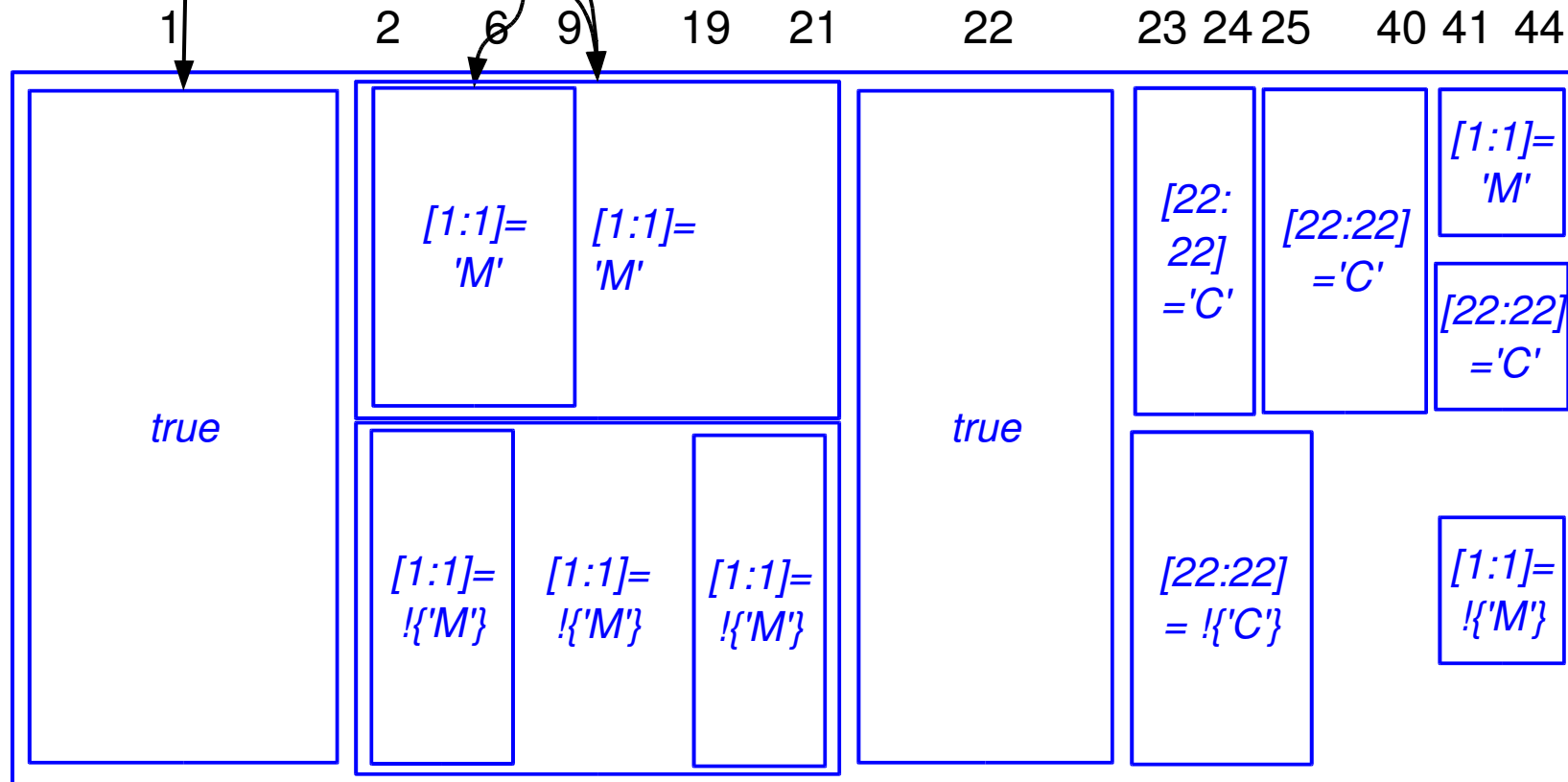
LOCATION-TYPE='M' ►
LOCATION-DETAILS[1:8]@/1/ →
MERCHANT-ID@/16/



Step 2: Compute cuts at each data source

CARD-TRANSACTION-REC

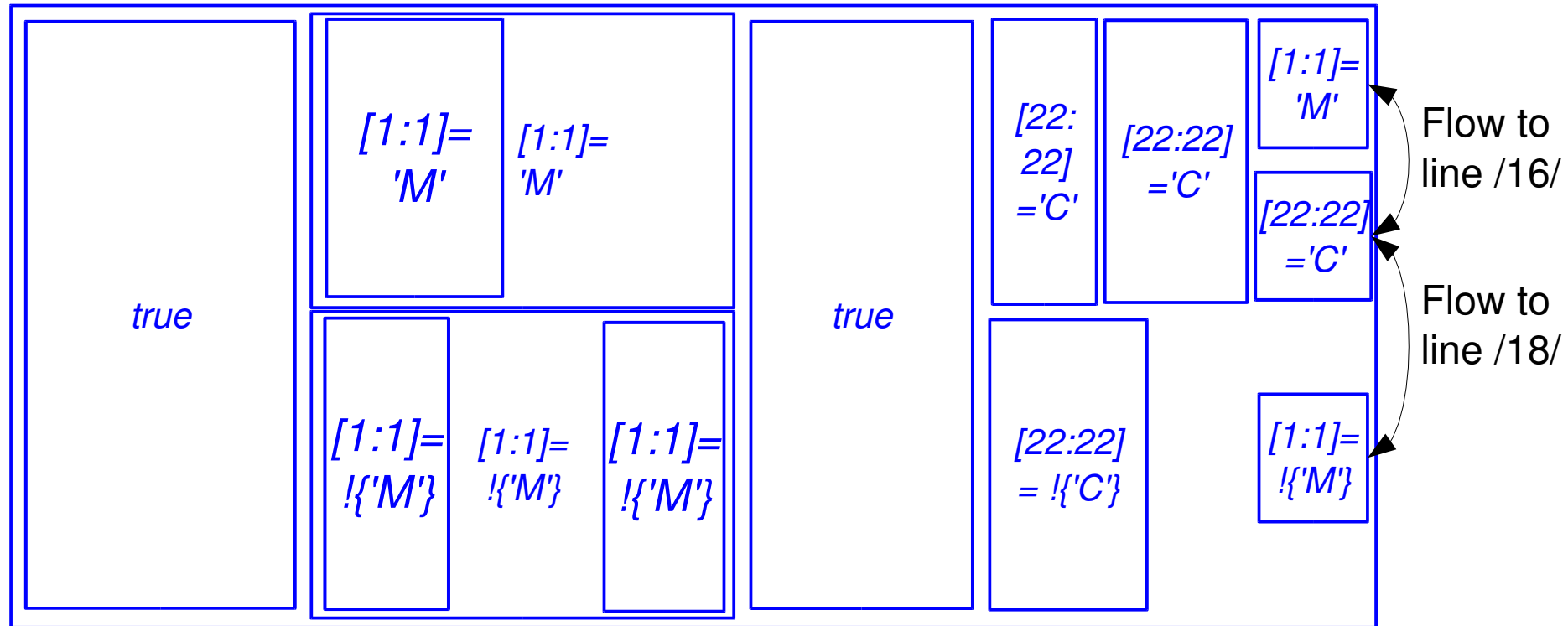
1. true ▶ LOCATION-TYPE@/1/ → LOCATION-TYPE@/2/
2. (LOCATION-TYPE='M') ▶ LOCATION-DETAILS@/1/ → LOCATION-DETAILS@/3/
1. (LOCATION-TYPE='M') ▶ LOCATION-DETAILS@/1/ → MERC-Details@/3/
2. (LOCATION-TYPE='M') ▶ LOCATION-DETAILS[1:8]@/1/ → MERCHANT-ID@/16/



Step 3: Organize cuts as tree

- There is intuitively a containment relation among cuts. Formalization:
 - c_i 's range “wider” than c_j 's range *and* c_i 's predicate “broader” than c_j 's predicate $\Rightarrow c_i$ contains c_j
- We broaden predicates of certain cuts such that
 - 1) Containment imposes a *tree* structure (not a DAG)
 - *Allows generation of a single-inheritance class hierarchy*
 - 2) Between any two siblings c_j and c_k there is no overlap;
i.e.:
 - *Either* their ranges are non-overlapping, *or* their predicates are non-overlapping

Illustration of Step 3

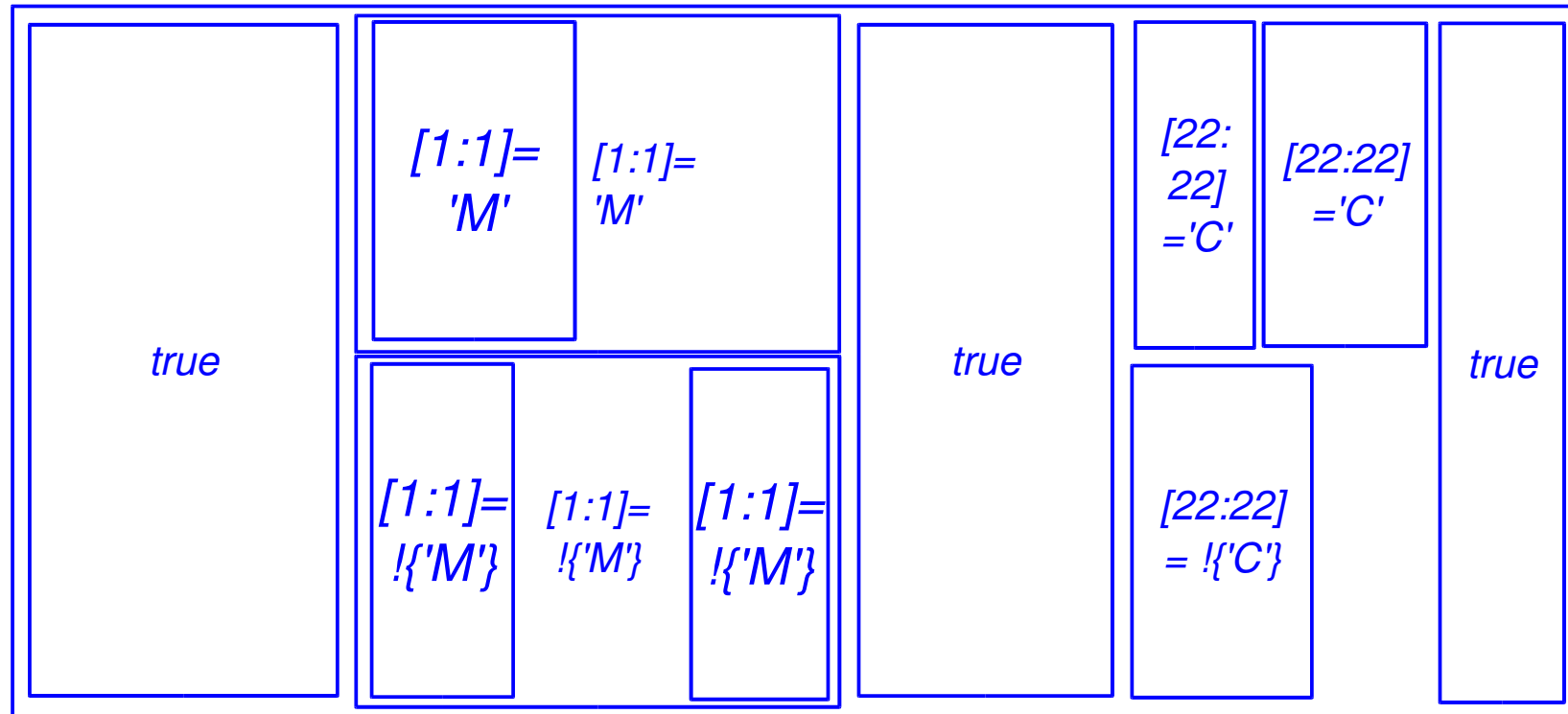


1) Cuts already form a tree structure. Good.

2) However, we have overlap problem!

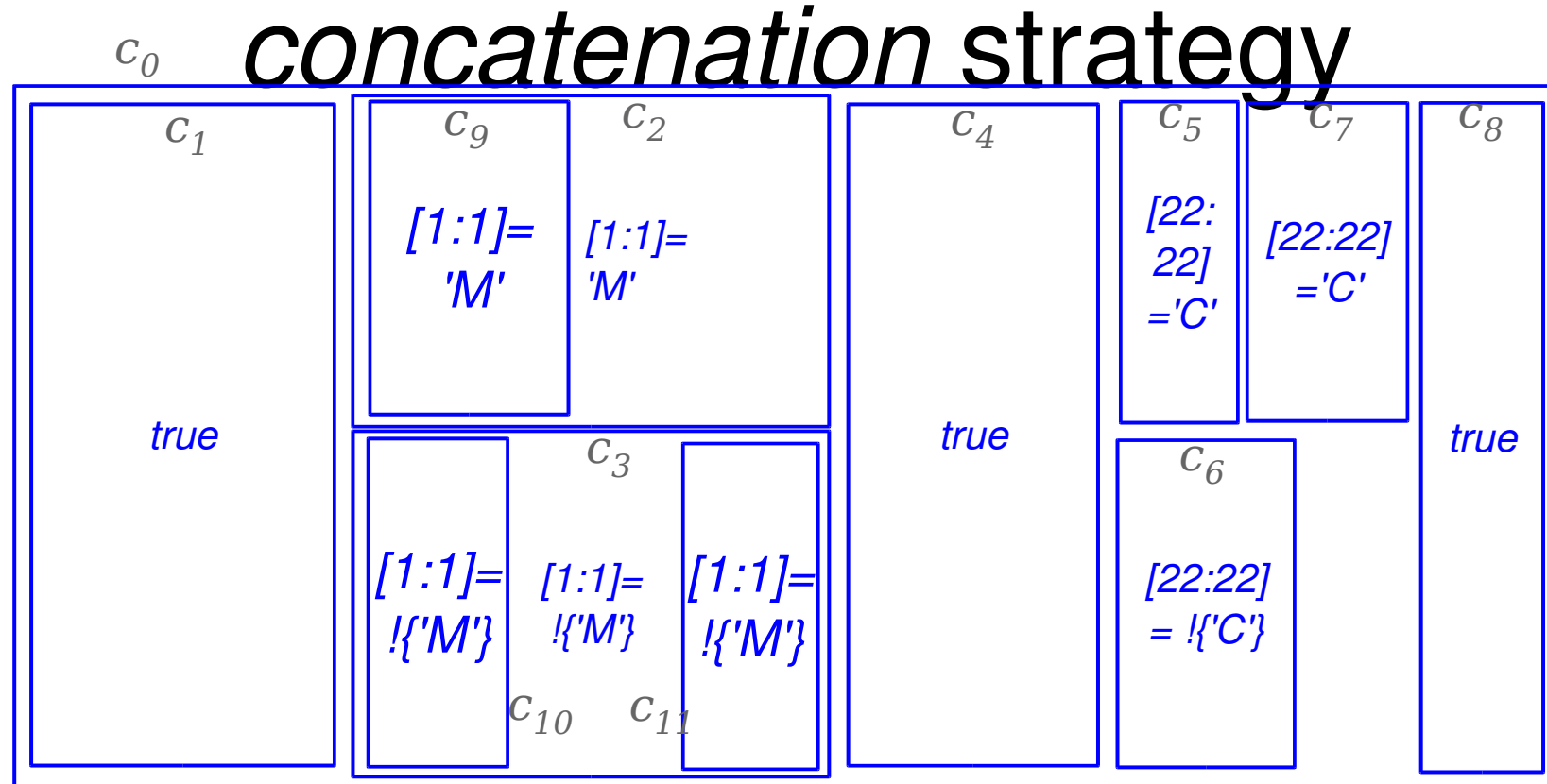
- Intuitively, two overlapping cuts \Rightarrow both flow into some variable reference;
- We would like a unique cut to flow into each variable ref.

Illustration of Step 3



Merge the three cuts.
That is, take logical
disjunction of their
predicates.

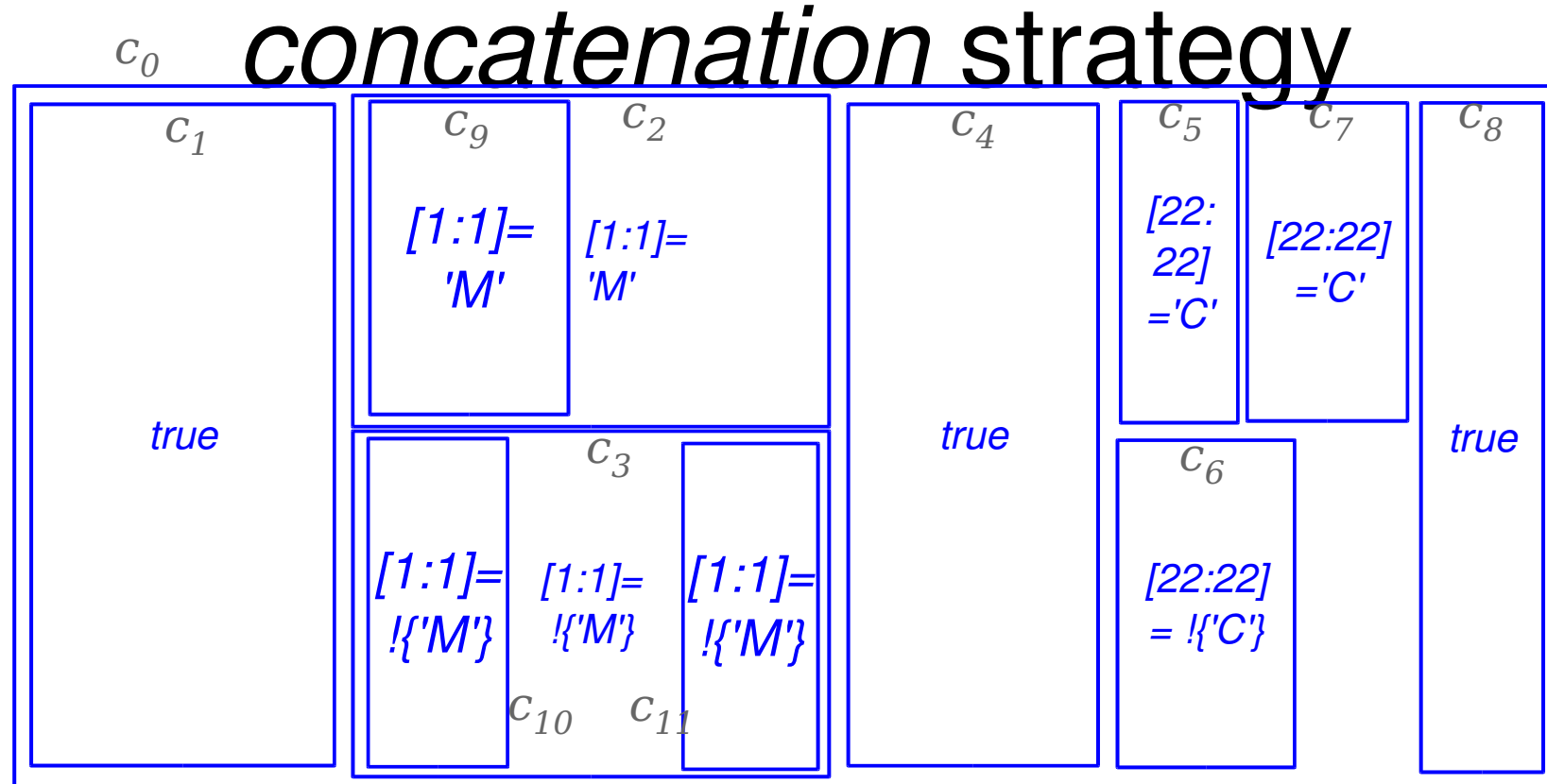
Generating a class hierarchy:



Approach: Turn each cut into a class, and each edge into a *field-of* relation.

- Class c_0 {f1: c_1 , f2: c_2 , f3: c_3 , ..., f8: c_8 }, Class c_1 {}, ..., Class c_8 {}

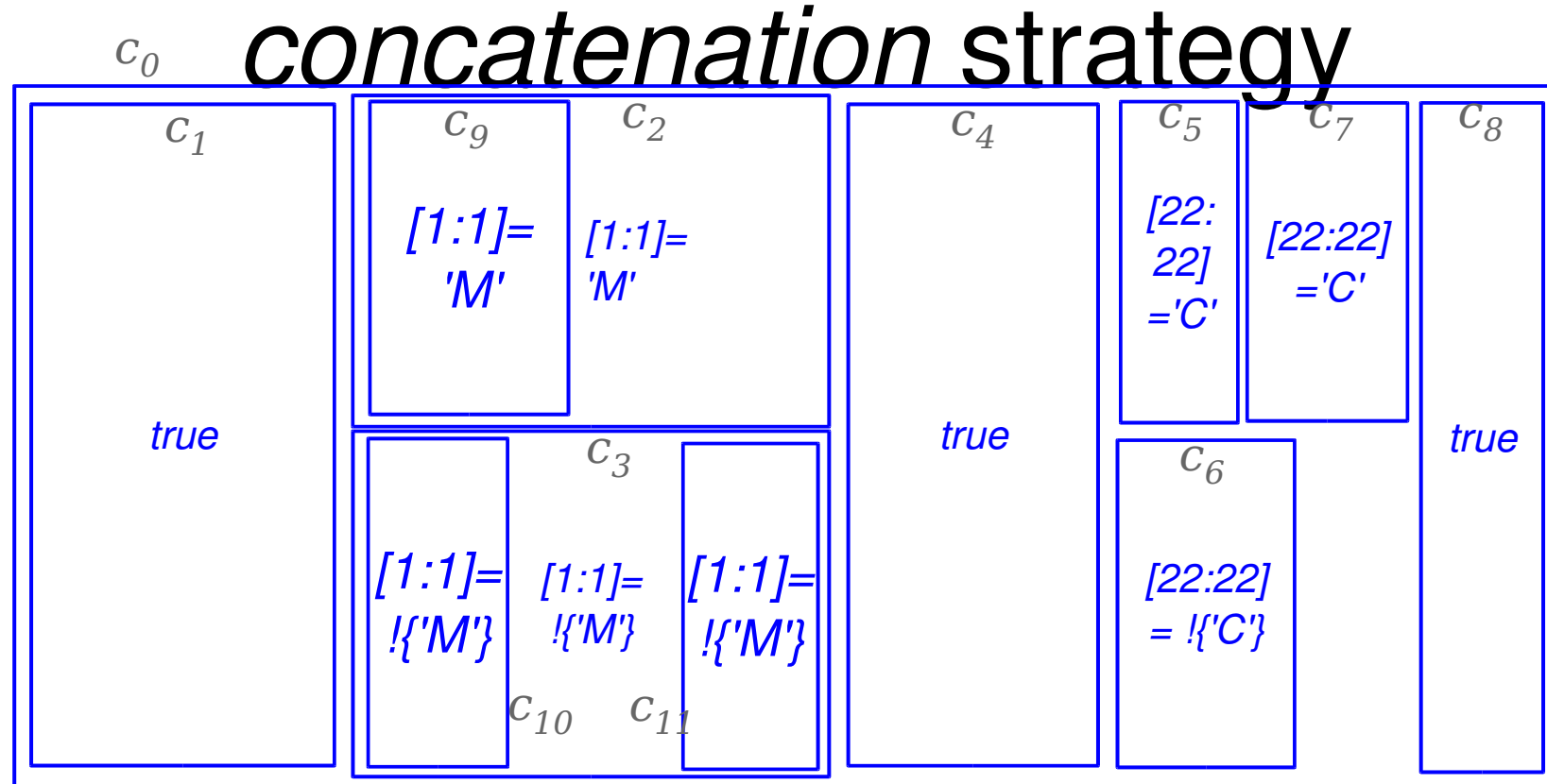
Generating a class hierarchy:



Approach: Turn each cut into a class, and each edge into a *field-of* relation.

- Class c_0 {f1: c_1 , f2: c_2 , f3: c_3 , ..., f8: c_8 }, Class c_1 {}, ..., Class c_8 {}
- However, predicates are *lost in translation*, hence loss of precision: fields f2 and f3 ought not to co-exist!

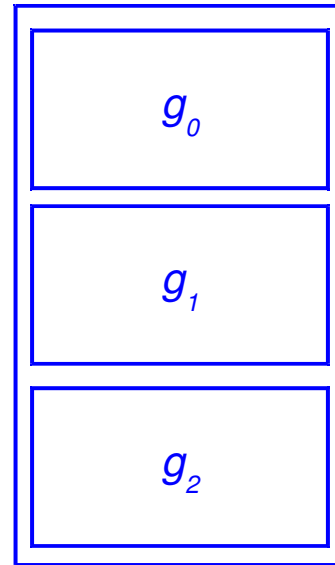
Generating a class hierarchy:



Approach: Turn each cut into a class, and each edge into a *field-of* relation.

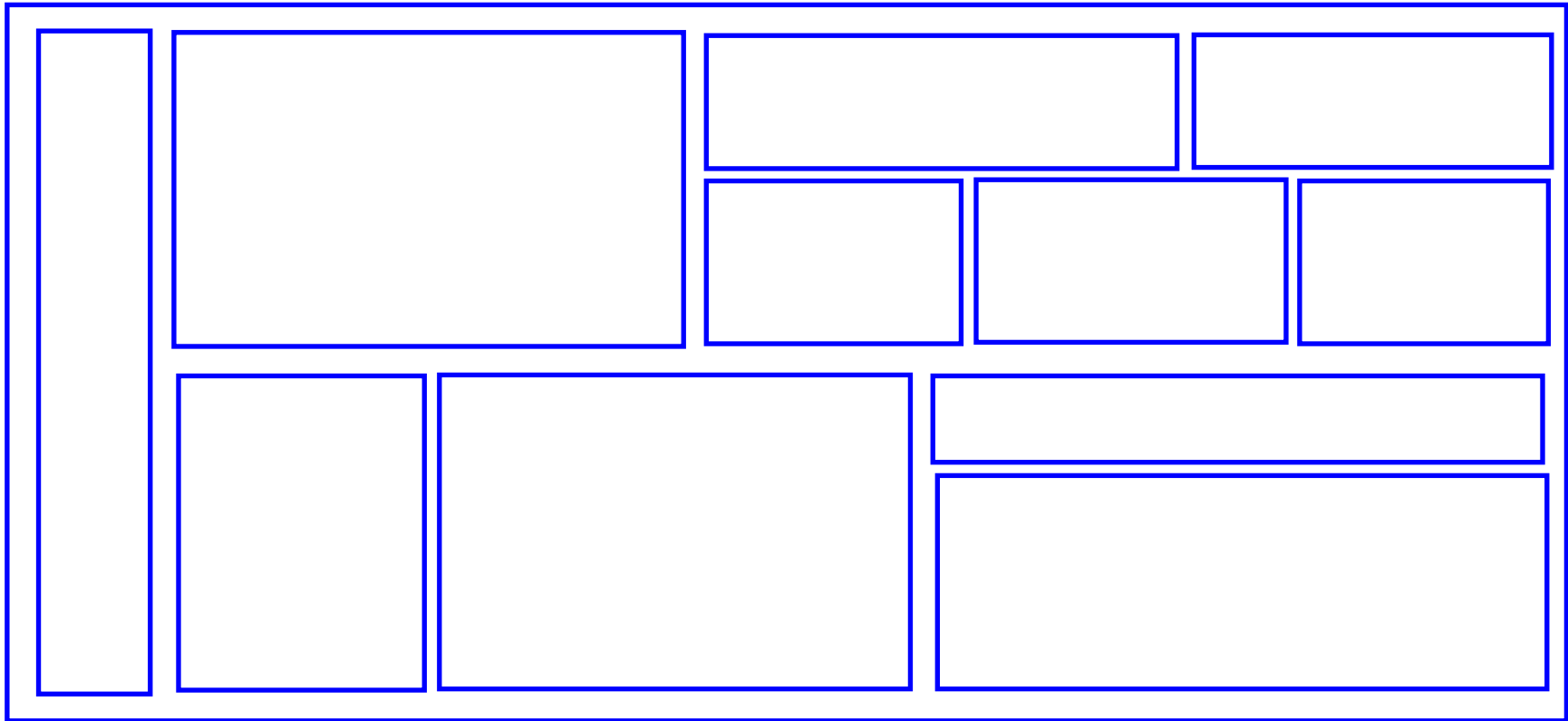
- Class c_0 {f1: c_1 , f2: c_2 , f3: c_3 , ..., f8: c_8 }, Class c_1 {}, ..., Class c_8 {}
- However, predicates are *lost in translation*, hence loss of precision: fields f2 and f3 ought not to co-exist!
- No loss of precision when when all children have the same guard

Vertical partitioning



- Applicable *only when* all children have mutually disjoint predicates
 - parent corresponds to a base class
 - children correspond to derived classes

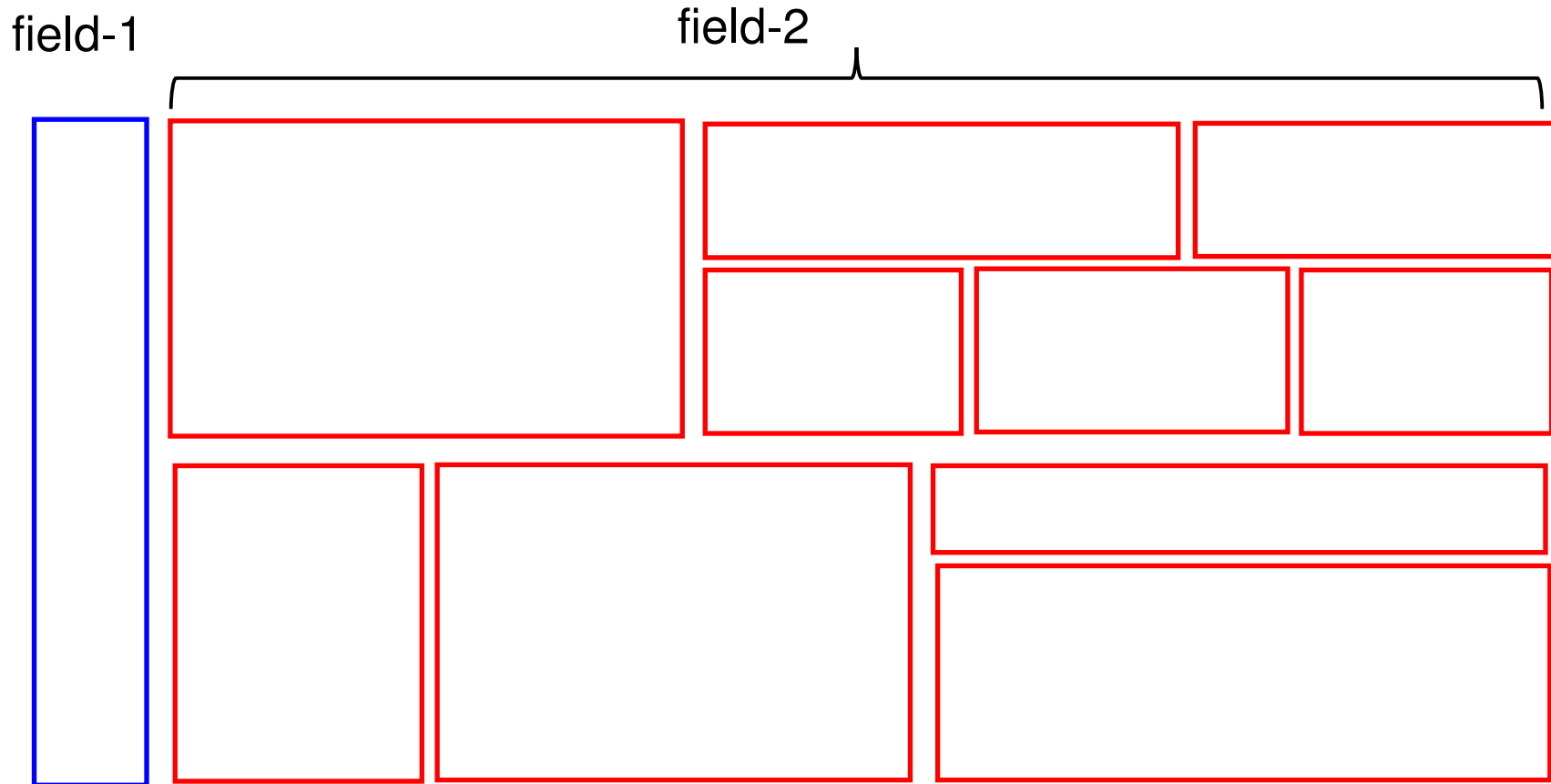
Step 4: Factoring cut-structure tree



Generalized Horizontal Partitioning

- Add edges between boxes with disjoint guards
- each connected component == a field

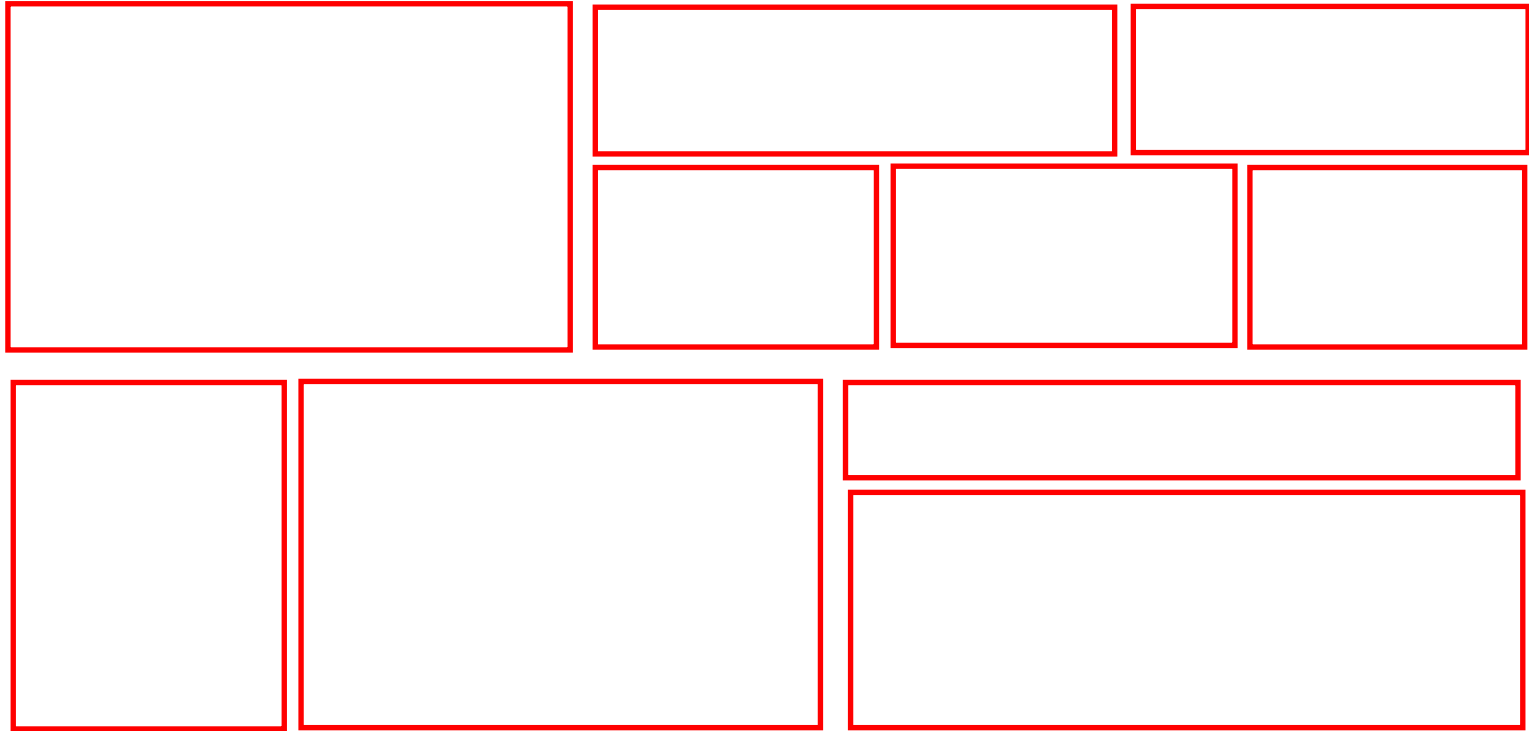
Step 4



Generalized Horizontal Partitioning

- Add edges between boxes with disjoint guards
- each connected component == a field

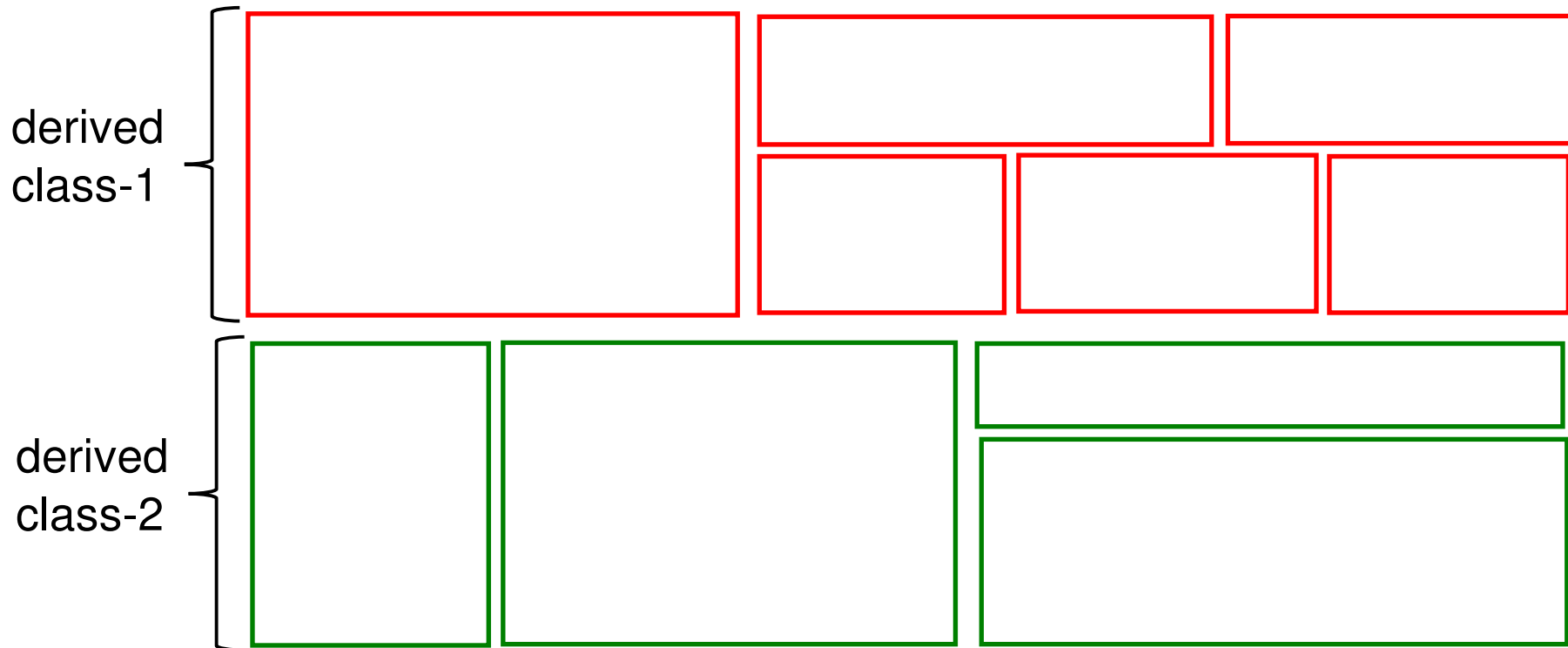
Step 4



Generalized Vertical Partitioning

- Add edges between boxes with overlapping guards
- each connected component == a derived class

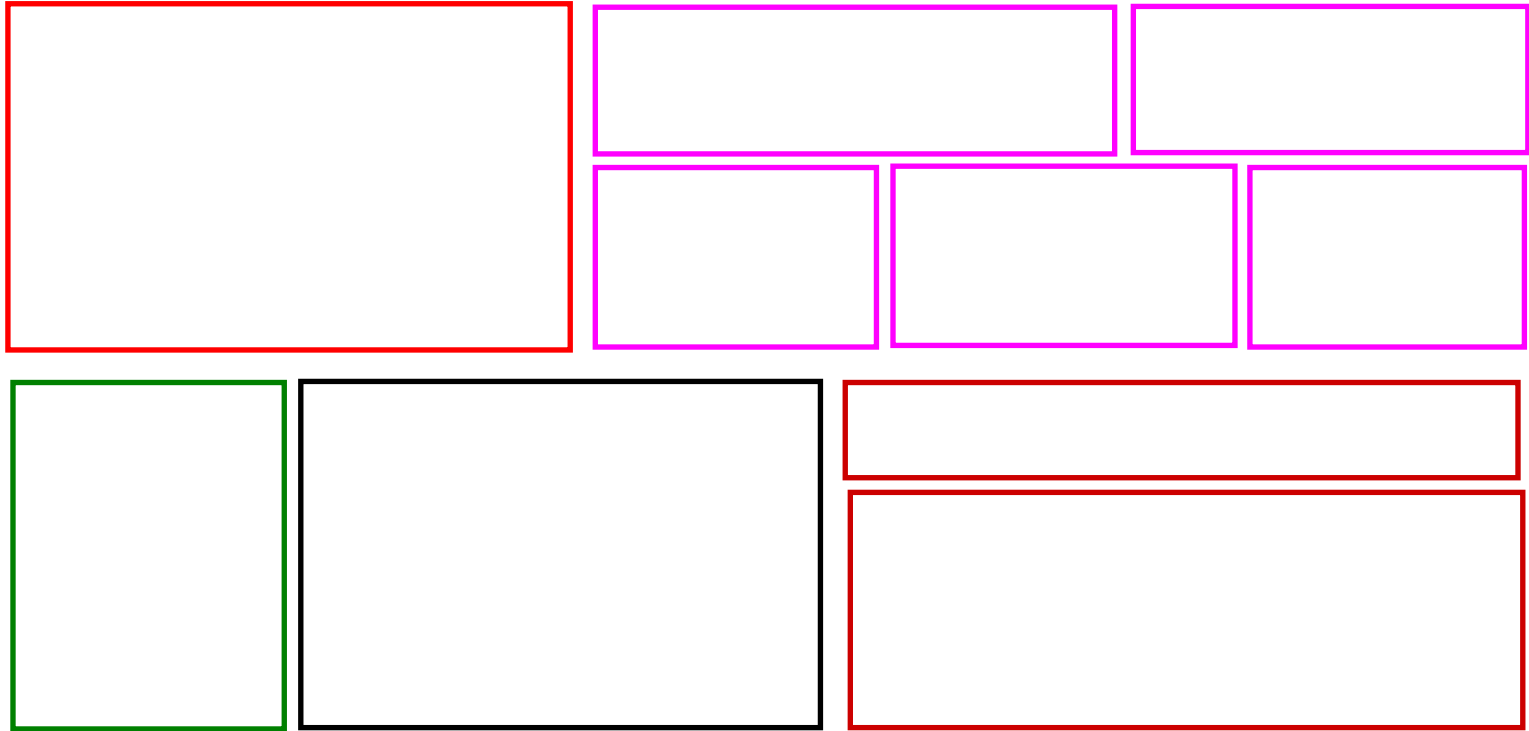
Step 4



Generalized Vertical Partitioning

- Add edges between boxes with overlapping guards
- each connected component == a derived class

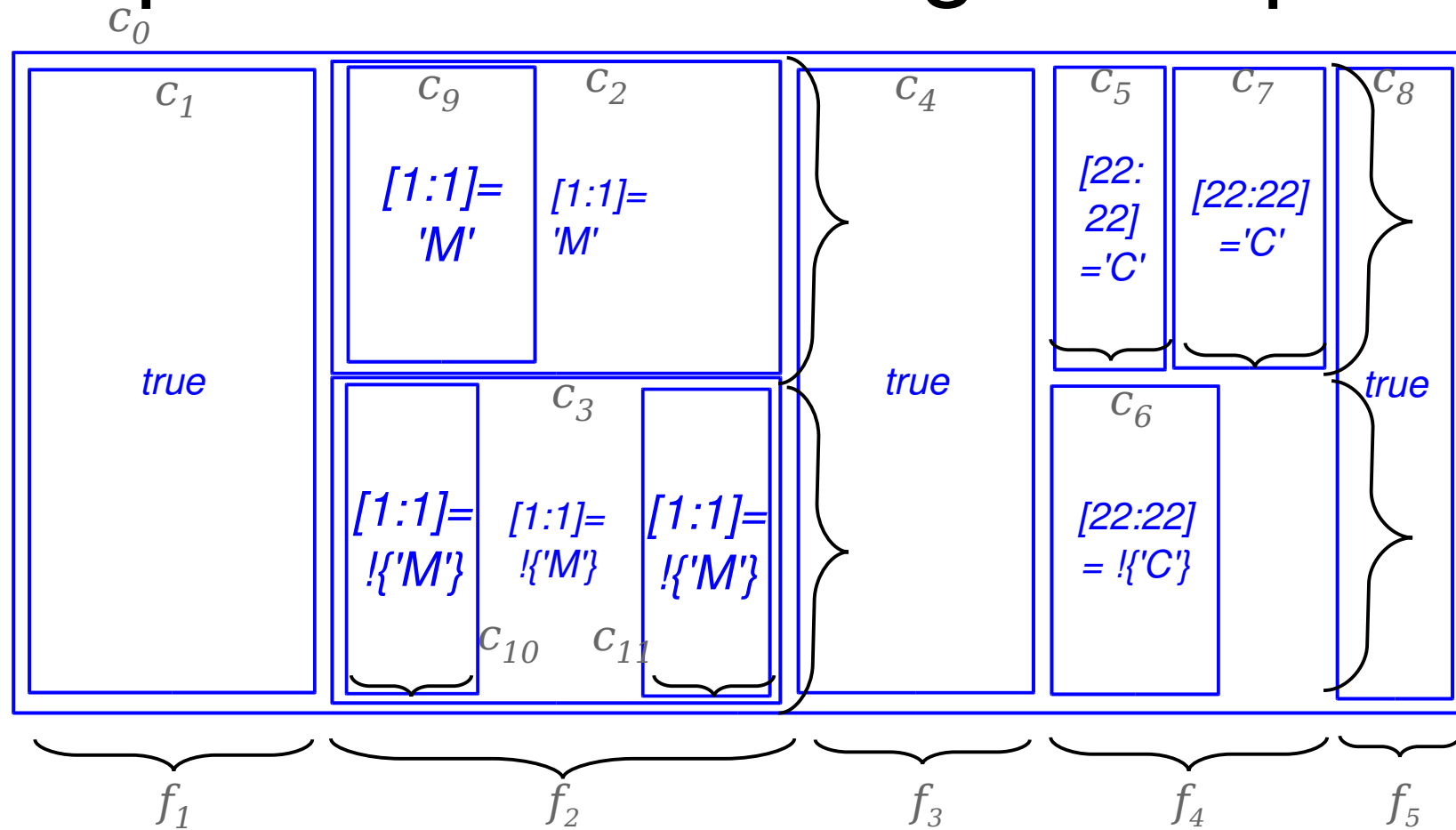
Step 4



Generalized Horizontal Partitioning

- Add edges between boxes with disjoint guards
- each connected component == a field

Step 4 on the running example ...



```

01 CARD-TRANSACTION-REC.
  05 LOCATION-TYPE PIC X.
  05 LOCATION-DETAILS PIC X(20).
  05 CARD-INFO PIC X(19).
  05 AMT PIC X(4).

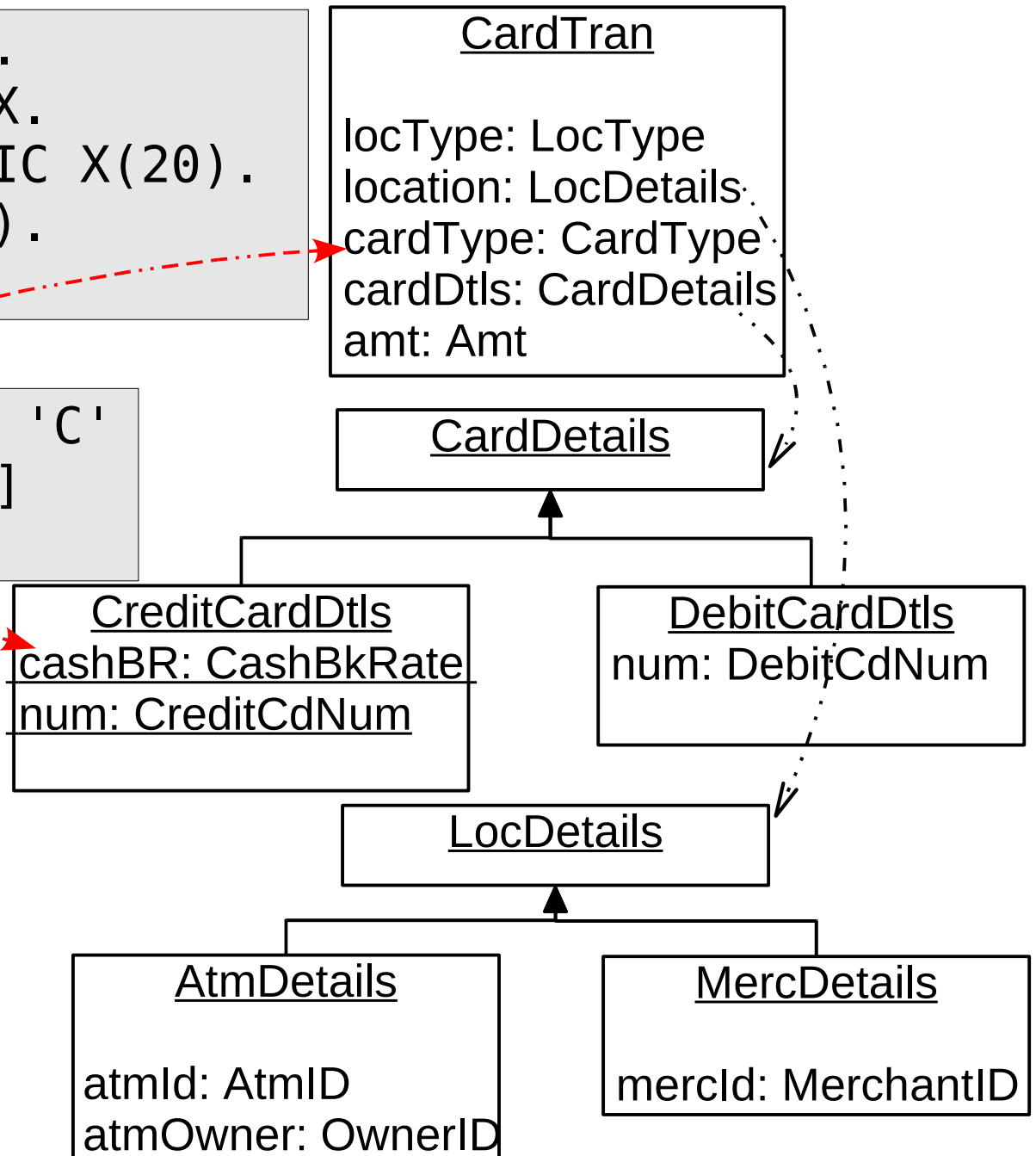
```

```

/7/ IF CARD-INFO[1:1] = 'C'
/8/ MOVE CARD-INFO[2:3]
    TO CASHBACK-RATE

```

- How can we say if a given OO model is correct for a given program?
- Executing the program using an altered data representation as suggested by the OOM does not affect the observable behavior of the program.
 - See [ICSE '06] for details



Details of Step 1: Computing guarded dependences

- **guard** ► **source** → **target**
 - **source** is a pair **memory range @ program-pt**
 - **target** is similar (however, we restrict ourself to variable dereference sites)
 - **guard** is a predicate on the state at **source program-point**.
- when **guard** is true, value at **source** may reach **target** (via some sequence of copies)

Guarded dependence analysis

- Guarded dependences
 - capture transitive data-dependences
 - capture conditions under which dependence is manifested
- Parametric guarded dependence computation
 - parameterized by abstraction for guards
 - can be computed in polynomial time for simple (common type of) guards

Transfer functions (without guards)

Statement S	$\alpha_{pi}[S] : 2^{\mathcal{D}_{pi}} \rightarrow 2^{\mathcal{D}_{pi}}$
WRITE Y^d	$\lambda \text{Out}. \text{Out} \cup \{ Y \rightsquigarrow d \}$
READ Y^d	$\lambda \text{Out}. \{ Y \rightsquigarrow d \} \cup$ $\{ r \rightsquigarrow t \mid r \rightsquigarrow t \in \text{Out} \text{ and } r \not\subseteq Y \}$
MOVE X^{dX} TO Y^{dY} where $Y = [y_1 : y_2]$ and $X = [x_1 : x_2]$	$\lambda \text{Out}. \{ X \rightsquigarrow dX, Y \rightsquigarrow dY \} \cup$ $\{ r \rightsquigarrow t \mid r \rightsquigarrow t \in \text{Out} \wedge r \not\subseteq Y \} \cup$ $\{ r' \rightsquigarrow t \mid r \rightsquigarrow t \in \text{Out}, r \subset Y,$ and $r' = r - y_1 + x_1 \}$
ASSUME pred	$\lambda \text{Out}. \text{Out} \cup \{ r \rightsquigarrow d \mid \text{pred contains}$ a data-reference d to a range $r \}$

Backward dataflow analysis. Dataflow fact is: set of memory-range X variable-reference-site. Meet operation: set union.

Transfer functions (with guards)

$$\alpha_{eps}[t](\text{Out}) = \{true \triangleright r' \rightsquigarrow d' \mid r' \rightsquigarrow d' \in \alpha_{pi}[t]\{\}\} \cup \\ \{\alpha_g[t](g) \triangleright r' \rightsquigarrow d' \mid g \triangleright r \rightsquigarrow d \in \text{Out}, \\ r' \rightsquigarrow d' \in \alpha_{pi}[t]\{r \rightsquigarrow d\}\}$$

$\alpha_g[t](g)$ = weakest pre-condition semantics; i.e.,
broadest condition *before* statement t that implies g
after statement t .

Ensuring polynomial-time analysis

- Atomic predicates
 - variable = constant | variable \notin set-of-constants
 - $x = 1, y \neq 2, z \notin \{1,2\}$
- Each guard is
 - a conjunction of atomic predicates
 - (at most one per variable)
- Use Map(memory-range X variable-reference-site, guard) as dataflow fact domain, instead of memory-range X variable-reference-site X guard.

Algorithm 2 : Contributions

- An efficient approach to infer OO data models from weakly-typed programs
- Inferred models are provably compact and *correct*
- Prototype implementation, and manual examination of results
- Therefore, is a sound basis for program understanding, migration, and transformation

Related work

- Canfora et al. [SEKE 96]
- O'Callahan and Jackson [ICSE 97]
- van Deursen and Moonen [WCRE 98, ...]
- Eidorff et al. [POPL 99]
- Ramalingam, Field, and Tip [POPL 99]
- Balakrishnan and Reps [CC 04]

- Distinguishing attributes of our work:
 - path-sensitive analysis
 - semantic correctness criterion of inferred OO model