

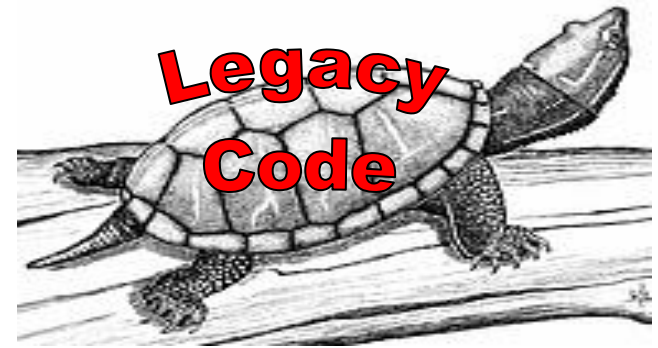
When less is more

G. Ramalingam

Microsoft Research India

The Dark Side of Programming

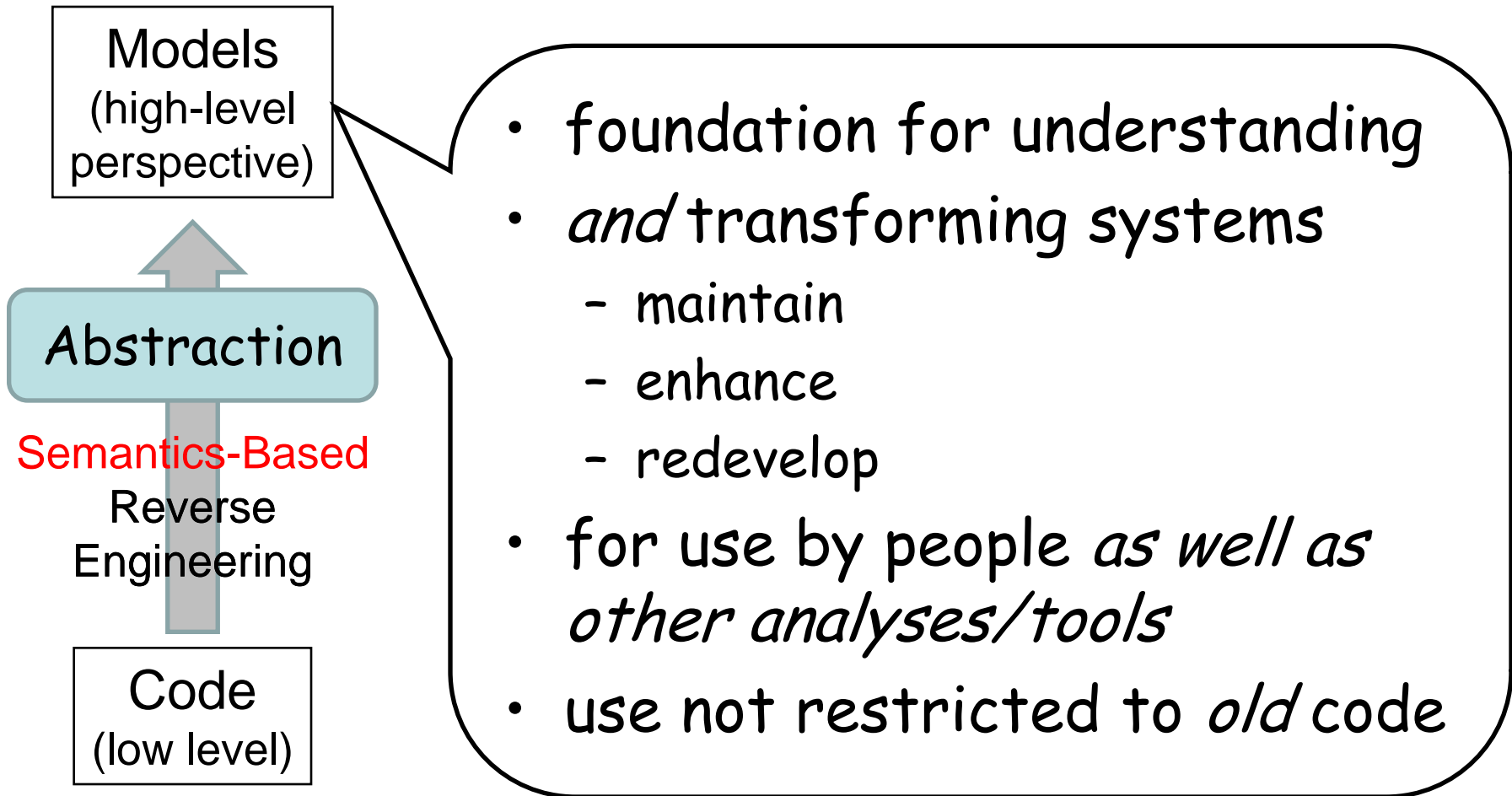
- Common scenario
 - huge existing legacy code base
 - building on top of existing code
 - transforming existing code
 - integrating legacy systems
- Legacy code can be surprisingly hard to work with
 - lack of documentation and understanding of existing code
- Need tools to help understand (legacy) code



Towards program-understanding

- Stumbling blocks
 - large and complex code-base
 - low-level code
 - implementation details
 - tangling of multiple concerns
 - scattering (of a logical aspect all over code)
- Understanding requires models that
 - are high level
 - focus on some desired aspect
 - eliminate irrelevant details

Reverse Engineering: From Code To Models



1. Mining Library Usage Rules

(Joint work with
Kanika Nema, Kapil Vaswani, Venkatesh-Prasad,
Microsoft Research India)

- Libraries typically have usage constraints
 - data preconditions
 - input parameter p must be non-null
 - parameter x must point to array of size $> y$
 - temporal constraints
 - call $m1(p)$ before calling $m2(p)$
 - do not call $m3(p)$ after calling $m4(p)$
- Goal: Recover (part of) API contract from large code base that uses the API

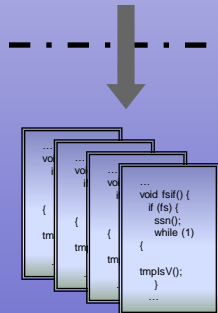
Abstracting Client-Library Interaction: What & Why

Library Side Concerns

Clients



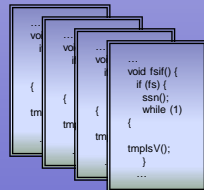
- Which clients are likely to break?
- Why?
- How can we ensure compatibility?



Library



evolution



Updated Library

Client Side Concerns

How do I use this?

Development/
Maintenance



Verification tools maturing.
How do we feed them?

Large Code Base



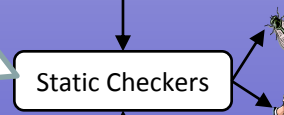
API Usage
Rules/Patterns

Static Checkers

Code

Where could the problem be?

Debugging



2. Static Analysis of Binaries

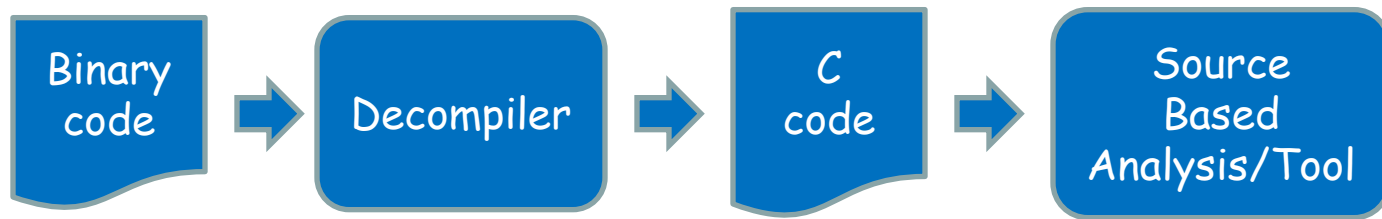
(Joint work with:

Pushkar Tripathi, V.L. Subrahmanyam, Venkatesh-Prasad,
Microsoft Research India)

- Many useful analyses and tools exist
 - e.g., Static Driver Verifier (SLAM)
- Most target source languages, such as C
- But sometimes only binaries are available
 - third party code (drivers)
 - potentially malicious code

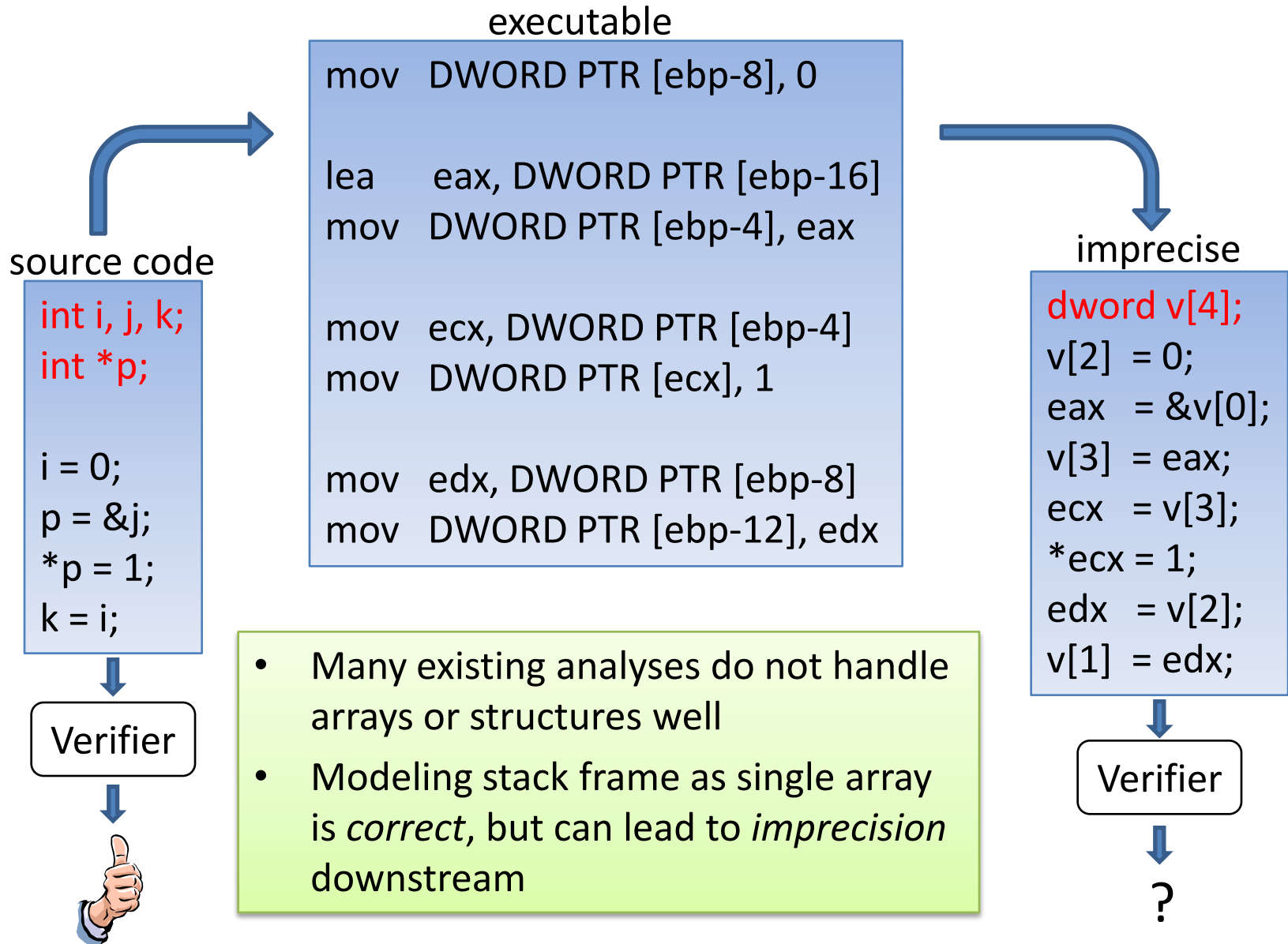
The Problem

- Decompile binary code to C ...
- ... permitting reuse of source-code based analyses and tools

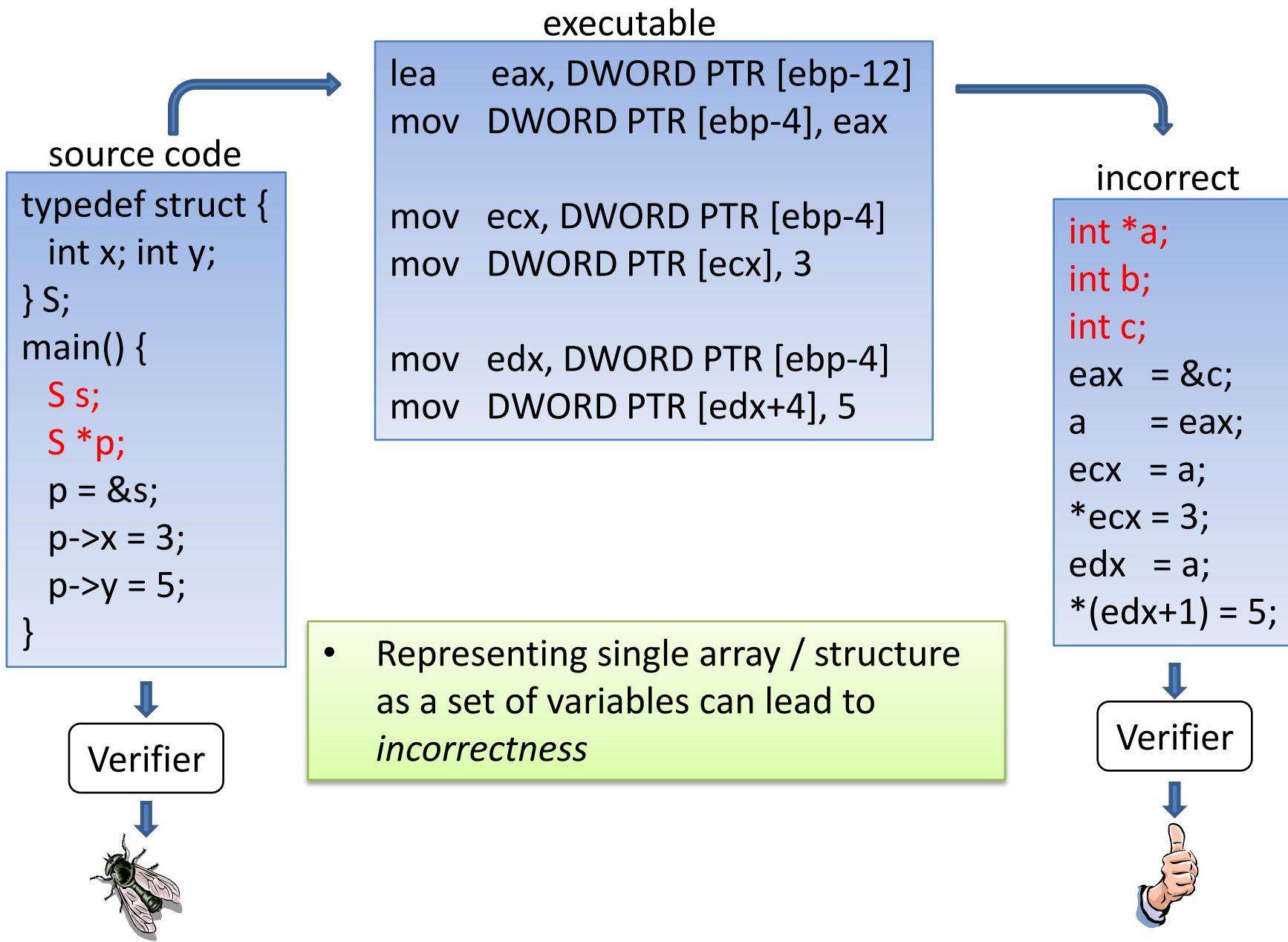


- Naïve transliteration will not be very effective
 - Existing analyses exploit higher level abstractions, such as variables and types
- Problem: Inferring *variables* and *types*

Variable Identification: Issues



Variable Identification: Issues



3. Maintenance of Legacy Code

(Joint work with: Raghavan Komondoor, John Field, Satish Chandra, Saurabh Sinha, IBM Research)

- Focus:
 - legacy business applications, written in (weakly-typed) languages such as Cobol
 - assist in common legacy transformation and maintenance tasks

Motivating Example

- Field Expansion
 - a frequently recurring scenario
 - e.g., expand "account number" from 8 to 10 digits
 - Y2K: well-known instance of field expansion



Get Ready for the New ISBN!

- [NISO Home](#)
- [About NISO](#)
- [NISO Standards](#)
- [NISO Registration Process](#)
- [Creating Standards](#)
- [Become a Member](#)
- [The Standards World](#)
- [International](#)

The new 13-digit ISBN has been approved and plans are underway to transition to the new number industry-wide, world-wide by January 1, 2007. Find out how the expansion of the ISBN from 10-digits to 13-digits will impact your business and operations:

- [Information for Publishers](#)
- [A Librarian's ISBN FAQ](#)
- [More ISBN Links](#)

Information for Publishers

As of January 1, 2007, all book and book-related products must carry 13-digit ISBNs.

All 10-digit ISBNs in circulation will have the 3-digit EAN prefix "978" added (which currently represents the book industry). This 13-digit

An Example Cobol Program

```
01 input-record  
   05 ir-trans-code pic x(1).  
   05 ir-amount pic x(8).  
   05 ir-account-number pic x(8).  
   05 ir-data pic x(8).  
01 withdrawal-record  
   05 wi-amount pic x(8).  
   05 wi-data pic x(8).  
01 account-record  
   05 ar-account-number pic x(8).  
   05 ar-user-name pic x(8).  
   05 ar-data pic x(8).  
...  
MOVE ar-account-number TO ir-account-number  
IF ir-trans-code = 'w' THEN  
   MOVE ir-data TO withdrawal-record  
...
```

Cobol

- no user-defined types
- lack of typing constraints
- primitive type: string (sequence of bytes)
- facility to define structured variables
- memory overlay (like *union* in C)
- declarative info: no semantic guarantee

Is this an account-number?

Logical Data Model Inference

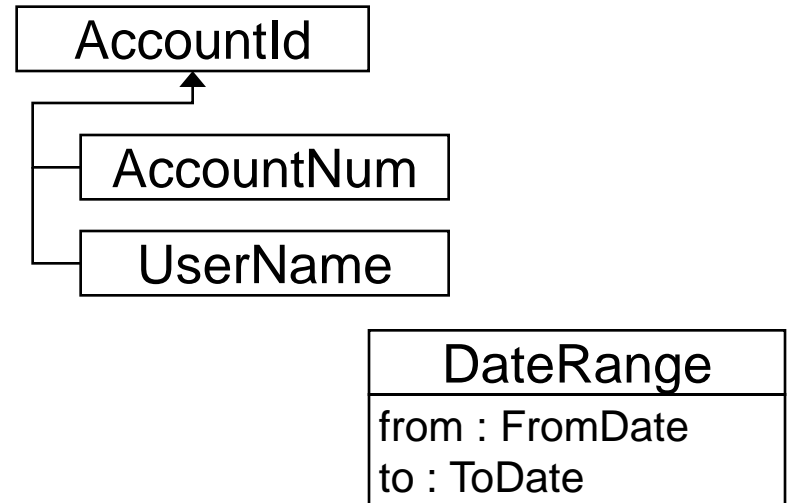
- Problem: Reverse engineer a *logical data model* of a given (legacy) program

The Inference Problem

Source Program

```
...  
01 account-rec.  
   05 ar-acc-num pic x(8).  
   05 ar-user-name pic x(8).  
   05 ar-data pic x(5).  
...  
MOVE ar-acc-num TO ir-acc-num.  
...  
MOVE ir-data TO ...  
...
```

Object-Oriented Data Model



- class definitions
- inheritance
- typed fields

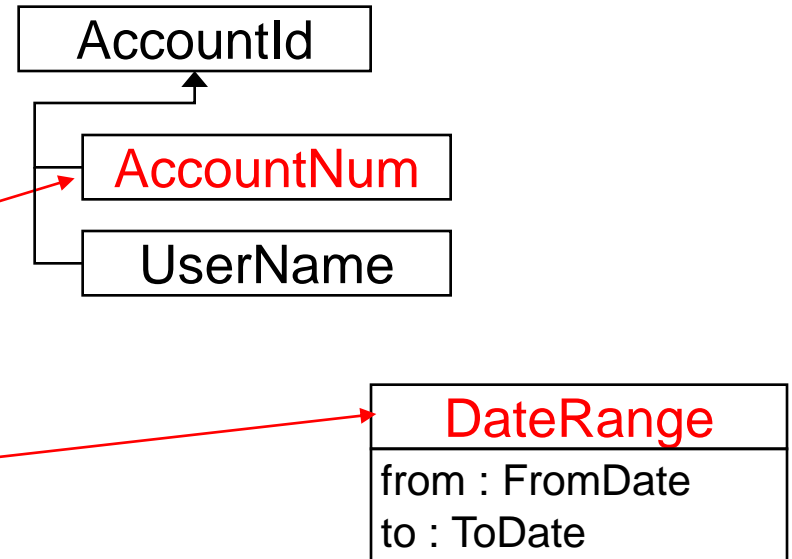
- Infer
 - Object-Oriented Data Model
 - + Links
 - connect a source program to an OO model

Object-Oriented Data Model + Links (LOOM)

Source Program

```
...  
01 account-rec.  
   05 ar-acc-num pic x(8).  
   05 ar-user-name pic x(8).  
   05 ar-data pic x(5).  
...  
MOVE ar-acc-num TO ir-acc-num.  
...  
MOVE ir-data TO ...  
...
```

Object-Oriented Data Model



• Infer

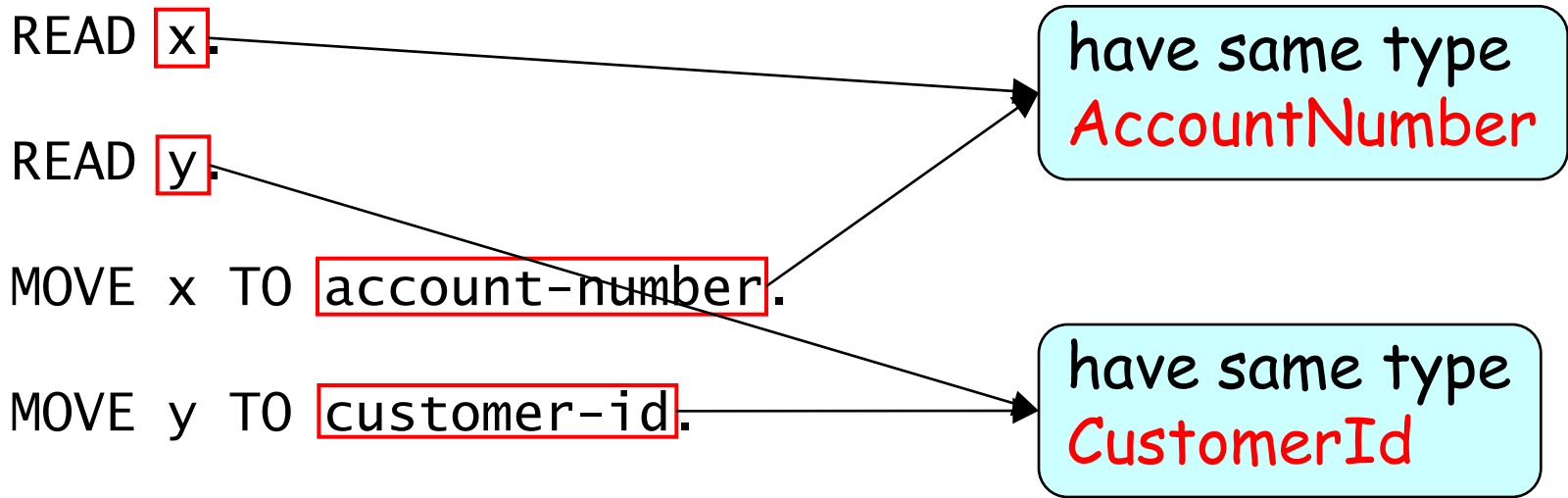
- Object-Oriented Data Model
- + Links

- connect a source program to an OO model
- map variable occurrences to types (& more)

The Inference Algorithm

Basic Idea:

1. Infer Type Equivalence



Basic Idea:

2. Infer Record Structure

READ **R**.

MOVE R[1:8] TO **account-number**

MOVE R[9:18] TO **customer-id**

```
class R {  
  acctNum: AcctNum;  
  custId: CustId;  
}
```

```
class AcctNum {}
```

```
class CustId {}
```

Basic Idea:

3. Infer Subtypes

READ **R**.

MOVE R[1:8] TO account-number.

MOVE R[9:18] TO customer-id.

IF (R[19:19] = 'D') THEN

... **process deposit** ...

ELSE

... **process withdrawal** ...

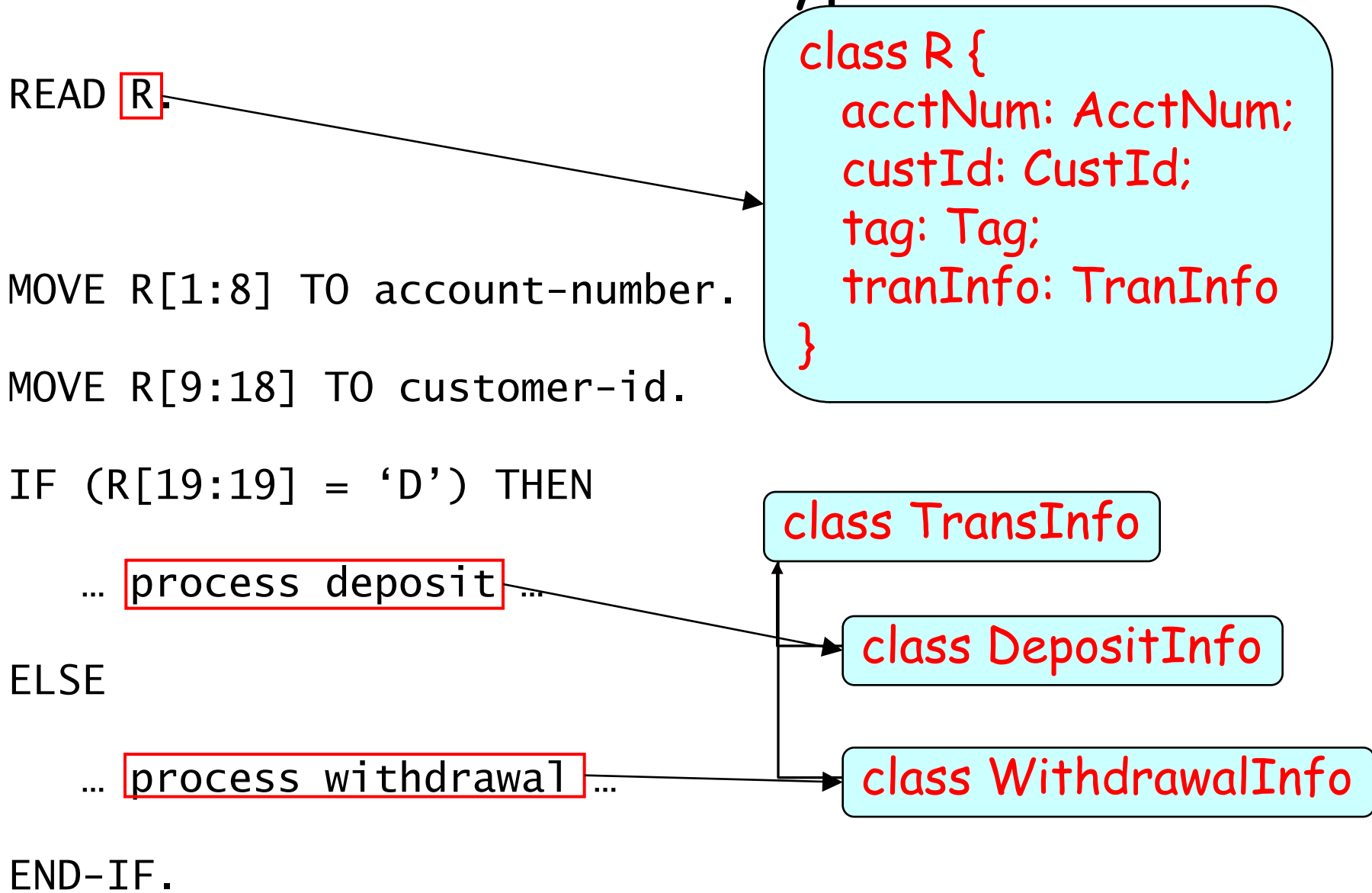
END-IF.

```
class R {  
  acctNum: AcctNum;  
  custId: CustId;  
  tag: Tag;  
  tranInfo: TranInfo  
}
```

class TransInfo

class DepositInfo

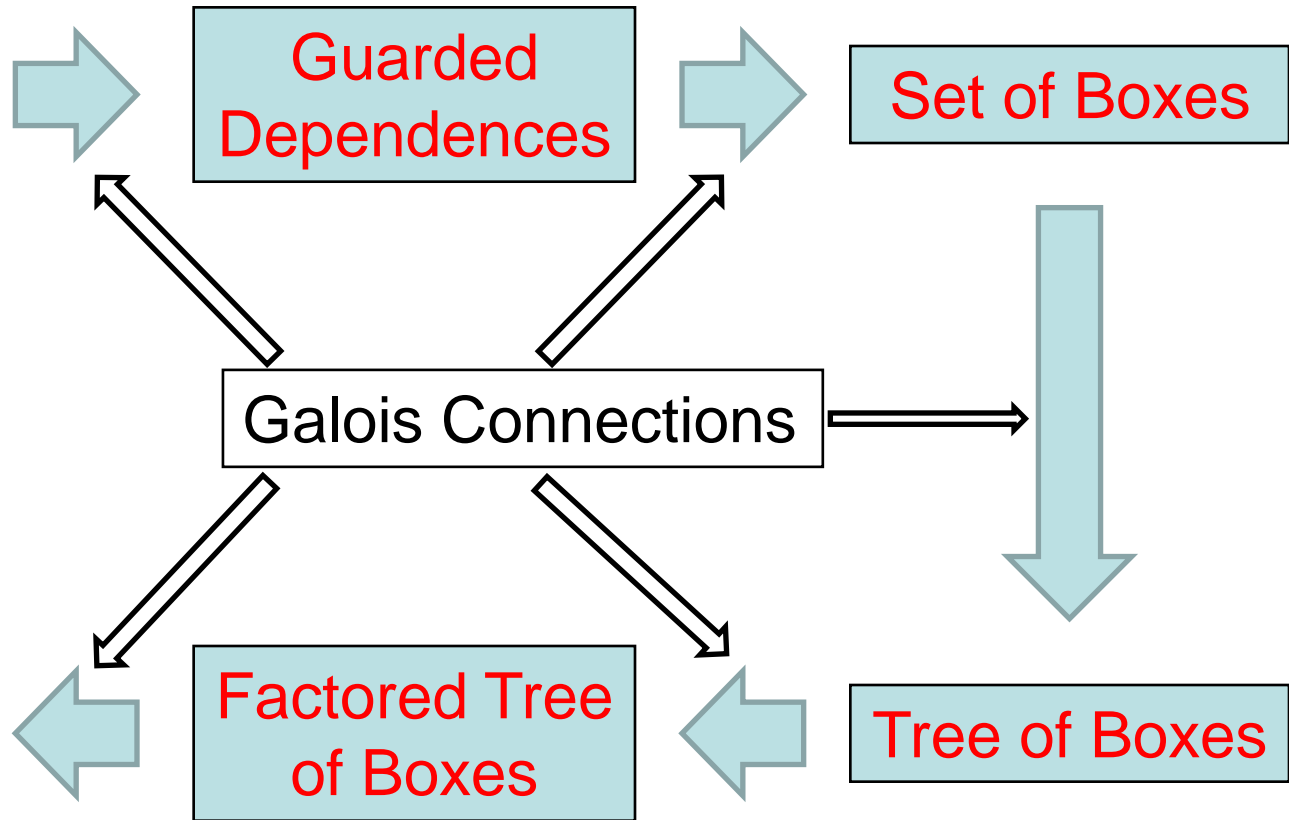
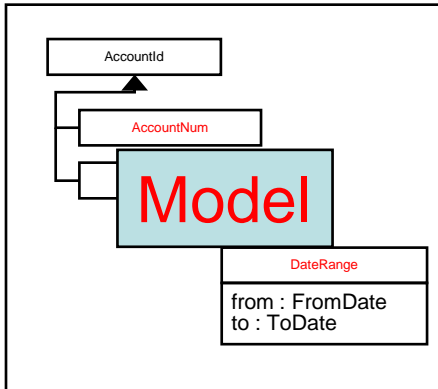
class WithdrawalInfo



Algorithm Overview

```
01 input-record.  
05 ir-trans-code pic x(1).  
05 ir-user-name pic x(8).  
05 ir-acc-num redefines ir-user-  
name.  
05 ir-data pic x(12).  
01 account-rec.  
05 ar-acc-num pic x(6).  
05 ar-acc-name pic x(20).  
05 ar-acc-type pic x(1).  
01 withdraw-info.  
05 wi-acc-num pic x(6).  
05 wi-amount pic x(12).  
05 wi-date pic x(8).  
MOVE ar-acc-num TO ir-acc-num.  
IF ir-trans-code = 'w' THEN  
MOVE ir-data TO withdrawal-info  
...
```

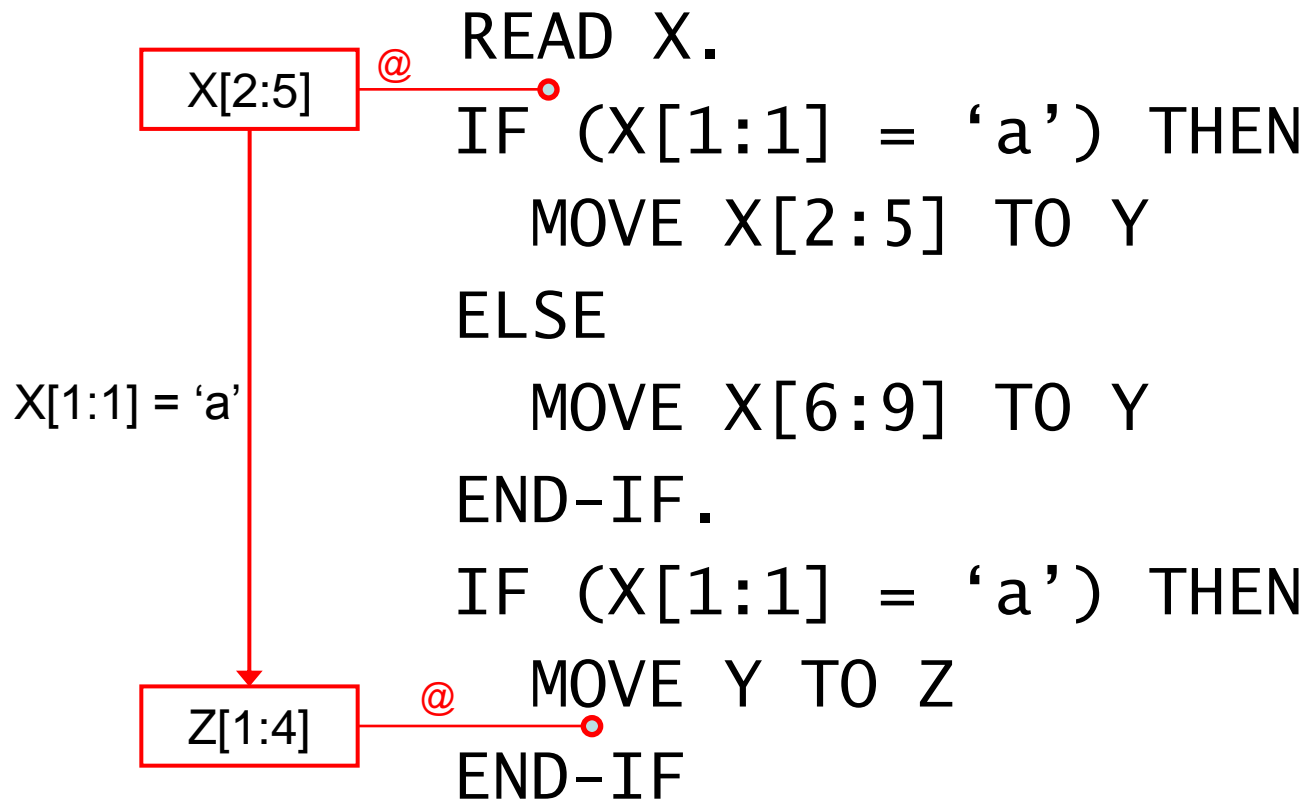
Code



Guarded Dependence Definition

- $\text{guard} \Rightarrow \text{source} \rightarrow \text{target}$
 - source is a pair $\text{range} @ \text{program-pt}$
 - target is similar
 - 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)
- example:
 - $(x = 'a') \Rightarrow y @ \text{pp1} \rightarrow z @ \text{pp2}$
 - $([1:1] = 'a') \Rightarrow [2:5] @ \text{pp1} \rightarrow [6:9] @ \text{pp2}$

Guarded Dependence Example

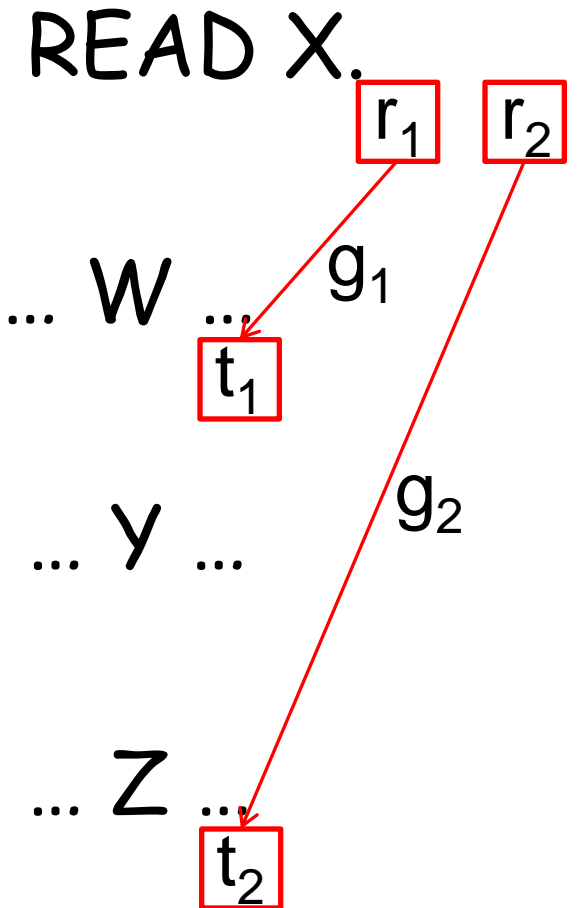


Guarded Dependence Analysis

- *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

Step 2

From Dependences to Boxes

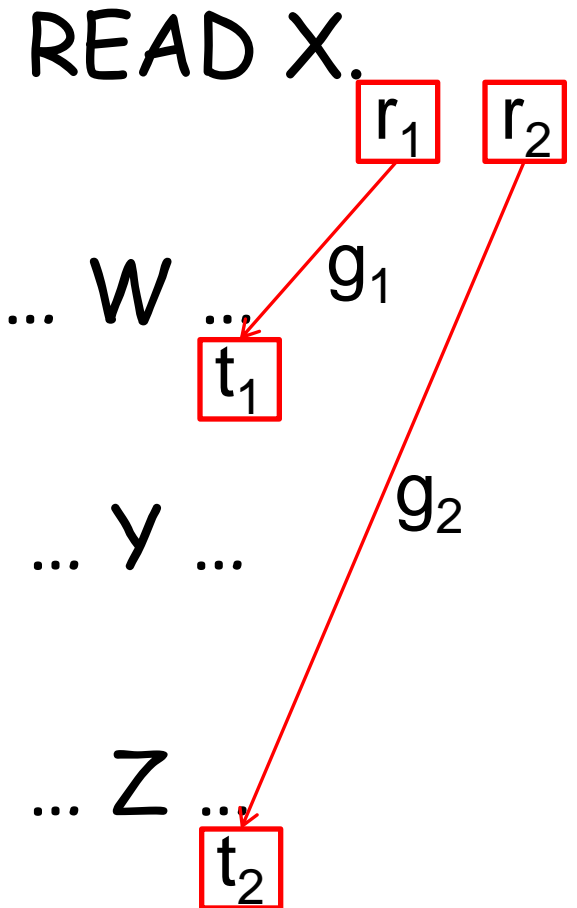


for each data-source:

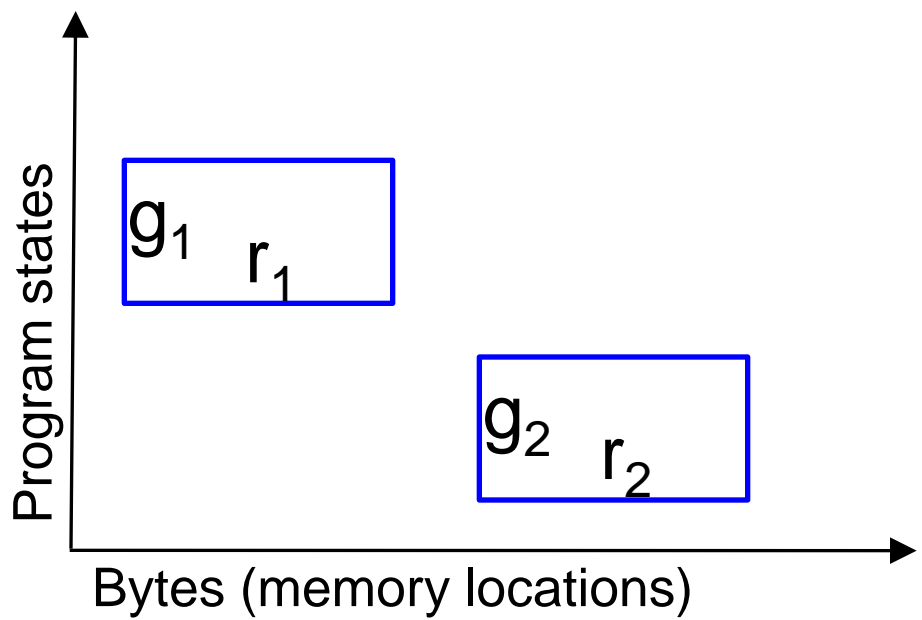
- consider all guarded dependences originating at that source and ending at data-references
- accumulate the set of all **(guard, source-range)** pairs from these dependences

Step 2

From Dependences to Boxes



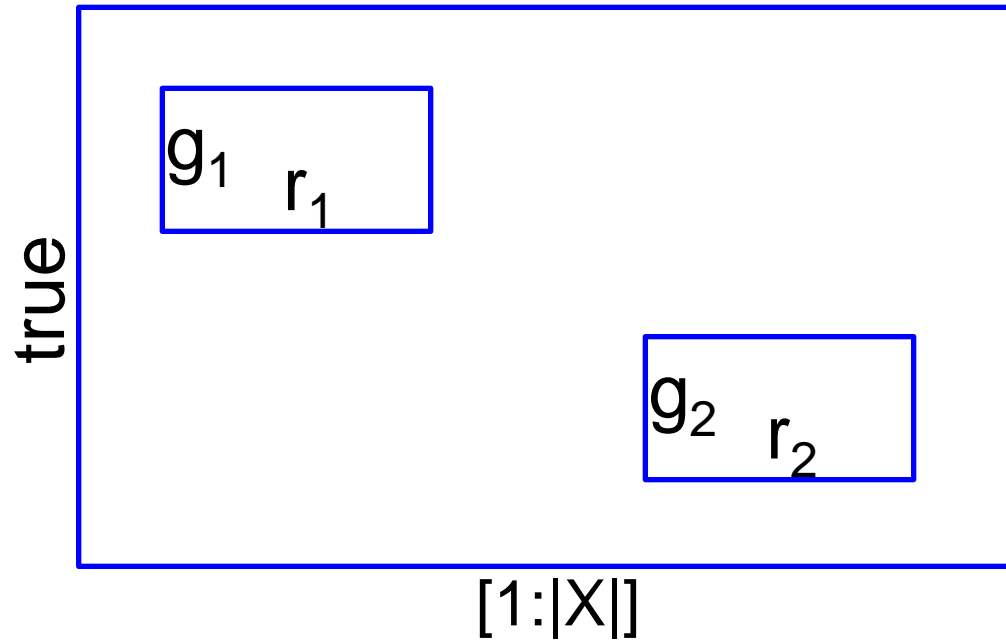
$\{ (g_1, r_1), (g_2, r_2) \}$



Step 3

From a Set of Boxes ...

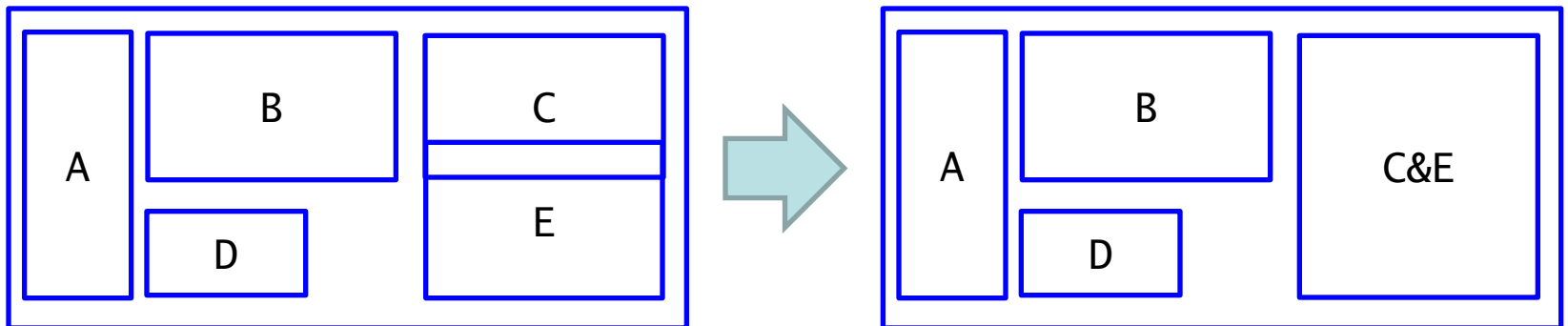
READ X.



Step 3

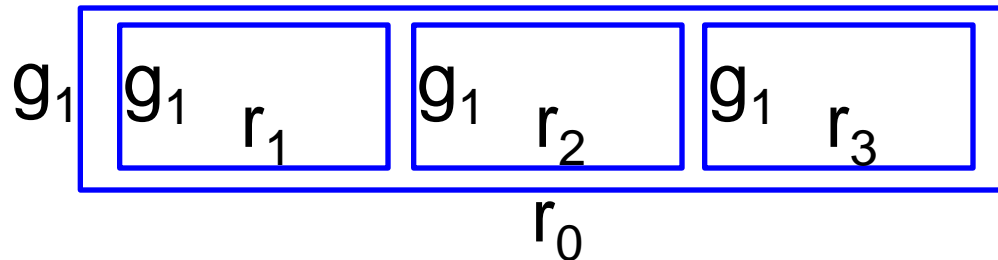
... to a Tree of Boxes

- Ensure boxes are properly nested
 - any two boxes must be disjoint, or one must properly contain the other
 - ensure property by widening box vertically as necessary
 - if proper nesting can not be ensured, the algorithm halts (with failure)



Step 4: Factoring the Tree

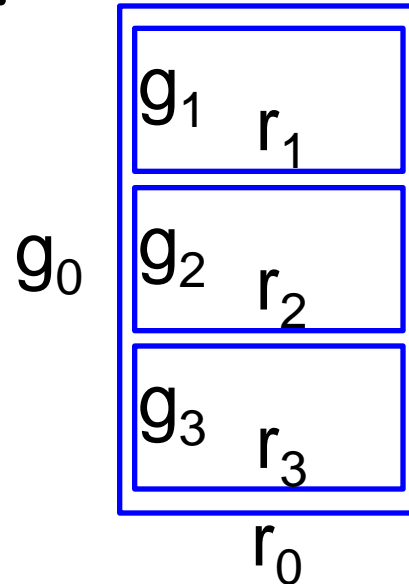
Special Case 1



- All children have same guard as parent
 - parent corresponds to a record (class)
 - children correspond to fields
- (Horizontal Partitioning)

Step 4: Factoring the Tree

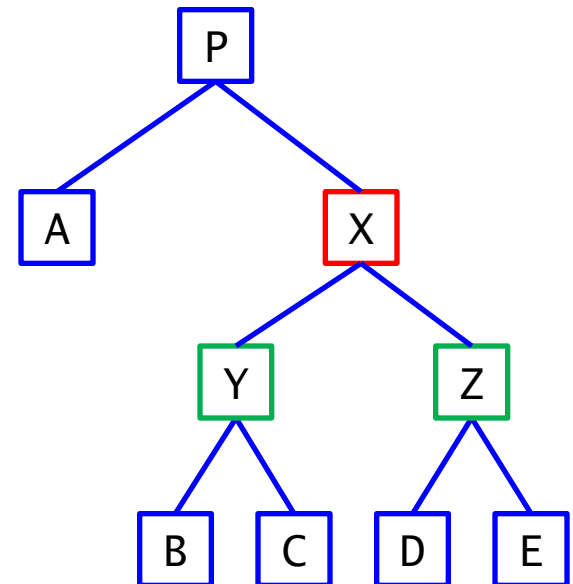
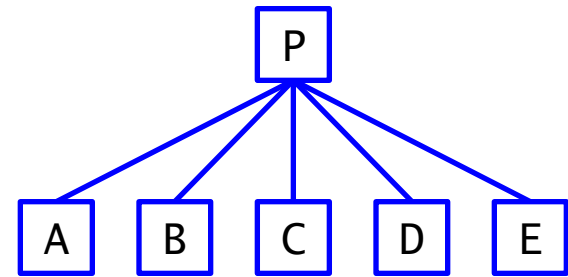
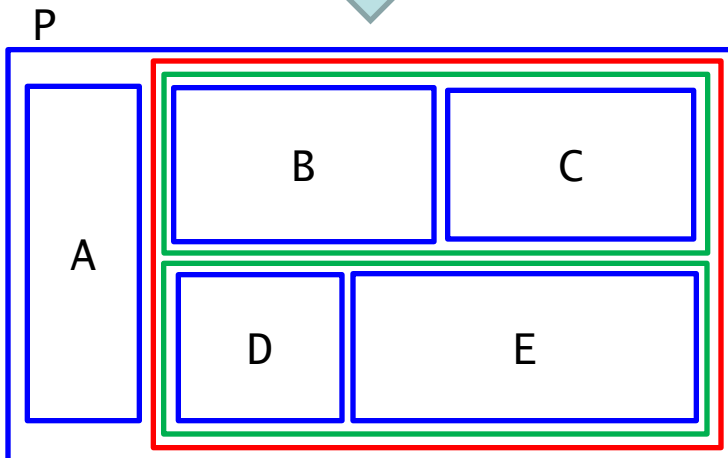
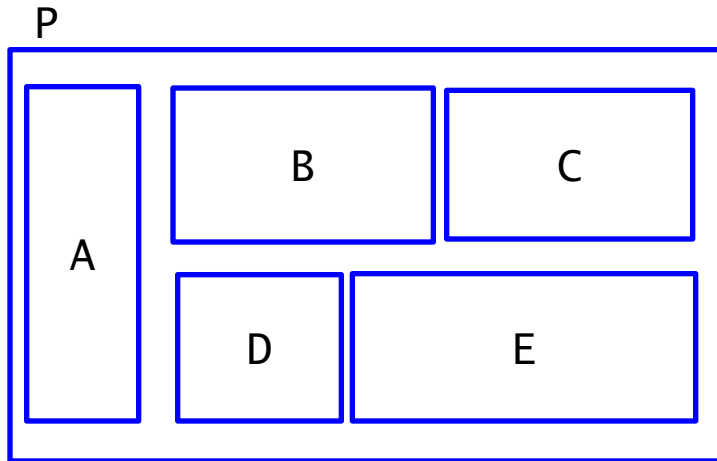
Special Case 2



- All children have mutually disjoint predicates
 - parent corresponds to a base class
 - children correspond to derived classes
- (Vertical Partitioning)

Step 4: Factoring the Tree

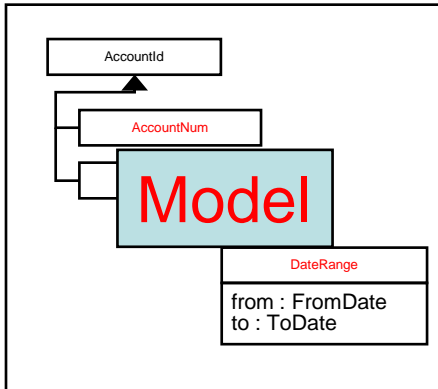
The General Case



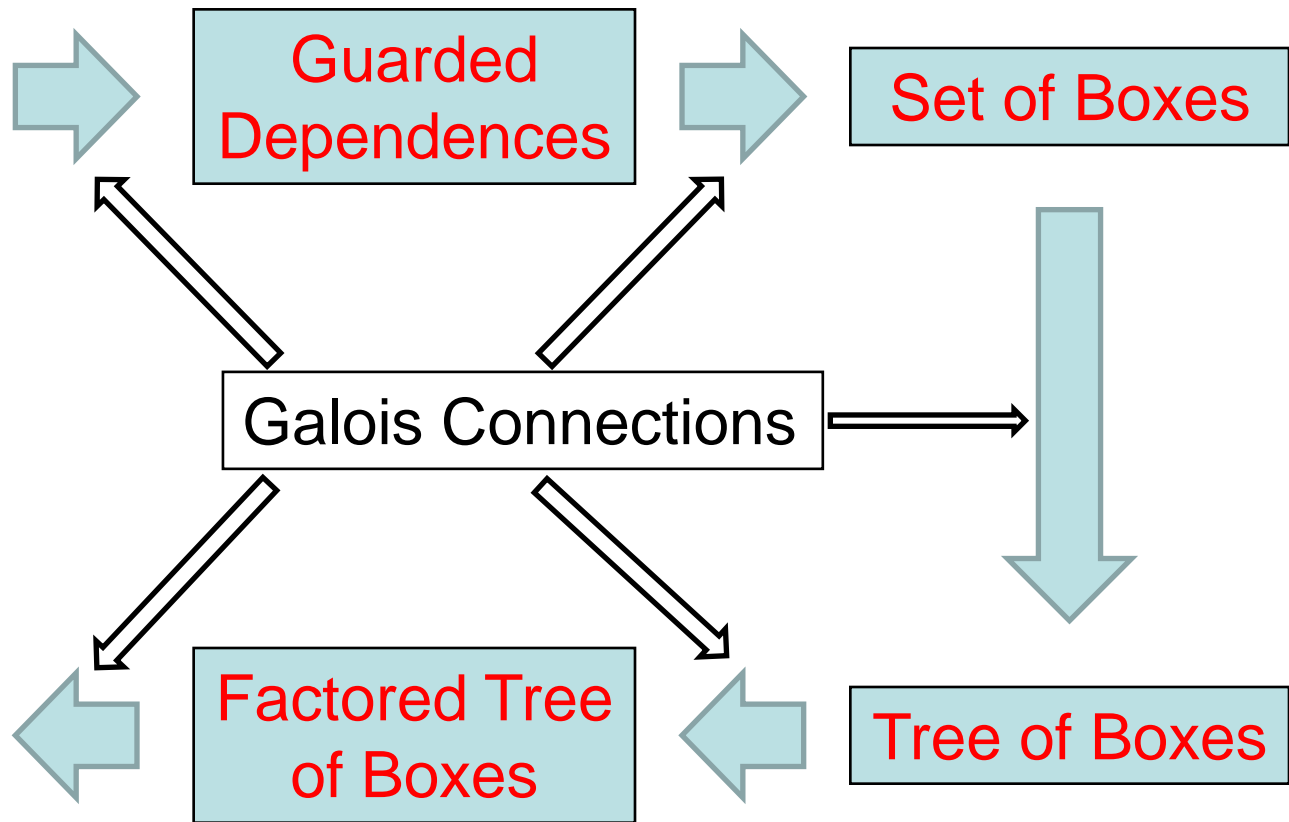
Algorithm Overview

```
01 input-record.  
05 ir-trans-code pic x(1).  
05 ir-user-name pic x(8).  
05 ir-acc-num redefines ir-user-  
name.  
05 ir-data pic x(12).  
01 account-rec.  
05 ar-acc-num pic x(6).  
05 ar-acc-name pic x(20).  
05 ar-acc-type pic x(3).  
01 withdraw-info.  
05 wi-acc-num pic x(6).  
05 wi-amount pic x(12).  
05 wi-date pic x(8).  
MOVE ar-acc-num TO ir-acc-num.  
IF ir-trans-code = 'w' THEN  
MOVE ir-data TO withdrawal-info  
...
```

Code



Model



Guarded Dependences

Set of Boxes

Galois Connections

Factored Tree of Boxes

Tree of Boxes

thank you!