

# Refinement-Based CFG Reconstruction from Unstructured Programs

Sébastien Bardin, Philippe Herrmann, Franck Védrine

CEA LIST  
(Paris, France)

## Automatic analysis of executable files

- recent research field [Codesurfer/x86, SAGE, Jakstab, Osmose, etc.]
- many promising applications (COTS, mobile code, malware, etc.)

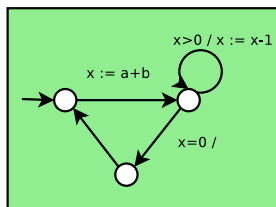
## A key issue : Control-Flow Graph (CFG) reconstruction

- prior to any other static analysis (SA)
- must be safe : otherwise, other SA unsafe
- must be precise : otherwise, other SA imprecise

## This talk is about CFG reconstruction (from executable files)

- safe and precise technique
- based on abstraction-refinement

## Model



## Source code

```
int foo(int x, int y) {  
    int k = x;  
    int c = y;  
    while (c > 0) do {  
        k++;  
        c--;}  
    return k;  
}
```

## Assembly

```
_start:  
    load A 100  
    add B A  
    cmp B 0  
    jle label  
  
label:  
    move @100 B
```

## Executable

```
ABFFF780BD70696CA101001BDE45  
145634789234ABFFE678ABDCF456  
5A2B4C6D009F5F5D1E0835715697  
145FEDBCADACBDAD459700346901  
3456KAHA305G67H345BFFADECAD3  
00113456735FFD451E13AB080DAD  
344252FFAADBDA457345FD780001  
FFF22546ADDAE989776600000000
```

# Binary code analysis is useful !

## Always available

- commercial off-the-shelf software
- mobile code (including malware)
- third-party certification

## Faithful

- optimising compilers and security
- optimising compilers and safety
- What You See Is Not What You eXecute [Reps 04,05]

## Very precise

- worst case execution time, memory consumption, etc.

# BUT binary code analysis is difficult ...

... i.e. more difficult than usual source-code analysis

---

## Low-level semantic of data

- machine arithmetic, bitvector operations
- systematic usage of untyped memory (stack)

## Low-level semantic of control

- no clear distinction between data and control
- no clean encapsulation of procedure calls
- dynamic jumps (goto R0)

No easy (syntactic) recovery of the Control Flow Graph (CFG)

## Diversity of architectures and instruction sets

- each ISA contains dozen of instructions
- lots of engineering work

# BUT binary code analysis is difficult ...

... i.e. more difficult than usual source-code analysis

## Low-level semantic of data

- machine arithmetic, bitvector operations
- systematic usage of untyped memory (stack)

## Low-level semantic of control

- no clear distinction between data and control
- no clean encapsulation of procedure calls
- dynamic jumps (goto R0)

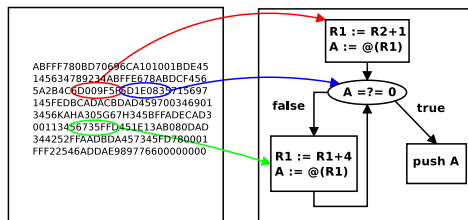
No easy (syntactic) recovery of the Control Flow Graph (CFG)

## Diversity of architectures and instruction sets

- each ISA contains dozen of instructions
- lots of engineering work

## Input

- an executable file, i.e. an array of bytes
- the address of the initial instruction
- a basic decoder :  $\text{exec f.} \times \text{address} \mapsto \text{instruction} \times \text{size}$



Output : CFG of the program

Successor addresses are often syntactically known

- `< addr : move a b >` →
- `< addr : goto 100 >` →
- `< addr : ble 100 >` →



## Successor addresses are often syntactically known

- `⟨ addr : move a b ⟩` → successor at `addr+size`
- `⟨ addr : goto 100 ⟩` → successor at `100`
- `⟨ addr : ble 100 ⟩` → successors at `100` and `addr+size`

Successor addresses are often syntactically known

- `⟨ addr : move a b ⟩` → successor at `addr+size`
- `⟨ addr : goto 100 ⟩` → successor at `100`
- `⟨ addr : ble 100 ⟩` → successors at `100` and `addr+size`

But not always : successors of `⟨addr : goto a ⟩`?

Successor addresses are often syntactically known

- `< addr : move a b >` → successor at `addr+size`
- `< addr : goto 100 >` → successor at `100`
- `< addr : ble 100 >` → successors at `100` and `addr+size`

But not always : successors of `<addr : goto a >`?

Dynamic jump is the enemy !

Dynamic jumps are pervasive : introduced by compilers

- `switch`, function pointers, virtual methods, etc.

... current industrial practise ...

---

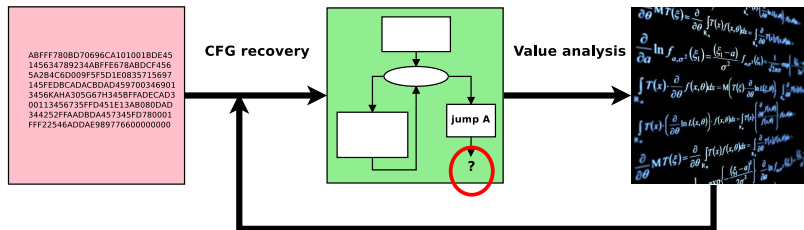
## Linear sweep decoding [brute force]

- decode instructions at each code address
- miss every “dynamic” edge of the CFG
- may still miss instructions [too optimistic hypothesisises]

## Recursive traversal

- decode recursively from entry point, stop on dynamic jump
- miss large parts of CFG

## VA and CFG reconstruction must be interleaved



**Very difficult to get precise :** imprecision on jumps  $\rightarrow$  extra propagation on false targets  $\rightarrow$  more imprecision on value analysis  $\rightarrow$  possibly more imprecision on jumps  $\rightarrow$  ...

## CodeSurfer/x86 [Balakrishnan-Reps 04,05,07,...]

- abstract domain : strided intervals (+ affine relationships)
- imprecise : abstract domain not suited to sets of jump targets (arbitrary values from compiler)
- in practice many false targets

## Jakstab [Kinder-Veith 08,09,10]

- abstract domain : sets of bounded cardinality ( $k$ -sets)
- precise when the bound  $k$  is well-tuned
- not robust to the parameter  $k$  : possibly inefficient if  $k$  too large, very imprecise if  $k$  not large enough

## Key observations

- k-sets are the only domain well-suited to **precise** CFG reconstruction
  - for most programs, only a few facts need to be tracked precisely to resolve dynamic jumps
  - **good candidate for abstraction-refinement**
- 

## Contribution [VMCAI 2011]

- A refinement-based approach to safe CFG reconstruction
- An implementation and a few experiments
- The technique is **safe, precise, robust** and **reasonably efficient**

**Unstructured Programs** :  $P = (L, V, A, T, l_0)$  where

- $L \subseteq \mathbb{N}$  finite set of code addresses
- $V$  finite set of program variables,  $A$  finite set of arrays
- $T$  maps code addresses to instructions
- $l_0$  initial code address
- instructions : assignments  $v := e$  and  $a[e_1] := e_2$ , static jumps `goto l`, branching instructions `ite(cond, l1, l2)`, dynamic jumps `cgoto(v)`

---

**Problem** : compute an invariant of  $P$  such that no dynamic target evaluates to  $\top$ , or fail

- do not fail “too often”
- do not add “too many” false targets



# Sketch of the procedure

abstract domain = k-sets

k-set cardinality bounds are local to each location

- gain efficiency through loss of precision
- still a global bound  $K_{max}$  over local bounds

procedure : propagate forward until a dynamic target expression evaluates to  $\top$ , then try to refine the domain to avoid this  $\top$  value

- domain refinement = increase some k-set cardinality bounds
- if no domain update then fail, else restart propagation with new domains

## For each target evaluating to $\top$

- follows backward data dependencies
- only interested in  $\top$ -values (other locations are safe until now)
- only interested in correcting **initial causes of precision loss**

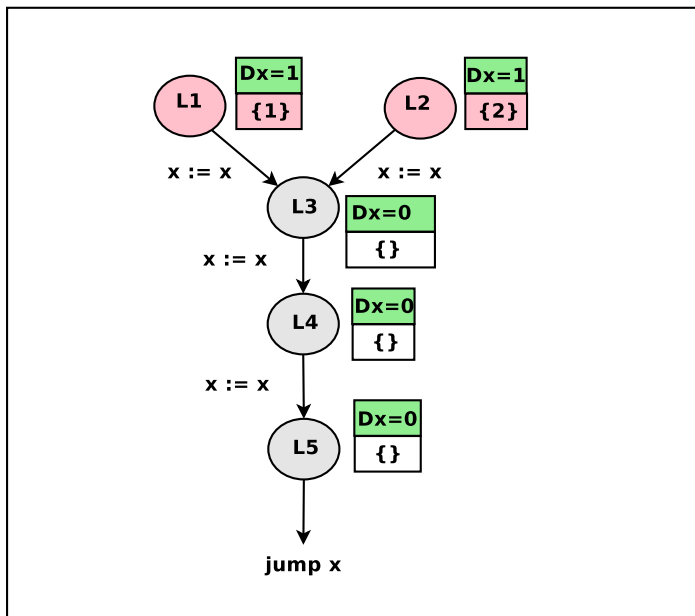
## Finding the initial causes of precision loss

- add **tags** to  $\top$ -values, recording origin :  $\top, \top_{init}, \top_{\langle c_1, \dots, c_n \rangle}$
- initial causes of precision loss are of the form  $\top_{init}, \top_{\langle c_1, \dots, c_n \rangle}$

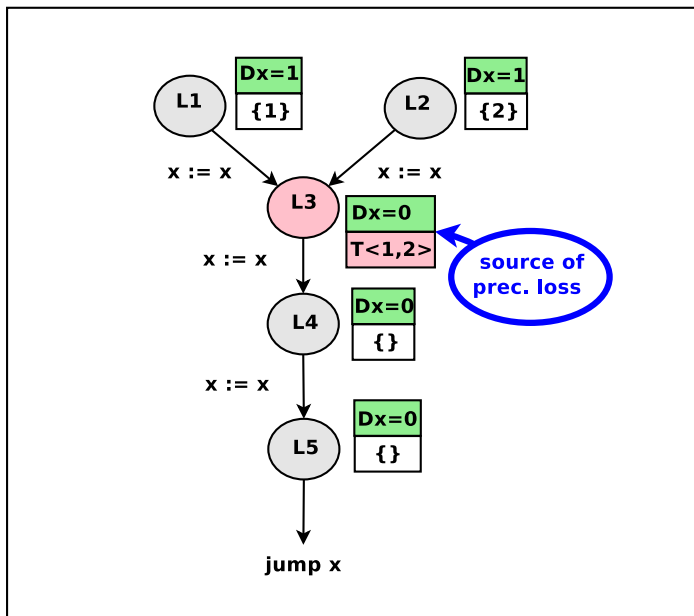
## How to correct

- $\top_{init}$  cannot be avoided
- $\top_{\langle c_1, \dots, c_n \rangle}$  may be avoided if  $n \leq Kmax$  (set local bound to  $n$ )

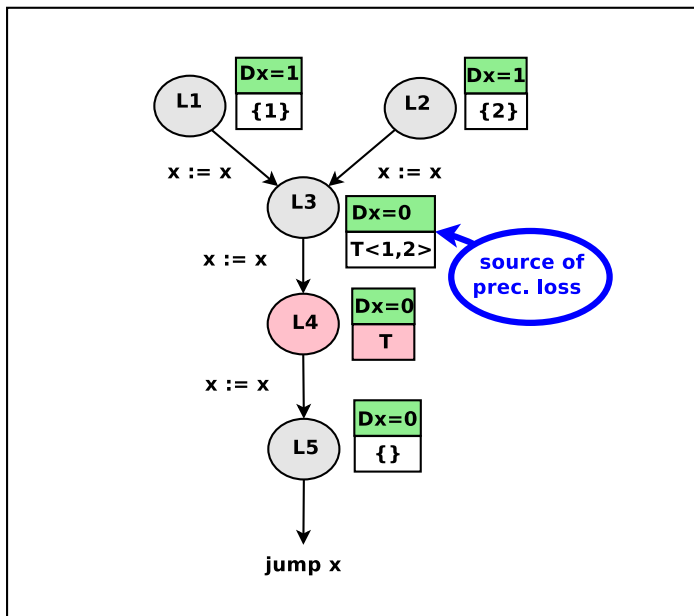
# Example



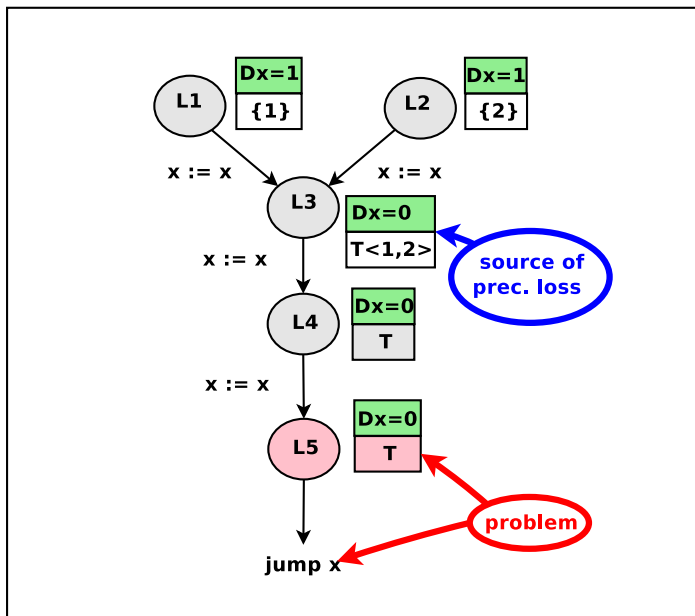
# Example



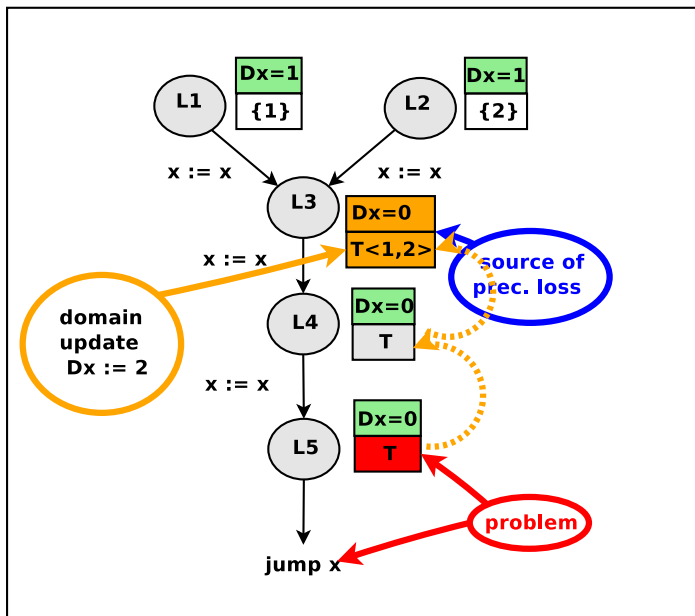
# Example



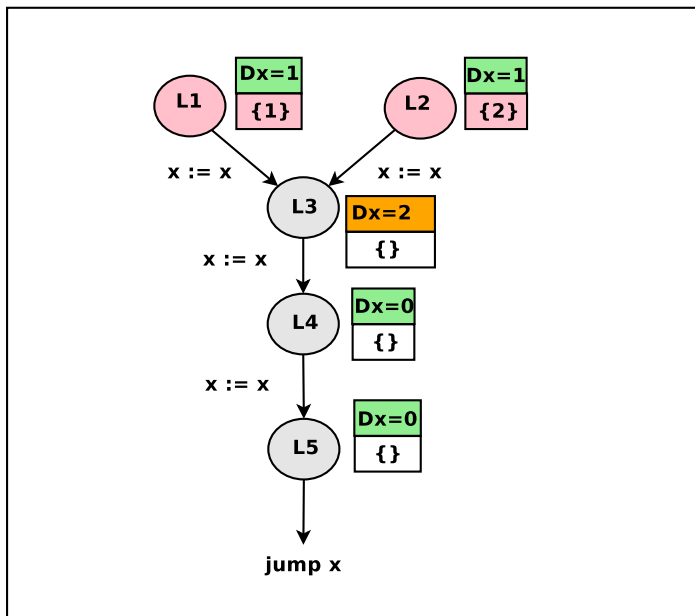
# Example



# Example

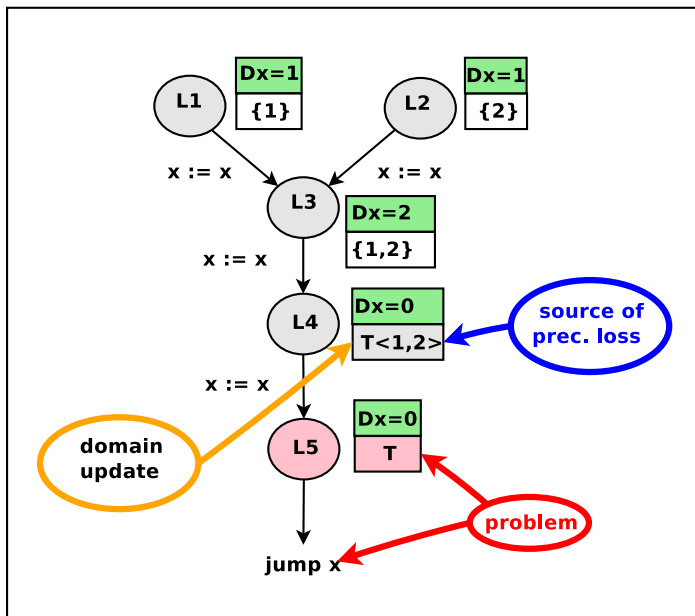


# Example

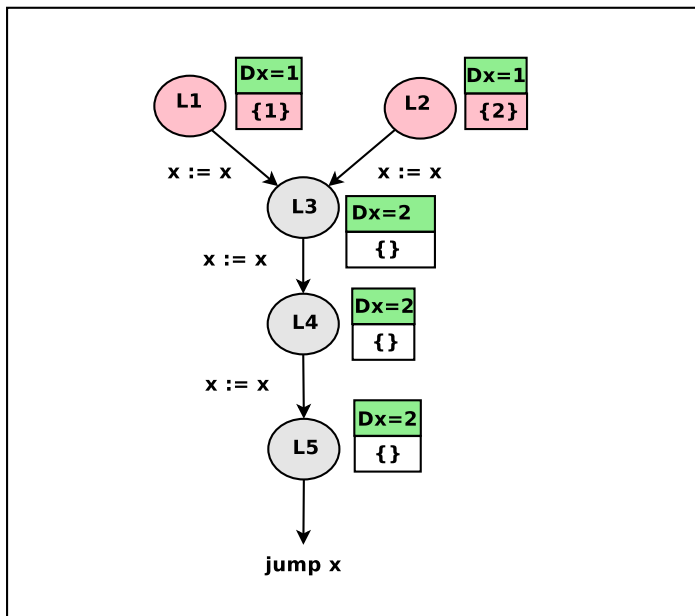




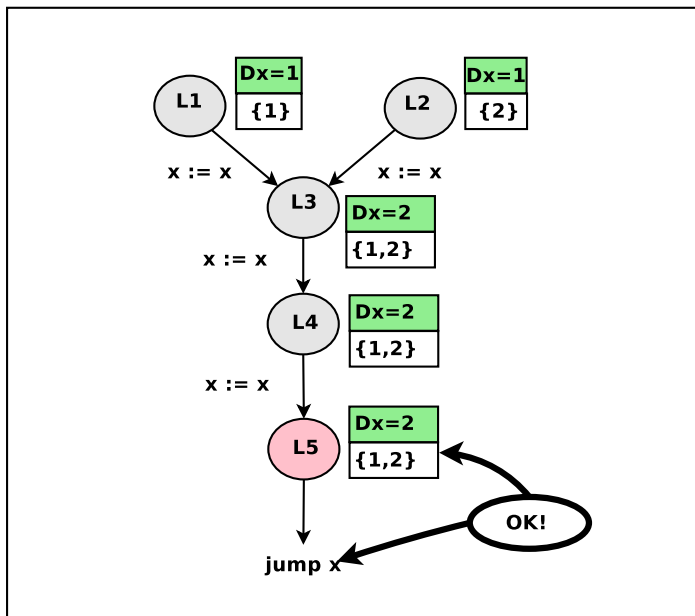
# Example



# Example



# Example



Too much refinement = inefficiency

## A journal of the forward propagation

- record observed feasible branches / alias / dynamic targets
  - prune backward data dependencies accordingly
- 

## Two possible failure policies during refinement

- optimistic : fails only when no local domain is corrected
- pessimistic : fails as soon as one “cause of precision loss” cannot be corrected

**Soundness** : returns either FAIL or an invariant such that no jump target evaluates to  $\top$

**Complexity** : polynomial number of refinements

**Precision** : perfect relative precision for a non trivial subclass of programs (see next)

# What about precision ?

**Relative completeness (RC)** : PaR is relatively complete if  $\text{PaR}(P, Kmax)$  returns successfully when the forward k-set propagation with parameter  $Kmax$  does

---

**Bad news** : no RC in the general case

- mainly because of control dependencies

**Good news** : RC for a non trivial subclass of programs

- non deterministic branches [new : only feasible branches]
- guarded aliases
- restricted class of operators :  $+$ ,  $-$ ,  $\times k$  ok, but not  $\times$
- RC even for the procedure with “pessimistic failure”

**Implementation** : CFG reconstruction from 32-bit PowerPC (PPC)

**Bench** : Safety critical program from Sagem

- 32 kloc, 51 dynamic jumps, up to 16 targets a jump
- 

## Results

- **precision** : resolve every jump, only 7% of false targets  
( standard program analysis cannot recover better than between 400% and 4000% of false targets )
- **robustness** : results independent of  $K_{max}$  (if large enough)
- **locality** : tight value of max- $k$ , low value of mean- $k$

Terminates in 18 min [ $\leq 5$  min now]

- ok for a preliminary implementation
- already sufficient for some industrial application
- however (as expected) procedure inlining is an issue

1x - 3x faster than adequate k-set propag

3x - 5x faster than iterated k-set propag

- we expected more gap
- lots of redundant work from one refinement step to the other
- can probably be improved



We investigate safe CFG reconstruction from executable files

## Results

- an original refinement-based procedure
- safe, precise, robust and reasonably efficient
- both theoretical and empirical evidence

## Future work

- improve efficiency [inlining, redundant work]
- experiments on non-critical programs [dynamic alloc]
- ultimate goal : executables coming from large C++ programs

## Relative completeness : why it does not work (general case)

let us suppose  $K_{max} = 1$

1.  $x := 1$ , goto 2
2. if  $x == 1$  then goto 3 else goto 4
3.  $t := 100$ , goto 5
4.  $t := 200$ , goto 5      // dead code
5. jump t

## Relative completeness : why it does not work (general case)

let us suppose  $K_{max} = 1$

1.  $x := 1$ , goto 2 //  $x = \top$
2. if  $x == 1$  then goto 3 else goto 4 //  $x = \{1\}$
3.  $t := 100$ , goto 5 //  $x = \{1\}$
4.  $t := 200$ , goto 5 //  $x = \perp$  // dead code
5. jump t //  $t = \{100\}$

Forward propagation with  $K_{max} = 1$  succeeds.

## Relative completeness : why it does not work (general case)

let us suppose  $K_{max} = 1$

1.  $x := 1$ , goto 2 //  $x = T$
2. if  $x == 1$  then goto 3 else goto 4 //  $x = T$
3.  $t := 100$ , goto 5 //  $x = T$
4.  $t := 200$ , goto 5 //  $x = T$  // dead code
5. jump t //  $t = T_{\langle 100, 200 \rangle}$

Forward propagation with  $K_{max} = 1$  succeeds.

Our procedure fails :

- believes that  $(5, t)$  can take at least values  $\{100, 200\}$
- do not notice that else branch infeasible

## Relative completeness : why it works (restricted class)

- $\text{KSET}(k)$  is as precise as  $\text{KSET}(K_{\max})$ , as long as there is no  $\top$ -cast
- loss of relative precision happens only because of  $\top$ -cast
- ⇒ on the restricted subclass, as long as no alias / jump evaluates to  $\top$ ,  $\text{KSET}(k)$  and  $\text{KSET}(K_{\max})$  computes the same proper k-sets
- ⇒ same aliases and same dynamic targets (if proper k-sets)

Actually, more powerful than RC ...