# Binary Program Rewriting with Diablo

## Bjorn De Sutter

## Ghent University

PLDI06 , June 06, 2006

# Credits

- Bruno De Bus
- Dominique Chanet
- Ludo Van Put
- Matias Madou
- Bertrand Anckaert
- Koen De Bosschere

# Overview



- Some background
- Diablo
  - Extensibility
  - Retargetability
  - Reliability

# Overview



- INTRODUCTION (45 min)
- DATASTRUCTURES (1 hr)
- ANALYSES AND TRANSFORMATIONS (45 min)
- BACKENDS (30 min)

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

p. 4

# Program Development



programming → **Source code** → compiler → **Assembly code** → assembler → **Code 1** / **Data 1** / **Meta 1** → linker → **Executable**

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

UNIVERSITEIT GENT

# Linking

**Linker**

| | | |
|---|---|---|
| Code 1 | | Code 1 |
| Code 2 | | Code 2 |
| Code 3 | | Code 3 |
| | | Data 1 |
| Data 1 | | Data 2 |
| Data 2 | | Data 3 |
| Data 3 | | |

Meta 1

Meta 2

Meta 3

**start print stop**

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

p. 6

UNIVERSITEIT GENT

# Program Development

programming → **Source code**

**Source code** → compiler

compiler → **Assembly code**

**Assembly code** → assembler

assembler → **Code 1** / **Data 1** / **Meta 1**

**Code 1** / **Data 1** / **Meta 1** → linker

linker → **Executable**

**Executable** → Squeeze++

Squeeze++ → **Compacted Executable**

UNIVERSITEIT GENT

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

p. 7

# Additional optimization opportunities?

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
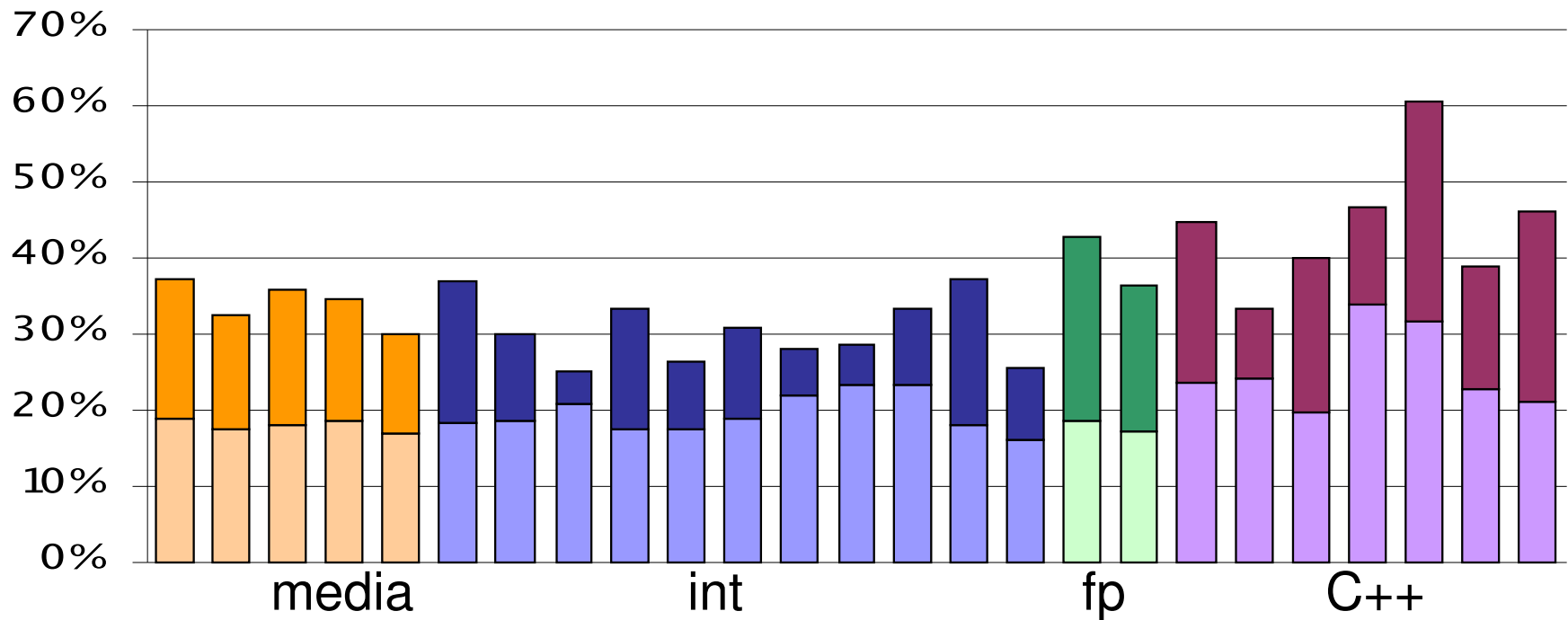Engineering Sciences Faculty – Electronics and Information Systems Department

UNIVERSITEIT GENT

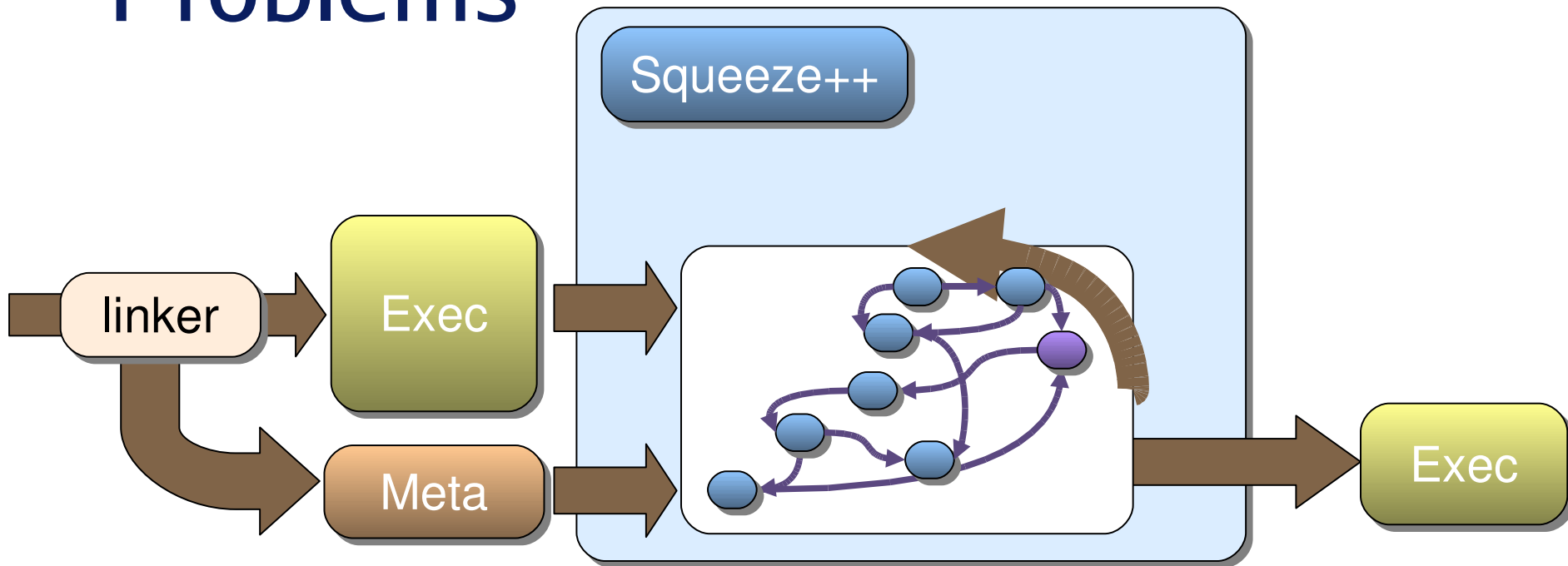# Results (Squeeze++)

Code size reduction obtained with optimization and code abstraction



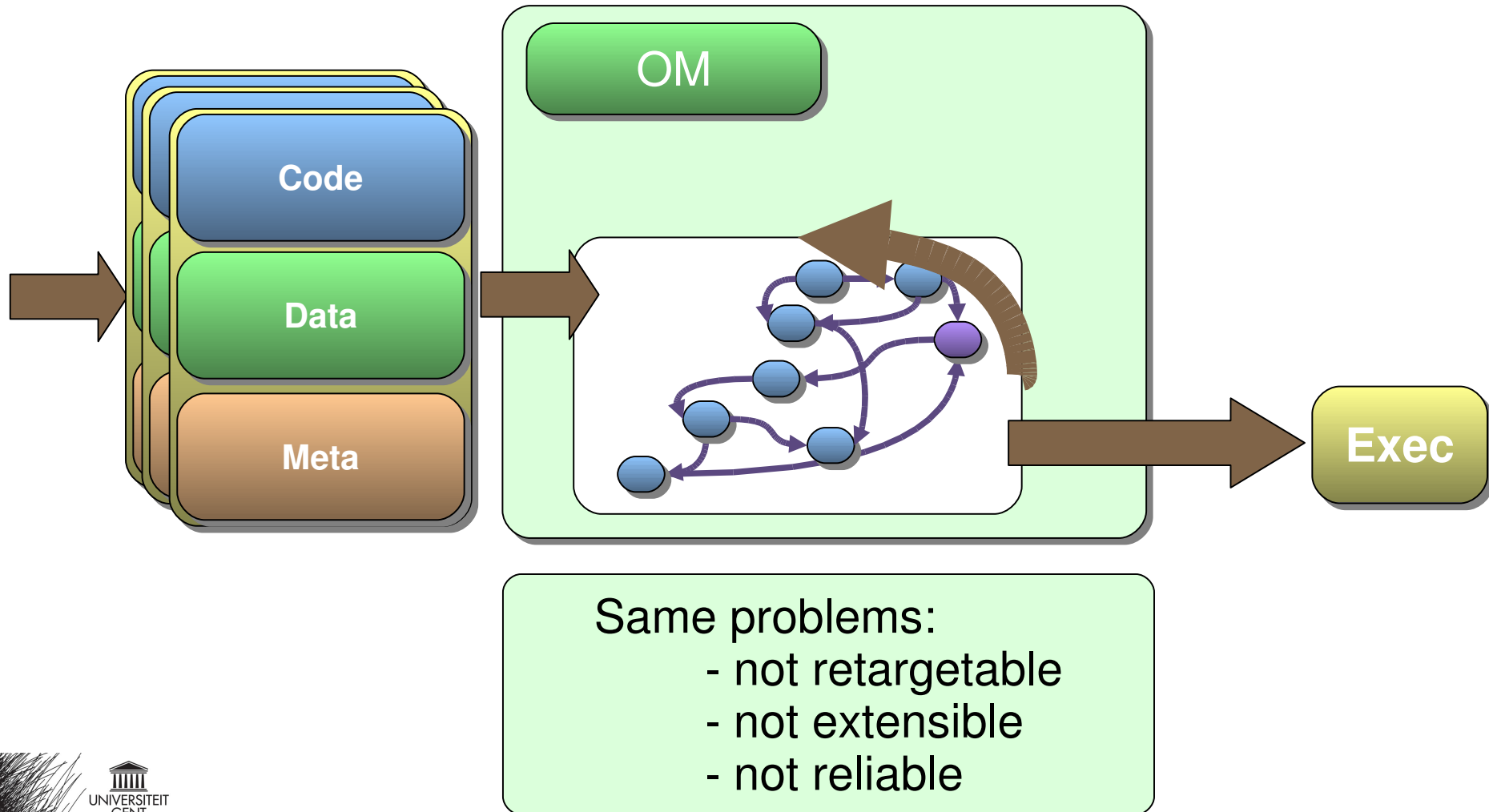[De Sutter, De Bus and De Bosschere, ACM TOPLAS, Sept 2005]

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

UNIVERSITEIT
GENT

# Problems

Squeeze++
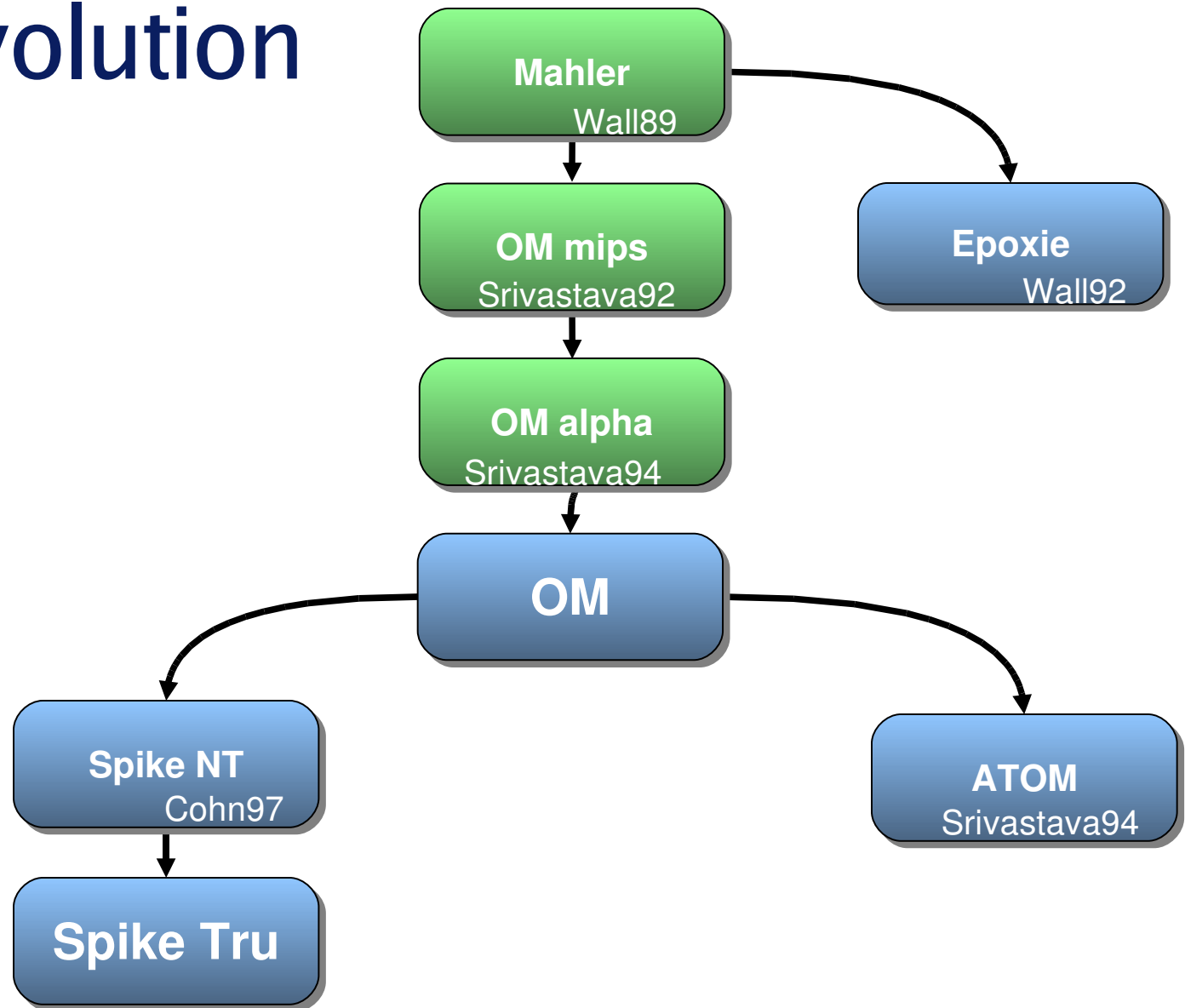
linker → Exec

linker → Meta

Exec

Only for the Alpha architecture.
Only for compaction/optimization
Small change implies days of debugging
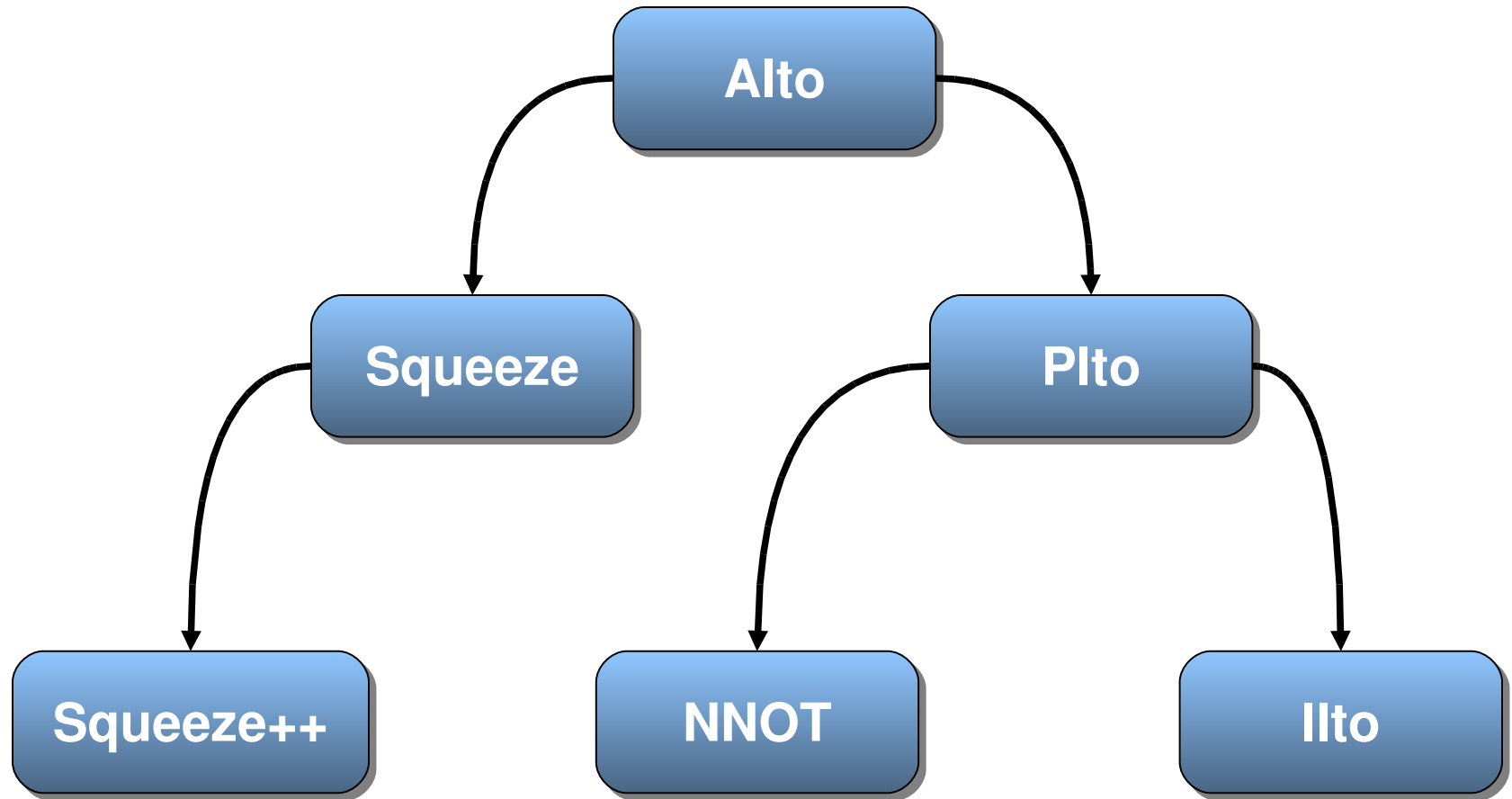
Not retargetable
Not extensible
Not reliable

UNIVERSITEIT GENT

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

p. 10

# Other rewriters?

**OM**

**Code**

**Data**

**Meta**

**Exec**

Same problems:
- not retargetable
- not extensible
- not reliable

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

UNIVERSITEIT GENT

# OM's evolution

```
Mahler
Wall89
   │
   ├──────────────────────────► Epoxie
   │                             Wall92
   ▼
OM mips
Srivastava92
   │
   ▼
OM alpha
Srivastava94
   │
   ▼
  OM
 ┌──┴──┐
 ▼     ▼
Spike NT    ATOM
Cohn97      Srivastava94
 │
 ▼
Spike Tru
```

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

UNIVERSITEIT
GENT

# Alto's Evolution

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

# static binary rewriting is useful...

## Applications

- Optimization, compaction
- Instrumentation
- Obfuscation

- Program understanding, visualisation
- Debugging
- ...

# but it is a bit problematic...

## Problems

- Not retargetable
- Not extensible
- Not reliable

UNIVERSITEIT GENT

# Overview

- Some background
- Diablo
  - Extensibility
  - Retargetability
  - Reliability

UNIVERSITEIT
GENT

# Pre-link or post-link?
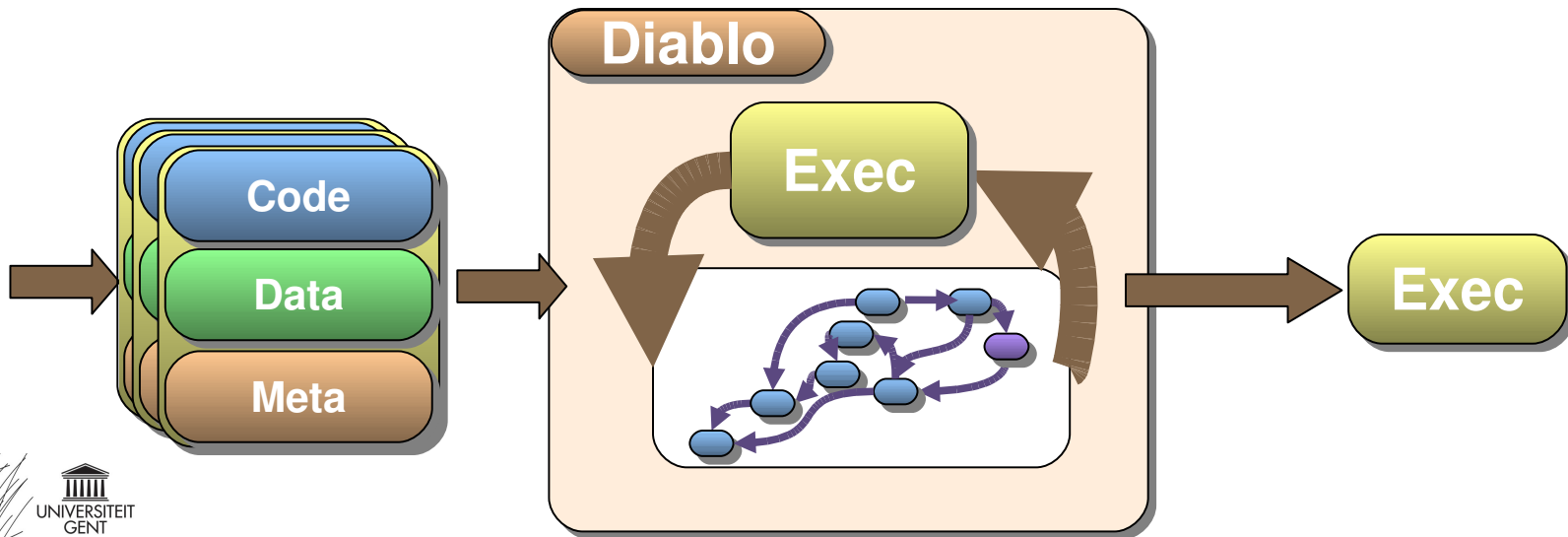
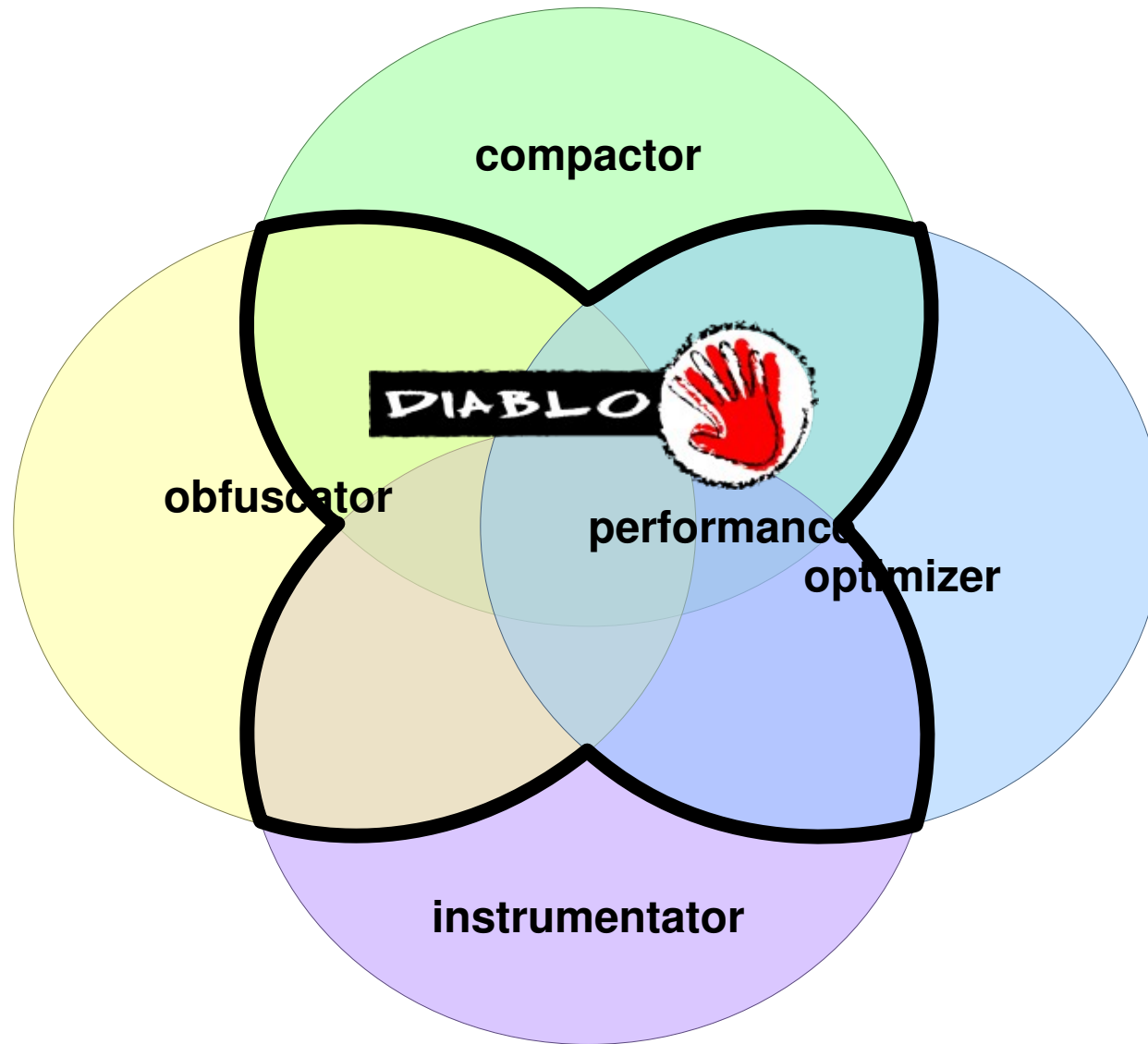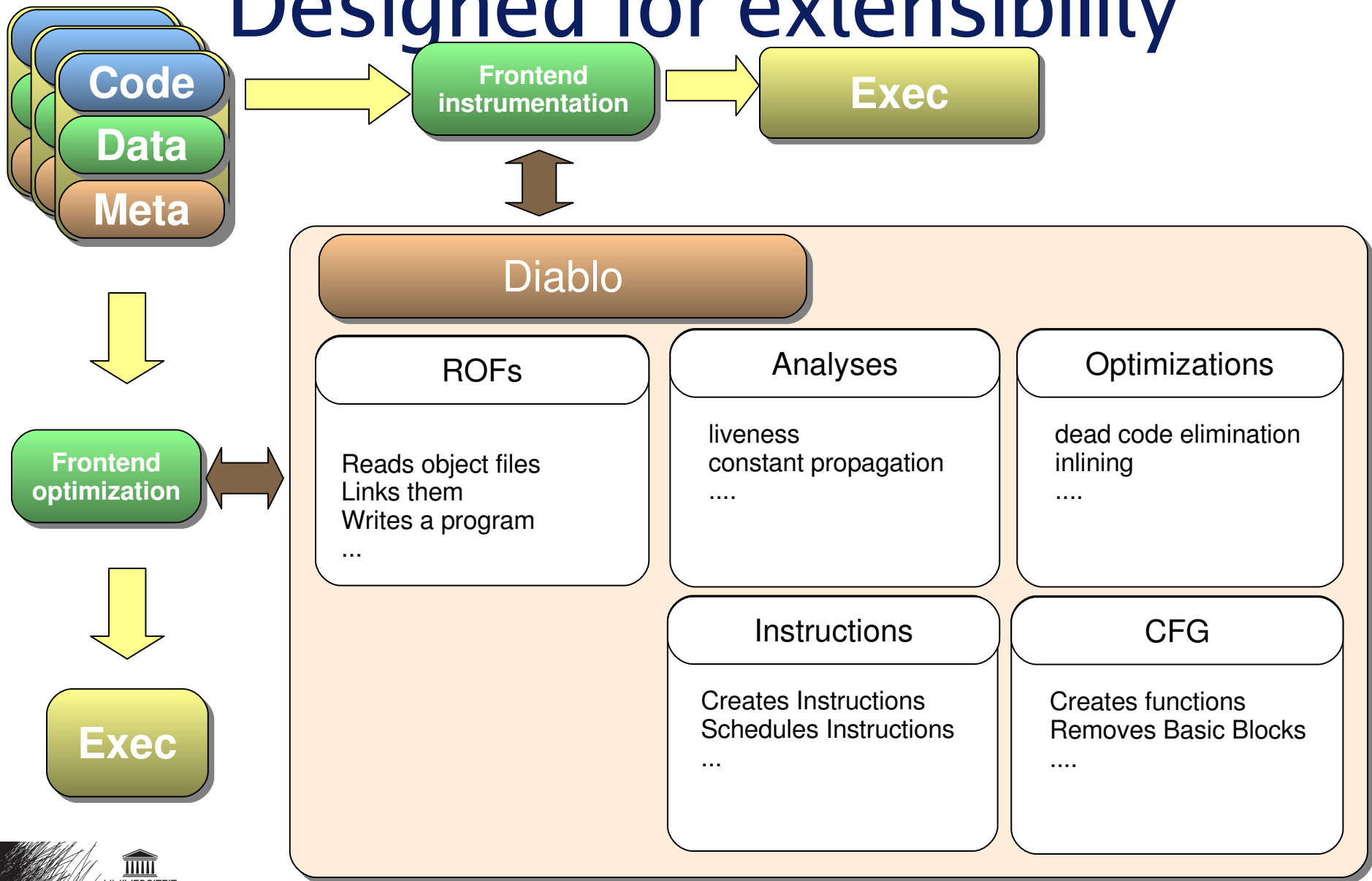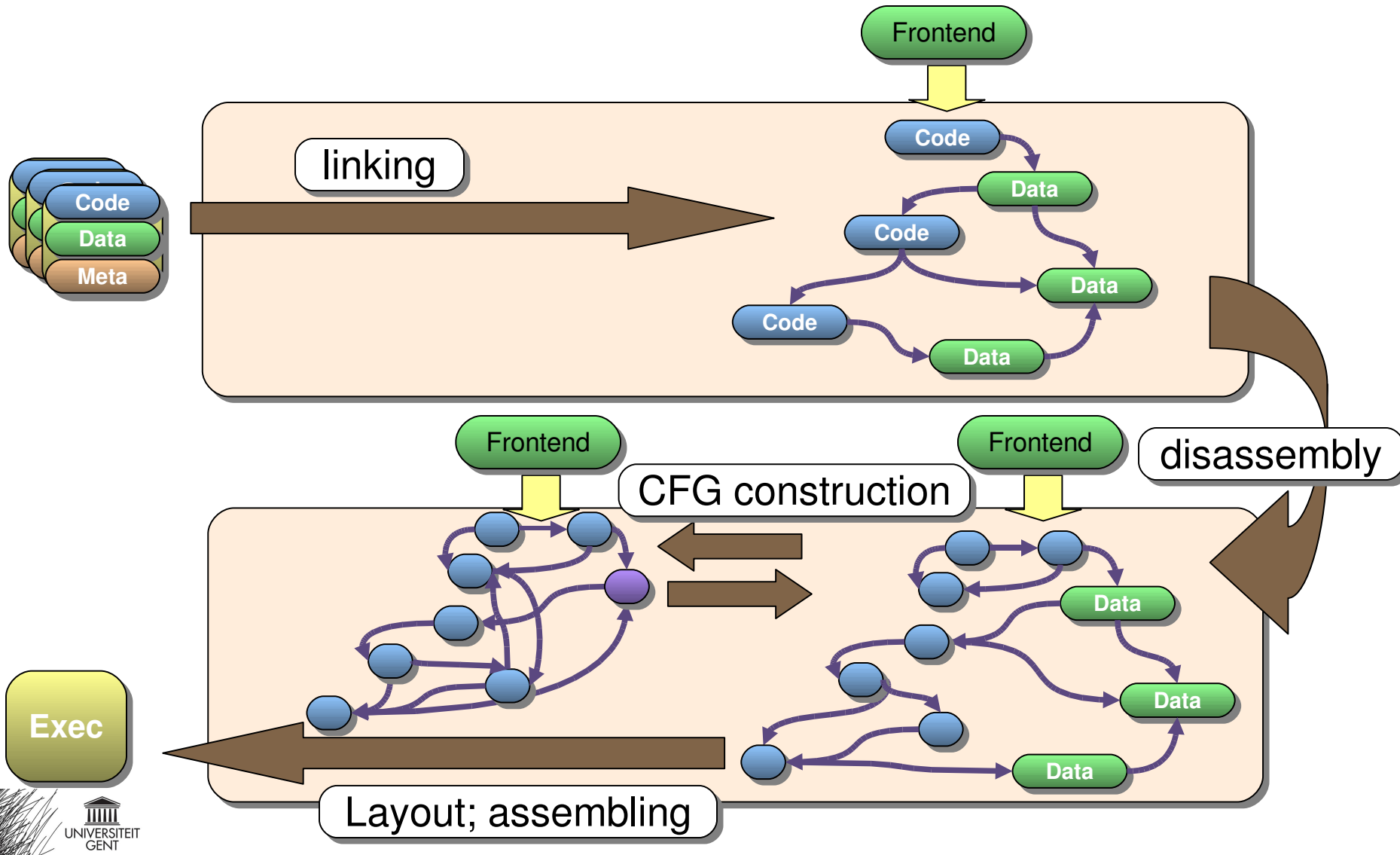| Pre-link | Link - time | Post-link |
|---|---|---|
| **+ More meta information**<br>= more aggressive transformations<br><br>**- No program overview** | **+ More meta information**<br>**+ Whole-program overview**<br><br>**- More implementation work** | **- Less meta information**<br>= more conservative transformations<br><br>**+ Whole-program overview** |

**Code**

**Data**

**Meta**

**Diablo**

**Exec**

**Exec**

# Extensibility – The problem

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

# Designed for extensibility

**Code**
**Data**
**Meta**

**Frontend instrumentation** → **Exec**

**Frontend optimization**

**Exec**

## Diablo

### ROFs

Reads object files
Links them
Writes a program
...

### Analyses

liveness
constant propagation
....

### Optimizations

dead code elimination
inlining
....

### Instructions

Creates Instructions
Schedules Instructions
...

### CFG

Creates functions
Removes Basic Blocks
....

UNIVERSITEIT GENT

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

p. 18

# Operation at different levels

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

# Extensibility

**Diablo** — **Application optimization and compaction frontend** — LCTES'04

**FIT** — **Instrumentation frontend** — PASTE'04

**Stilo** — **Steganography frontend** — ICISC'04

**kDiablo** — **Linux kernel specialization frontend** — LCTES'05

**Lancet** — **Interactive binary program editor** — PASTE'05

**Loco** — **Interactive program obfuscator** — PEPM'05

UNIVERSITEIT GENT

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

p. 20
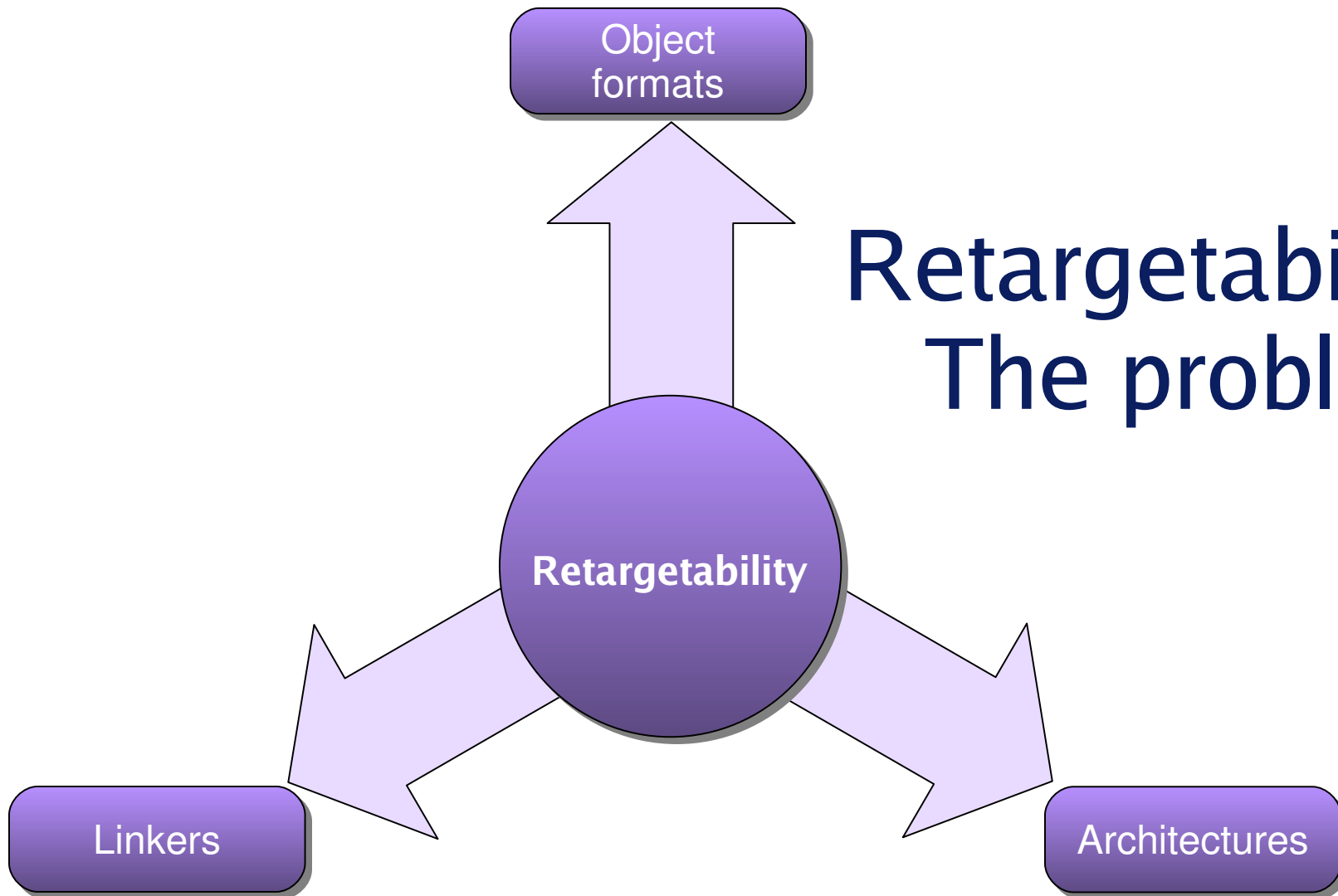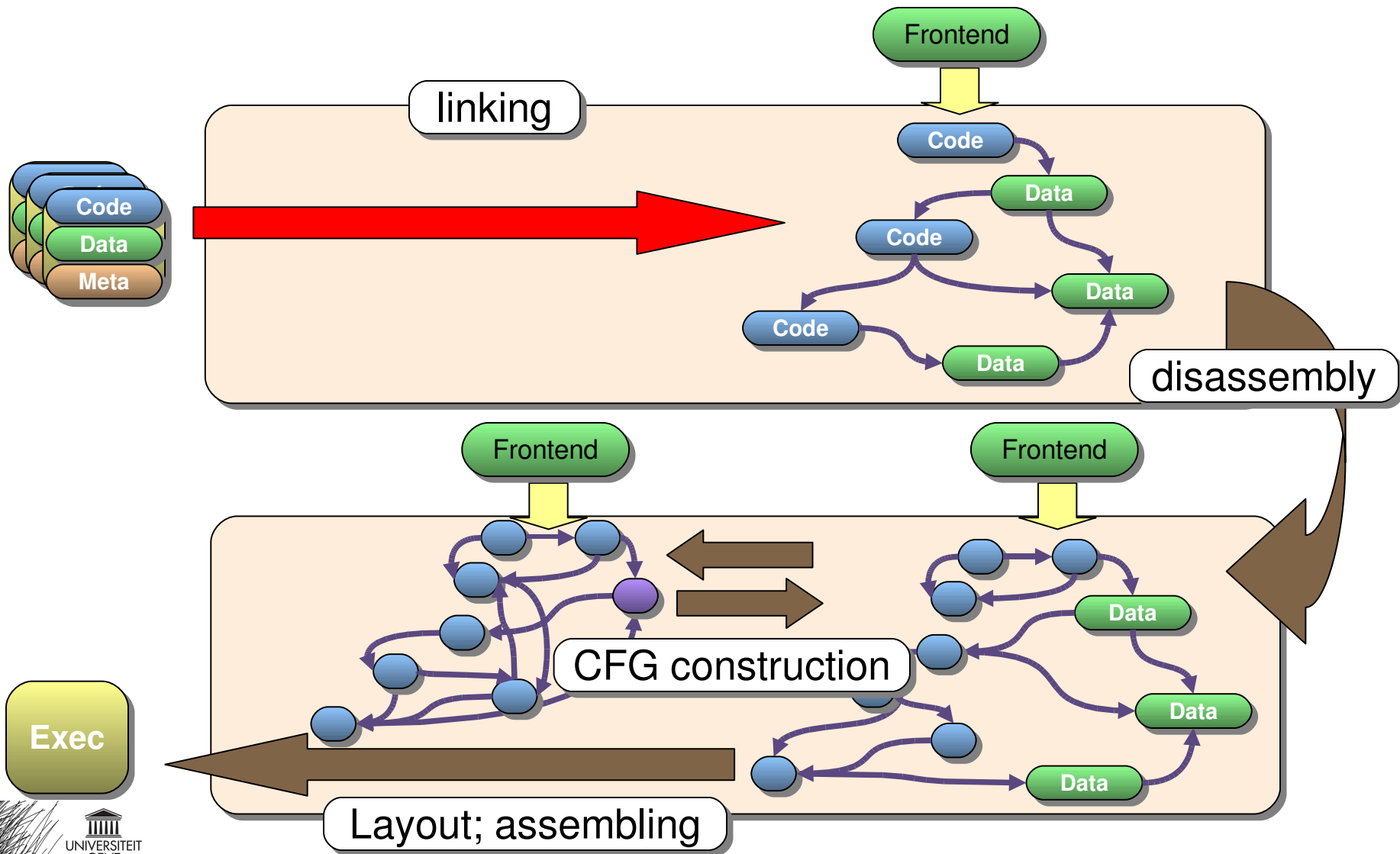
# Overview



- Some background
- Diablo
  - Extensibility
  - Retargetability
  - Reliability

Retargetability - The problem

Object formats

Retargetability

Linkers

Architectures

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

p. 22

UNIVERSITEIT GENT

# Operation at different levels

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

p. 23

# Retargeting multiple ROF & Linkers

Binary Pr... Diablo - Bjorn D...
Engineer... Electronics ar... Department

**Object file backend**

**Linker map & linker script**

linking

Code

Data

Meta

Code

Data

Meta

Code

Data

Code

Data

Code

Data

disassembly

**Architecture Backend**

CFG construction

Data

Data

Data

**Exec**

Layout; assembling

**Exe backend**

**Layout script & layout code**

UNIVERSITEIT GENT

# Major implemented backends

**ARM** — Diablo, FIT, kDiablo, Lancet — LCTES'04

**x86** — Diablo, FIT, Stilo, kDiablo, Lancet, Loco — LCTES'05

**IA64** — Diablo — Europar'04

**MIPS32** — Diablo — ESA'04

**Alpha** — Diablo, FIT

UNIVERSITEIT GENT

# Overview

- Some background
- Diablo
  - Extensibility
  - Retargetability
  - Reliability

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

p. 26

Reliability – The Problem

CFG construction

Conservative for safety

Reliability

Aggressive, precise enough to be useful

Relocation

Data flow analysis

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

UNIVERSITEIT GENT

# CFG Construction

## Disassembling

**Potential problems:**

- Differentiate data from code
- Detect self-modifying code
- Detect unrewritable code

**Solutions:**

- Section information
- Symbols annotate data in code (ARM ABI)
- Self-modifying code in data:
    no problem at this point
- True self-modifying code:
    look at system calls and protection

## Conservatively modelling control flow

**Potential problems:**

- Indirect control flow transfers
- Code that is treated as data
- Unrealizable paths (procedures)

**Solutions:**

- Use relocation information:
    identifies computable addresses
- Use pattern matching:
    identifies known address computations
- Use knowledge on compiler-generated code

UNIVERSITEIT
GENT

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

p. 28

# Detecting Data

```
$code   0x0080:   mov r2, 0x0a0
        0x0084:   cmp r1, $0
        0x0088:   jl 0x0b4
        0x008c:   cmp r1, 5
        0x0090:   jge 0x0b4
        0x0094:   add r1, r2, r1
        0x0098:   ldr r1, [r1]
        0x009c:   jmp r1
$data   0x00a0:   0x00000120
        0x00a4:   0x0000012c
        0x00a8:   0x000000d0
        0x00ac:   0x00000248
        0x00b0:   0x00000210
$code   0x00b4:   mov r3, r5
```

**Solution:**
 add mapping symbols

UNIVERSITEIT GENT

# Detecting Control Flow Targets

```
$code   0x0080:    mov r2, 0x0a0
        0x0084:    cmp r1, $0
        0x0088:    jl 0x0b4
        0x008c:    cmp r1, 5
        0x0090:    jge 0x0b4
        0x0094:    add r1, r2, r1
        0x0098:    ldr r1, [r1]
        0x009c:    jmp r1
$data   0x00a0:    0x00000120
        0x00a4:    0x0000012c
        0x00a8:    0x000000d0
        0x00ac:    0x00000248
        0x00b0:    0x00000210
$code   0x00b4:    mov r3, r5
        ...
        0x00d0:    add r4, r6, r6
        ...
        0x0120:    ldr r4, [r5]
```

**Direct control flow:**
trivial

**Indirect control flow:**
only to code-addresses that are targets of relocations!

**Problem:**
what about *unrelocated* computations on code-addresses?

UNIVERSITEIT GENT
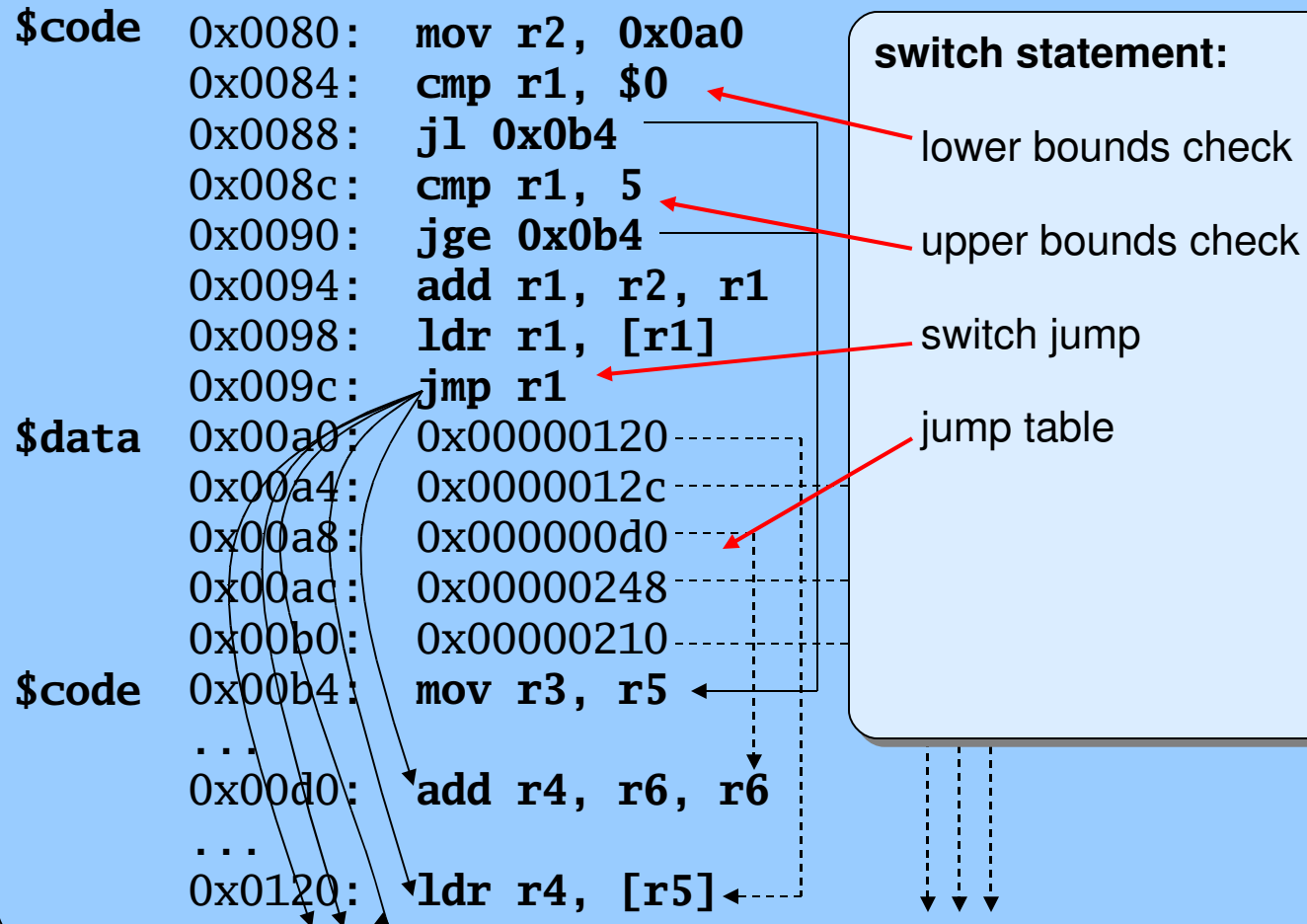
# Control Flow: Pattern Matching

```
$code   0x0080:   mov r2, 0x0a0
        0x0084:   cmp r1, $0
        0x0088:   jl 0x0b4
        0x008c:   cmp r1, 5
        0x0090:   jge 0x0b4
        0x0094:   add r1, r2, r1
        0x0098:   ldr r1, [r1]
        0x009c:   jmp r1
$data   0x00a0:   0x00000120
        0x00a4:   0x0000012c
        0x00a8:   0x000000d0
        0x00ac:   0x00000248
        0x00b0:   0x00000210
$code   0x00b4:   mov r3, r5
        ...
        0x00d0:   add r4, r6, r6
        ...
        0x0120:   ldr r4, [r5]
```

**Use pattern matching to improve accuracy of control flow graph:**
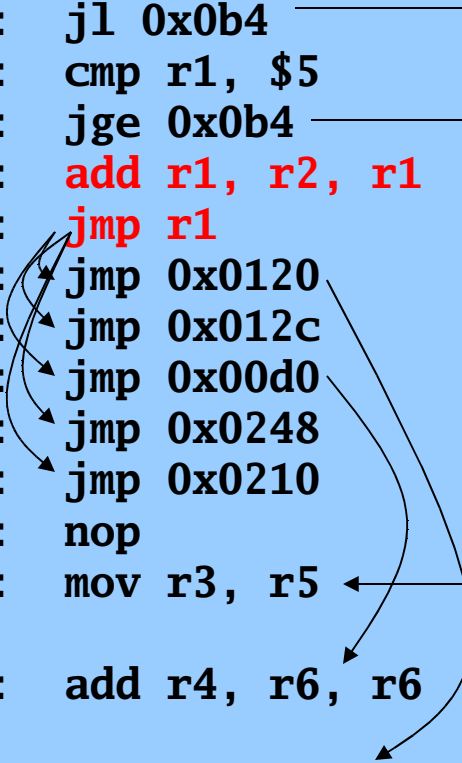disallow computations on code-addresses that are not part of a recognized pattern

UNIVERSITEIT GENT

# Pattern Matching: example

```
$code    0x0080:    mov r2, 0x0a0
         0x0084:    cmp r1, $0
         0x0088:    jl 0x0b4
         0x008c:    cmp r1, 5
         0x0090:    jge 0x0b4
         0x0094:    add r1, r2, r1
         0x0098:    ldr r1, [r1]
         0x009c:    jmp r1
$data    0x00a0:    0x00000120
         0x00a4:    0x0000012c
         0x00a8:    0x000000d0
         0x00ac:    0x00000248
         0x00b0:    0x00000210
$code    0x00b4:    mov r3, r5
         ...
         0x00d0:    add r4, r6, r6
         ...
         0x0120:    ldr r4, [r5]
```

**switch statement:**

lower bounds check

upper bounds check

switch jump

jump table

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

UNIVERSITEIT GENT

# Pattern Matching: example 2

```
$code   0x0080:    mov r2, 0x09c
        0x0084:    cmp r1, $0
        0x0088:    jl 0x0b4
        0x008c:    cmp r1, $5
        0x0090:    jge 0x0b4
        0x0094:    add r1, r2, r1
        0x0098:    jmp r1
        0x009c:    jmp 0x0120
        0x00a0:    jmp 0x012c
        0x00a4:    jmp 0x00d0
        0x00a8:    jmp 0x0248
        0x00ac:    jmp 0x0210
        0x00b0:    nop
        0x00b4:    mov r3, r5

        ...
        0x00d0:    add r4, r6, r6

        ...
        0x0120:    ldr r4, [r5]
```

**switch statement 2:**
address table is replaced by a series of direct jumps to the switch cases.

**unrecognized pattern!**

**Solution:**
add pattern to Diablo

UNIVERSITEIT GENT

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

# Procedure Calls and Returns

**ARM indirect procedure call:**

```
mov r14, pc
mov pc, r2
```

**ARM procedure return:**

```
mov pc, r14
```

or
```
ldr pc, [r13], #4
```
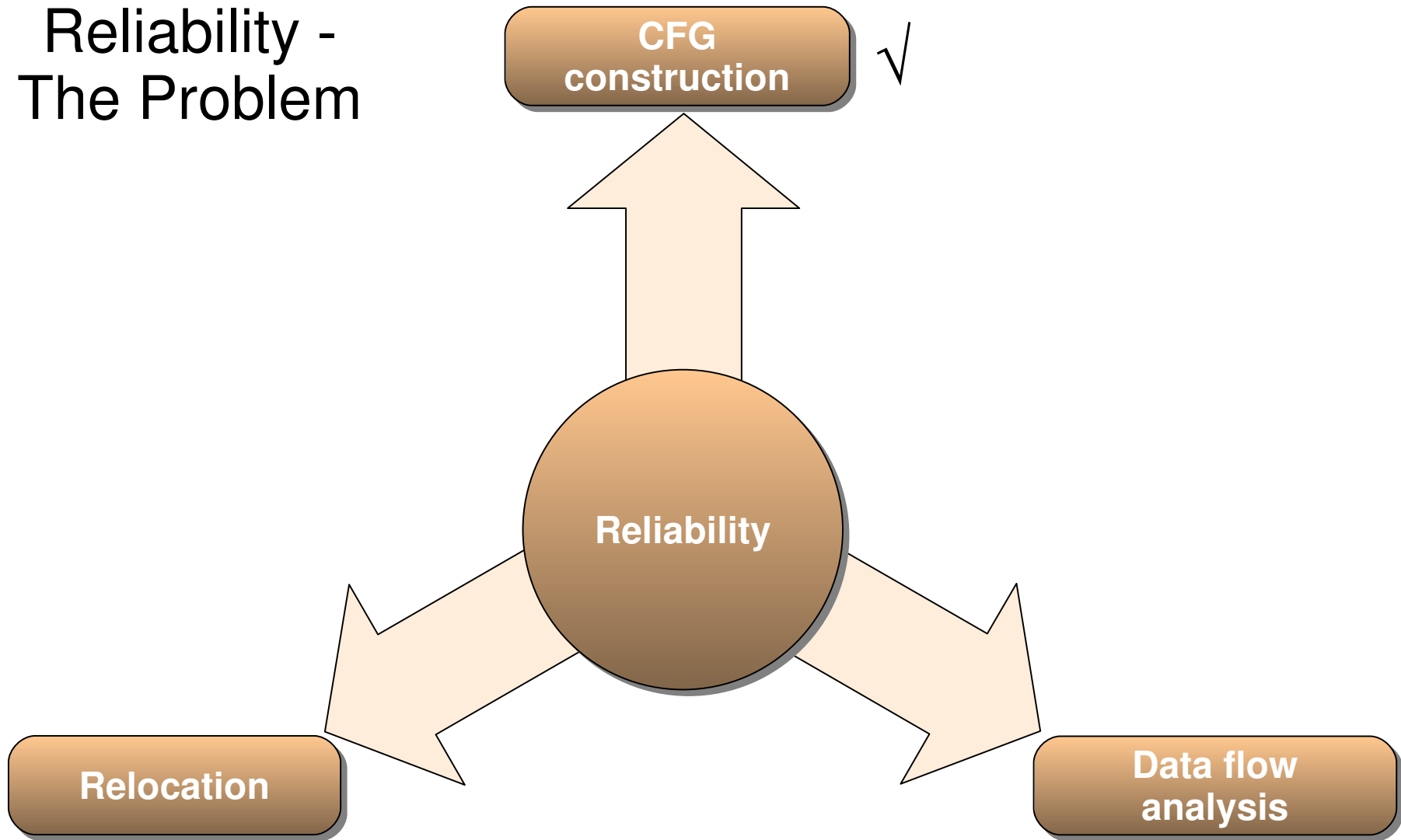
or
```
ldmia r13, {r4–r7,r15}!
```

**Function calls and returns are often just "special" indirect jumps:**
not recognizing them makes the flow graph much too conservative

**Solution:**
use pattern matching to recognize them:
- rely on ABI
- rely on compiler conventions

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

p. 34

UNIVERSITEIT
GENT

# Reliability - The Problem



**CFG construction** √

**Reliability**

**Relocation**

**Data flow analysis**

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

p. 35

UNIVERSITEIT GENT

# Data flow analyses

## Stack analysis

**Problem**

- Difficult to analyse
- Necessary to improve precision
- Especially for C++-like languages
  (calls through function pointers)

**Solution**

- rely on calling-conventions
- use symbol information
- use mapping symbols
- use source code information
- use stack unwind information

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

p. 36

UNIVERSITEIT
GENT

# Calling convention adherence

**A.c**

```
extern int B(int x);

int A(int x)
{
        return B(x);
}
```

**B.c**

```
int B(int x)
{
        return x * 2;
}

int C(int x)
{
        return B(x)*2;
}
```

B is unknown

**call to B respects
calling conventions**

C is known
but *A is unknown*
**B respects
calling convention**

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

p. 37

UNIVERSITEIT
GENT

# Calling convention adherence

**B.c**

```
static int B(int x)
{
        return x * 2;
}

int C(int x)
{
        return B(x)*2;
}
```

C is known
*no unknown callers of B*
**B does not need to respect
calling convention**

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

p. 38

UNIVERSITEIT
GENT

# Calling convention adherence

**A.c**

```
extern int B(int x);

int A(int x)
{
        int y;
        asm("
            movl %ecx, x
            call B
            movl y, %ecx
        ");
        return y;
}
```

**B.s**
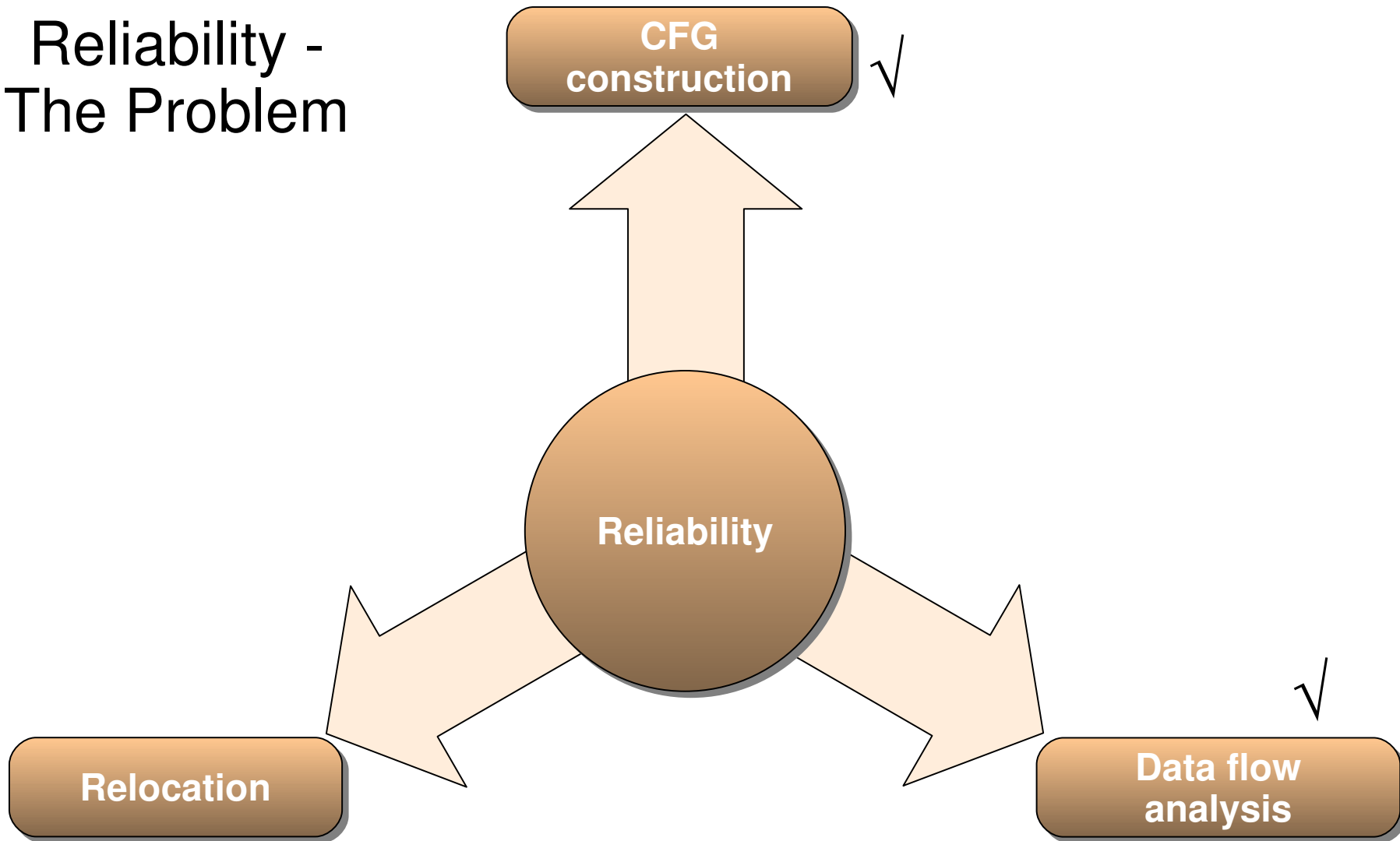
```
B:

shl %ecx, #1
ret
```

even though B is global, the programmer has control over all call sites
**B does not need to adhere to calling conventions**

**solution:**
identify assembler code through *mapping symbols* (for inline assembler) and object file header info

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

p. 39

UNIVERSITEIT
GENT

# Reliability - The Problem

CFG construction √

Reliability

Relocation

Data flow analysis √

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

p. 40

UNIVERSITEIT GENT

# Relocation

**Producing binary program again**

**Problem**

- How to write a correct program?
- How to layout data?
- How to update pointers?
- How to update addresses?

**Observation**
- Most "strange" requirements come from linker manipulations

**Solution**
- make relocations expressive
- make relocations first class objects
- let transformations update relocations
- use linker scripts

UNIVERSITEIT
GENT

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

p. 41

# Overview



- Some background
- Diablo
  - Extensibility
  - Retargetability
  - Reliability (no, really now)

UNIVERSITEIT GENT

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

**Reliability - The Real Problem**

CFG construction

Conservative for safety

Aggressive, precise enough to be useful

Reliability

Relocation

Data flow analysis

UNIVERSITEIT GENT

**CFG construction**

**Conservative for safety**

**Reliability - The Real Problem**

**Aggressive, precise enough to be useful**

**Reliability**

**Solution: limit most conservative assumptions to parts of program**

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

p. 44

# Limit imprecision to some parts

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

# Limit imprecision to some parts

Binary Program Rewriting with Diablo - Bjorn De Sutter – 2006-6-11
Engineering Sciences Faculty – Electronics and Information Systems Department

# What program parts?

- Sections from object files
  - only refer to each other via symbols
  - special code addresses identified by relocations
  - extend relocations where necessary
    - no relaxation
    - annotate PIC code with relocations if necessary
    - mark data
    - ...

UNIVERSITEIT GENT

# When/why does this work?

- Under separate compilation
  - Partial-separate compilation
  - Compiler-generated code only, not manually-written assembler
- Compiler needs to maintain conventions
- Assembly writers do not know compiler-generated code
  - Because multiple compiler versions are available
- Whenever imprecision could become viral, the linker (rewriter) is informed!