Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# A generic approach to the definition of low-level components for multi-architecture binary analysis

Cédric Valensi

PhD advisor: William Jalby

University of Versailles Saint-Quentin-en-Yvelines, France
Exascale Computing Research Center, France
LRC ITACA, France

July 2nd, 2014

UNIVERSITÉ DE
VERSAILLES
SAINT-QUENTIN-EN-YVELINES

**LRC IT@CA**

Exascale ∞
computing research

Introduction
Multi architecture support
Disassembly of binary files
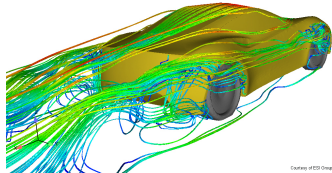Binary rewriting
Conclusion

# High Performance Computing

## Supercomputers

- Front-line of the computing capacity
- Multiprocessor systems
- Current top speed 33 Petaflop/s

## Applications

- Physical simulations
- Natural resources exploration
- Molecular modeling
- Weather forecasts



天河二号

TH-2 High Performance Computer System



Courtesy of ESI Group

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

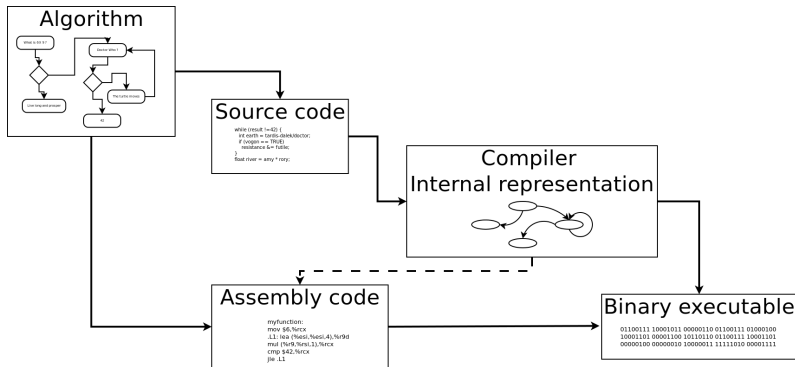# Performance analysis for optimisation

## Optimising performance of HPC applications

- Optimise use of processors in terms of speed and power
- Pinpoint bottlenecks
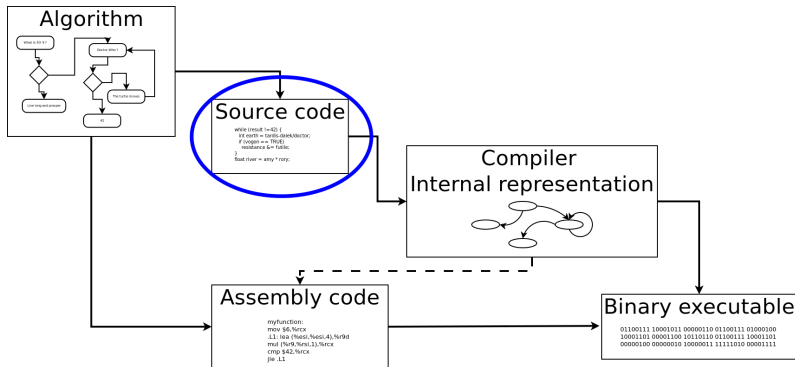- Estimate gain from improvements

## Performance analysis

- Static or dynamic
- Instrumentation
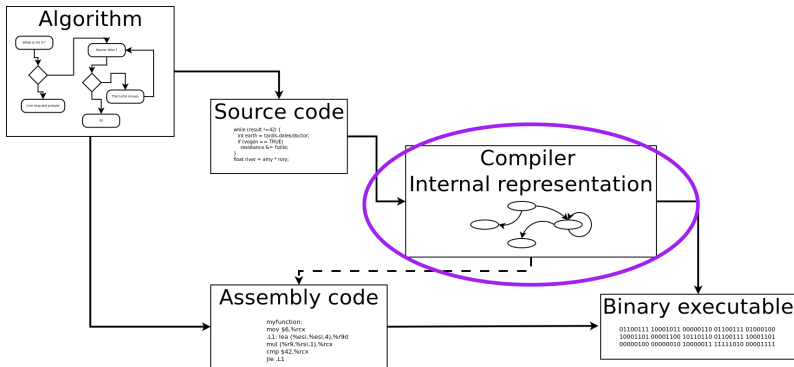- Possible in all steps of the design process

Introduction
Multi architecture support
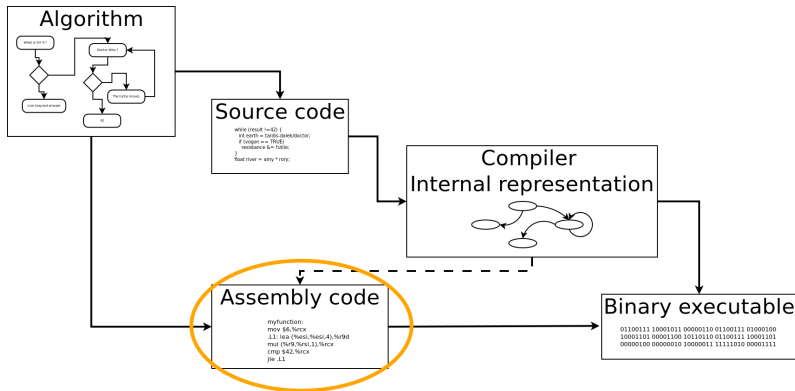Disassembly of binary files
Binary rewriting
Conclusion

# Steps of an application design process

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Steps of an application design process

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Steps of an application design process

**Introduction**
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Steps of an application design process

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Steps of an application design process

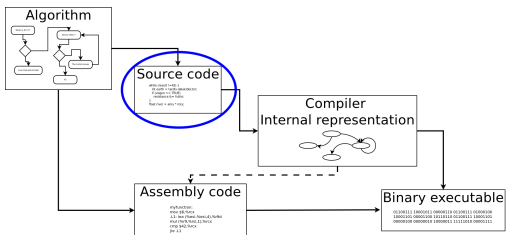Introduction
Multi architecture support
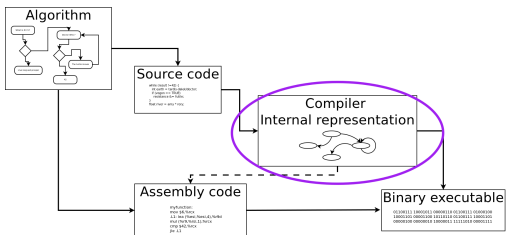Disassembly of binary files
Binary rewriting
Conclusion

# Performance analysis levels



## Source code

- Knowledge of source language
- Requires access to source files
- Compilation may perform complex transformations
- Instrumenting at the source level may modify these transformations

Introduction
Multi architecture support
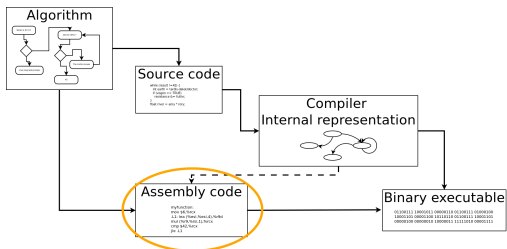Disassembly of binary files
Binary rewriting
Conclusion

# Performance analysis levels



## Compiler Internal Representation

- More accurate
- Requires access to compiler internals
- Requires intrusion into compilation process
- Ineffective for code written in assembly

Introduction
Multi architecture support
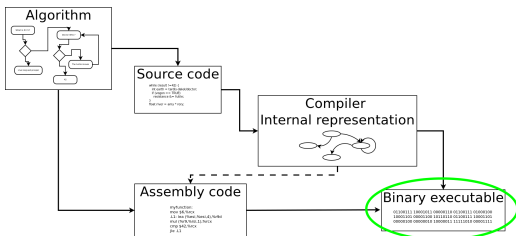Disassembly of binary files
Binary rewriting
Conclusion

# Performance analysis levels



## Assembly analysis

- Closer to the actual executable
- Not available by default
- Requires intrusion into compilation process

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Performance analysis levels



## Binary analysis

- "What you see is what you run"
- Allows to retrieve additional information
- More complex

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Challenges of binary analysis

## Dependent on the architecture

- Multiple architectures may be used by a single application
- Binary architectures evolve frequently

## Static Analysis

- Requires disassembly of binary code

## Instrumentation

- Requires static or dynamic patching
- Extensive changes can be needed

**Introduction**
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

## Contribution

### Low level binary encoder and decoder

- Able to support multiple architectures
- Minimised implementation workload

### Usage in analysis context

- Customisable behaviour
- Unified output format
- Acceptable performance
- Static analysis and instrumentation

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Outline

1. Introduction

2. Multi architecture support

3. Disassembly of binary files

4. Binary rewriting

5. Conclusion

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Objectives

### Generic encoder and decoder

- Multi-architecture support
- Customisable output and behaviour
- Reduced implementation workload

### Challenges

- Complex binary coding rules
- Coding rules and assembly vary significantly between architectures
- Avoid hard coding

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an Intel 64 instruction

0x66 49 89 4C 90 20 <=> mov %r9, 0x20(%eax,%edx,4)

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an Intel 64 instruction

0x66 49 89 4C 90 20 <=> mov %r9, 0x20(%eax,%edx,4)

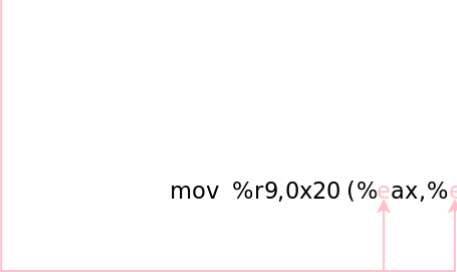01100111   0100 1100   10001001   01001100   10 010000   00100000
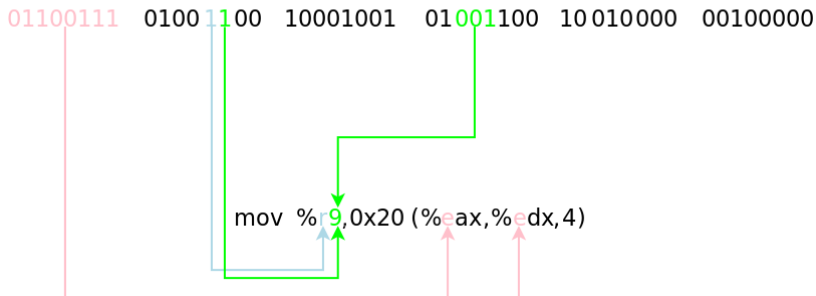
mov %r9,0x20 (%eax,%edx,4)

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an Intel 64 instruction

0x66 49 89 4C 90 20 <=> mov %r9, 0x20(%eax,%edx,4)

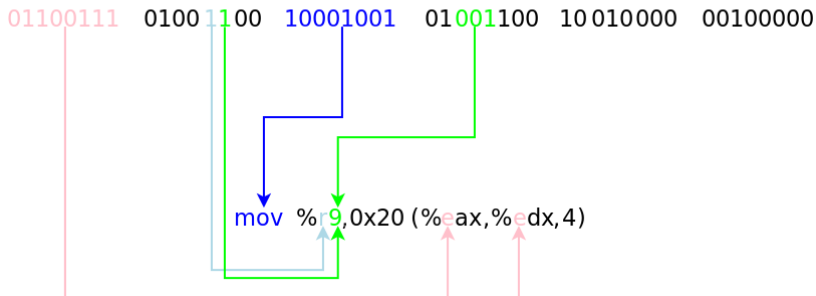01100111  0100 1100  10001001  01001100  10 010000  00100000

mov %r9,0x20 (%eax,%edx,4)

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an Intel 64 instruction



0x66 49 89 4C 90 20 <=> mov %r9, 0x20(%eax,%edx,4)

01100111   0100 1100   10001001   01001100   10 010000   00100000

mov  %r9,0x20 (%eax,%edx,4)

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an Intel 64 instruction



0x66 49 89 4C 90 20 <=> mov %r9, 0x20(%eax,%edx,4)

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an Intel 64 instruction



0x66 49 89 4C 90 20 <=> mov %r9, 0x20(%eax,%edx,4)

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an Intel 64 instruction

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an Intel 64 instruction



0x66 49 89 4C 90 20 <=> mov %r9, 0x20(%eax,%edx,4)

01100111   0100 1100   10001001   01001100   10010000   00100000

mov %r9,0x20(%eax,%edx,4)

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an Intel 64 instruction

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an Intel 64 instruction

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an Intel 64 instruction



0x66 49 89 4C 90 20 <=> mov %r9, 0x20(%eax,%edx,4)

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an Intel 64 instruction



0x66 49 89 4C 90 20 <=> mov %r9, 0x20(%eax,%edx,4)

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an ARM instruction

0x15 2D 40 05 <=> strne r4, [sp, #-5]!

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an ARM instruction

0x15 2D 40 05 <=> strne r4, [sp, #-5]!

00010101 00101101 01000000 00000101

strne   r4,[sp,#-5]!

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an ARM instruction

0x15 2D 40 05 <=> strne r4, [sp, #-5]!

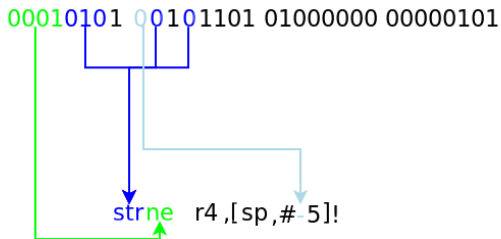0001 0101 00101101 01000000 00000101

strne r4,[sp,#-5]!

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an ARM instruction

0x15 2D 40 05 <=> strne r4, [sp, #-5]!

0001**0**101 00**1**01101 01000000 00000101

strne r4,[sp,#-5]!

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an ARM instruction

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an ARM instruction

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an ARM instruction

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an ARM instruction

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an ARM instruction

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Requirements

### Ensuring agnosticism with regard to architecture

- Unified representation of an architecture encoding rules
- Decorrelation of decoding from post parsing actions
- Same representation to generate encoder and decoder

### Remaining close to the documentation format

- Handling exclusions and restricted cases
- Possibility of fields with no fixed value

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Using a context-free grammar formalism

## Advantages

- Allows to decorrelate the encoding rules from the actions performed
- Decoder implemented as the corresponding parser
- Multiple possible uses for the decoder
- Encoder built from the same grammar

## Challenges

- Grammars usually operate at the character level
- Using a bit by bit parsing would be inefficient
- Lookahead challenged by instructions of variable sizes

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Standard notions

### Context free grammars

- Symbols associated to list of productions
- A production contains terminal and nonterminal symbols
- Terminal symbols have no production
- Semantic actions associated to productions

### LR parsers

- Processing left to right
- Bottom-up matching
- Implemented as finite state automata
- Shift and reduction states

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Our algorithm for parser generation

## New principles

- Bits can have a fixed or unfixed value
- Terminals are defined as groups of bits
- A state represents the matching of bits anywhere in the production
- Transitions over terminals can include bits ahead of the parsing step
- Shift/reduce states are authorised

## Parser execution

- Processing left to right
- Terminals containing less unfixed bits are tested first

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Context free grammar

```
%token <2> d
Start:
A 00
|B 01
;
A:
C
|0111
;
B:
0111
|d 11
;
C:
00 d
|0000
;
```

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Context free grammar

```
%token <2> d
Start:
A 00
|B 01
;
A:
C
|0111
;
B:
0111
|d 11
;
C:
00 d
|0000
;
```

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Context free grammar

```
%token <2> d
Start:
A 00
|B 01
;
A:
C
|0111
;
B:
0111
|d 11
;
C:
00 d
|0000
;
```

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Context free grammar

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion
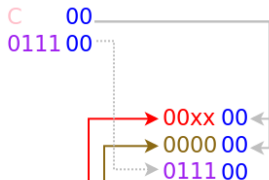
# Example: Context free grammar

```
%token <2> d
Start:
A 00
|B 01
;
A:
C
|0111
;
B:
0111
|d 11
;
C:
00 d
|0000
;
```

00xx 00
0000 00
0111 00

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Context free grammar



```
%token <2> d
Start:
A 00
|B 01
;
A:
C
|0111
;
B:
0111
|d 11
;
C:
00 d
|0000
;
```

```
00xx 00
0000 00
0111 00
0111 01
xx11 01
```

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Context free grammar

```
%token <2> d
Start:
A 00
|B 01
;
A:
C
|0111
;
B:
0111
|d 11
;
C:
00 d
|0000
;
```

```
00xx 00
0000 00
0111 00
0111 01
xx11 01
```

Introduction
**Multi architecture support**
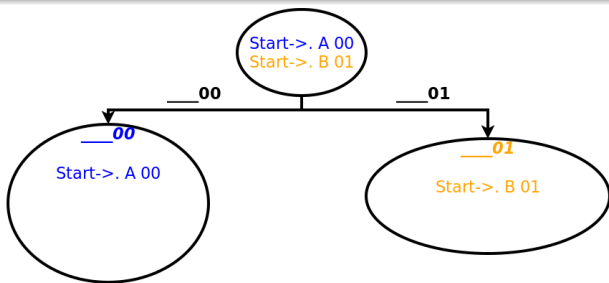Disassembly of binary files
Binary rewriting
Conclusion

# Example: Parser Generation from grammar

```
%token <2> d
Start:
A 00
| B 01
;
A:
C
| 0111
;
B:
0111
| d 11
;
C:
00 d
| 0000
;
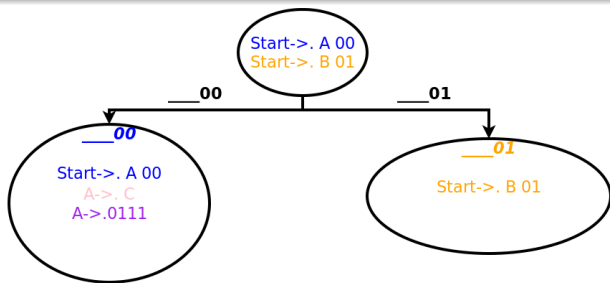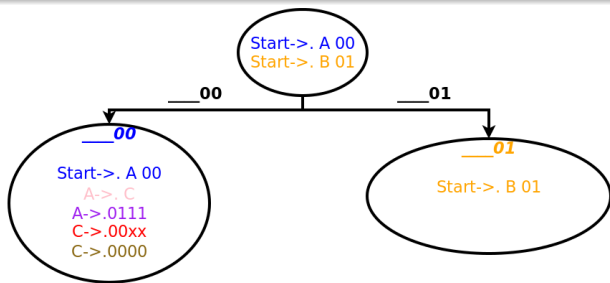```

Start->. A 00
Start->. B 01

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion
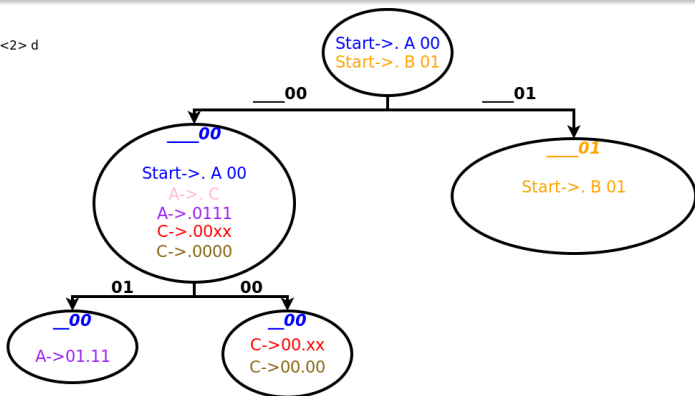
# Example: Parser Generation from grammar



```
%token <2> d
Start:
A 00
|B 01
;
A:
C
|0111
;
B:
0111
|d 11
;
C:
00 d
|0000
;
```

Start->. A 00
Start->. B 01

___00        ___01

___00
Start->. A 00

___01
Start->. B 01

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Parser Generation from grammar

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Parser Generation from grammar

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Parser Generation from grammar

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Parser Generation from grammar

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Parser Generation from grammar

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Parser Generation from grammar

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Parser Generation from grammar

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

## Encoder generation

### Building an encoder from the same grammar file

- Semantic actions are redefined as matching functions
- Input tentatively matched over all productions of nonterminals
- Shortest productions are matched first
- Nonterminals in a matching production are recursively matched
- Resulting encoder algorithm corresponds to a top-down parser

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoder algorithm

```
%token <2> d
Start:
A 00  #[ S_ACT1($1) ]#
| B 01 #[ S_ACT2($1) ]# ;
A:
C     #[ A_ACT1($1) ]#
| 0111 #[ A_ACT2()   ]#;
B:
0111 #[ B_ACT1()   ]#
| d 11 #[ B_ACT2($1) ]# ;
C:
00 d   #[ C_ACT1($1) ]#
| 0000 #[ C_ACT2() ]#   ;
```

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoder algorithm

```
%token <2> d
Start:
A 00  #[ S_ACT1($1) ]#
| B 01 #[ S_ACT2($1) ]# ;
A:
C      #[ A_ACT1($1) ]#
| 0111 #[ A_ACT2()   ]#;
B:
0111 #[ B_ACT1()   ]#
| d 11 #[ B_ACT2($1) ]#;
C:
00 d   #[ C_ACT1($1) ]#
| 0000 #[ C_ACT2() ]#  ;
```

INPUT

Introduction
**Multi architecture support**
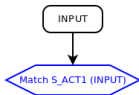Disassembly of binary files
Binary rewriting
Conclusion
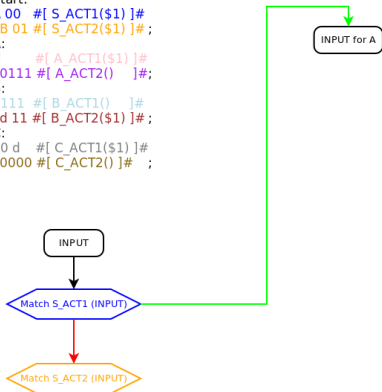
# Example: Encoder algorithm

```
%token <2> d
Start:
A 00  #[ S_ACT1($1) ]#
| B 01 #[ S_ACT2($1) ]# ;
A:
C     #[ A_ACT1($1) ]#
| 0111 #[ A_ACT2()   ]#;
B:
0111  #[ B_ACT1()    ]#
| d 11 #[ B_ACT2($1) ]# ;
C:
00 d   #[ C_ACT1($1) ]#
| 0000 #[ C_ACT2() ]#   ;
```

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoder algorithm

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoder algorithm

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoder algorithm

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoder algorithm

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoder algorithm

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoder algorithm

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoder algorithm

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoder algorithm

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoder algorithm



```
%token <2> d
Start:
A 00  #[ S_ACT1($1) ]#
| B 01 #[ S_ACT2($1) ]# ;
A:
C     #[ A_ACT1($1) ]#
| 0111 #[ A_ACT2()   ]#;
B:
0111 #[ B_ACT1()   ]#
| d 11 #[ B_ACT2($1) ]# ;
C:
00 d   #[ C_ACT1($1) ]#
| 0000 #[ C_ACT2() ]#  ;
```
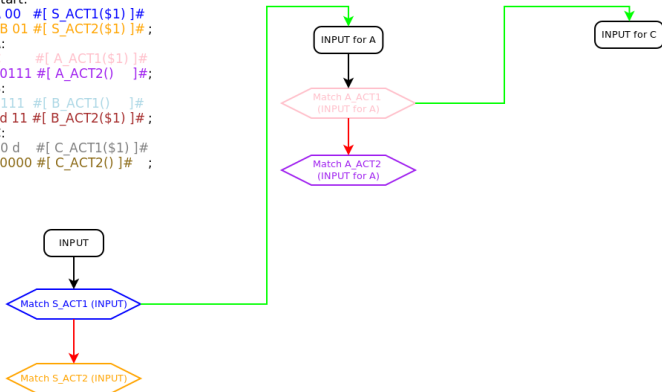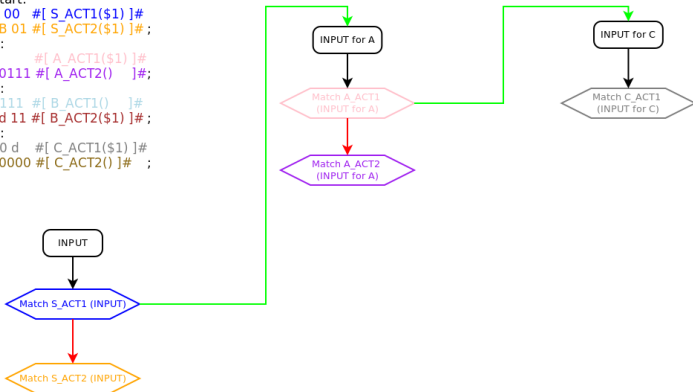
Introduction
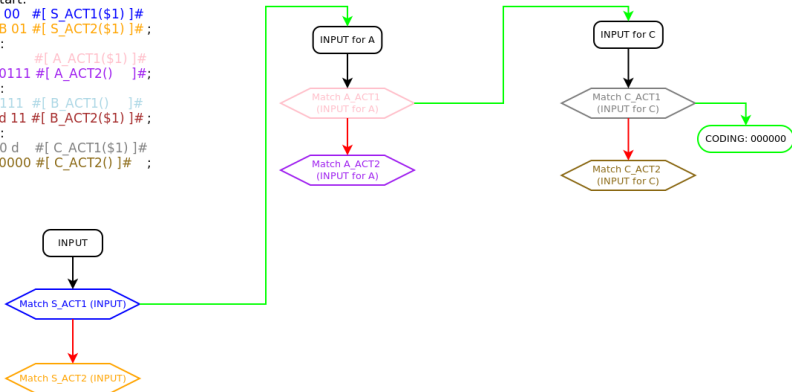**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoder algorithm

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoder algorithm

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Validation

### MINJAG

- Uses a context-free grammar describing the architecture
- Grammar generated from architecture documentation through simple transformations
- Generates the code for decoder and encoder from the same grammar
- Functional tool used in a production context
- Tested over Intel 64, Intel Xeon Phi coprocessor and ARM

Introduction
**Multi architecture support**
Disassembly of binary files
Binary rewriting
Conclusion

# Characteristics of implemented architectures

| Architecture | Intel 64 | Intel Xeon Phi | ARM |
|---|---|---|---|
| Lines in instruction list | 2,398 | 1,194 | 1,512 |
| Lines in grammar | 6,082 | 3,082 | 1,491 |
| Reduction states | 5,950 | 2,406 | 1,625 |
| Shift states | 4,019 | 1,468 | 2,916 |
| Shift/reduce states | 2 | 2 | 6 |
| Total states | 9,971 | 3,876 | 4,547 |

Introduction
Multi architecture support
**Disassembly of binary files**
Binary rewriting
Conclusion

1 Introduction

2 Multi architecture support

3 Disassembly of binary files

4 Binary rewriting

5 Conclusion

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Challenges of disassembly

## Binary code is not intended to be read

- No constraints on the code as long as the program can be executed
- No separation between instructions
- Instructions may be of varying sizes

## Specific examples

- Interleaved foreign bytes
- Overlapping instructions
- Obfuscated code or binary format
- Self rewriting code

Introduction
Multi architecture support
**Disassembly of binary files**
Binary rewriting
Conclusion

# Example: Interleaved foreign bytes

| Binary code | Corresponding assembly instructions |
|---|---|
| 67 4C 89 4C 90 20 | mov %r9,0x20(%eax,%edx,4) |
| EB 04 | jmp <+4 bytes> |
| *80* | *Alignment byte, never executed* |
| 4C 89 C8 | mov %r9,%rax |
| F2 0F 10 EE | movsd %xmm6, %xmm5 |
| F2 0F 10 F4 | movsd %xmm4, %xmm6 |

Correct disassembly up to that point

Disassembly

| | | |
|---|---|---|
| 67 4C 89 4C 90 20 | mov %r9,0x20(%eax,%edx,4) | Instructions correctly disassembled |
| EB 04 | jmp <+4 bytes> | |
| 80 4C 89 C8 F2 | or $0xf2,-0x38(%rcx,%rcx,4) | Erroneous instructions |
| 0F 10 EE | movups %xmm6, %xmm5 | |
| F2 0F 10 F4 | movsd %xmm4, %xmm6 | Instructions correctly disassembled |

Mistaking the alignment byte for the beginning of the next instruction

Realignement of the parser on a valid boundary

Introduction
Multi architecture support
**Disassembly of binary files**
Binary rewriting
Conclusion

# Disassembly algorithms

### Linear sweep

- Decoding one instruction after another
- Errors when encountering interleaved foreign bytes
- Vulnerability to obfuscation methods
- Faster disassembly

### Recursive traversal

- Decoding following the actual execution of the program
- Resists to some obfuscation techniques
- Finding the destination of a branch can be difficult
- Slower disassembly

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Linear Sweep vs Recursive Traversal

Linear sweep

011001110100110010001000101001100
1001000000100000 1110101100000011
010010011000100111001110 10010000
1111001000001111000100000111011110

Recursive traversal

011001110100110010001000101001100
1001000000100000 1110101100000011
010010011000100111001110 10010000
1111001000001111000100000111011110

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Linear Sweep vs Recursive Traversal

Linear sweep

Recursive traversal

0110011101001100100010010100110
1001000000100000 1110101100000011
0100100110001001110011110 10010000
11110010000011110001000011101110

→ 00: mov %r9,0x20(%eax,%edx,4)

0110011101001100100010010100110
1001000000100000 1110101100000011
0100100110001001110011110 10010000
11110010000011110001000011101110

→ 00: mov %r9,0x20(%eax,%edx,4)

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Linear Sweep vs Recursive Traversal

Linear sweep

Recursive traversal

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Linear Sweep vs Recursive Traversal



Linear sweep

```
0110011010011001000100101001100
1001000000100000 1110101100000011
010010011000100111001110 10010000
11110010000011110001000011101110
```

```
00: mov %r9,0x20(%eax,%edx,4)
06: jmp <0C> #+4 bytes
08: mov %rcx,%r14
```

Recursive traversal

```
0110011010011001000100101001100
1001000000100000 1110101100000011
010010011000100111001110 10010000
11110010000011110001000011101110
```

```
00: mov %r9,0x20(%eax,%edx,4)
06: jmp <0C> #+4 bytes
```

```
0C: movsd %xmm6,%xmm5
```

Introduction
Multi architecture support
**Disassembly of binary files**
Binary rewriting
Conclusion

## Our constraints

### Disassembler intended to be used by analysis tools

- Retrieve all possible available information from the file
- Architecture independent output format
- Possibility to add customisable additional information
- Acceptable performance in terms of speed and accuracy

Introduction
Multi architecture support
**Disassembly of binary files**
Binary rewriting
Conclusion

# Our disassembling algorithm

## General execution

- Linear sweep parsing
- Extraction of executable code from binary format
- Retrieval of labels and debug information if present

## Additional processing

- Resolving destination of direct branches
- Associating labels and debug information to instructions
- Post parsing actions to fill additional information
- Detection of unreachable instructions
- Identification of dubious disassembled data

Introduction
Multi architecture support
**Disassembly of binary files**
Binary rewriting
Conclusion

# Implementation: the $\mathrm{MADRAS}$ disassembler

### Multi Architecture **Disassembler**, Rewriter and ASsembler

- Relies on $\mathrm{MINJAG}$ for source code of decoder
- Processes binaries using the $\mathrm{ELF}$ format used by Unix and Linux
- Disassembler available for Intel 64, Xeon Phi coprocessor and ARM
- Base component of the $\mathrm{MAQAO}$ framework

Introduction
Multi architecture support
**Disassembly of binary files**
Binary rewriting
Conclusion

# Performance tests

## Protocol

- Comparison between MADRAS and hard coded disassemblers
- Disassembling SPEC benchmarks and test files
    - Size of executable code varying between 1 and 23 MBytes
    - Executables compiled for Intel 64 and Xeon Phi coprocessor
- Speed measured as disassembled instructions per second

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Disassembler performance on Intel 64 files

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Disassembler performance on Xeon Phi files

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Parallel disassembler performance



Intel 64 files



Xeon Phi files

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Disassembler accuracy

Introduction
Multi architecture support
Disassembly of binary files
**Binary rewriting**
Conclusion

1 Introduction

2 Multi architecture support

3 Disassembly of binary files

4 Binary rewriting

5 Conclusion

Introduction
Multi architecture support
Disassembly of binary files
**Binary rewriting**
Conclusion

# Instrumentation

### Retrieving information during execution

- Monitoring memory usage
- Value profiling

### Dynamic: Performed during execution

- Monitoring code execution using a supervising thread
- Invoking functions under specified conditions
- Modifying the image loaded in memory

### Static: Modifying the executable file

- Probe insertion
- Instructions modification

Introduction
Multi architecture support
Disassembly of binary files
**Binary rewriting**
Conclusion

# Binary rewriting

## Static instrumentation

- No recompilation needed
- No overhead from instrumentation process
- No additional requirements for execution

## Binary rewriting allows other modifications to the program

- Deleting or adding instructions to test their overall impact
- Modifying variables defined in the file

Introduction
Multi architecture support
Disassembly of binary files
**Binary rewriting**
Conclusion

# Challenges of binary rewriting

## Patched file must remain valid

- Preservation of the structure of the binary file
- Preservation of the control flow
- Preservation of data environment

## Executables are not intended to be modified

- All references are fixed
- No relocation tables
- Addresses can appear as immediate operands

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Example of patching pitfalls



```
00:   48 8B 04 25 0E 00 00 00    mov  $0x14, %rax

08:   FF E0                       jmp  *%rax

0A:   83 45 F8 01                 add  $1, -8(%rbp)

0E:   83 7D F8 01                 cmp $1, -8(%rbp)

12:   7E F6                       jle   0A

14:   B8 00 00 00 00              mov  $0, %eax
```

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Example of patching pitfalls

**E8 xx xx xx xx    callq  &lt;myfunc&gt;**

| 00: | 48 8B 04 25 0E 00 00 00 | mov  $0x14, %rax |
| 08: | FF E0 | jmp  *%rax |
| 0A: | 83 45 F8 01 | add  $1, -8(%rbp) |
| 0E: | 83 7D F8 01 | cmp $1, -8(%rbp) |
| 12: | 7E F6 | jle   0A |
| 14: | B8 00 00 00 00 | mov  $0, %eax |

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Example of patching pitfalls

Introduction
Multi architecture support
Disassembly of binary files
**Binary rewriting**
Conclusion

# Binary rewriting algorithm

### Block relocation

- The code to be modified is moved in a new section in the executable
- Code moved at the basic block level
- Use of trampolines if the patching site is too small

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Code relocation

Original code

Introduction
Multi architecture support
Disassembly of binary files
**Binary rewriting**
Conclusion

# Code relocation

## Original code



Modification site

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Code relocation

Original code



Basic block
surrounding
the site

Introduction
Multi architecture support
Disassembly of binary files
**Binary rewriting**
Conclusion

# Code relocation

Original code



Branch instruction

Introduction
Multi architecture support
Disassembly of binary files
**Binary rewriting**
Conclusion

# Code relocation



Original code

Added code section

Relocated block

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Code relocation

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Code relocation



Original code

Added code section

Return branch

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Trampolines

Original code

Modification site

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Trampolines

Original code

Basic block
surrounding
the site

Introduction
Multi architecture support
Disassembly of binary files
**Binary rewriting**
Conclusion

# Trampolines



Original code

Trampoline
block

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Trampolines



Original code

Added code section

Moved trampoline block

Introduction
Multi architecture support
Disassembly of binary files
**Binary rewriting**
Conclusion

# Trampolines

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Trampolines



Original code            Added code section

Branch

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Trampolines



Original code                    Added code section

Moved block

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Trampolines

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Trampolines



Original code      Added code section

Return branch

Introduction
Multi architecture support
Disassembly of binary files
**Binary rewriting**
Conclusion

# Implementation: the MADRAS patcher

## Multi Architecture Disassembler, **Rewriter and ASsembler**

- Relies on MINJAG for source code of assembler
- Processes binaries under the ELF format used by Unix and Linux
- Available for Intel 64 and Xeon Phi coprocessor

## A production tool

- C API
- Back end of the MAQAO Instrumentation Language (MIL)
- Used by the DECAN module

Introduction
Multi architecture support
Disassembly of binary files
**Binary rewriting**
Conclusion

## Patcher features

### Code insertion

- Insertion of calls to functions from external or static libraries
- Insertion of lists of assembly instructions

### Conditions

- Possibility to set conditions on the execution of an inserted code
- Possibility to specify code to execute if such a condition is not met

### Other features

- Modification or deletion of instructions
- Insertion of global variables usable by inserted code

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Using MADRAS API to insert a function call

```
void insert(char* in,char* lib,char* fct,uint addr,char* out) {
 //Disassemble the file and inits the modifications
 elfdis_t* madras = madras_disass_file(in);
 madras_modifs_init(madras, STACK_SHIFT, 512);
 //Adds a function call at the given address
 insert_t* ifct = madras_fctcall_new(madras, fct, lib, addr, 0);
 //Adds the given address as an immediate parameter
 madras_fctcall_addparam_imm(madras, ifct, addr, 0);
 //Commit changes
 madras_modifs_commit(madras,out);
 //Terminates the madras structure
 madras_terminate(madras);
}
```

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Interface with the MAQAO Instrumentation Language

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Performance of code patched by MIL

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
**Conclusion**

## Contributions

### Generic representation of binary encoding rules

- Unified format
- Use of the same grammar for encoder and decoder generation
- Validated for the Intel and ARM architectures
- Implemented as the functional tool $\mathrm{MINJAG}$

### Disassembly

- Easier updates of architecture specific code
- Performance comparable to existing hard coded tools
- Customisable output

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
**Conclusion**

## Contributions

### Patching

- Fine granularity offering wide range of options
- Patched code has similar or better performance than existing tools

### MADRAS

- Functional tool
- Standalone implementation of the whole disassembly and instrumentation chain
- Handling of multiple architectures from a single executable
- Integral component of the MAQAO framework
- Used by the DECAN module

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
**Conclusion**

# Future work

### General

- Implement additional architectures
- Support additional binary file formats

### Generic encoder and decoder

- Generic meta language for representing instruction lists
- Extensions allowing to specialise generated parser

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
**Conclusion**

# Future work

## Disassembler

- Improve accuracy through use of recursive traversal
- Detection of switch tables
- Improve speed
- Parallel disassembly
- Application to domains outside performance analysis

## Patcher

- Improve safety of patching
- Update of indirect branch destinations

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Q/A

Thank you for your attention!

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

Additional slides

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Example: Encoding of an ARM instruction

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

## Example of grammar for binary definition

```
%token <3,b> reg
%%
Start: template ;
template: Legacy3 Insn #[ FULLINSN_L3PREFIX($1,$2) ]#
        | Insn #[ FULLINSN($1) ]# ;
MemModRM: 00 reg RMSIB_00 #[ OPRS_REG_MEM($1,$2) ]#
        | 01 reg RMSIB_01 #[ OPRS_REG_MEM($1,$2) ]#
        | 10 reg RMSIB_10 #[ OPRS_REG_MEM($1,$2) ]# ;
RegModRM: 11 reg RMSIB_11 #[ OPRS_REG_REG($1,$2) ]# ;
Insn: 00010000 RegModRM #[ INSN(ADC,
        REG(GEN8b,R,$1),REG(GEN8b,RW,$1)) ]#
    | REX 00010000 MemModRM #[ INSN(ADC,
        REG(GEN8b,R,$1,$2),MEM(MEM8b,RW,$1,$2)) ]# ;
```

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

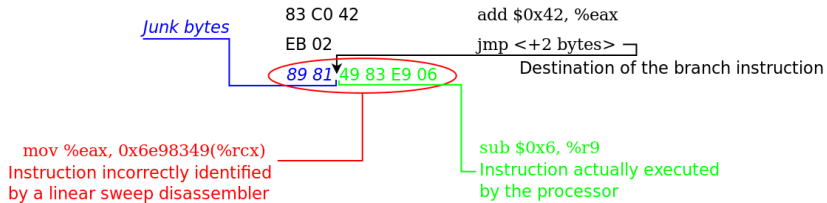# Overlapping instructions

Destination of the branch instruction

F3 AB                    rep stos

48 FF C1                 inc %rcx

48 83 F9 7F              cmp $127,%rcx

75 F6                    jne <-10 bytes>

The first iteration of the loop
will execute instruction rep stos
Later iterations will skip the F3 (rep)
prefix and execute only the stos
instruction

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Obfuscated code



*Junk bytes*

83 C0 42    add $0x42, %eax

EB 02     jmp <+2 bytes>

*89 81* 49 83 E9 06 Destination of the branch instruction

mov %eax, 0x6e98349(%rcx)
Instruction incorrectly identified
by a linear sweep disassembler

sub $0x6, %r9
Instruction actually executed
by the processor

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Performance tests

### Disassemblers

- objdump
- XED
- udis86
- distorm
- ndisasm

### Disassembly modes

- Print only mode for comparison against objdump and XED
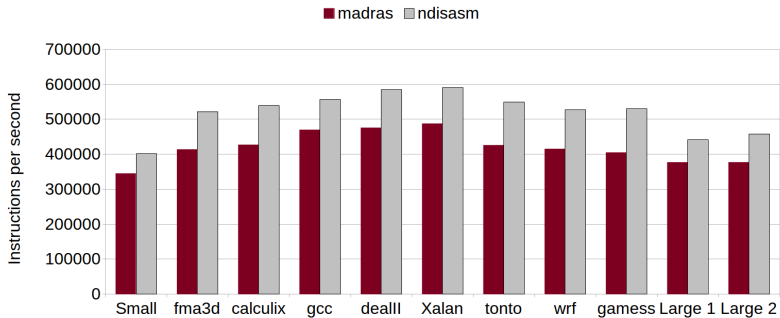- Without parsing of the binary file against udis86 and distorm

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Intel 64 files used for the disassembler performance tests

| File | File size (MByte) | Code size (MByte) | Description |
|---|---|---|---|
| Small | 0,96 | 0,96 | Test file |
| fma3d | 3,78 | 1,75 | SPEC2001 |
| calculix | 5 | 2,31 | SPEC2006 |
| gcc | 9,02 | 2,56 | SPEC2006 |
| dealII | 60,94 | 2,83 | SPEC2006 |
| Xalan | 130,64 | 3,46 | SPEC2006 |
| tonto | 33,27 | 5,81 | SPEC2006 |
| wrf | 19,52 | 6,83 | SPEC2006 |
| gamess | 18,2 | 10,55 | SPEC2006 |
| Large 1 | 11,95 | 11,94 | Test file |
| Large 2 | 23,22 | 23,22 | Test file |

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Xeon Phi files used for the disassembler performance tests

| File | File size (Mb) | Code size (Mb) | Description |
|---|---|---|---|
| equake | 0,12 | 0,05 | SPEC2001 |
| art | 0,21 | 0,12 | SPEC2001 |
| ammp | 0,84 | 0,44 | SPEC2001 |
| swim | 0,96 | 0,57 | SPEC2001 |
| wupwise | 0,96 | 0,66 | SPEC2001 |
| mgrid | 0,95 | 0,68 | SPEC2001 |
| applu | 1,03 | 0,71 | SPEC2001 |
| apsi | 2,61 | 1,72 | SPEC2001 |
| galgel | 2,84 | 2,08 | SPEC2001 |
| fma3d | 4,62 | 2,35 | SPEC2001 |

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Disassembler performance on Intel 64 files

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Parallel disassembler performance

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Performance of patched code

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Performance of instrumentation

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# Performance of patched code

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# MADRAS overall architecture

Introduction
Multi architecture support
Disassembly of binary files
Binary rewriting
Conclusion

# MAQAO Framework