# Analysis of Binaries using SAT/SMT

Jörg Brauer

RWTH Aachen University

brauer@embedded.rwth-aachen.de

with Andy King and Thomas Reinbacher

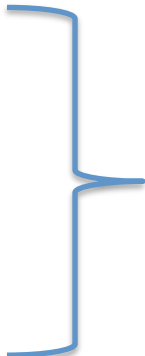01.03.2012 @ InBot, Aachen

# The Ultimate Goal

Rather: What I want to get a PhD for

- Given: Binary program

- Goal: Compute over-approximations of all registers & memory locations

- Approach: SAT & SMT solving

- Applications: Jump-target recovery, and many more

# Example #1

```
XOR  r0  r1
XOR  r1  r0
XOR  r0  r1
```

Swaps contents of two registers without involving a third

# Example #1

Introduce input-output variables and consider abstract domain of two-variable equalities:

$$(\mathbf{r0} = \mathbf{r1}), (\mathbf{r0} = \mathbf{r0}'), (\mathbf{r0} = \mathbf{r1}'), \dots, \top$$

```
XOR  r0  r1
XOR  r1  r0
XOR  r0  r1
```

Swaps contents of two registers without involving a third

# Example #1
## Traditional Abstraction

$(\mathbf{r0} = \mathbf{r1})$

XOR r0 r1

$\top$

XOR r1 r0

$\top$

XOR r0 r1

$\top$

# Example #1
# Block-Wise Abstraction

$(\mathbf{r0} = \mathbf{r1})$

XOR r0 r1

⊤

XOR r1 r0

⊤

XOR r0 r1

⊤

$(\mathbf{r0} = \mathbf{r1})$

XOR r0 r1

XOR r1 r0

XOR r0 r1

$(\mathbf{r0} = \mathbf{r1}) \wedge (\mathbf{r0} = \mathbf{r1'}) \wedge (\mathbf{r1} = \mathbf{r0'})$

# Lesson #1

- Reasoning about blocks rather than instructions can increase precision

  - **Problem**: Blocks are program-dependent, whereas instructions are not

  - **Solution**: Automatic techniques to compute abstractions for each block in a program

# Example #2

SBC  R2  R2  } Subtract with carry, result is either 00000000 or 11111111

XOR  R0  R0  } Reset register and program status word at once

# Example #2

`SBC R2 R2` } Subtract with carry, result is either 00000000 or 11111111

`XOR R0 R0` } Reset register and program status word at once

Need to handle such cases (and many more) to get precise analysis results

# Lesson #2

- Reasoning about binaries involves thinking about many nifty cases and obfuscated code
- Manual abstraction requires a lot of experience and engineering work
- Automatic abstraction deobfuscates the binary, represents the semantics as is
  - Nice for arithmetic obfuscations
  - There is no need to be smart here

# Example #3

```
AND R0 #15
AND R1 #15
XOR R0 R1
ADD R0 R1
```

$$0 \leq \mathbf{r0'} \leq 30$$
$$0 \leq \mathbf{r1'} \leq 15$$
$$0 \leq \mathbf{r0'} + \mathbf{r1'} \leq 45$$
$$0 \leq \mathbf{r0'} - \mathbf{r1'} \leq 15$$

```
XOR R0 R1
XOR R1 R0
XOR R0 R1
```

$$\mathbf{r0'} = \mathbf{r1}$$
$$\mathbf{r1'} = \mathbf{r0}$$

# Lesson #3

- Need to combine different abstract domains that can model different properties
  - Ranges vs. equalities
  - Doing this by hand is a horrible task
- Automatic abstraction computes different abstractions for free
  - No need to be busy here!

# Lessons Learnt

1. Block-wise abstraction improves precision
2. Automatic abstraction takes care of the evil cornercases
3. Automatic abstraction supports a variety of abstract domains for free

That's enough blurb, let's go technical!

# Approach

- Specify semantics for each instruction once and for all
- Combine different semantics to form a basic block
- Apply SSA conversion within the block
- Use decision procedure to compute abstractions for variety of domains
  - Equalities, affine relations, congruences, intervals, value sets, etc.
  - SAT/SMT solvers are great to reason about bit-vectors

# Conversion Into SSA

```
INC R0                    R0_1 := INC R0
MOV R1 R0                 R1_1 := R0_1
LSL R1                    R1_2 := LSL R1_1
SBC R1 R1        =>       R1' := SBC R1_2 R1_2
EOR R0 R1                 R0_2 := EOR R0_1 R1'
SUB R0 R1                 R0' := SUB R0_2 R1'
```

$R0_1 := INC\ R0$

$R1_1 := R0_1$

$R1_2 := LSL\ R1_1$

$R1' := SBC\ R1_2\ R1_2$

$R0_2 := EOR\ R0_1\ R1'$

$R0' := SUB\ R0_2\ R1'$

# Conversion Into Logic

```
R0₁ := INC R0
R1₁ := R0₁
R1₂ := LSL R1₁
R1' := SBC R1₂ R1₂
R0₂ := EOR R0₁ R1'
R0' := SUB R0₂ R1'
```

$$\begin{aligned}
& \sigma_{\text{INC}}(\mathbf{r0}_1, \mathbf{r0}) \\
\wedge\ & \sigma_{\text{MOV}}(\mathbf{r1}_1, \mathbf{r0}_1) \\
\wedge\ & \sigma_{\text{LSL}}(\mathbf{r1}_2, \mathbf{r1}_1) \\
\wedge\ & \sigma_{\text{SBC}}(\mathbf{r1'}, \mathbf{r1}_2, \mathbf{r1}_2) \\
\wedge\ & \sigma_{\text{EOR}}(\mathbf{r0}_2, \mathbf{r0}_1, \mathbf{r1'}) \\
\wedge\ & \sigma_{\text{SUB}}(\mathbf{r0'}, \mathbf{r0}_2, \mathbf{r1'})
\end{aligned}$$

$$\varphi$$

# A Variety of Abstractions

- Intervals
- Value sets
- Affine equalities
- Octagons

# Interval Abstraction

- Suppose $\mathbf{r0} \in [-10, 20]$ on input

- Put $\varphi' = \varphi \wedge (\mathbf{r0} \in [-10, 20])$

- What's the upper bound of $\mathbf{r0}'$ on output?

# Interval Abstraction

- Suppose $\mathbf{r0} \in [-10, 20]$ on input
- Put $\varphi' = \varphi \wedge (\mathbf{r0} \in [-10, 20])$
- What's the upper bound of $\mathbf{r0}'$ on output?

$-128$            $0$           $127$

$\varphi' \wedge \neg \mathbf{r0}'[7]$ ?

It's somewhere in this range

# Interval Abstraction

- Suppose $\mathbf{r0} \in [-10, 20]$ on input

- Put $\varphi' = \varphi \wedge (\mathbf{r0} \in [-10, 20])$

- What's the upper bound of $\mathbf{r0}'$ on output?

$-128$                 $0$                 $127$

$0$             $127$

$\varphi' \wedge \neg \mathbf{r0}'[7]$ ?

$\varphi' \wedge \neg \mathbf{r0}'[7] \wedge \mathbf{r0}'[6]$ ?

# Interval Abstraction

- Suppose $\mathbf{r0} \in [-10, 20]$ on input

- Put $\varphi' = \varphi \wedge (\mathbf{r0} \in [-10, 20])$

- What's the upper bound of $\mathbf{r0}'$ on output?

$-128$             $0$            $127$

$\varphi' \wedge \neg\mathbf{r0}'[7]$ ?     $0$       $127$

$\varphi' \wedge \neg\mathbf{r0}'[7] \wedge \mathbf{r0}'[6]$ ?     $0$     $63$

$\varphi' \wedge \neg\mathbf{r0}'[7] \wedge \mathbf{r0}'[6] \wedge \mathbf{r0}'[5]$ ?

# Interval Abstraction

- Suppose $\mathbf{r0} \in [-10, 20]$ on input
- Put $\varphi' = \varphi \wedge (\mathbf{r0} \in [-10, 20])$
- What's the upper bound of $\mathbf{r0}'$ on output?

$-128$            $0$          $127$

$\varphi' \wedge \neg\mathbf{r0}'[7]$ ?      $0$      $127$

$\varphi' \wedge \neg\mathbf{r0}'[7] \wedge \mathbf{r0}'[6]$ ?      $0$      $63$

$\varphi' \wedge \neg\mathbf{r0}'[7] \wedge \mathbf{r0}'[6] \wedge \mathbf{r0}'[5]$ ?      $0$      $31$

$\varphi' \wedge \neg\mathbf{r0}'[7] \wedge \mathbf{r0}'[6] \wedge \mathbf{r0}'[5] \wedge \mathbf{r0}'[4]$ ?
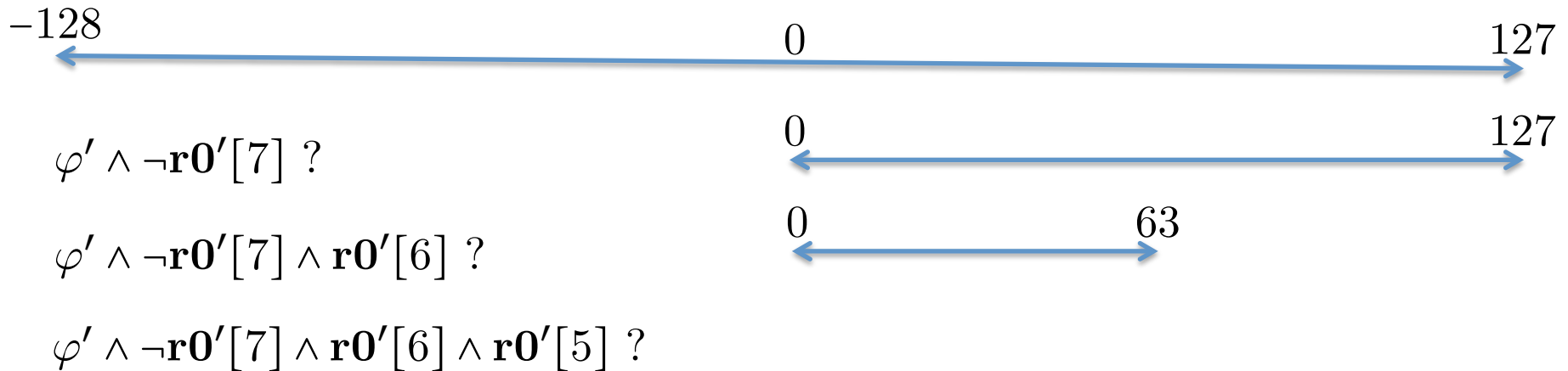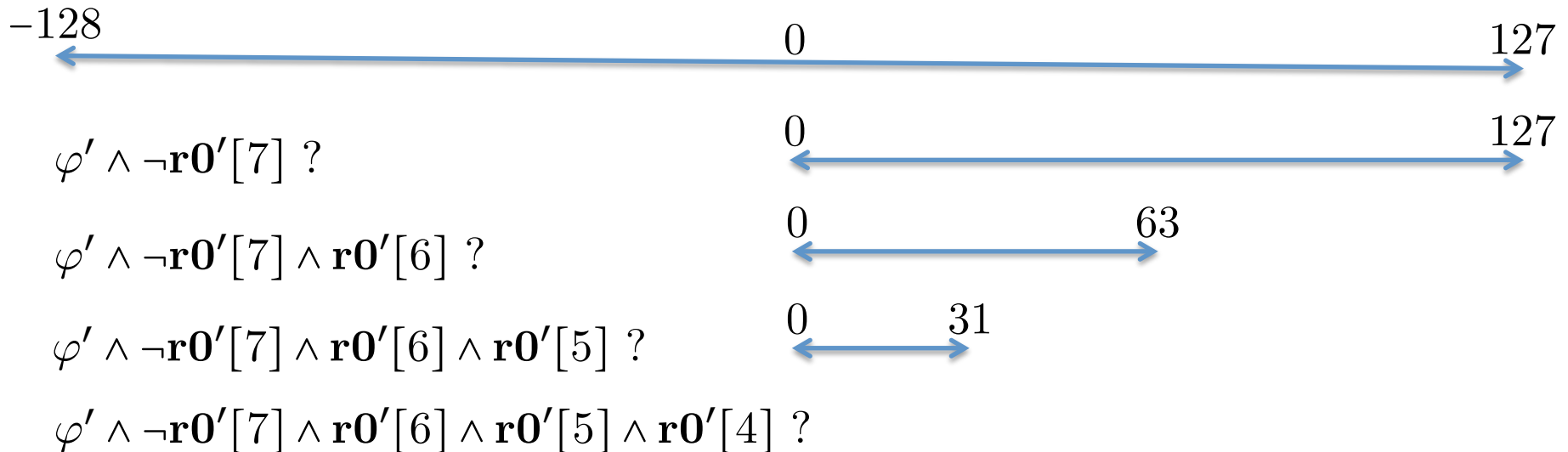
# Interval Abstraction

- Suppose $\mathbf{r0} \in [-10, 20]$ on input
- Put $\varphi' = \varphi \wedge (\mathbf{r0} \in [-10, 20])$
- What's the upper bound of $\mathbf{r0}'$ on output?

$-128$          $0$          $127$

$\varphi' \wedge \neg\mathbf{r0}'[7]$ ?    $0$       $127$

$\varphi' \wedge \neg\mathbf{r0}'[7] \wedge \mathbf{r0}'[6]$ ?    $0$      $63$

$\varphi' \wedge \neg\mathbf{r0}'[7] \wedge \mathbf{r0}'[6] \wedge \mathbf{r0}'[5]$ ?    $0$    $31$

$\varphi' \wedge \neg\mathbf{r0}'[7] \wedge \mathbf{r0}'[6] \wedge \mathbf{r0}'[5] \wedge \mathbf{r0}'[4]$ ?    $16$   $31$
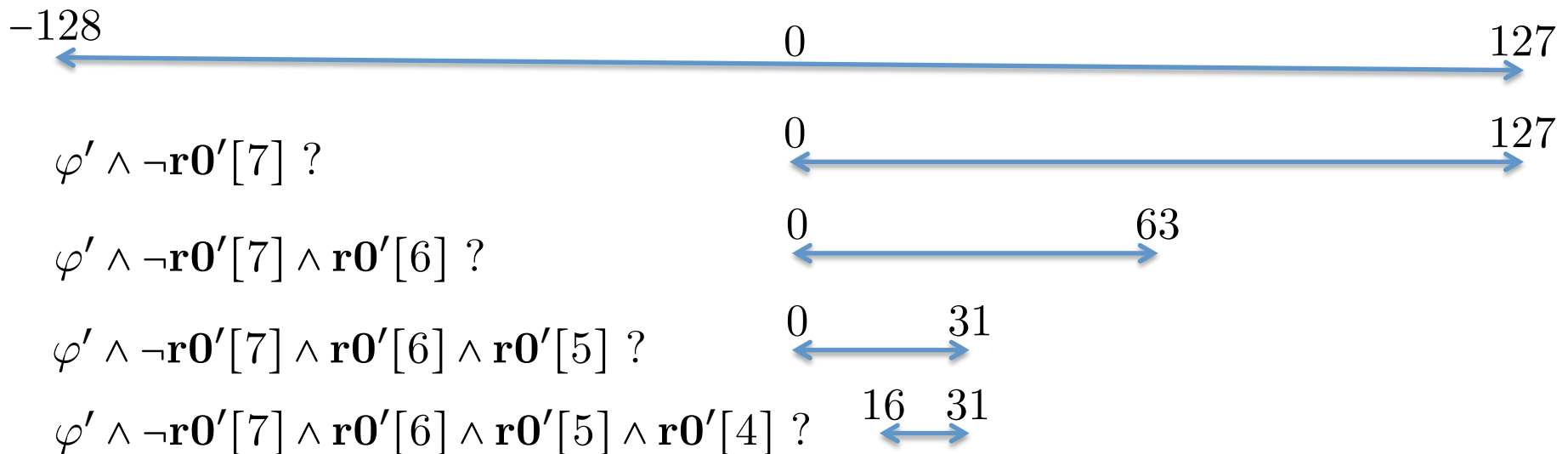
# Interval Abstraction

- Suppose $\mathbf{r0} \in [-10, 20]$ on input

- Put $\varphi' = \varphi \wedge (\mathbf{r0} \in [-10, 20])$

- What's the upper bound of $\mathbf{r0}'$ on output?

- Eventually: $\mathbf{r0}' \leq 21$

- Likewise, minimization gives $\mathbf{r0}' \geq 0$

# Literature: Interval Abstraction

- Codish, Lagoon & Stuckey: Logic Programming with Satisfiability (TPLP'08)

- Barrett & King: Range and Set Abstraction using SAT (NSAD'10)

- Brauer, King & Kowalewski: Range Analysis of Microcontroller Code using Bit-Level Congruences (FMICS'10)

- Brauer & King: Transfer Function Synthesis without Quantifier Elimination (ESOP'11)

# Value-Set Abstraction

- Suppose $\mathbf{r0} \in \{-10, 4, 20\}$ on input

- Put $\varphi' = \varphi \wedge (\mathbf{r0} \in \{-10, 4, 20\}$

- What's the value set of $\mathbf{r0}'$ on output?

- Enumerate models of $\varphi'$ using incremental SAT

# Value-Set Abstraction

- Suppose $\mathbf{r0} \in \{-10, 4, 20\}$ on input

- Put $\varphi' = \varphi \wedge (\mathbf{r0} \in \{-10, 4, 20\}$

- What's the value set of $\mathbf{r0}'$ on output?

- Enumerate models of $\varphi'$ using incremental SAT

$-128$       $0$    5       $127$

# Value-Set Abstraction

- Suppose $\mathbf{r0} \in \{-10, 4, 20\}$ on input

- Put $\varphi' = \varphi \wedge (\mathbf{r0} \in \{-10, 4, 20\}$

- What's the value set of $\mathbf{r0}'$ on output?

- Enumerate models of $\varphi'$ using incremental SAT

$-128$        $0$   5      $127$

- Consider $\varphi' \wedge (\mathbf{r0}' \neq 5)$

# Value-Set Abstraction

- Suppose $\mathbf{r0} \in \{-10, 4, 20\}$ on input

- Put $\varphi' = \varphi \wedge (\mathbf{r0} \in \{-10, 4, 20\}$

- What's the value set of $\mathbf{r0'}$ on output?

- Enumerate models of $\varphi'$ using incremental SAT

$-128$ ←————————————— $0$ ⑤ ——— ㉑ ————————————→ $127$

- Consider $\varphi' \wedge (\mathbf{r0'} \neq 5) \wedge (\mathbf{r0'} \neq 21)$

# Value-Set Abstraction

- Suppose $\mathbf{r0} \in \{-10, 4, 20\}$ on input

- Put $\varphi' = \varphi \wedge (\mathbf{r0} \in \{-10, 4, 20\}$

- What's the value set of $\mathbf{r0}'$ on output?

- Enumerate models of $\varphi'$ using incremental SAT

$-128$                     $0$   5   9   21               $127$

- $\varphi' \wedge (\mathbf{r0}' \neq 5) \wedge (\mathbf{r0}' \neq 21) \wedge (\mathbf{r0}' \neq 9)$ is unsat

# Value-Set Abstraction
## The Other Direction

- Suppose $\mathbf{r0}' \in \{5, 9, 21\}$
- Put $\varphi' = \varphi \wedge (\mathbf{r0}' \in \{5, 9, 21\})$

$-128$ 

-22  -10  -6  4  8  20

5  9  21

$127$

- We apply the same algorithm backwards
  - Over-approximation of inputs we had before
  - Also works for intervals

# Literature: Value-Set Abstraction

- Barrett & King: Range and Set Abstraction using SAT (NSAD'10)

- Reinbacher & Brauer: Precise Control Flow Reconstruction using Boolean Logic (EMSOFT'11)

- Brauer, King & Kriener: Existential Quantification as Incremental SAT (CAV'11)

# Affine Equalities

- Consider again

$$\varphi \quad = \quad \begin{cases} & \sigma_{\text{INC}}(\mathbf{r0}_1, \mathbf{r0}) \\ \wedge & \sigma_{\text{MOV}}(\mathbf{r1}_1, \mathbf{r0}_1) \\ \wedge & \sigma_{\text{LSL}}(\mathbf{r1}_2, \mathbf{r1}_1) \\ \wedge & \sigma_{\text{SBC}}(\mathbf{r1}', \mathbf{r1}_2, \mathbf{r1}_2) \\ \wedge & \sigma_{\text{EOR}}(\mathbf{r0}_2, \mathbf{r0}_1, \mathbf{r1}') \\ \wedge & \sigma_{\text{SUB}}(\mathbf{r0}', \mathbf{r0}_2, \mathbf{r1}') \end{cases}$$

restricted to normal operation of `INC` and `SUB` and overflow of `LSL`, denoted $\psi$

# Affine Equalities

- Goal: compute affine equality that describes relation between $r0$ on input and $r0'$ on output

- Natural choice: represent affine equalities as matrices

# Affine Hull

Iteration #1

- Pass $\psi$ to a solver

- Gives model $\quad m_1 = (\mathbf{r0} = -4 \wedge \mathbf{r0}' = 3)$

- Represent solution as matrix

$$\left[ \begin{array}{cc|c} 1 & 0 & -4 \\ 0 & 1 & 3 \end{array} \right] \quad \longleftarrow \quad \text{Focus on this row}$$

- Put $\psi' = \psi \wedge (\mathbf{r0}' \neq 3)$

# Affine Hull
Iteration #2

- Pass $\psi' = \psi \wedge (\mathbf{r0}' \neq 3)$ to a solver

- Gives model $m_2 = (\mathbf{r0} = -5 \wedge \mathbf{r0}' = 4)$

- Represent solution as matrix $\begin{bmatrix} 1 & 0 & | & -5 \\ 0 & 1 & | & 4 \end{bmatrix}$

- Now compute

$$\begin{bmatrix} 1 & 0 & | & -4 \\ 0 & 1 & | & 3 \end{bmatrix} \sqcup \begin{bmatrix} 1 & 0 & | & -5 \\ 0 & 1 & | & 4 \end{bmatrix} = \begin{bmatrix} 1 & 1 & | & -1 \end{bmatrix}$$

# Affine Hull

$$\begin{bmatrix} 1 & 0 & \Big| & -4 \\ 0 & 1 & \Big| & 3 \end{bmatrix} \sqcup \begin{bmatrix} 1 & 0 & \Big| & -5 \\ 0 & 1 & \Big| & 4 \end{bmatrix} \quad = \quad \begin{bmatrix} 1 & 1 & \Big| & -1 \end{bmatrix} \longleftarrow$$

- Put $\psi'' = \psi \wedge (\mathbf{r0} + \mathbf{r1} \neq -1)$

- Then, $\psi''$ is unsatisfiable

- Hence, $\mathbf{r0}' = -\mathbf{r0} - 1$ is the optimal affine abstraction of $\psi$

# Literature: Affine Relations

- Karr: Affine Relationships among Variables of a Program (Acta Informatica'76)

- Müller-Olm & Seidl: A Note on Karr's Algorithm (ICALP'04)

- Müller-Olm & Seidl: Analysis of Modular Arithmetic (ACM TOPLAS'07)

- Brauer & King: Automatic Abstraction for Intervals using Boolean Formulae (SAS'10)

# Combined Analyses

```
INC R0
MOV R1 R0
LSL R1
SBC R1 R1
EOR R0 R1
SUB R0 R1
```

$$(\mathbf{r0} = 127) \qquad \Rightarrow \qquad (\mathbf{r0'} = -128)$$
$$(-128 \leq \mathbf{r0} \leq -2) \qquad \Rightarrow \qquad (\mathbf{r0'} = -\mathbf{r0} - 1)$$
$$(-1 \leq \mathbf{r0} \leq 126) \qquad \Rightarrow \qquad (\mathbf{r0'} = \mathbf{r0} + 1)$$

```
R0' := abs(R0 + 1)
```

# Control Flow Reconstruction

[0x003]
$\{r_7 : 1-3, 101-103\}$
$\{flags : C := false\}$

MOV 0x08 ↤ R7
MOV A ↤ 0x08
CLR A
SUBB A, #N_HANDL
JC C:0x00F

$\{r_8 : 1-3, 101-103\}$
$\{r_A : 251-253, 95-97\}$
$\{flags : C := true, false\}$

$c := 1$

$c := 0$

[0x00C]
$\{r_A : 95-97\}$
$\{flags : C := false\}$
MOV R7 ↤ #C_FAIL
RET
$\{r_7 : \#C\_FAIL\}$

[0x00F]
$\{r_{0x08} : 1-3\}$
$\{flags : C := true\}$
MOV R7 ↤ 0x08
MOV A ↤ R7
MOV B ↤ #0x03
MUL AB
ADD A, #0x26
MOV DPL ↤ A
CLR A
ADDC A, #0x00
MOV DPH ↤ A
AJMP C:0x038
$\{r_{dph} : 0\}$
$\{r_{dpl} : 41, 44, 47\}$

[0x038]
$\{r_{dph} : 0\}$
$\{r_{dpl} : 41, 44, 47\}$
CLR A
MOVC A ↤ @(A+DP)
MOV R3 ↤ A
MOV A ↤ #0x01
MOVC A ↤ @(A+DP)
MOV R2 ↤ A
MOV A ↤ #0x02
MOVC ↤ @(A+DP)
MOV R1 ↤ A
AJMP C:0x100
$\{r_3 : 255\}$
$\{r_2 : 0\}$
$\{r_1 : 105, 110, 115\}$

[0x100]
$\{r_2 : 0\}$
$\{r_1 : 105, 110, 115\}$
MOV DPH ↤ R2
MOV DPL ↤ R1
CLR A
$\{r_{dph} : 0\}$
$\{r_{dpl} : 105, 110, 115\}$
$\{r_A : 0\}$

[0x106]
$\{r_{dph} : 0\}$
$\{r_{dpl} : 105, 110, 115\}$
$\{r_A : 0\}$
**iJMP @ (A+DP)**

abstracted jump targets
hdr2()✓,hdr3()✓,hdr4()✓

# Experimental Results

| Binary Program | | | | | $\overrightarrow{\mathcal{F}}$ interpreter | | | $\overrightarrow{\mathcal{F}} + \overleftarrow{\mathcal{B}}$ interpreter | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Name | Compiler | $loc_C$ | $instr_B$ | JT | RT | FT | Time | RS | $k$ | RT | FT | Time |
| Single row input | Keil | 80 | 67 | 6 | 2401 | 2395 | 2.6 | 2 | 2 | 6 | – | 3.32 |
|  | Sdcc |  | 52 |  | 460 | 454 | 2.4 | 2 | 2 | 6 | – | 2.0 |
| Keypad | Keil | 113 | 113 | 9 | 3844 | 3835 | 3.49 | 4 | 2 | 9 | – | 4.33 |
|  | Sdcc |  | 80 |  | 1508 | 1499 | 3.08 | 4 | 2 | 9 | – | 2.57 |
| Communication Link | Keil | 111 | 164 | 8 | 6889 | 6881 | 4.56 | 2 | 2 | 8 | – | 4.37 |
|  | Sdcc |  | 118 |  | 84 | 76 | 3.38 | 2 | 2 | 8 | – | 4.29 |
| Task Scheduler | Keil | 81 | 105 | 5 | >1000 | >995 | >5m | 17 | 2 | 5 | – | 14.03 |
|  | Sdcc |  | 97 |  |  |  |  | 23 | 2 | 5 | – | 10.23 |
| Switch Case | Keil | 82 | 166 | 19 | >5000 | >4981 | >5m | 94 | 2 | 19 | – | 17.49 |
|  | Sdcc |  | 180 |  | 3304 | 3285 | 2.31 | 6 | 2 | 38 | 19 | 2.6 |
| Emergency Stop | Keil | 138 | 150 | 9 | 768 | 759 | 2.8 | 2 | 2 | 9 | – | 2.6 |
|  | Sdcc |  | 141 |  | 256 | 247 | 2.9 | 2 | 2 | 9 | – | 3.1 |

| | | | | | |
|---|---|---|---|---|---|
| $loc_C$ | ... | Lines of C code | FT | ... | Number of recovered false targets |
| $instr_B$ | ... | Number of assembly instructions | RS | ... | Number of refinement steps applied |
| JT | ... | Number of jump targets | $k$ | ... | Backtracking depth |
| RT | ... | Number of recovered targets | Time | ... | Analysis time in seconds |

# So as to not cause offense

- Reps, Sagiv & Yorsh: Symbolic Implementation of the Best Transformer (VMCAI'04)
- Regehr & Reid: HOIST – A System For Automatically Deriving Static Analyzers for Embedded Systems (ASPLOS'04)
- Monniaux: Automatic Modular Abstractions for Linear Constraints (POPL'09)
- Monniaux: Automatic Modular Abstractions for Template Numerical Constraints (LMCS'10)
- Brauer & King: Automatic Abstraction for Intervals using Boolean Formulae (SAS'10)
- Brauer, King & Kowalewski: Range Analysis of Microcontroller Binary Code using Bit-Level Congruences (FMICS'10)
- Brauer & King: Transfer Function Synthesis without Quantifier Elimination (ESOP'11)
- Reinbacher & Brauer: Precise Control Flow Reconstruction using Boolean Logic (EMSOFT'11)

# Conclusion

- We advocate automatic abstraction as opposed to manual design
- SAT/SMT solvers can easily solve thousands of structured problems per second
- All techniques rely on the same encoding of the semantics
  - Solving for different abstract domains is slightly different
  - Can be put into uniform framework, but it is more efficient the way we put it

Excellent