

Program verification

Scaling analyses: 2 case studies

Laure Gonnord and David Monniaux

University of Lyon / LIP

November 17, 2015

Plan

Symbolic range analysis

- Motivation and big picture

- Overview

- Technical context LLVM

- Symbolic Range Analysis

- A bit on the other analyses

- Experimental results

- Conclusion

Symbolic pointer analysis

- Range Analysis

- Pointer Range Analysis

- Experimental results

Static Single Information form (SSI)



Inspiration : OOPSLA'14 slides:

VALIDATION OF MEMORY ACCESSES THROUGH SYMBOLIC ANALYSES

Henrique Nazare
Izabela Maffra
Willer Santos
Leonardo Oliveira
Fernando Quintão
Laure Gonnord

DCC
DEPARTAMENTO DE
CIÊNCIA DA COMPUTAÇÃO



UF *m* G

dip



<http://homepages.dcc.ufmg.br/~fernando/publications/presentations/OOPSLA14.pdf>

dip

Goal : Safety

Prove that (some) memory accesses are safe:

```
int main() {  
    int v[10];  
    v[0] = 0; ✓  
    return v[20]; ✗  
}
```

- ▶ Fight against bugs and overflow attacks.

Contributions (OOPSLA'14)

- ▶ A technique to prove that (some) memory accesses are safe :
 - ▶ Less need for additional guards.
 - ▶ Based on abstract interpretation.
 - ▶ Precision and cost compromise.
- ▶ Implemented in LLVM-compiler infrastructure :
 - ▶ Eliminate 50% of the guards inserted by AddressSanitizer
 - ▶ SPEC CPU 2006 17% faster



Our key insight : Symbolic (parametric) Analyses

```
int main(int argc, char** a) {
  char* p = malloc(argc);
  int i = 0;
  while (i < argc) {
    p[i] = 0;
    i++;
  }
  return 0;
}
```

$W(p) = [0, \text{argc} - 1]$.

$R(i) = [0, \text{argc} - 1]$

► $R(i) \subseteq W(p)$ thus $p[i]$ is **safe**.



A bit on sanitizing memory accesses

Different techniques : but all have an overhead.

Ex : Address Sanitizer

- ▶ Shadow every memory allocated : 1 byte \rightarrow 1 bit (allocated or not).
- ▶ Guard every array access : check if its shadow bit is valid.
 - ▶ slows down SPEC CPU 2006 by 25%
- ▶ We want to **remove these guards**.



Green Arrays : overview 1/2

```
1. int main(int argc, char** argv) {
2.     int size = argc + 1;
3.     char* buf = malloc(size);
4.     unsigned index = 0;
5.     scanf("%u", &index);
6.     if (index < argc) {
7.         buf[index] = 0;
8.     }
9.     return index;
10. }
```

Any address
from buf + 0
to buf + argc
is safe!

Inside the
branch index is
at least 0 and
at most argc-1

We know that
"argc - 1" is
less than argc

As long as
we do not
have integer
overflows!

Green Arrays : overview 2/2

Symbolic Range Analysis:

finds the lower and upper values that variables can assume

Any address from $\text{buf} + 0$ to $\text{buf} + \text{argc}$ is safe!

Symbolic Region Analysis:

finds the lower and upper values that a pointer can address

Inside the branch index is at least 0 and at most argc-1

Integer Overflow Analysis:

Which arithmetic operations can overflow?

We know that " $\text{argc} - 1$ " is less than argc

As long as we do not have integer overflows!

A bit on LLVM

LLVM is a **compiler infrastructure**:

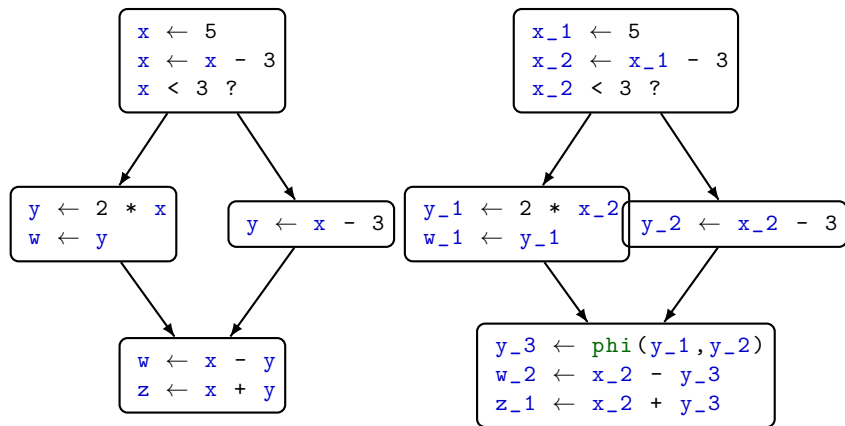


- ▶ Open source
- ▶ Various frontends (C, C++, Fortran)
- ▶ Various code generators (x86, ...)

Writing optimisations is easier :

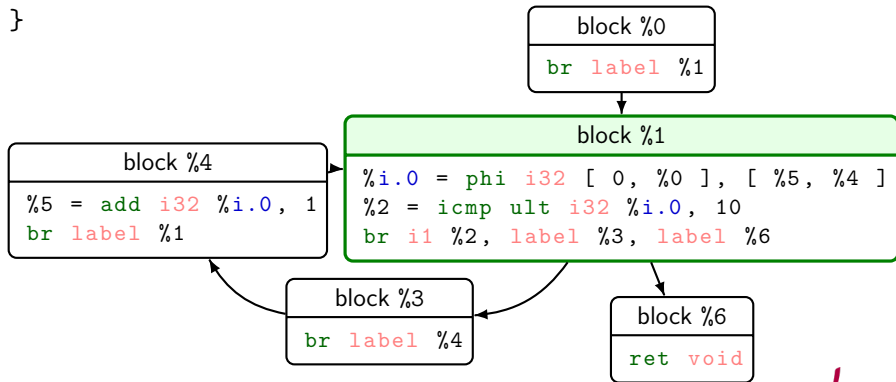
- ▶ A unique IR (**intermediate representation**)
- ▶ C++ iterators (functions, blocks, ...)

LLVM representation : SSA form

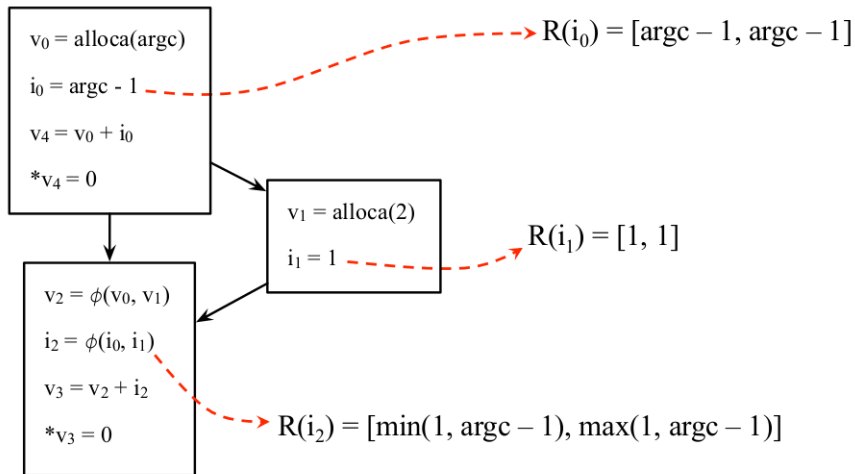


LLVM representation: a toy example

```
void simple_loop_constant() {  
    for(unsigned i=0; i<10; i++) {  
        // Do nothing  
    }  
}
```



Symbolic Ranges (SRA): Goal



SRA on SSA form: a sparse analysis

- ▶ An abstract interpretation-based technique.
- ▶ Very similar to classic range analysis.
- ▶ One abstract value (R) **per variable**: sparsity.
- ▶ Easy to implement (simple algorithm, simple data structure).



SRA on SSA form: constraint system

$$v = \bullet \Rightarrow R(v) = [v, v]$$

$$v = o \Rightarrow R(v) = R(o)$$

$$v = v_1 \oplus v_2 \Rightarrow R(v) = R(v_1) \oplus' R(v_2)$$

$$v = \phi(v_1, v_2) \Rightarrow R(v) = R(v_1) \sqcup R(v_2)$$

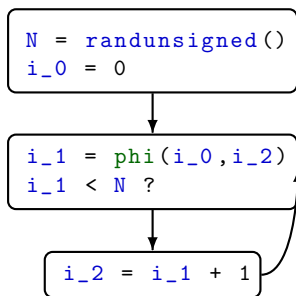
$$\text{other instructions} \Rightarrow \emptyset$$

\oplus' : abstract effect of the operation \oplus on two intervals.

\sqcup : convex hull of two intervals. \blacktriangleright All these operation are performed symbolically thanks to **GiNaC**



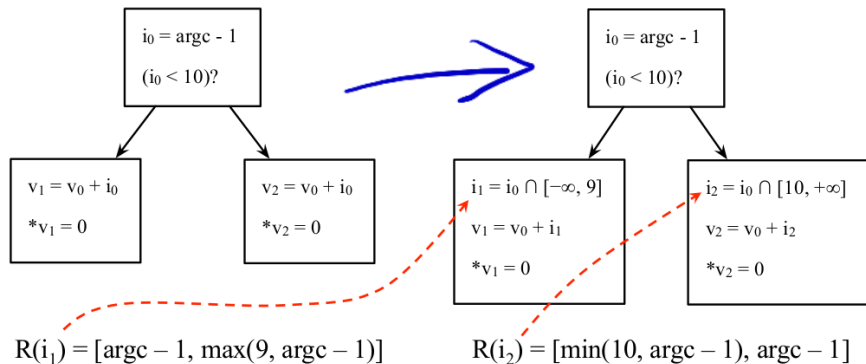
SRA on SSA form: an example



- ▶ $R(i_0) = [0, 0]$
- ▶ $R(i_1) = [0, +\infty]$
- ▶ $R(i_2) = [1, +\infty]$

Improving precision of SRA : live-range splitting

1/2



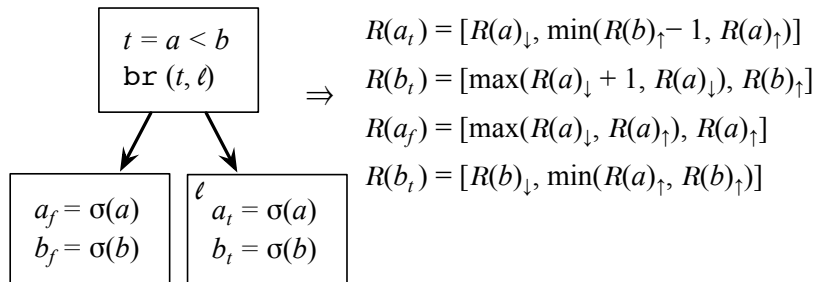
► e-SSA form.



Improving precision of SRA : live-range splitting

2/2

Rule for live-range splitting :



► All simplifications are done by GiNaC.

SRA + live-range on an example

```
N = randunsigned()
i_0 = 0
```

```
i_1 = phi(i_0, i_2)
i_1 < N ?
```

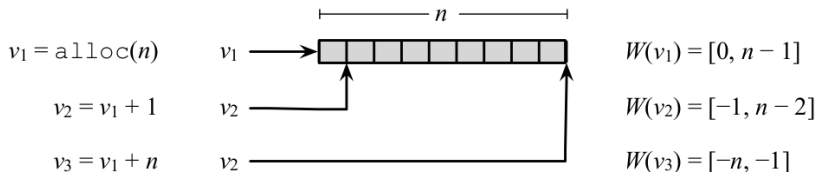
```
i_t = sigma(i_1)
i_2 = i_t + 1
```

- ▶ $R(i_0) = [0, 0]$
- ▶ $R(i_1) = [0, N]$

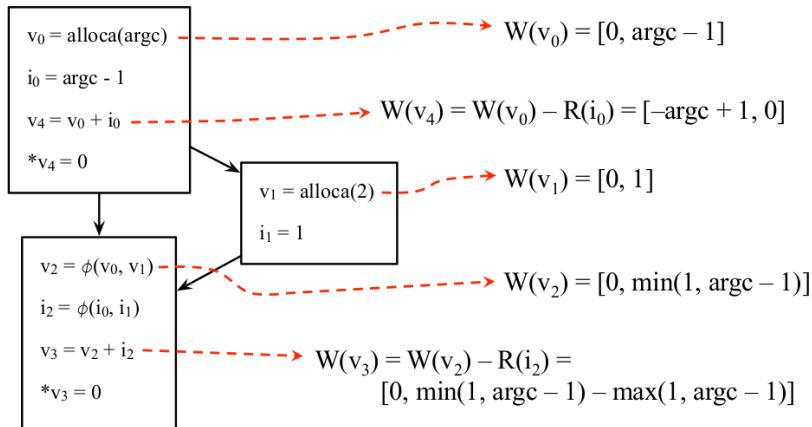
$$R(i_t) = [R(i_1) \downarrow, \min(N - 1, R(i_1) \uparrow)]$$

Symbolic regions 1/2: Goal

Compute (an underapproximation of) the range of **valid accesses** from base pointers:

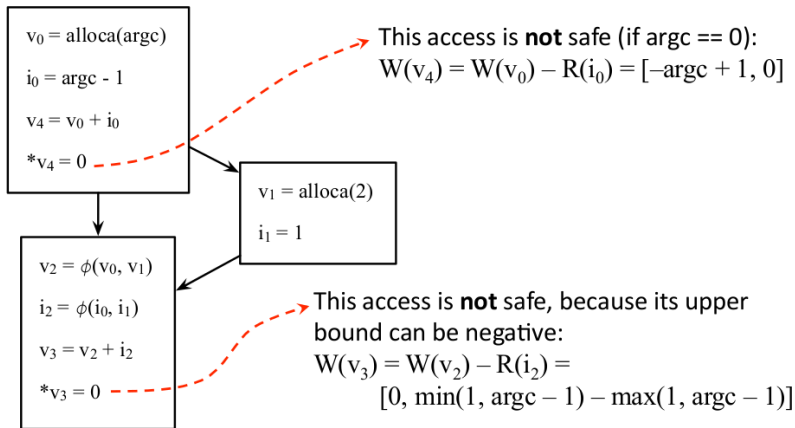


Symbolic regions 2/2: An example



Safety: result

If $0 \in W(p)$, then $*p$ is **safe**, else **DK**



Overflows 1/2

```
int main(int argc, char** argv) {  
    int index = argc + 1;  
    int size = index * index;  
    char* buf = malloc(size);  
    return buf[index];  
}
```

Because we manipulate symbols, "argc + 1 < (argc + 1) * (argc + 1)" only in the absence of integer overflows

index * index
may wrap
around.

Do you know
what malloc
will return?

Overflows 2/2

- We find every arithmetic operation that may influence memory **allocation** or memory **indexing**.

```
int main(int argc, char** argv) {  
    int index = argc + 1;  
    int size = index * index;  
    char* buf = malloc(size);  
    return buf[index];  
}
```



We find them
via program
slicing

- We **instrument the code** to detect overflows.

Experimental setup

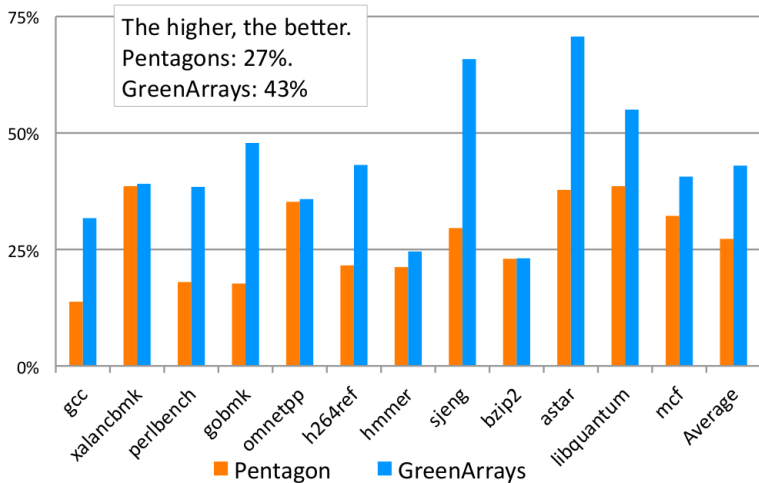
- **Implementation:** LLVM + AddressSanitizer
- **Benchmarks:** SPEC CPU 2006 + LLVM test suite
- **Machine:** Intel(R) Xeon(R) 2.00GHz, with 15,360KB of cache and 16GB of RAM
- **Baseline:** Pentagons
 - Abstract interpretation that combines "less-than" and "integer ranges".[†]

```
int i = 0;
unsigned j = read();
if (...)
    i = 9;
if (j < i)
    ...
```

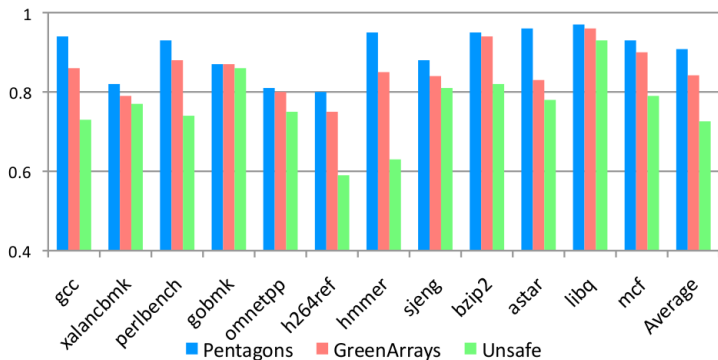
$P(j) = (\text{less than } \{i\}, [0, 8])$

[†]: Pentagons: A weakly relational abstract domain for the efficient validation of array accesses, 2010, Science of Computer Programming

Percentage of bound checks removed



Runtime improvement



The lower the bar, the faster. Time is normalized to AddressSanitizer without bound-check elimination. Average speedup: Pentagons = 9%. GreenArrays = 16%.



In the paper (OOPSLA'14)

A complete formalisation of all the analyses :

- ▶ Concrete and abstract semantics.
- ▶ Safety is proved.
- ▶ Interprocedural analysis.
- ▶ <https://code.google.com/p/ecosoc/>

Remaining question : improving precision of the symbolic range analysis ?



Plan

Symbolic range analysis

Motivation and big picture

Overview

Technical context LLVM

Symbolic Range Analysis

A bit on the other analyses

Experimental results

Conclusion

Symbolic pointer analysis

Range Analysis

Pointer Range Analysis

Experimental results

Static Single Information form (SSI)



Credit : M. Maalej. Accepted to CGO'16.



Goal + Contribution

Goal:

- ▶ Optimizing languages with pointers;
- ▶ Solving pointer arithmetic, disambiguating pointers;
- ▶ Low cost analysis.

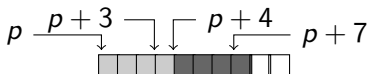
Contribution:

- ▶ Combine alias analysis with range analysis;
- ▶ Speed up;



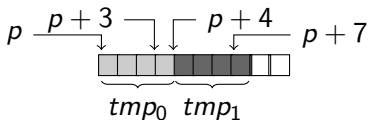
Motivating example

```
void fill_array (char* p) {  
    int i ;  
    for (i = 0; i < 4; i++)  
        p[i] = 0 ;  
    for (i = 4; i < 8; i++)  
        p[i] = 2 × i ;  
}
```

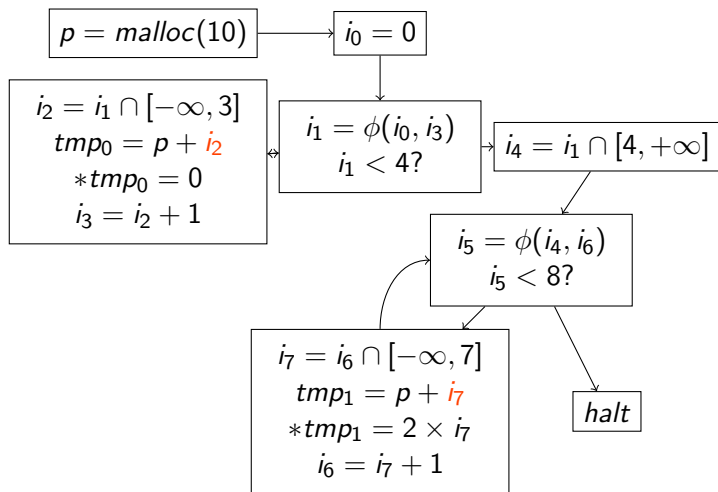


Motivating example

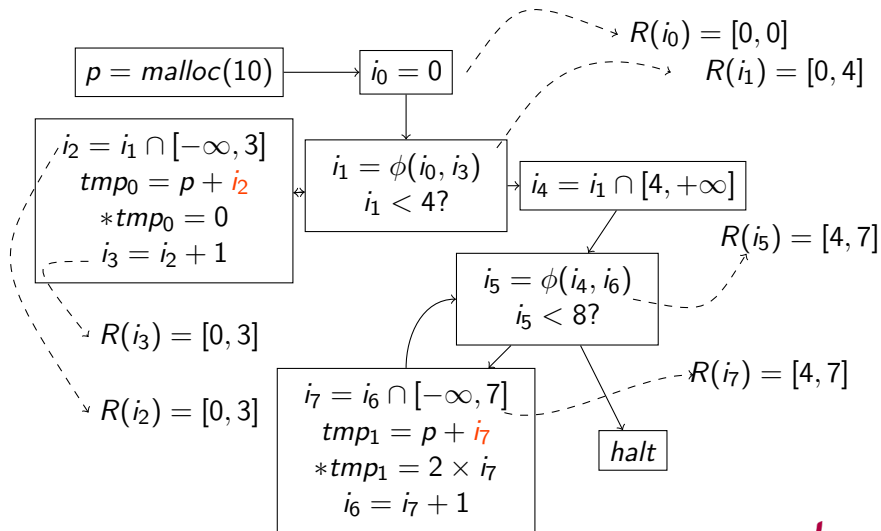
```
void fill_array (char* p) {  
    int i ;  
    for (i = 0; i < 4; i++)  
        p[i] = 0 ; // char* tmp0 = p + i ; *tmp0 = 0  
    for (i = 4; i < 8; i++)  
        p[i] = 2 × i ; // char* tmp1 = p + i ; *tmp0 = 2 × i  
}
```



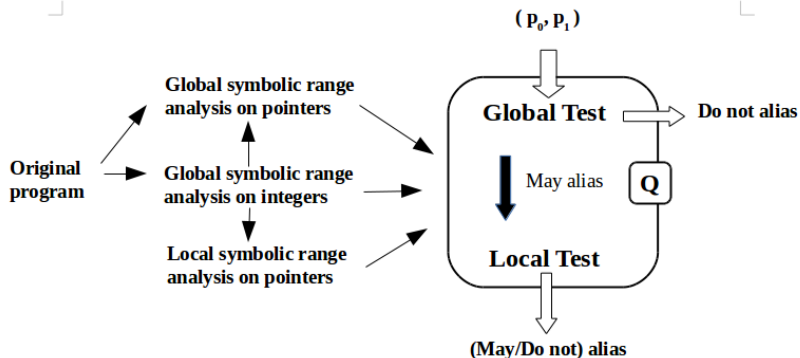
Range Analysis on e-SSA (recall)



Range Analysis on e-SSA (recall)



Pointer Range Analysis



Abstract Analysis

$$\begin{array}{l} p_i \rightsquigarrow loc_i + [l_i, u_i] \\ p_j \rightsquigarrow loc_j + [l_j, u_j] \end{array} \quad i = j \text{ or } i \neq i$$

Property (overapprox):

If $(loc_i = loc_j \text{ and } [l_i, u_i] \cap [l_j, u_j] = \emptyset)$

Then p_i and p_j **do not alias** Else **may alias** .



Global Pointer Range Analysis

Let n be the number of program sites where memory is allocated. We associate pointers with tuples of size n : $(\text{SymbRanges} \cup \perp)^n$: $\text{GR}(p) = (p_0, \dots, p_{n-1})$.

Notation

$$\blacktriangleright \text{GR}(p) = \{loc_i + p_i, loc_j + p_j, \dots\}$$

Constraint System:

$$j : p = \text{malloc}(v) \Rightarrow \text{GR}(p) = (\perp, \dots, [0, 0]_j, \dots)$$

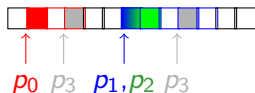
$$v = v_1 \Rightarrow \text{GR}(v) = \text{GR}(v_1)$$

$$\begin{array}{l} q = p + c \\ \text{with } c \text{ scalar} \end{array} \Rightarrow q_i = \begin{cases} \perp & \text{if } p_i = \perp \\ p_i + R(c) & \text{else} \end{cases}$$

$$q = \phi(p^1, p^2) \Rightarrow \text{GR}(q) = \text{GR}(p^1) \sqcup \text{GR}(p^2)$$



Global Pointer Range Analysis



$GR(p_0) = loc_0 + [0, 0]$

$GR(p_1) = loc_1 + [0, 0]$

$GR(p_2) = loc_1 + [0, 1]$

$GR(p_3) = \{loc_0 + [2, 2], loc_1 + [3, 3]\}$

```
 $p_0 = \text{malloc}(3) ;$ 
```

```
 $p_1 = \text{malloc}(5) ;$ 
```

```
if (...)  $p_2 = p_1 ;$ 
```

```
else  $p_2 = p_1 + 1 ;$ 
```

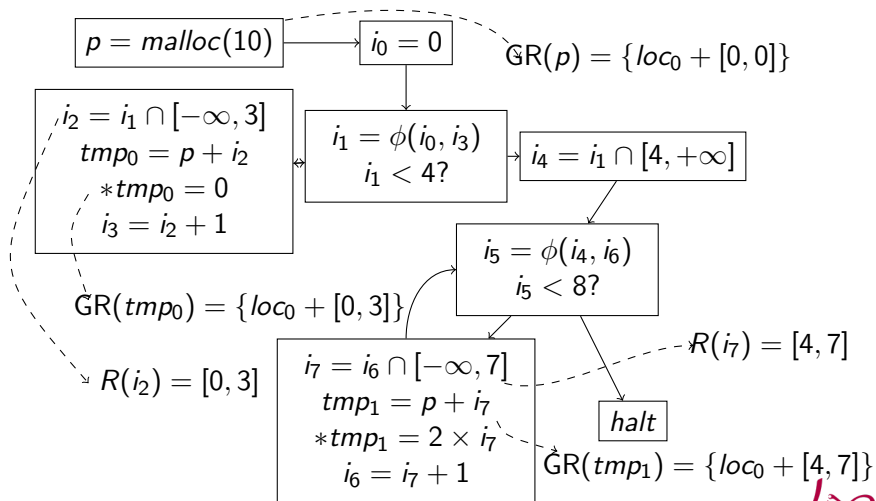
```
if (...)  $p_3 = p_0 + 2 ;$ 
```

```
else  $p_3 = p_1 + 3$ 
```



Global Pointer Range Analysis

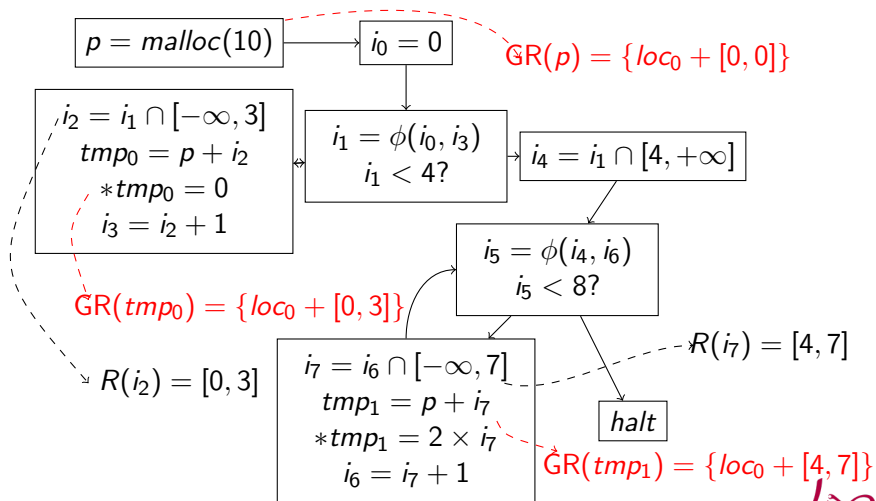
Example:



lip

Global Pointer Range Analysis

Example:



Local Pointer Range Analysis

Motivation:

$a_0 = \text{malloc}(N) ; \rightsquigarrow \text{GR}(a_0) = \{\text{loc}_0 + [0, 0]\}$

if (...)

$a_1 = a_0 + 1 ; \rightsquigarrow \text{GR}(a_1) = \{\text{loc}_0 + [1, 1]\}$

$a_2 = a_1 ; \rightsquigarrow \text{GR}(a_2) = \{\text{loc}_0 + [1, 1]\}$

else $a_2 = a_0 ; \rightsquigarrow \text{GR}(a_2) = \{\text{loc}_0 + [0, 1]\}$

$a_3 = a_2 + 1 ; \rightsquigarrow \text{GR}(a_3) = \{\text{loc}_0 + [1, 2]\}$

$a_4 = a_2 + 2 ; \rightsquigarrow \text{GR}(a_4) = \{\text{loc}_0 + [2, 3]\}$

$$[1, 2] \cup [2, 3] \neq \emptyset$$



Local Pointer Range Analysis

Motivation:

$a_0 = \text{malloc}(N)$; \rightsquigarrow $\text{GR}(a_0) = \{\text{loc}_0 + [0, 0]\}$

if (...)

$a_1 = a_0 + 1$; \rightsquigarrow $\text{GR}(a_1) = \{\text{loc}_0 + [1, 1]\}$

$a_2 = a_1$; \rightsquigarrow $\text{GR}(a_2) = \{\text{loc}_0 + [1, 1]\}$

else $a_2 = a_0$; \rightsquigarrow $\text{GR}(a_2) = \{\text{loc}_0 + [0, 1]\}$ \rightsquigarrow $\text{LR}(a_2) = \{\text{loc}_1 + [0, 0]\}$

$a_3 = a_2 + 1$; \rightsquigarrow $\text{GR}(a_3) = \{\text{loc}_0 + [1, 2]\}$ \rightsquigarrow $\text{LR}(a_3) = \{\text{loc}_1 + [1, 1]\}$

$a_4 = a_2 + 2$; \rightsquigarrow $\text{GR}(a_4) = \{\text{loc}_0 + [2, 3]\}$ \rightsquigarrow $\text{LR}(a_4) = \{\text{loc}_1 + [2, 2]\}$

$$[1, 2] \cup [2, 3] \neq \emptyset$$



Local Pointer Range Analysis

Constraint System:

$$p = \text{malloc}(v) \quad \Rightarrow \quad \text{LR}(p) = \text{NewLocs}() + [0, 0]$$

with v scalar

$$j : q = \phi(p_1, p_2) \quad \Rightarrow \quad \text{LR}(q) = \text{loc}_j + [0, 0]$$

with $\text{loc}_j = \text{NewLocs}()$



Experimental setup

- ▶ Implementation : LLVM 3.5
- ▶ Benchmarks : LLVM test suite + Micro benchmarks + PtrDist + Prolangs + MallocBench
- ▶ Machine : Intel i7-4770K, 8GB of memory Ubuntu 14.04.2

Experimental results

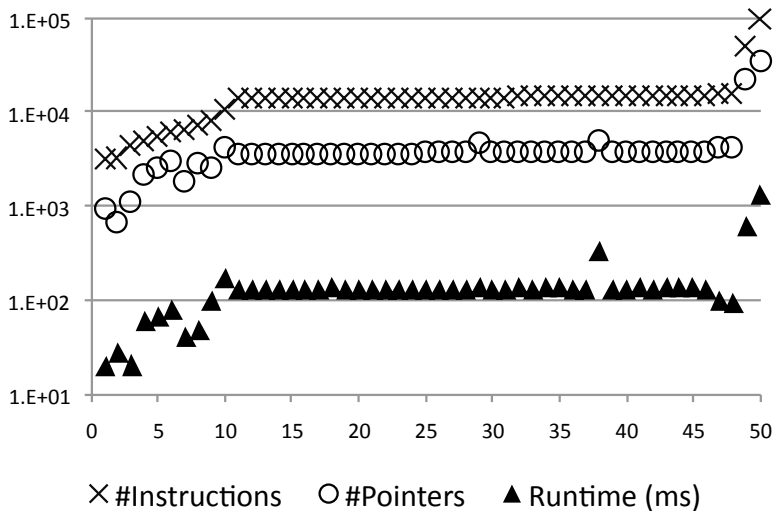
	#Queries	scev	basic	rbaa	rbaa + basic
Total %	7,243,418	2.8	17.8	16.9	22.1

- ▶ #Queries : number of pair of pointers.
- ▶ scev : scalar evolution based alias-analysis.
- ▶ basic : -O3 LLVM analysis (global + local pointers).
- ▶ rbaa : range based alias analysis.

Answering queries: number of pairs that do not alias.



Experimental results



Conclusion

This analysis scales well!



Plan

Symbolic range analysis

- Motivation and big picture

- Overview

- Technical context LLVM

- Symbolic Range Analysis

- A bit on the other analyses

- Experimental results

- Conclusion

Symbolic pointer analysis

- Range Analysis

- Pointer Range Analysis

- Experimental results

Static Single Information form (SSI)



Credit : F. Peirera and Fabrice Rastello. (Acaces 2015)



Goal + Contribution

Goal:

- ▶ Static Analyses that scale.
- ▶ Static but precise.

Contribution:

- ▶ A generic framework.
- ▶ A general way to solve the problem.





Introduction

- Data-flow analysis: discover facts (information) that are true about a program. Bind to *Variables* \times *ProgramPoints*.
- Static Single Information (SSI) property: IR such that information of a variable invariant along its whole live-range
- ϕ -functions split live-ranges where reaching definitions collide: SSA fulfills SSI property for constant analysis. Not for class inference (backward from uses).
- Extended SSA: SSI property for forward analysis flowing from definitions and conditional tests.
- SSU: SSI property for backward analysis flowing from uses

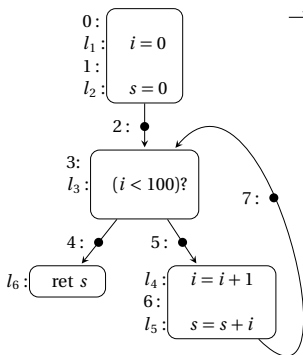
Can we generalize?

Sparse Analysis

Non-relational (dense) analysis: bind information to pairs *Variables* × *ProgPoints*

```

i = 0;
s = 0;
while (i < 100)
    i = i + 1;
    s = s + i;
ret
    
```



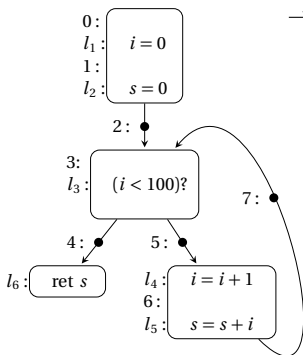
prog. point	[i]	[s]
0	⊥	⊥
1	[0, 0]	⊥
2	[0, 0]	[0, 0]
3	[0, 100]	[0, +∞[
4	[100, 100]	[0, +∞[
5	[0, 99]	[0, +∞[
6	[0, 100]	[0, +∞[
7	[0, 100]	[0, +∞[

Sparse Analysis

Range Analysis: $[v]^p$ intervals of possible values variable v might assume at program point p

```

i = 0;
s = 0;
while (i < 100)
    i = i + 1;
    s = s + i;
ret
    
```



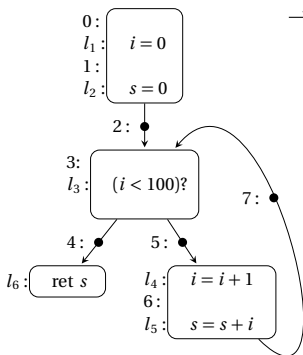
prog. point	$[i]$	$[s]$
0	\top	\top
1	$[0, 0]$	\top
2	$[0, 0]$	$[0, 0]$
3	$[0, 100]$	$[0, +\infty[$
4	$[100, 100]$	$[0, +\infty[$
5	$[0, 99]$	$[0, +\infty[$
6	$[0, 100]$	$[0, +\infty[$
7	$[0, 100]$	$[0, +\infty[$

Sparse Analysis

Redundancies: e.g. $[i]^1 = [i]^2$; because identity transfer function for $[i]$ from 1 to 2.

```

i = 0;
s = 0;
while (i < 100)
    i = i + 1;
    s = s + i;
ret
    
```



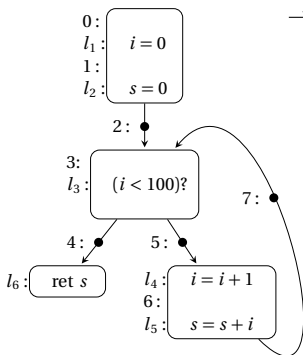
prog. point	$[i]$	$[s]$
0	\top	\top
1	$[0, 0]$	\top
2	$[0, 0]$	$[0, 0]$
3	$[0, 100]$	$[0, +\infty[$
4	$[100, 100]$	$[0, +\infty[$
5	$[0, 99]$	$[0, +\infty[$
6	$[0, 100]$	$[0, +\infty[$
7	$[0, 100]$	$[0, +\infty[$

Sparse Analysis

Sparse data-flow analysis: shortcut identity transfer functions by grouping contiguous program points bound to identities into larger regions

```

i = 0;
s = 0;
while (i < 100)
    i = i + 1;
    s = s + i;
ret
    
```



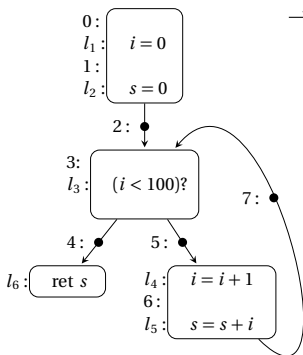
prog. point	$[i]$	$[s]$
0	\top	\top
1	$[0, 0]$	\top
2	$[0, 0]$	$[0, 0]$
3	$[0, 100]$	$[0, +\infty[$
4	$[100, 100]$	$[0, +\infty[$
5	$[0, 99]$	$[0, +\infty[$
6	$[0, 100]$	$[0, +\infty[$
7	$[0, 100]$	$[0, +\infty[$

Sparse Analysis

Sparse data-flow analysis: replace all $[v]^p$ by $[v]$ ($\forall v, p \in \text{live}(v)$); propagate along def-use chains.

```

i = 0;
s = 0;
while (i < 100)
    i = i + 1;
    s = s + i;
ret
    
```



prog. point	$[i]$	$[s]$
0	\top	\top
1	$[0, 0]$	\top
2	$[0, 0]$	$[0, 0]$
3	$[0, 100]$	$[0, +\infty[$
4	$[100, 100]$	$[0, +\infty[$
5	$[0, 99]$	$[0, +\infty[$
6	$[0, 100]$	$[0, +\infty[$
7	$[0, 100]$	$[0, +\infty[$

Partitioned Lattice per Variable Problems

Partitioned Lattice per Variable (PLV) Problem

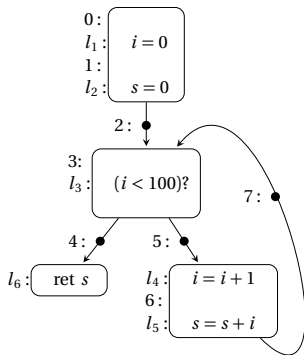
- program variables: v_i ; program points: p ; lattice: \mathcal{L}
- abstract state associated to prog. point p : x^p
- transfer function associated with $s \in \text{preds}(p)$: $F^{s,p}$
- constraint system: $x^p = x^p \wedge F^{s,p}(x^s)$ (or eq. $x^p \sqsubseteq F^{s,p}(x^s)$)

The corresponding Max. Fixed Point (MFP) problem is a PLV problem iff $\mathcal{L} = \mathcal{L}_{v_1} \times \dots \times \mathcal{L}_{v_n}$ where each \mathcal{L}_{v_i} is the lattice associated with v_i i.e. $x^s = ([v_1]^s, \dots, [v_n]^s)$. Thus $F^{s,p} = F_{v_1}^{s,p} \times \dots \times F_{v_n}^{s,p}$ and $[v_i]^p = [v_i]^p \wedge F_{v_i}^{s,p}([v_1]^s, \dots, [v_n]^s)$.

Partitioned Lattice per Variable Data-Flow Problem

Range analysis

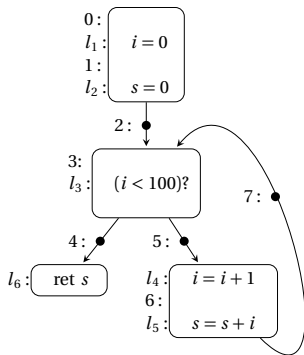
- $[i]^0 = [i]^0 \wedge F_i^{r,0}([i]^r, [s]^r)$
- $[i]^1 = [i]^1 \wedge F_i^{l_1}([i]^0, [s]^0)$
- $[i]^2 = [i]^2 \wedge F_i^{l_2}([i]^1, [s]^1)$
- $[i]^3 = [i]^3 \wedge F_i^{2,3}([i]^2, [s]^2)$
- $[i]^3 = [i]^3 \wedge F_i^{7,3}([i]^7, [s]^7)$
- $[i]^4 = [i]^4 \wedge F_i^{\bar{l}_3}([i]^3, [s]^3)$
- $[i]^5 = [i]^5 \wedge F_i^{l_3}([i]^3, [s]^3)$
- $[i]^6 = [i]^6 \wedge F_i^{l_4}([i]^5, [s]^5)$
- $[i]^7 = [i]^7 \wedge F_i^{l_5}([i]^6, [s]^6)$



Partitioned Lattice per Variable Data-Flow Problem

Range analysis

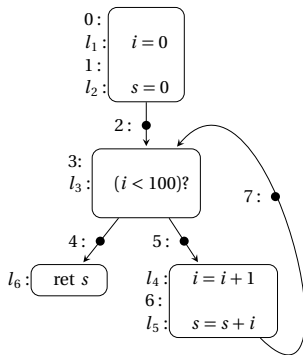
- $[i]^0 = [i]^0 \cup F_i^{r,0}([i]^r, [s]^r)$
- $[i]^1 = [i]^1 \cup F_i^{l_1}([i]^0, [s]^0)$
- $[i]^2 = [i]^2 \cup F_i^{l_2}([i]^1, [s]^1)$
- $[i]^3 = [i]^3 \cup F_i^{2,3}([i]^2, [s]^2)$
- $[i]^3 = [i]^3 \cup F_i^{7,3}([i]^7, [s]^7)$
- $[i]^4 = [i]^4 \cup F_i^{\bar{l}_3}([i]^3, [s]^3)$
- $[i]^5 = [i]^5 \cup F_i^{l_3}([i]^3, [s]^3)$
- $[i]^6 = [i]^6 \cup F_i^{l_4}([i]^5, [s]^5)$
- $[i]^7 = [i]^7 \cup F_i^{l_5}([i]^6, [s]^6)$



Partitioned Lattice per Variable Data-Flow Problem

Range analysis

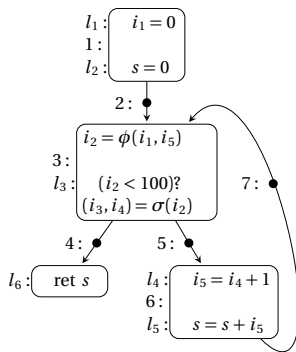
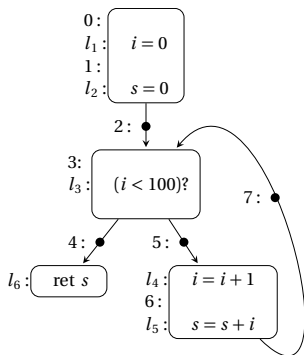
- $[i]^0 = [i]^0$
- $[i]^1 = [i]^1 \cup [0, 0]$
- $[i]^2 = [i]^2 \cup [i]^1$
- $[i]^3 = [i]^3 \cup [i]^2$
- $[i]^3 = [i]^3 \cup [i]^7$
- $[i]^4 = [i]^4 \cup ([i]^3 \cap [100, +\infty[)$
- $[i]^5 = [i]^5 \cup ([i]^3 \cap]-\infty, 99])$
- $[i]^6 = [i]^6 \cup ([i]^5 + 1)$
- $[i]^7 = [i]^7 \cup [i]^6$



The Static Single Information Property

SSIfy (forward)

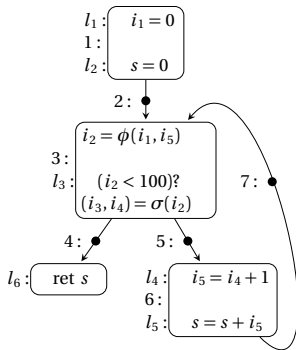
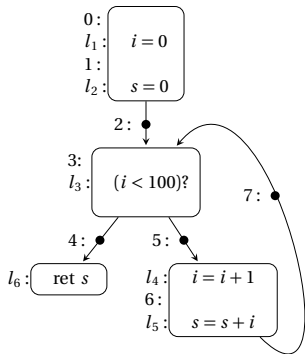
Modify the code (split live-ranges) without modifying its semantic s.t. fullfils SSI property



The Static Single Information Property

SPLIT

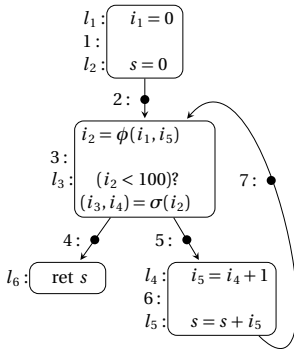
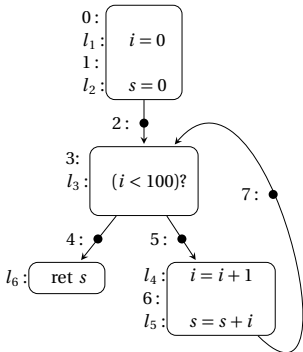
if s unique pred. of $p \in \text{live}(v)$ and such that $F_v^{s,p} \neq \lambda x. \top$ is non-trivial, then s should contain a definition of v



The Static Single Information Property

SPLIT

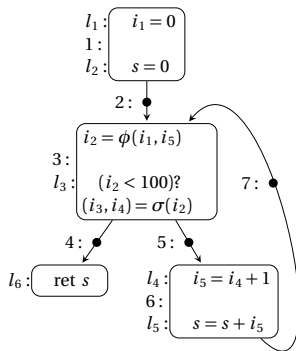
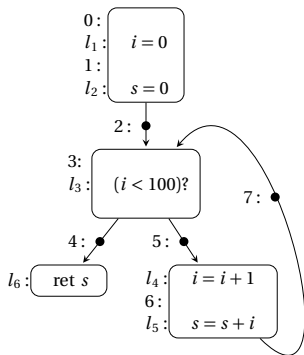
if s and t two preds of p such that $F_v^{s,p}(Y) \neq F_v^{t,p}(Y)$ (Y a MFP solution), then there must be a ϕ -function at entry of p



The Static Single Information Property

INFO

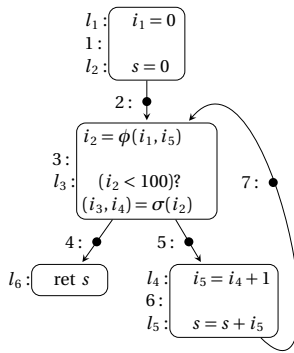
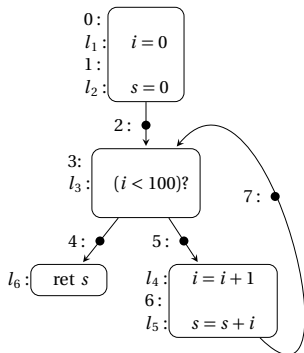
if $F_v^{s,p} \neq \lambda x. \top$,
then $v \in \text{live}(p)$



The Static Single Information Property

VERSION

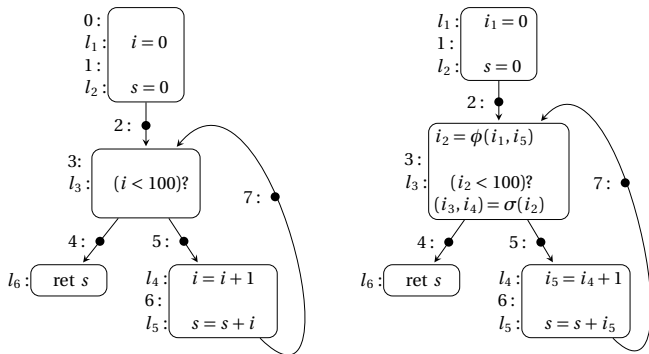
for each variable v , $\text{live}(v)$ is a connected component of the CFG.



The Static Single Information Property

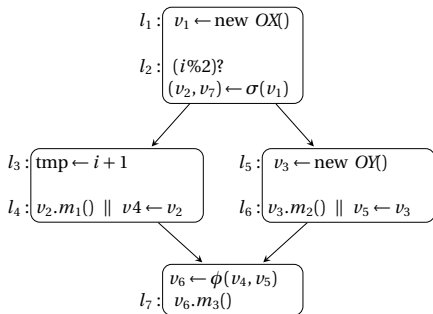
LINK

if F_v^{inst} depends on some $[u]^s$,
then $inst$ should contain an use of u live-in at $inst$.



Special instructions used to split live ranges

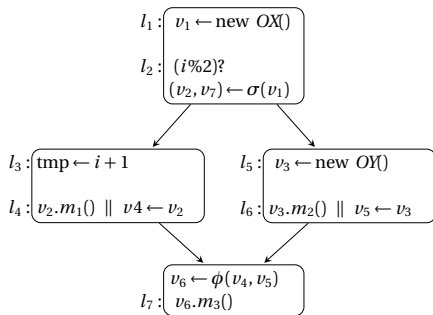
Interior nodes (unique predecessor, unique successor)

$$inst \parallel v_1 = v'_1 \parallel \dots \parallel v_m = v'_m$$


Special instructions used to split live ranges

joins (multiple predecessors, one successor)

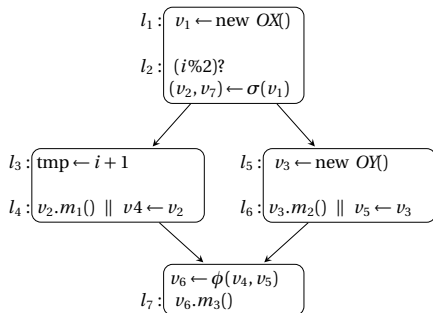
ϕ -functions



Special instructions used to split live ranges

branch points (one predecessor, multiple successors)

$$(l^1 : v_1^1, \dots, l^q : v_1^q) = \sigma(v_1) \quad || \quad \dots \quad || \quad (l^1 : v_m^1, \dots, l^q : v_m^q) = \sigma(v_m)$$



Propagating Information Forwardly and Backwardly

Dense constrained system

$$[v]^p = [v]^p \wedge F_v^{s,p}([v_1]^s, \dots, [v_n]^s)$$

Sparse SSI constrained system

$$[v] = [v] \wedge G_v^i([a], \dots, [z]) \text{ where } a, \dots, z \text{ are used (resp. defined) at } i$$

Proof

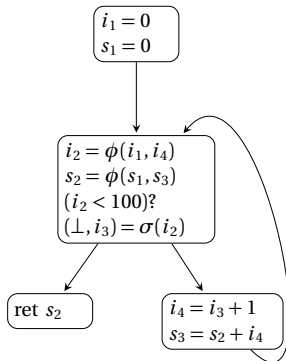
- coalesce all $[v]^p$ such that $v \in \text{live}(p)$ into $[v]$;
replace all $[v]^p$ such that $v \notin \text{live}(p)$ by \top
- for each instruction *inst* with uses $a \dots z$, let
 $G_v^i([a], \dots, [z]) = F_v^i([v_1], \dots, [v_n])$
- remove redundancies

Propagating Information Forwardly and Backwardly

Backward propagation engine under SSI

```
1  function back_propagate(transfer_functions  $\mathcal{G}$ )
2      worklist =  $\emptyset$ 
3      foreach  $v \in \text{vars}$ :  $[v] = \top$ 
4      foreach  $i \in \text{insts}$ : worklist +=  $i$ 
5      while worklist  $\neq \emptyset$ :
6          let  $i \in \text{worklist}$ ; worklist -=  $i$ 
7          foreach  $v \in i.\text{uses}()$ :
8               $[v]_{\text{new}} = [v] \wedge G_v^i([i.\text{defs}()])$ 
9              if  $[v] \neq [v]_{\text{new}}$ :
10                 worklist +=  $v.\text{defs}()$ 
11                  $[v] = [v]_{\text{new}}$ 
```

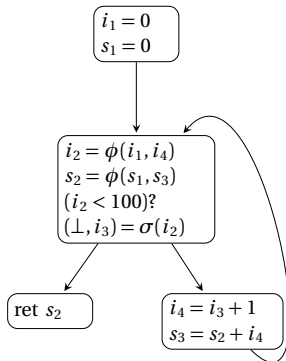

Examples of sparse data-flow analyses



Range analysis (forward from defs & conds)

$[i_1] \cup = [0, 0]$	\emptyset
$[s_1] \cup = [0, 0]$	\emptyset
$[i_2] \cup = [i_1] \cup [i_4]$	\emptyset
$[s_2] \cup = [s_1] \cup [s_3]$	\emptyset
$[i_3] \cup = ([i_2] \cap] - \infty, 99])$	\emptyset
$[i_4] \cup = ([i_3] + 1)$	\emptyset
$[s_3] \cup = ([s_2] + [i_4])$	\emptyset

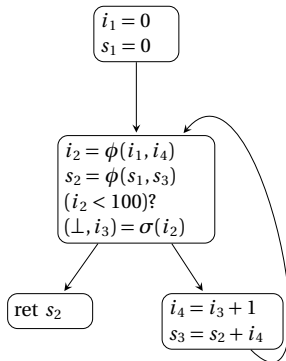
Examples of sparse data-flow analyses



Range analysis (forward from defs & conds)

$[i_1] \cup = [0, 0]$	\emptyset
$[s_1] \cup = [0, 0]$	\emptyset
$[i_2] \cup = [i_1] \cup [i_4]$	\emptyset
$[s_2] \cup = [s_1] \cup [s_3]$	\emptyset
$[i_3] \cup = ([i_2] \cap] - \infty, 99])$	\emptyset
$[i_4] \cup = ([i_3] + 1)$	\emptyset
$[s_3] \cup = ([s_2] + [i_4])$	\emptyset

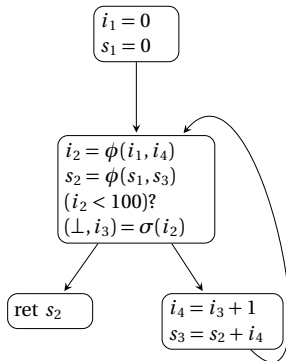
Examples of sparse data-flow analyses



Range analysis (forward from defs & conds)

$[i_1] \cup = [0, 0]$	$[0, 0]$
$[s_1] \cup = [0, 0]$	\emptyset
$[i_2] \cup = [i_1] \cup [i_4]$	\emptyset
$[s_2] \cup = [s_1] \cup [s_3]$	\emptyset
$[i_3] \cup = ([i_2] \cap]-\infty, 99])$	\emptyset
$[i_4] \cup = ([i_3] + 1)$	\emptyset
$[s_3] \cup = ([s_2] + [i_4])$	\emptyset

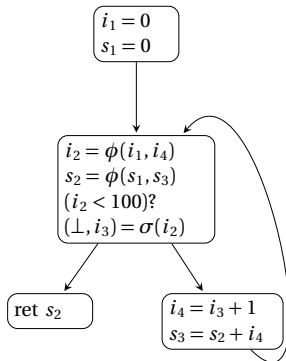
Examples of sparse data-flow analyses



Range analysis (forward from defs & conds)

$[i_1] \cup = [0, 0]$	$[0, 0]$
$[s_1] \cup = [0, 0]$	\emptyset
$[i_2] \cup = [i_1] \cup [i_4]$	\emptyset
$[s_2] \cup = [s_1] \cup [s_3]$	\emptyset
$[i_3] \cup = ([i_2] \cap]-\infty, 99])$	\emptyset
$[i_4] \cup = ([i_3] + 1)$	\emptyset
$[s_3] \cup = ([s_2] + [i_4])$	\emptyset

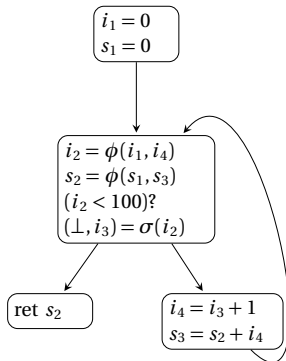
Examples of sparse data-flow analyses



Range analysis (forward from defs & conds)

$[i_1] \cup = [0, 0]$	$[0, 0]$
$[s_1] \cup = [0, 0]$	$[0, 0]$
$[i_2] \cup = [i_1] \cup [i_4]$	\emptyset
$[s_2] \cup = [s_1] \cup [s_3]$	\emptyset
$[i_3] \cup = ([i_2] \cap]-\infty, 99])$	\emptyset
$[i_4] \cup = ([i_3] + 1)$	\emptyset
$[s_3] \cup = ([s_2] + [i_4])$	\emptyset

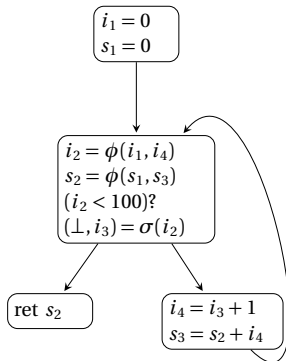
Examples of sparse data-flow analyses



Range analysis (forward from defs & conds)

$[i_1] \cup = [0, 0]$	$[0, 0]$
$[s_1] \cup = [0, 0]$	$[0, 0]$
$[i_2] \cup = [i_1] \cup [i_4]$	\emptyset
$[s_2] \cup = [s_1] \cup [s_3]$	\emptyset
$[i_3] \cup = ([i_2] \cap]-\infty, 99])$	\emptyset
$[i_4] \cup = ([i_3] + 1)$	\emptyset
$[s_3] \cup = ([s_2] + [i_4])$	\emptyset

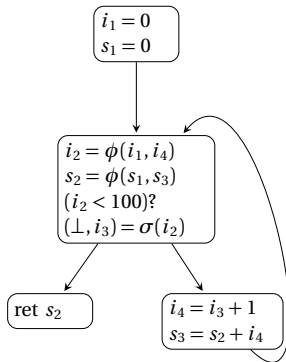
Examples of sparse data-flow analyses



Range analysis (forward from defs & conds)

$[i_1] \cup = [0, 0]$	$[0, 0]$
$[s_1] \cup = [0, 0]$	$[0, 0]$
$[i_2] \cup = [i_1] \cup [i_4]$	$[0, 0]$
$[s_2] \cup = [s_1] \cup [s_3]$	\emptyset
$[i_3] \cup = ([i_2] \cap]-\infty, 99])$	\emptyset
$[i_4] \cup = ([i_3] + 1)$	\emptyset
$[s_3] \cup = ([s_2] + [i_4])$	\emptyset

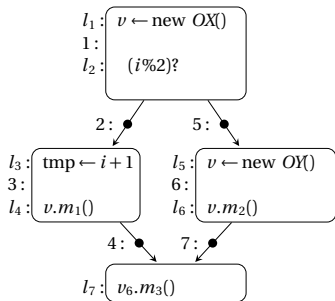
Examples of sparse data-flow analyses



Range analysis (forward from defs & conds)

$[i_1] \cup = [0, 0]$	$[0, 0]$
$[s_1] \cup = [0, 0]$	$[0, 0]$
$[i_2] \cup = [i_1] \cup [i_4]$	$[0, 100]$
$[s_2] \cup = [s_1] \cup [s_3]$	$[0, +\infty[$
$[i_3] \cup = ([i_2] \cap]-\infty, 99])$	$[0, 99]$
$[i_4] \cup = ([i_3] + 1)$	$[1, 100]$
$[s_3] \cup = ([s_2] + [i_4])$	$[1, +\infty[$

Examples of sparse data-flow analyses

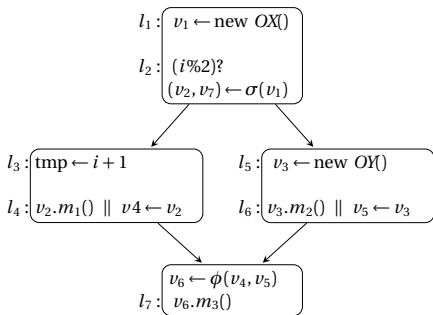


Class inference (backward from uses)

<i>prog. point</i>	$[v]$
1	$\{m_1, m_3\}$
2	$\{m_1, m_3\}$
3	$\{m_1, m_3\}$
4	$\{m_3\}$
5	\top
6	$\{m_2, m_3\}$
7	$\{m_3\}$

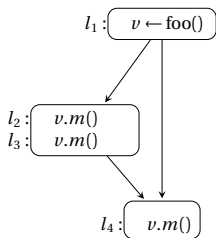
Examples of sparse data-flow analyses

Class inference (backward from uses)



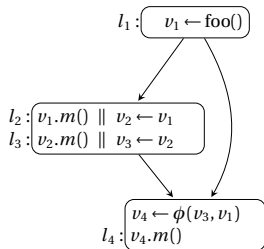
$[v_7]$	\top
$[v_6] \cup = \{m_3\}$	$\{m_3\}$
$[v_5] \cup = [v_6]$	$\{m_3\}$
$[v_4] \cup = [v_6]$	$\{m_3\}$
$[v_2] \cup = (\{m_1\} \wedge [v_4])$	$\{m_1, m_3\}$
$[v_3] \cup = (\{m_2\} \wedge [v_4])$	$\{m_2, m_3\}$
$[v_1] \cup = [v_2]$	$\{m_1, m_3\}$
$[v_1] \cup = [v_7]$	$\{m_1, m_3\}$

Examples of sparse data-flow analyses



Null pointer (forward from defs & uses)

Examples of sparse data-flow analyses



Null pointer (forward from defs & uses)

$$\begin{array}{l|l}
 [v_1] \wedge = 0 & 0 \\
 [v_2] \wedge = \emptyset & \emptyset \\
 [v_3] \wedge = \emptyset & \emptyset \\
 [v_4] \wedge = ([v_3] \wedge [v_1]) & 0
 \end{array}$$



Splitting strategy

Live range splitting strategy $\mathcal{P}_v = I_\uparrow \cup I_\downarrow$

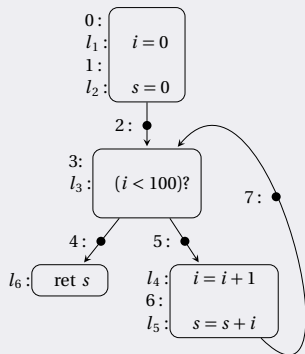
I_\downarrow : set of points i with forward direction

I_\uparrow : set of points i with backward direction

- 1 function SSIfy(var v , Splitting_Strategy \mathcal{P}_v)
- 2 split(v , \mathcal{P}_v)
- 3 rename(v)
- 4 clean(v)

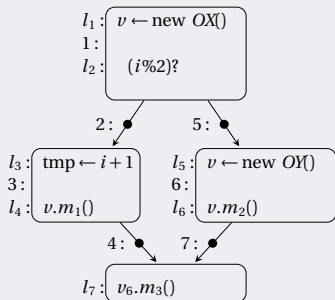
Splitting strategy

Range analysis: $\mathcal{P}_i = \{l_1, \text{Out}(l_3), l_4\}_\downarrow$



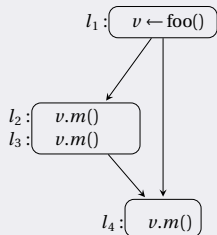
Splitting strategy

Class inference: $\mathcal{P}_v = \{l_4, l_6, l_7\} \uparrow$



Splitting strategy

Null pointer: $\mathcal{P}_v = \{l_1, l_2, l_3, l_4\} \downarrow$



Splitting strategy

Client	Splitting strategy \mathcal{P}
Alias analysis, reaching definitions cond. constant propagation	$Defs_{\downarrow}$
Partial Redundancy Elimination	$Defs_{\downarrow} \cup LastUses_{\uparrow}$
ABCD, taint analysis, range analysis	$Defs_{\downarrow} \cup Out(Conds)_{\downarrow}$
Stephenson's bitwidth analysis	$Defs_{\downarrow} \cup Out(Conds)_{\downarrow} \cup Uses_{\uparrow}$
Mahlke's bitwidth analysis	$Defs_{\downarrow} \cup Uses_{\uparrow}$
An's type inference, Class inference	$Uses_{\uparrow}$
Hochstadt's type inference	$Uses_{\uparrow} \cup Out(Conds)_{\uparrow}$
Null-pointer analysis	$Defs_{\downarrow} \cup Uses_{\downarrow}$



Splitting live ranges

- Split live range of v at each $p \in \mathcal{P}_v$
- Split live range where the information collide (join set $\mathcal{J}(I_\downarrow)$ and split set $\mathcal{S}(I_\uparrow)$)
- Iterated dominance frontier $DF^+(S) = \mathcal{J}(S \cup \{r\})$ can be computed efficiently (as opposed to $\mathcal{J}(S)$)
- Iterated post dominance frontier $pDF^+(S) = \mathcal{J}(S \cup \{r\})$ for the reverse CFG

function split(var v , Splitting_Strategy $\mathcal{P}_v = I_\downarrow \cup I_\uparrow$)

$[I_\downarrow \cup \text{In}(DF^+(I_\downarrow))]$



Splitting live ranges

- Split live range of v at each $p \in \mathcal{P}_v$
- Split live range where the information collide (join set $\mathcal{J}(I_{\downarrow})$ and split set $\mathcal{S}(I_{\uparrow})$)
- Iterated dominance frontier $DF^+(S) = \mathcal{J}(S \cup \{r\})$ can be computed efficiently (as opposed to $\mathcal{J}(S)$)
- Iterated post dominance frontier $pDF^+(S) = \mathcal{J}(S \cup \{r\})$ for the reverse CFG

function split(var v , Splitting_Strategy $\mathcal{P}_v = I_{\downarrow} \cup I_{\uparrow}$)

$$[I_{\downarrow} \cup \text{In}(DF^+(I_{\downarrow}))] \cup [I_{\uparrow} \cup \text{Out}(pDF^+(I_{\uparrow}))]$$



Splitting live ranges

- Split live range of v at each $p \in \mathcal{P}_v$
- Split live range where the information collide (join set $\mathcal{J}(I_\downarrow)$ and split set $\mathcal{S}(I_\uparrow)$)
- Iterated dominance frontier $DF^+(S) = \mathcal{J}(S \cup \{r\})$ can be computed efficiently (as opposed to $\mathcal{J}(S)$)
- Iterated post dominance frontier $pDF^+(S) = \mathcal{J}(S \cup \{r\})$ for the reverse CFG

function $\text{split}(\text{var } v, \text{Splitting_Strategy } \mathcal{P}_v = I_\downarrow \cup I_\uparrow)$

$$\mathcal{P}_v \cup \text{In} [DF^+(I_\downarrow \cup [I_\uparrow \cup \text{Out}(pDF^+(I_\uparrow))])]$$



Variable Renaming

function rename(var v)

- traverses the CFG along topological order
- give a unique version to each definition of v
- stack the versions that dominates the current program point
- rename each use of v with the version of immediately dominating definition



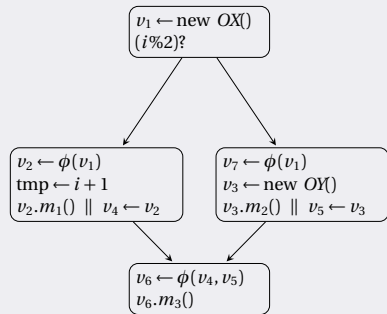
Dead and Undefined Code Elimination

clean(var v)

- actual instructions: instructions originally in the code
- SSA graph: nodes are instructions; edges are def-use chains
- active instructions: instructions connected to an actual instruction
- simple traversal of the SSA graph from actual instructions that mark active ones
- remove non-active instructions (inserted phi and sigma functions)

Implementation Details

Implementing σ -functions



SSI destruction

